



HAL
open science

PAC-Bayesian Bounds and Beyond : Self-Bounding Algorithms and New Perspectives on Generalization in Machine Learning

Paul Viallard

► **To cite this version:**

Paul Viallard. PAC-Bayesian Bounds and Beyond : Self-Bounding Algorithms and New Perspectives on Generalization in Machine Learning. Machine Learning [cs.LG]. Université Jean Monnet - Saint-Etienne, 2022. English. NNT : 2022STET0057 . tel-04496162

HAL Id: tel-04496162

<https://theses.hal.science/tel-04496162v1>

Submitted on 8 Mar 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



N°d'ordre NNT : 2022LYSES057

THÈSE de DOCTORAT DE L'UNIVERSITÉ JEAN MONNET SAINT-ÉTIENNE

Membre de la COMUE de LYON

École Doctorale N° 488
Sciences Ingénierie Santé

Spécialité / discipline de doctorat :
Informatique

Soutenue publiquement le 07/12/2022, par :
Paul Viillard

PAC-Bayesian Bounds and Beyond: Self-Bounding Algorithms and New Perspectives on Generalization in Machine Learning

Devant le jury composé de :

Gribonval, Rémi	Directeur de recherche	ENS Lyon, Inria	Président
Canu, Stéphane	Professeur	INSA de Rouen	Rapporteur
Ralaivola, Liva	VP Research	Criteo AI Lab	Rapporteur
Tommasi, Marc	Professeur	Université de Lille, Inria	Examineur
Habrador, Amaury	Professeur	Université de Saint-Etienne	Directeur
Germain, Pascal	Professeur adjoint	Université Laval	Co-encadrant
Morvant, Emilie	Maître de conférences	Université de Saint-Etienne	Co-encadrante

ACKNOWLEDGEMENTS

First of all, I would like to thank Liva and Stéphane for reviewing my thesis (in a short time). Moreover, I would like to thank Marc and Rémi for accepting to be, respectively, the examiner and the president; I am very honored to have you on my jury!

I would like to sincerely thank my supervisors, who were always here and patient with me; I believe that I was very lucky! Thank you, Amaury, for proposing this PhD subject to me and for always having a critical viewpoint on my work. In the meantime, I always felt free to work where and when I want; thank you. Thank you, Pascal, for your guidance when writing the papers. Our discussions were always interesting, and I learned a lot. Visiting Laval University has not been possible with the pandemic. . . I hope that it will be possible later! Thank you, Emilie, for being there literally every day for me. I really enjoyed talking about science and personal things with you. Thank you for showing me (for a long time) how to write a scientific paper and how to do presentations. Moreover, thank you for letting me the chance to perform lectures (which was hard but so interesting). Lastly, thanks, Emilie and Cécile, for making me discover Kung Fu (which I gave up but enjoyed a lot)! Thanks, Valentina and Rémi, for having contributed together. This thesis would not have looked like this without you.

This PhD journey has also been full of fun moments and scientific discussions among the permanent and non-permanent staff of the lab. Generally, thank you all¹ for the good mood in the lab that I really enjoyed during these years. More specifically, thanks, Léo and Jules, for introducing me to the “ASEC” parties that were awesome with you! I also want to thank the “ASEC” team (*i.e.*, Aubin, Damien, Jules, and, Léo) for the good nights. Unfortunately, the pandemic² stopped, for a long time, these good moments. . .

Moreover, this thesis marks the end of all those school years. I would have never bet on myself about writing a “book”, especially in English! This could never have been done if I had stopped studying early; thanks, mom and dad, for your support during these years. More generally, I want to thank my family, who have all been supportive during these moments of hard work with the projects, the exams, and during the PhD.

I want to thank some special friends that have been important to me since primary,

¹The list is too long to be enumerated!

²Thanks, Brrmmaou, Titi, and Bibi, for sleeping next to me when I was working; it was motivating!

middle, high school, and the university. Thanks, Jules, for all these years: I would never have been the same without you! Thank you for making me discover the Loire gorges along with hiking since middle school. Our philosophical discussions in this area are so memorable! Your recklessness helped me a lot to push my limits, *e.g.*, during the urban and cave explorations. Thank you for accompanying me in my first climbing and paragliding! Thanks, Fleurian, for these school years: our hiking and climbing sessions were so cool (and usually with Jules)! Moreover, thank you for making me discover a new sport while writing my thesis: karting. I hope that there will be other road trips with the Dyane car! Thanks, Benoit, for sharing my passion for computer science and coding together for all these years. For instance, we coded an optimization algorithm in high school! Likewise, thanks, Angélique, for sharing your interest in computer science with me since high school. I hope that we will continue to hike or drink tea at the Oriol Tower! Thanks, Guillaume, for these wonderful student years spent with you since the first year of “licence”! We shared many parties, met our respective friends, and had good times. Moreover, thank you for doing sports with me, such as climbing, ice skating, and hiking (during classes). These good moments were accompanied by a lot of work: for you for being my project partner for almost five years! I also appreciated our first internship together at the same desk and our NeurIPS paper! It is funny where life leads us. I would like to thank Coralie for asking me for information about horse riding during my first internship. Since then, I have ridden again and learned a lot about horses! Your restlessness (and your desire to always do something) has been valuable, especially when avoiding thinking about the thesis.

Last but not least, thank you, Chloé, for your unwavering support over these past two years. It was precious, notably for the difficult moment like the submissions and the thesis writing! Thank you for removing (sometimes) my nose from the grindstone to enjoy life (with you). I really hope to continue horse riding, hiking, traveling, or having common projects.

The further backward you look, the further forward you can see.

WINSTON CHURCHILL

Science is a bit like the joke about the drunk who is looking under a lamppost for a key that he has lost on the other side of the street, because that's where the light is. It has no other choice.

NOAM CHOMSKY

CONTENTS

Contents	6
List of Figures	9
List of Theorems	12
List of Notations	15
Preamble	17
List of Publications	25
International Conference	25
International Workshop	25
National Conference	25
Research Report	26
I Background	27
1 An introduction to Statistical Learning Theory	29
1.1 Introduction	30
1.2 Hypothesis Selection	40
1.3 Generalization Bounds	43
1.4 Conclusion and Summary	53
2 The PAC-Bayesian Theory and the Majority Vote	55
2.1 Introduction	56
2.2 PAC-Bayesian Majority Votes	56
2.3 PAC-Bayesian Bounds	64
2.4 Disintegrated PAC-Bayesian Bounds	75
2.5 Conclusion and Summary	78
II PAC-Bayesian Majority Vote: Theory and Self-bounding Algorithms	79
3 PAC-Bayesian Theory for the Robust Majority Vote	81

3.1	Introduction	82
3.2	Adversarially Robust Majority Vote	83
3.3	Adversarially Robust PAC-Bayes	87
3.4	Experimental Evaluation on Differentiable Decision Trees	95
3.5	Conclusion and Summary	98
4	Self-Bounding Algorithms for the Majority Vote	101
4.1	Introduction	102
4.2	Setting	103
4.3	State of the Art: PAC-Bayesian Bounds for the Majority Vote	104
4.4	Contribution: Algorithms based on the PAC-Bayesian C-Bounds	108
4.5	Experiments	114
4.6	Conclusion and Summary	123
5	Toward a Stochastic Majority Vote	125
5.1	Introduction	126
5.2	The Stochastic Majority Vote	129
5.3	From a PAC-Bayesian Bound to an Algorithm	133
5.4	Experiments	138
5.5	Conclusion and Summary	141
III Beyond PAC-Bayesian Bounds: From Disintegration to Novel Bounds		147
6	On the Practical uses of the Disintegrated Bounds	149
6.1	Introduction	150
6.2	Setting and PAC-Bayesian Bounds	151
6.3	Disintegrated PAC-Bayesian Theorems	152
6.4	The Disintegration in Action	158
6.5	Experiments with Neural Networks	162
6.6	Perspectives for the Majority Vote	171
6.7	Summary and Conclusion	172
7	Generalization Bounds with Complexity Measures	175
7.1	Introduction	176
7.2	Preliminaries	177
7.3	Integrating Arbitrary Complexities in Generalization Bounds	179
7.4	Using Arbitrary Complexities in Practice	184
7.5	Comparison with the Generalization Bounds of the Literature	192
7.6	Conclusion and Summary	197

IV Conclusion and Perspectives	199
Conclusion and Perspectives	201
V Appendix	205
A Some Mathematical Tools	207
A.1 JENSEN's Inequality	207
A.2 MARKOV's Inequality	207
A.3 2nd Order MARKOV's Inequality	208
A.4 CHEBYSHEV-CANTELLI Inequality	209
A.5 HÖLDER's Inequality	209
B Appendix of Chapter 2	213
C Appendix of Chapter 3	245
D Appendix of Chapter 4	265
E Appendix of Chapter 5	285
F Appendix of Chapter 6	303
G Appendix of Chapter 7	347
References	367

LIST OF FIGURES

Preamble

1	Illustration of the Supervised Classification Setting in Machine Learning . . .	18
2	Example of the Majority Vote's Prediction	20
3	Example of a Majority Vote with Three Voters	21

Chapter 1

1.1	Representation of Statistical Learning	30
1.2	Schematic Representation of a Data Distribution	32
1.3	Example of Decision Stump	33
1.4	Example of Decision Tree	34
1.5	Plot of the Exclusive Or Function	36
1.6	Examples of Linear Classifiers	37
1.7	Example of Neural Network	37
1.8	Details on the Neural Network of Figure 1.7	38
1.9	Plot of Different Losses	39
1.10	Representation of the Underfitting, Overfitting and the Generalization . . .	42
1.11	Illustration of the VC-dimension for Linear Classifiers in 2-dimension	46
1.12	Illustration of the Rademacher Complexity in Multi-class Classification . . .	47
1.13	Schematic Representation of the Uniform Stability	49
1.14	Schematic Representation of the Algorithmic Robustness	51

Chapter 2

2.1	Illustration of the Margin of the Majority Vote	59
2.2	Plots that Summarize the Relationship Between the Surrogates	65
2.3	Illustration of the Functions \underline{k}_l and \overline{k}_l	71
2.4	Illustration of the Tightness of \underline{k}_l and \overline{k}_l	72

Chapter 3

3.1	Illustration of the Adversarial Examples	82
3.2	Visualization of the impact of the TV term in Equation (3.9)	98
3.3	Visualization of the risk and bound values for "Defense=Attack"	99

Chapter 4

4.2	Comparison Between the Test Risks and the Bounds (1/6)	115
4.3	Comparison Between the Test Risks and the Bounds (2/6)	116
4.4	Comparison Between the Test Risks and the Bounds (3/6)	117
4.5	Comparison Between the Test Risks and the Bounds (4/6)	118
4.6	Comparison Between the Test Risks and the Bounds (5/6)	119
4.7	Comparison Between the Test Risks and the Bounds (6/6)	120

Chapter 5

5.1	Comparison of the Minimization of the Surrogates on Moons	126
5.2	Comparison of the Self-bounding Algorithms on Moons	128
5.3	Examples of Probability Density Functions for the Dirichlet Distribution	130
5.4	Plot of the Digamma Function and its Derivative	134
5.5	Plot of the Performance of Algorithms 5.1 and 5.2 on Moons (1/2)	140
5.6	Plot of the Performance of Algorithms 5.1 and 5.2 on Moons (2/2)	141
5.7	Comparison of the Risks and the Bound Values (1/4)	142
5.8	Comparison of the Risks and the Bound Values (2/4)	143
5.9	Comparison of the Risks and the Bound Values (3/4)	144
5.10	Comparison of the Risks and the Bound Values (4/4)	145

Chapter 6

6.1	Evolution of the Bound Values in Terms of the Split Ratio (1/2)	165
6.2	Evolution of the Bound Values in Terms of the Split Ratio (2/2)	166
6.3	Comparisons of the PAC-Bayesian Bounds and the Disintegrated Bounds	167
6.4	Comparison of the Disintegrated Bounds and the Test Risks	168
6.5	Evolution of the Bound Values in Terms of the Variance σ^2	170

Chapter 7

7.1	Tightness of Equations (7.7) and (7.8) on MNIST	188
7.2	Tightness of Equations (7.7) and (7.8) on FashionMNIST	189
7.3	Influence of the Parameter α	190
7.4	Influence of the Depth/Width	191

LIST OF ALGORITHMS

Chapter 1

1.1	Example of Decision Stump	33
1.2	Example of Decision Tree	34
1.3	Empirical Risk Minimization	41
1.4	Structural Risk Minimization	42

Chapter 2

2.1	Compute $\bar{\text{kl}}(q \tau)$ resp. $\underline{\text{kl}}(q \tau)$ through the bisection method	72
-----	--	----

Chapter 3

3.1	Average Adversarial Training with Guarantee	94
-----	---	----

Chapter 4

4.1	Minimization of Equation (4.2) by Stochastic Gradient Descent	108
4.2	Minimization of Equation (4.3) by Stochastic Gradient Descent	111
4.3	Minimization of Equation (4.4) by Stochastic Gradient Descent	113

Chapter 5

5.1	Approximating the Stochastic Risk	131
5.2	Computing Exactly the Stochastic Majority Vote's Risk	133
5.3	Minimization of Theorem 5.3.1's Bound	137
5.4	Minimization of Corollary 5.3.1's Bound	137

Chapter 7

7.1	Stochastic MALA	185
-----	---------------------------	-----

LIST OF THEOREMS

Chapter 1

1.1.1 Definition (Learning sample)	31
1.1.2 Definition (Loss function)	38
1.1.3 Definition (Empirical Risk)	39
1.1.4 Definition (True Risk)	40
1.3.1 Definition (PAC Generalization Bound)	43
1.3.2 Definition (Uniform Convergence Bound)	43
1.3.1 Theorem (Generalization Bound for Finite \mathbb{H})	44
1.3.3 Definition (VAPNIK-CHEVONENKIS (VC) Dimension)	45
1.3.2 Theorem (VC-Dimension-based Generalization Bounds)	45
1.3.4 Definition (Rademacher Complexity)	46
1.3.3 Theorem (Rademacher-complexity-based Generalization Bound)	47
1.3.5 Definition (Algorithmic-dependent Generalization Bound)	48
1.3.6 Definition (Uniform Stability)	49
1.3.4 Theorem (Stability-based Bounds)	50
1.3.7 Definition (Robustness)	50
1.3.5 Theorem (Robustness-based Bounds)	51
1.3.6 Theorem (PAC-Bayesian Bound of MAURER (2004))	52

Chapter 2

2.2.1 Definition (Majority Vote)	56
2.2.2 Definition (Risk of the Majority Vote)	58
2.2.3 Definition (Margin of the Majority Vote)	58
2.2.4 Definition ($\frac{1}{2}$ -Margin of the Majority Vote)	59
2.2.5 Definition (Gibbs Risk)	60
2.2.1 Theorem (Risk Upper Bound Based on the Gibbs Risk)	61
2.2.6 Definition (Joint Error)	61
2.2.2 Theorem (Risk Upper Bound Based on the Joint Error)	62
2.2.7 Definition (Disagreement)	62
2.2.3 Theorem (The C-Bound)	63
2.2.4 Theorem (Relationship between Theorems 2.2.1 to 2.2.3)	64
2.3.1 Definition (PAC-Bayesian Generalization Bound)	66
2.3.1 Theorem (General PAC-Bayesian Bound of GERMAIN <i>et al.</i> (2009))	67
2.3.2 Theorem (PAC-Bayesian Bound of MCALLESTER (2003))	67

2.3.3 Theorem (PAC-Bayesian Bound of CATONI (2007))	68
2.3.2 Definition (KL Divergence Between Two Bernoulli Distributions)	69
2.3.4 Theorem (PAC-Bayesian Bound of SEEGER (2002))	70
2.3.3 Definition (Inverting Functions of $\text{kl}()$)	70
2.3.1 Proposition (DONSKER-VARADHAN Variational Representation)	73
2.3.5 Theorem (General PAC-Bayesian Bound of BÉGIN <i>et al.</i> (2016))	74
2.4.1 Definition (Disintegrated PAC-Bayesian Generalization Bound)	75
2.4.1 Theorem (General Disintegrated Bound of RIVASPLATA <i>et al.</i> (2020))	76
2.4.2 Theorem (Disintegrated Bound of CATONI (2007))	76
2.4.3 Theorem (Disintegrated Bound of BLANCHARD and FLEURET (2007))	77

Chapter 3

3.2.1 Definition (True/Empirical Adversarial Risk)	84
3.2.2 Definition (Averaged Adversarial Risk)	86
3.2.3 Definition (Averaged-Max Adversarial Risk)	87
3.3.1 Proposition (Relations Between the Averaged Adversarial Risks)	87
3.3.2 Proposition (Classical and Averaged Adversarial Risks)	88
3.3.1 Definition (Surrogates on the Averaged Adversarial Risks)	89
3.3.1 Theorem (Upper Bounds on the Surrogates)	89
3.3.2 Theorem (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$)	90
3.3.3 Theorem (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$)	90
3.3.1 Corollary (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$)	92
3.3.2 Corollary (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$)	92

Chapter 4

4.3.1 Theorem (PAC-Bayesian Bound Based on the Gibbs Risk)	104
4.3.2 Theorem (PAC-Bayesian Bound Based on the Joint Error)	105
4.3.3 Theorem (PAC-Bayesian C-Bound of ROY <i>et al.</i> (2016))	106
4.3.4 Theorem (PAC-Bayesian C-Bound of GERMAIN <i>et al.</i> (2015))	107
4.3.5 Theorem (PAC-Bayesian C-Bound of LACASSE <i>et al.</i> (2006))	107
4.4.1 Theorem (Reformulation of LACASSE <i>et al.</i> 's PAC-Bayesian C-Bound)	111

Chapter 5

5.2.1 Definition (Risks of the stochastic majority vote)	129
5.2.2 Definition (Dirichlet Distribution)	130
5.2.1 Lemma (Computation of the Stochastic Risk)	131
5.2.1 Corollary (Closed-form Solution of the Stochastic Risks)	132
5.3.1 Theorem (PAC-Bayesian Bound for Stochastic Majority Votes)	133

5.3.2 Theorem (PAC-Bayesian bound with data-dependent voters)	135
5.3.1 Corollary (PAC-Bayesian bound with data-dependent voters)	136

Chapter 6

2.3.5 Theorem (General PAC-Bayesian Bound of BÉGIN <i>et al.</i> (2016))	151
2.4.1 Definition (Disintegrated PAC-Bayesian Generalization Bound)	152
6.3.1 Theorem (General Disintegrated PAC-Bayes Bound)	153
6.3.1 Corollary (Extreme Cases of Theorem 6.3.1)	154
2.4.1 Theorem (General Disintegrated Bound of RIVASPLATA <i>et al.</i> (2020))	155
6.3.2 Theorem (Parametrizable Disintegrated PAC-Bayes Bound)	156
6.3.1 Proposition (Optimal Bound of Theorem 6.3.2)	157
6.4.1 Corollary (Instantiation of Theorem 6.3.1 for Neural Networks)	159
6.4.2 Corollary (Instantiation of Known Bounds for Neural Networks)	159
6.4.3 Corollary (PAC-Bayesian Bound for Stochastic Neural Networks)	161
6.6.1 Corollary (Instantiation of Theorem 2.4.1 to Stochastic Majority Votes)	171
6.6.2 Corollary (Instantiation of Theorem 6.3.1 to Stochastic Majority Votes)	172

Chapter 7

2.4.1 Definition (Disintegrated PAC-Bayesian Generalization Bound)	178
2.4.1 Theorem (General Disintegrated Bound of RIVASPLATA <i>et al.</i> (2020))	178
7.3.1 Definition (Generalization Bound with Complexity Measures)	180
7.3.1 Theorem (Generalization Bound with Complexity Measures)	182
7.3.1 Corollary (Practical Generalization Bound with Complexity Measures)	182
7.5.1 Proposition (Set-theoretic view of Definition 7.3.1)	193
1.3.2 Definition (Uniform Convergence Bound)	193
7.5.2 Proposition (Set-theoretic View of Uniform Convergence Bounds)	194
7.5.1 Corollary (Uniform Convergence Bound from Theorem 7.3.1)	194
1.3.5 Definition (Algorithmic-dependent Generalization Bound)	195
7.5.3 Proposition (Set-theoretic View of Algorithmic Dependent Bounds)	196
7.5.2 Corollary (Algorithmic-dependent Bound from Theorem 7.3.1)	196

Appendix

A.1.1 Theorem (JENSEN's Inequality)	207
A.2.1 Theorem (MARKOV's Inequality)	207
A.3.1 Theorem (2nd Order MARKOV's Inequality)	208
A.4.1 Theorem (CHEBYSHEV-CANTELLI Inequality)	209
A.5.1 Lemma (YOUNG's Inequality)	209
A.5.1 Theorem (HÖLDER's Inequality)	210

LIST OF NOTATIONS

General

a	A scalar (integer or real)
\mathbf{a}	A vector
\mathbf{A}	A matrix
\mathbb{A}, \mathcal{A}	A set
\mathbb{R}	The set of real numbers
\mathbb{R}_*	The set of real numbers excluding 0
\mathbb{R}_*^+	The set of positive real numbers excluding 0
\mathbb{N}	The set of natural numbers
\mathbb{N}_*	The set of natural numbers excluding 0
$\text{card}(\cdot)$	The cardinal of a set
a_i	i -th element of the vector \mathbf{a}

Statistical Learning Theory

\mathcal{X}	Set of d -dimensional inputs ($\subseteq \mathbb{R}^d$)
\mathcal{Y}	Set of labels
\mathbf{x}	A real-valued input $\mathbf{x} \in \mathcal{X}$
y	A label $y \in \mathcal{Y}$ associated to the input \mathbf{x}
\mathcal{D}	Unknown data distribution on $\mathcal{X} \times \mathcal{Y}$
\mathcal{D}^m	Unknown data distribution on the m -samples, <i>i.e.</i> , on $(\mathcal{X} \times \mathcal{Y})^m$
\mathcal{S}	Learning sample $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ drawn from \mathcal{D}^m
\mathcal{S}	Uniform distribution on \mathcal{S}
\mathbb{T}	Test set drawn from \mathcal{D}^m

\mathcal{T}	Uniform distribution on \mathbb{T}
\mathbb{H}	The set of hypotheses
h	A hypothesis $h \in \mathbb{H}$
$\ell(\cdot, \cdot)$	Loss function
$R_{\mathcal{D}'}^\ell(h)$	Risk of the hypothesis $h \in \mathbb{H}$ w.r.t. the loss function $\ell(\cdot)$ on \mathcal{D}'
$R_{\mathcal{D}'}(h)$	Risk of the hypothesis $h \in \mathbb{H}$ w.r.t. the 0-1 loss on \mathcal{D}'

Probability Theory

$\mathbb{E}_{X \sim \mathcal{X}}[\cdot]$	The expectation w.r.t. the random variable $X \sim \mathcal{X}$
$\mathbb{P}_{X \sim \mathcal{X}}[\cdot]$	The probability w.r.t. the random variable $X \sim \mathcal{X}$
$\mathbb{I}[a]$	Indicator function; returns 1 if a is true and 0 otherwise
$\mathbb{M}(\mathbb{H})$	Set of Probability densities w.r.t. the reference measure on \mathbb{H}
ρ	Posterior distribution $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H}
π	Prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H}
$\text{KL}(\rho \pi)$	Kullback-Leibler (KL) divergence between ρ and π
$D_\alpha(\rho \pi)$	Rényi Divergence between ρ and π
$\text{Uni}(\mathbb{A})$	Uniform distribution on \mathbb{A}
$\text{Dir}(\boldsymbol{\alpha})$	Dirichlet distribution of parameters $\boldsymbol{\alpha} \in \mathbb{R}_*^+$

Majority Vote

$\text{MV}_\rho(\cdot)$	The majority vote classifier
$m_\rho(\mathbf{x}, y)$	Majority vote's margin for the example $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$
$\widehat{m}_\rho(\mathbf{x}, y)$	$\frac{1}{2}$ -Margin for the example $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$
$\text{sign}(a)$	Sign function; returns +1 if $a \geq 0$ and -1 otherwise
$r_{\mathcal{D}'}(\rho)$	Gibbs risk on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$
$e_{\mathcal{D}'}(\rho)$	Joint Error on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$
$d_{\mathcal{D}'}(\rho)$	Disagreement on the distribution \mathcal{D}' associated to the majority vote $\text{MV}_\rho(\cdot)$

PREAMBLE

Introduction

Statistical Machine Learning is a subfield of artificial intelligence at the intersection of computer science, statistics, and optimization that consists of a set of learning methods that learn mathematical models³ automatically to solve a task from a statistical perspective. We refer the reader to the textbook of RUSSELL and NORVIG (2020) for a general introduction to artificial intelligence and to BISHOP (2007), MOHRI *et al.* (2012), or SHALEV-SHWARTZ and BEN-DAVID (2014) for an introduction to machine learning.

Various tasks can be solved through these methods, such as image recognition, medical diagnosis, fraud detection, recommendation system, etc. These machine learning methods aim to find a model h belonging to a set \mathbb{H} that solves a given task. These methods assume that we have some data, *i.e.*, a set of examples, that are sufficiently representative of the task. Each example obtained from the task is generally composed of an input represented by some features and its corresponding output. Different types of output can be considered: the supervised regression setting uses real-valued output, while the supervised classification setting assumes that the outputs are categories (*a.k.a.* classes or labels). This thesis stands in the supervised classification setting. The supervised classification methods learn a model, called a classifier, that separates/classifies the inputs into different categories.⁴ For instance, in Figure 1, we illustrate an image classification task: it consists in predicting if an image contains a horse or a cat. More precisely, the input is an image, and the label is either “cat” (the red image) or “horse” (the blue image). From the examples in the data, the learned model h (the black line) separates the red images from the blue images: the classifier correctly predicts all the images from the learning sample.

One way to assess if a model performs well on the examples is to compute the probability that the model misclassifies an example in the available data; this quantity is called the *empirical risk*. However, the machine learning model may learn by heart the examples with an empirical risk at nearly 0. In this case, the model may be completely inefficient on new examples from the task, *i.e.*, unseen data; we say that the model *overfits* the data. To characterize if a model performs well on unseen examples from

³The machine learning models are also referred to as hypotheses in statistical learning theory.

⁴When the outputs are absent from the examples, unsupervised learning methods learn models that group (*i.e.* cluster) together similar inputs.

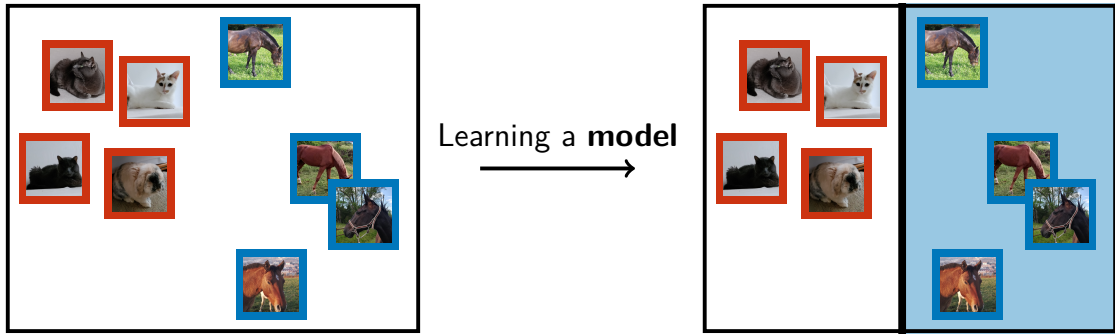


Figure 1. Illustration of the supervised classification setting in machine learning. Given some labeled examples (i.e., the image and its associated category), a model $h \in \mathbb{H}$ is learned, and then, once learned, it can be used to classify, possibly new examples. The images in the blue area are classified as “horse” while the white area corresponds to the images classified as “cat”.

the task, we can define the notion of *true risk*. This notion is the probability that the model misclassifies an example from the task (represented by an unknown distribution). To assess if the empirical risk is representative of the true risk for a given model h , we are interested in the *generalization gap* defined as

$$\text{Generalization Gap}(h) = \left| \text{True Risk}(h) - \text{Empirical Risk}(h) \right|.$$

The model overfits the data when the generalization gap is high (close to 1 in the worst case) while the empirical risk is close to 0. On the contrary, when the gap is close to 0, the empirical risk is a good approximation of the true risk. Hence, to obtain a model that performs well on a task, the generalization gap must be close to 0 and the empirical risk close to 0 as well. However, the gap is not computable because of the true risk since it relies on an unknown quantity: the underlying distribution of the task. Then, another strategy to assess the quality of the model is to consider computable upper bounds on the generalization gap called generalization bounds. The form of the first generalization bound, introduced by VAPNIK and CHERVONENKIS (1968, 1971, 1974), has the following form:

$$\text{For all model } h \in \mathbb{H}, \left| \text{True Risk}(h) - \text{Empirical Risk}(h) \right| \leq \text{Generalization Bound}(\mathbb{H}).$$

The generalization bounds are *probabilistic*, meaning that with high probability (over the examples sampled from the unknown distribution), the bound holds. They generally depend on the number of examples and a complexity term. This complexity term determines the potential of a model to overfit: the higher the complexity, the more

plausible the model overfits the data. Ideally, the upper bound decreases when the number of examples increases for a given finite fixed complexity term. After the seminal work of VAPNIK and CHERVONENKIS (1968, 1971, 1974), generalization bounds have been extended in several directions (see e.g., MCALLESTER, 1998; BARTLETT and MENDELSON, 2002; BOUSQUET and ELISSEEFF, 2002). By rearranging the terms, a bound (that is computable) on the true risk (uncomputable) can be deduced for all model $h \in \mathbb{H}$:

$$\text{True Risk}(h) \leq \text{Empirical Risk}(h) + \text{Generalization Bound}(\mathbb{H}). \quad (1)$$

This leads to a central point of this thesis: the possibility to derive algorithms minimizing a generalization bound, e.g., the right hand side of Equation (1). Algorithms minimizing a generalization bound are called self-bounding algorithms (FREUND, 1998). The advantage of minimizing a generalization bound is the capacity to directly control or have at least an influence on the evolution of the true risk. In particular, as a side effect, the minimization of the generalization bound allows us to control the overfitting phenomenon better.

One particular type of generalization bounds in which we are specifically interested comes from the PAC-Bayesian framework (SHAWE-TAYLOR and WILLIAMSON, 1997; MCALLESTER, 1998). This framework assumes that each model $h \in \mathbb{H}$ is associated with a positive weight $\rho(h)$ that forms a probability distribution over \mathbb{H} called the *posterior distribution* ρ . Based on this assumption, the PAC-Bayesian generalization bounds allow us to upper-bound the expected generalization gap; the form of the bounds is defined as

$$\text{Expectation}_{h \text{ sampled from } \rho} \left[\left| \text{True Risk}(h) - \text{Empirical Risk}(h) \right| \right] \leq \text{Generalization Bound}(\rho).$$

This framework allows to upper-bound the risk of a *stochastic* model which, for each input \mathbf{x} , (i) samples a new model $h \in \mathbb{H}$ from ρ and (ii) predicts the output of \mathbf{x} with $h(\mathbf{x})$. Actually, the risk of the stochastic model can be linked to the risk of a model in which we are particularly interested in this thesis: the majority vote; we provide an overview of such a model in Figure 2. The majority vote has a long history in science: CONDORCET (1785) started to explore mathematically voting systems. Famous machine learning models can be seen as a majority vote, such as linear classifiers, Support Vector Machine (GRAEPEL *et al.*, 2005), k -Nearest Neighbors (BELLET *et al.*, 2014), or neural networks vote (KAWAGUCHI *et al.*, 2017; VIALARD *et al.*, 2019). Some approaches that learn majority votes belong to the ensemble methods (DIETTERICH, 2000) aiming to combine supervised classifiers (called voters) to create an accurate model. For instance, bagging (BREIMAN, 1996), random forest (BREIMAN, 2001) and boosting (FREUND and SCHAPIRE, 1996) are famous examples of ensemble methods.

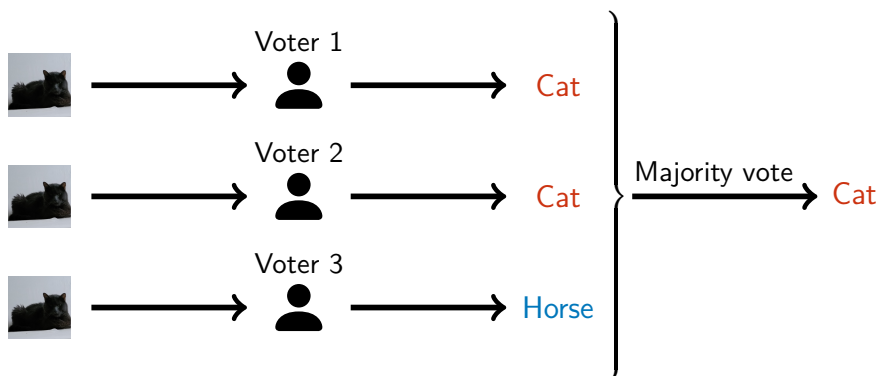


Figure 2. Example of the majority vote’s prediction: given an image, each voter outputs a label (“cat” or “horse”), and the majority vote gathers the results to output the majority label.

In these methods, the voters’ decisions are combined to obtain a better decision compared to the individual voters’ decisions (which can be weak). An important notion when combining different classifiers is the notion of diversity (DIETTERICH, 2000; KUNCHEVA, 2014). Indeed, when voters are *weak* and perform a bit better than random as in boosting (FREUND and SCHAPIRE, 1996), sufficiently diverse voters may improve the accuracy of the majority vote. We give in Figure 3 an example of the diversity’s importance. The combination can be done in very different ways depending on the methods: in bagging (BREIMAN, 1996) and random forest (BREIMAN, 2001), the voters’ predictions are only averaged while boosting (FREUND and SCHAPIRE, 1996) performs a weighted average. In this thesis, we consider a convex combination of the voters where each voter h is associated with the weight $\rho(h)$ encoding its importance in the majority vote. More formally, in the binary classification setting, each model h (*i.e.*, voter) belonging to \mathbb{H} predicts either the class -1 or the class $+1$. A weighted majority vote over the voters in \mathbb{H} applied on a given input \mathbf{x} is defined as:

$$\text{sign} \left[\sum_{h \in \mathbb{H}} \rho(h) h(\mathbf{x}) \right],$$

where $\text{sign}[a] = -1$ if $a < 0$ and $\text{sign}[a] = +1$ otherwise.

However, one drawback of the PAC-Bayesian theory is that it is not possible to bound the generalization gap of only one model h in \mathbb{H} . Hopefully, the *disintegrated* PAC-Bayesian bounds – introduced by BLANCHARD and FLEURET (2007) and CATONI (2007) – overcome this drawback. The bounds have the following form:

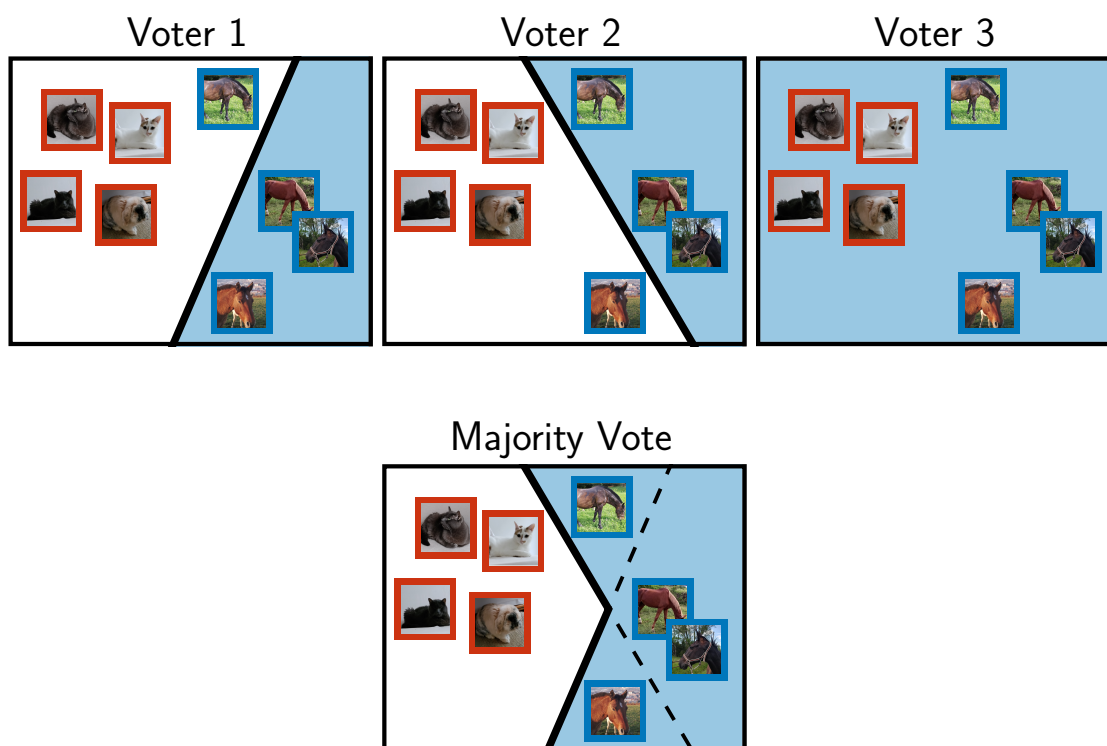


Figure 3. Example of a majority vote with three voters on the classification horse/cat classification task presented before. Each of the three voters makes some mistakes in the data. However, when the majority rule combines the voters, the final vote classifies all the data correctly. It is mainly because the three voters are diverse: they do not make the same mistakes, while the combination corrects the individual errors.

With high probability over the model h sampled from the posterior ρ ,

$$\left| \text{True Risk}(h) - \text{Empirical Risk}(h) \right| \leq \text{Generalization Bound}(\rho, h).$$

They allow us to obtain a bound on the generalization gap for a unique model h (sampled from the posterior distribution) that holds with high probability and which will serve as a basis for some contributions of this thesis.

Long Story Short

Motivations of this thesis. As discussed above, generalization bounds can be used to assess when machine learning models generalize, *i.e.*, when the empirical risk is repre-

sentative of the true risk. In this context, the PAC-Bayesian theory is adapted to upper-bound the generalization gap of models based on the majority vote or the stochastic classifier. However, in the PAC-Bayesian literature, few works propose to minimize a generalization bound to learn a machine learning model, such as, MASEGOSA *et al.* (2020). In the first series of contributions of this thesis, we develop new self-bounding algorithms for three settings. Firstly, we develop a new adversarial robustness setting tailored for the PAC-Bayesian and robustify majority votes after proving generalization bounds. Secondly, we minimize three particular PAC-Bayesian bounds on the majority vote’s risk that was considered difficult to optimize (MASEGOSA *et al.*, 2020). Finally, we introduce a stochastic version of the majority vote, *i.e.*, where the weights are assumed to be sampled from a probability distribution. The stochastic majority vote allows one to derive guarantees on majority-vote-based models. However, upper-bounding the generalization gap of a single classifier with the PAC-Bayesian theory is tedious and generally applicable only to certain classifiers such as the majority vote (see *e.g.*, LANGFORD and SHAWE-TAYLOR, 2002; GERMAIN *et al.*, 2009; LETARTE *et al.*, 2019). In our second series of contributions, we propose to overcome this drawback by considering the notion of disintegrated PAC-Bayesian bounds. Such bounds are able to provide generalization bounds for a single model. By leveraging this framework, we provide new bounds that are easily optimizable and allow us to derive new self-bounding algorithms. In our last contribution, we make use of this framework to develop a general way for incorporating arbitrary complexity measures in generalization bounds.

Outline of this thesis. This thesis is composed of three parts.

Part I is dedicated to the introduction of the field of statistical learning theory and the PAC-Bayesian theory.

- (i) Chapter 1 presents the general setting of this thesis. We introduce the notion of learning and solving a task with a statistical machine learning algorithm. Then, we introduce some machine learning models and some methods to learn them. Afterwards, we recall several classical generalization bounds, notably from VAPNIK and CHERVONENKIS (1974), that assess the quality of the obtained model for the chosen task.
- (ii) In Chapter 2, we mainly recall some results from the PAC-Bayesian framework. After a reminder about the majority vote, we recall different PAC-Bayesian bounds, which will serve as a basis for deriving new results in Part II. We also remind the first *disintegrated* PAC-Bayesian bounds, which are useful when we are interested in one model sampled from the posterior distribution.

Based on the PAC-Bayesian theory, Part II deals with our first series of contributions focusing on the derivation of self-bounding algorithms (FREUND, 1998) that minimize PAC-Bayesian bounds to obtain a majority vote with guarantees on the true risk.

- (i) Chapter 3 stands in the adversarial robustness setting (GOODFELLOW *et al.*, 2015): the goal is to make the majority vote robust to small changes/perturbations in the input. This setting is in contrast with the classical setting in machine learning, where no perturbations are applied to the input. To the best of our knowledge, we are the first to (i) formalize the robustness setting in the PAC-Bayesian framework and (ii) assess the robustness of the majority vote with this framework. We also derive a self-bounding algorithm that minimizes our new generalization bounds.
- (ii) In Chapter 4, we come back to the classical supervised classification setting. We introduce the minimization of PAC-Bayesian bounds on the majority vote risk's surrogate called the C-Bound (recalled in Chapter 2). Unlike the algorithms of the PAC-Bayesian literature, our learning algorithms better consider the voters' correlations.
- (iii) However, the self-bounding algorithms (including ours) do not fully exploit the diversity of the voters in general. Hence, after introducing the *stochastic* majority vote in Chapter 5, we develop a self-bounding learning algorithm to minimize the risk. It allows us to optimize the expected risk directly without requiring the use of a surrogate of the majority vote's risk.

The PAC-Bayesian theory, as considered in Part II, has a major drawback. While the majority vote's generalization abilities can be analyzed through the PAC-Bayesian theory, it becomes more difficult to analyze the generalization of a single voter chosen randomly according to the weights. Thanks to the disintegrated bounds (recalled in Chapter 2), we present two contributions in Part III that introduce self-bounding algorithms to choose a *single* classifier.

- (i) The disintegrated bounds of the literature are difficult to optimize (and to obtain self-bounding algorithms). Hence, in Chapter 6, we (i) derive new disintegrated PAC-Bayesian bounds (easier to optimize) and (ii) provide the first empirical study of self-bounding algorithms using these bounds. We also instantiate the bounds with neural networks and compare them with the PAC-Bayesian considered, *e.g.*, by DZIUGAITE and ROY (2017), ZHOU *et al.* (2019), and PÉREZ-ORTIZ *et al.* (2021).
- (ii) Chapter 7 offers a new viewpoint on generalization bounds by leveraging the disintegrated PAC-Bayesian framework. To the best of our knowledge, Chapter 7

introduces – for the first time – a way to include arbitrary complexity measures in generalization bounds. This work is another step toward the practical use of generalization bounds since the users can now include their complexity measures.

Finally, Part IV presents some perspectives and future works. Note that, for the sake of completeness and clarity, we provide in Appendix all the proofs; we give a hyperlink to the proof for each theorem, corollary, and proposition. Moreover, in order to reproduce the experiments and the figures, we provide the different source codes developed in the context of this thesis at

<https://github.com/paulviallard/PhDThesis>.

Context of this thesis

This thesis was carried out in the Data Intelligence team of the laboratoire Hubert Curien: a joint research unit (UMR 5516) affiliated with the French National Center for Scientific Research (CNRS); the Institut d’Optique Graduate School, and the Université Jean Monnet in Saint-Etienne, France. The french ANR (Agence Nationale de la Recherche) financially supported this thesis through the project APRIORI (A PAC-Bayesian RepresentatIOn LeaRnIng Perspective) ANR-18-CE23-0015. This research project was also the subject of a collaboration with Université Laval, Québec city, Canada.

LIST OF PUBLICATIONS

International Conference

PAUL VIALLARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. Self-bounding Majority Vote Learning Algorithms by the Direct Minimization of a Tight PAC-Bayesian C-Bound. *Machine Learning and Knowledge Discovery in Databases (ECML PKDD)*. (2021a).

PAUL VIALLARD, GUILLAUME VIDOT, AMAURY HABRARD, and EMILIE MORVANT. A PAC-Bayes Analysis of Adversarial Robustness. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021d).

VALENTINA ZANTEDESCHI, PAUL VIALLARD, EMILIE MORVANT, RÉMI EMONET, AMAURY HABRARD, PASCAL GERMAIN, and BENJAMIN GUEDJ. Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021).

International Workshop

PAUL VIALLARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. Interpreting Neural Networks as Majority Votes through the PAC-Bayesian Theory. *NeurIPS 2019 Workshop on Machine Learning with Guarantees*. (2019).

National Conference

PAUL VIALLARD, RÉMI EMONET, PASCAL GERMAIN, AMAURY HABRARD, EMILIE MORVANT, and VALENTINA ZANTEDESCHI. Intérêt des bornes désintégrées pour la généralisation avec des mesures de complexité. *Conférence sur l'Apprentissage automatique (CAp)*. (2022a).

VALENTINA ZANTEDESCHI, PAUL VIALLARD, EMILIE MORVANT, RÉMI EMONET, AMAURY HABRARD, PASCAL GERMAIN, and BENJAMIN GUEDJ. Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound. *Conférence sur l'Apprentissage automatique (CAp)*. (2022).

PAUL VIALLARD, PASCAL GERMAIN, and EMILIE MORVANT. Apprentissage de Vote de Majorité par Minimisation d'une C-Borne PAC-Bayésienne. *Conférence sur l'Apprentissage automatique (CAp)*. (2021b).

PAUL VIALLARD, PASCAL GERMAIN, and EMILIE MORVANT. Dérandomisation des Bornes PAC-Bayésiennes. *Conférence sur l'Apprentissage automatique (CAp)*. (2021c).

GUILLAUME VIDOT, PAUL VIALLARD, and EMILIE MORVANT. Une Analyse PAC-Bayésienne de la Robustesse Adversariale. *Conférence sur l'Apprentissage automatique (CAp)*. (2021).

PAUL VIALLARD, RÉMI EMONET, AMAURY HABRARD, EMILIE MORVANT, and PASCAL GERMAIN. Théorie PAC-Bayésienne pour l'apprentissage en deux étapes de réseaux de neurones. *Conférence sur l'Apprentissage automatique (CAp)*. (2020).

Research Report

PAUL VIALLARD, RÉMI EMONET, AMAURY HABRARD, EMILIE MORVANT, and VALENTINA ZANTEDESCHI. Generalization Bounds with Arbitrary Complexity Measures. *Submitted to ICLR 2023*. (2022b).

PAUL VIALLARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. A General Framework for the Disintegration of PAC-Bayesian Bounds. *Submitted to Machine Learning Journal*. (2022c).

PART I

Background

AN INTRODUCTION TO STATISTICAL LEARNING THEORY

1

Contents

1.1	Introduction	30
1.1.1	Representation of a Task	30
1.1.2	“Solving” a Task	32
1.2	Hypothesis Selection	40
1.2.1	Empirical Risk Minimization	41
1.2.2	Structural Risk Minimization	41
1.3	Generalization Bounds	43
1.3.1	Uniform Convergence Bounds	43
1.3.2	Algorithmic-dependent Generalization Bounds	48
1.3.3	PAC-Bayesian Generalization Bounds	52
1.4	Conclusion and Summary	53

Abstract

This chapter provides an overview of statistical learning theory. We introduce the main concepts and notations used in this thesis focusing on supervised learning. Moreover, we review the main theoretical frameworks allowing one to derive some guarantees on the quality of the learning process.

1.1 Introduction

Statistical learning encompasses a set of statistical methods that aims to automatically solve a task with the help of a computer from data. For example, the task can consist in identifying a digit in an image. In practice, the task is represented by a database (*a.k.a.* dataset or learning sample) containing a finite number of examples. More precisely, one example is composed of an input and its corresponding output (*a.k.a.* label). For the digit recognition problem, one example is made up of an image (the input) with its corresponding digit (the output).

To solve the task, we have to design a model that outputs the correct label given the input. This model is usually a potentially complicated mathematical function. From a human perspective, designing by hand a mathematical function that solves the task can be tedious. Instead, in statistical learning, we aim to find *automatically* such a function with an algorithm. This algorithm is usually called *learning algorithm* since the process of finding a function boils down to *learning* a function that solves the task. The setting of statistical learning and the learning process is illustrated in Figure 1.1; it is defined more formally in the rest of the section.

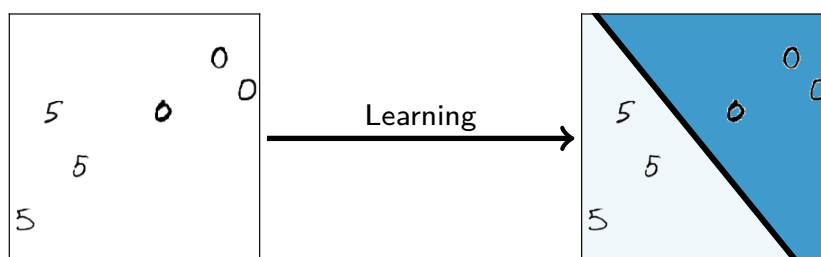


Figure 1.1. Rough representation of statistical learning for the digit recognition task. On the left, each input (the handwritten digits) is represented in 2 dimensions. On the right, we illustrate the notion of learning, *i.e.*, an algorithm finds a function (represented by the black line) predicting the digit 0 for the points belonging to the blue area and the digit 5 for those belonging to the white one. The black line is also called the decision boundary.

1.1.1 Representation of a Task

This thesis studies the classification problem from a supervised learning perspective. More formally, we consider a set of d -dimensional inputs $\mathcal{X} \subseteq \mathbb{R}^d$ and a set of outputs \mathcal{Y} (*a.k.a.* set of labels). In the binary classification setting, when $\mathcal{Y} = \{-1, +1\}$,

the input is classified either into the label -1 or $+1$. In the multi-class classification setting, when $\mathbb{Y} = \{1, 2, \dots, l\}$, the input is classified into $l \geq 2$ different labels. Given an input set \mathbb{X} and a label set \mathbb{Y} , we assume the task to be represented by an *unknown* function $h^* : \mathbb{X} \rightarrow \mathbb{Y}$. The unknown function h^* can be stochastic, meaning that $h^*(\mathbf{x})$ involves randomness, and thus different outputs $y \in \mathbb{Y}$ are probable for the same input $\mathbf{x} \in \mathbb{X}$. Moreover, some inputs $\mathbf{x} \in \mathbb{X}$ are more representative of the task, *i.e.*, some inputs may be more probable than others. To take these probabilities into account, the function h^* is replaced with an *unknown* distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$; the distribution \mathcal{D} represents the probability to sample a given input $\mathbf{x} \in \mathbb{X}$ and the output of the function h^* . More precisely, we assume that each couple $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$ (*a.k.a.* example) is a realization from this *unknown* data distribution \mathcal{D} on the set $\mathbb{X} \times \mathbb{Y}$. Even if the distribution \mathcal{D} is *unknown*, we usually assume that we have access to some examples that we hope to be sufficiently representative and sampled from \mathcal{D} . This set of examples is defined as follows.

Definition 1.1.1 (Learning sample). We define as *learning sample* (or *training set*) a set of m random variables *independent and identically distributed* (*i.i.d.*) following the distribution \mathcal{D} . We have

$$\mathbb{S} \triangleq \bigcup_{i=1}^m \{(\mathbf{x}_i, y_i)\} \triangleq \left\{ (\mathbf{x}_i, y_i) \right\}_{i=1}^m \subseteq (\mathbb{X} \times \mathbb{Y})^m,$$

where $\forall i \in \{1, \dots, m\}, (\mathbf{x}_i, y_i) \sim \mathcal{D}$ and $\mathbb{S} = \left\{ (\mathbf{x}_i, y_i) \right\}_{i=1}^m \sim \mathcal{D}^m$

The notation \mathcal{D}^m denotes the distribution of m examples following \mathcal{D} :

$$\mathcal{D}^m \left(\{(\mathbf{x}_i, y_i)\}_{i=1}^m \right) = \mathcal{D}^m(\mathbb{S}) \triangleq \prod_{i=1}^m \mathcal{D}((\mathbf{x}_i, y_i)).$$

In practice, the set of examples, organized in a dataset, can be created manually by a human and/or collected automatically by a computer. One famous dataset associated with the digit recognition task is the MNIST dataset from LECUN *et al.* (1998). For instance, the set of inputs \mathbb{X} is the set of gray-scale images of size 28×28 , and the set of labels $\mathbb{Y} = \{0, \dots, 9\}$ is the set of possible digits. The learning sample \mathbb{S} is composed of $m = 60,000$ pairs of images (input) and digit (label). A schematic representation of this task is presented in Figure 1.2.

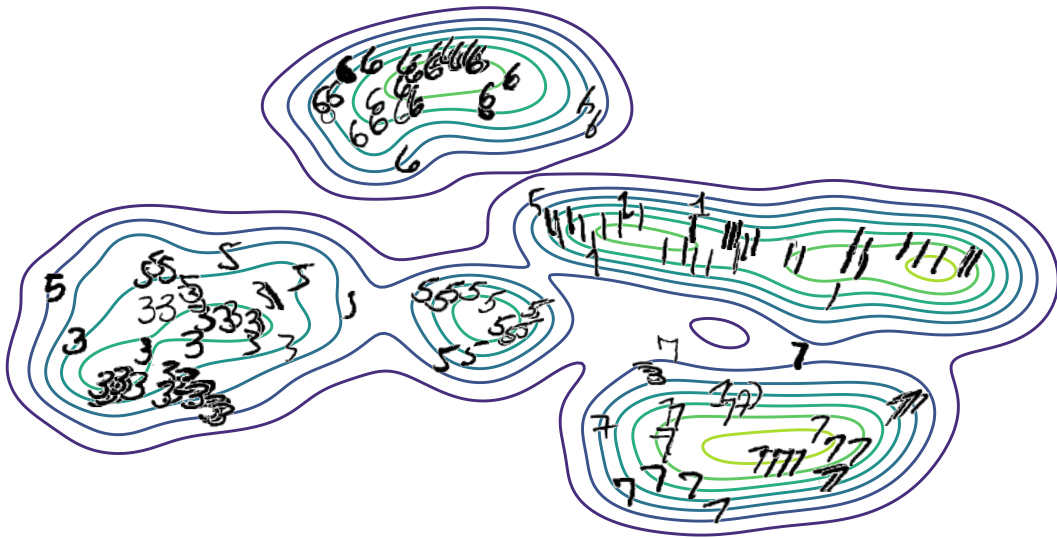


Figure 1.2. Schematic representation of the distribution \mathcal{D} for the digit recognition task. The density of the distribution \mathcal{D} is schematized with the contour lines where the purple resp. green is a low resp. high-density region. The examples (i.e., the handwritten digits) are sampled in high-density areas. A subset of the learning sample $\mathcal{S} \sim \mathcal{D}^m$ is represented.

Deducing a mathematical model that predicts, for all $i \in \{1, \dots, m\}$, the output y_i for the input \mathbf{x}_i is not necessarily trivial. To overcome this problem, statistical learning methods are developed to find *automatically* such a model that aims to solve the task by predicting the examples $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ correctly.

1.1.2 “Solving” a Task

The model found automatically by the statistical learning method is actually a function $h : \mathcal{X} \rightarrow \mathcal{Y}$ (called hypothesis in statistical learning theory) that takes an input $\mathbf{x} \in \mathcal{X}$ and outputs $h(\mathbf{x}) \in \mathcal{Y}$. This hypothesis h is selected by a learning algorithm among different candidates in a *hypothesis set* \mathbb{H} potentially infinite. As an illustration, we recall four types of hypotheses used in practice.

1.1.2.1 Examples of Hypotheses

The four types of hypotheses we recall are the decision stump (IBA and LANGLEY, 1992), the decision tree (BREIMAN *et al.*, 1984), the linear classifier, and the neural network. More precisely, these hypotheses are called *classifier* since they classify an input $\mathbf{x} \in \mathcal{X}$, i.e., assign a label to an input.

Decision Stump The decision stump (IBA and LANGLEY, 1992) is a classifier composed of one decision rule of the form “Is $x_i \leq \tau$?”, where x_i is the i -th component of the vector \mathbf{x} and $\tau \in \mathbb{R}$ is a threshold. More precisely, this *decision rule* (implemented as an if-condition) assigns a label to the input \mathbf{x} . For example, we present in Algorithm 1.1 a decision stump that classifies the input \mathbf{x} as $+1$ when $x_1 \leq 0.54$ and -1 otherwise. Its graphical representation and the associated decision are given in Figure 1.3.

Algorithm 1.1 Example of Decision Stump

Given: An input $\mathbf{x} \in \mathbb{R}^d$
if $x_1 \leq 0.54$ **then**
 return $+1$
else
 return -1

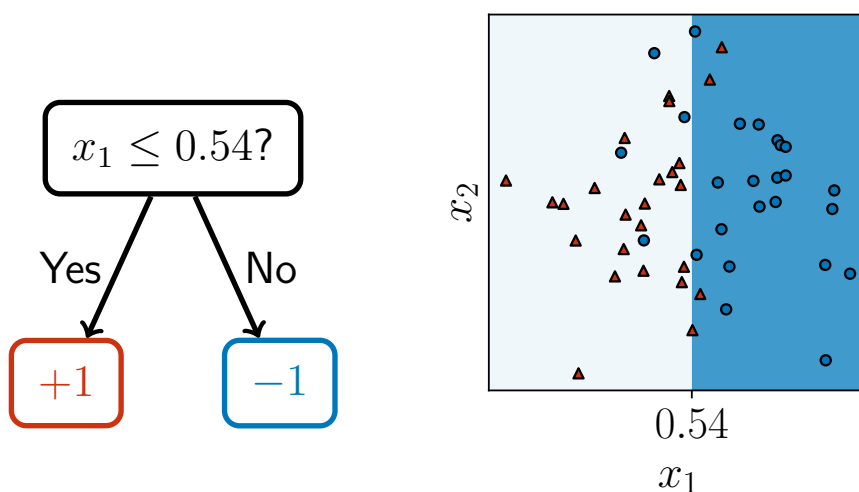


Figure 1.3. Example of decision stump with the decision rule “Is $x_1 \leq 0.54$?”. On the left, the decision of Algorithm 1.1 is shown with its binary tree representation. The right plot shows the decision boundary (in white and blue) of the given decision stump. Moreover, the red triangles resp. the blue dots are the examples in the learning sample \mathcal{S} with the label $+1$ resp. -1 .

To classify the input $\mathbf{x} \in \mathbb{X}$, the decision stump returns a label (by executing its associated algorithm). This algorithm is composed of only one if-condition that can be interpreted as a binary tree. This simple classifier can be complexified by adding more if-conditions to fit the data better. Indeed, a decision tree can be constructed

by considering nested if-conditions. Hence, the decision stump¹ is a particular case of a decision tree (BREIMAN *et al.*, 1984).

Decision Trees The *decision tree* (BREIMAN *et al.*, 1984) is a classifier that can be seen as a succession of *decision rules* to classify the label of the input $\mathbf{x} \in \mathbb{R}^d$. Algorithm 1.2 is an example of algorithm associated to the decision tree of Figure 1.4.

Algorithm 1.2 Example of Decision Tree

```

Given: An input  $\mathbf{x} \in \mathbb{R}^d$ 
if  $x_1 \leq 0.59$  then
  if  $x_2 \leq 0.63$  then
    return +1
  else
    return -1
else
  return -1

```

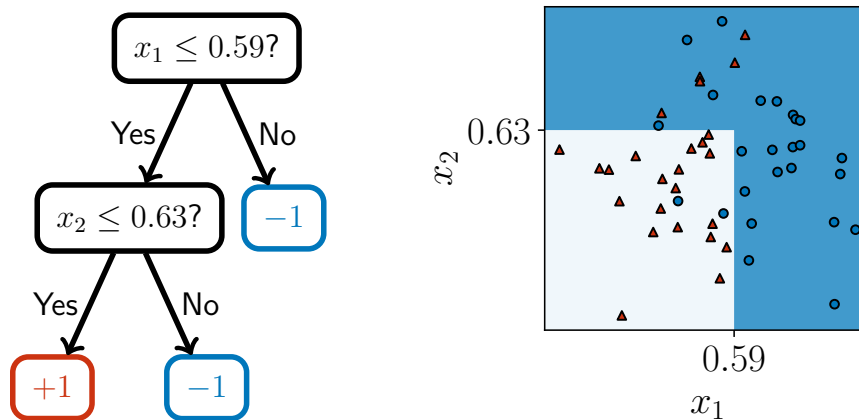


Figure 1.4. Example of decision tree composed of two decision rules “Is $x_1 \leq 0.59$?” and “Is $x_2 \leq 0.63$?”. The binary tree on the left is the graphical representation of the decision stump. The right plot is the stump’s decision boundary on the 2-dimensional data in \mathbb{S} (with the red triangles and the blue dots).

Numerous algorithms have been developed to infer decision trees (see *e.g.*, BREIMAN *et al.*, 1984; QUINLAN, 1986, 1993). One of the most popular is the CART algorithm (BREIMAN *et al.*, 1984): it is, *e.g.*, implemented in the well-known library

¹IBA and LANGLEY (1992) introduced the word “stump” for a one-level decision tree.

scikit-learn (PEDREGOSA *et al.*, 2011). This algorithm is summarized as follows. The algorithm tests the decision rules “Is $x_i \leq \tau$?” with different components x_i and thresholds τ . Given a new decision rule, the examples $(\mathbf{x}, y) \in \mathcal{S}$ are split into two groups: the ones that respect the rules and those that do not. Then, the algorithm selects the best decision rule according to a criterion and the two new groups. Finally, this step is repeated recursively to obtain the entire tree. In Part II, we learn a *majority vote* classifier which is a convex combination of the outputs of decision trees or decision stumps.

Linear classifier One of the most simple type of classifiers is the linear classifier. Given a function $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, this type of classifier is defined as

$$h(\mathbf{x}) = \sigma\left(\sum_{i=1}^d w_i x_i + b\right),$$

where $\mathbf{w} \in \mathbb{R}^d$ is the weight vector and $b \in \mathbb{R}$ is the bias that both need to be learned. The function $\sigma(\cdot)$ is called *activation* function notably in the context of neural networks (that we precise in the following). For instance, the Perceptron (MCCULLOCH and PITTS, 1943; ROSENBLATT, 1958) algorithm returns a linear classifier with the activation function $\sigma(\cdot)$ called threshold function defined as $\sigma(x) = 1$ if $x \geq 0$ and $\sigma(x) = 0$ otherwise. Moreover, the well-known Support Vector Machine (SVM) algorithm introduced by BOSER *et al.* (1992) and CORTES and VAPNIK (1995) learns a linear classifier with the sign as activation function, *i.e.*, $\sigma(x) = \text{sign}(x) = +1$ if $a \geq 1$ and -1 otherwise. We provide two 2-dimensional cases when the activation function is the sign function in Figure 1.6.

Neural networks. The Perceptron was abandoned when MINSKY and PAPERT (1972) show that this model cannot learn simple (boolean) functions like the exclusive or (see Figure 1.5). Hopefully, these restrictions are avoided when generalizing the Perceptron, and they regain interest in the 80s. Nevertheless, in the 2010s, these models have become popular when the neural network AlexNet (KRIZHEVSKY *et al.*, 2012) won the computer vision challenge “ImageNet Large Scale Visual Recognition Challenge”. The increase in popularity of such models has been helped by the existence of many programming frameworks such as Torch7 (COLLOBERT *et al.*, 2011), Tensorflow (ABADI *et al.*, 2015), Theano (AL-RFOU *et al.*, 2016), JAX (BRADBURY *et al.*, 2018; FROSTIG *et al.*, 2018), or PyTorch (PASZKE *et al.*, 2019). To learn more details about neural networks, we refer the reader to GOODFELLOW *et al.* (2016) for an extensive introduction. Actually, a neural network can be seen as a succession of linear classifiers: it is a composition of L linear classifiers $\mathbf{h}^{(1)}(), \dots, \mathbf{h}^{(L)}()$ called layer or module, *i.e.*,

$$h = \mathbf{h}^{(L)} \circ \dots \circ \mathbf{h}^{(2)} \circ \mathbf{h}^{(1)},$$

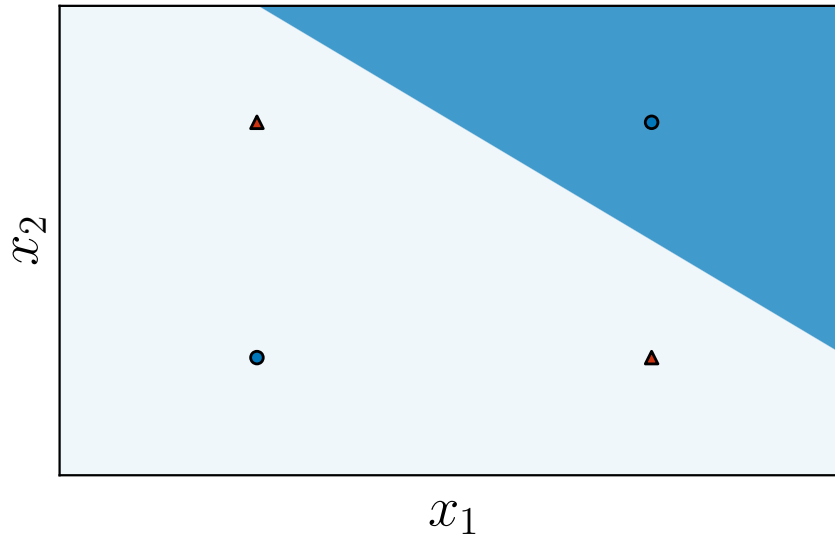


Figure 1.5. Plot of the decision boundary of a linear classifier learned on the learning sample $\mathbb{S} = \{([0, 0]^\top, -1), ([0, 1]^\top, +1), ([1, 1]^\top, -1), ([1, 0]^\top, +1)\}$ that represents the exclusive or function. As proved by MINSKY and PAPERT (1972), this model cannot predict correctly all the examples representing such a function.

where $f \circ g$ is the composition of the function $f()$ with $g()$. More precisely, given an activation function $\sigma^{(i)} : \mathbb{R}^{d^{(i)}} \rightarrow \mathbb{R}^{d^{(i)}}$, the i -th layer $\mathbf{h}^{(i)} : \mathbb{R}^{d^{(i-1)}} \rightarrow \mathbb{R}^{d^{(i)}}$ of the network is defined by

$$\mathbf{h}^{(i)}(\mathbf{x}) = \sigma^{(i)}(\mathbf{W}\mathbf{x} + \mathbf{b}),$$

where the matrix $\mathbf{W} \in \mathbb{R}^{d^{(i)} \times d^{(i-1)}}$ and the vector $\mathbf{b} \in \mathbb{R}^{d^{(i)}}$ are respectively the weights and the bias parameterizing the classifier that need to be learned. Note that, to respect the definition of the hypothesis $h : \mathbb{R}^d \rightarrow \mathbb{R}$, we have $d^{(0)} = d$ and $d^{(L)} = 1$.

An example of neural network is given in Figure 1.7. We see in this example that the succession of linear classifiers offers better expressiveness, *i.e.*, its decision boundary is not restricted to lines. Interestingly, these models produce new features of the original data in each layer. Indeed, in Figure 1.8, the i -th layer transforms all inputs $\mathbf{x} \in \mathbb{X}$ into $\mathbf{h}^{(i)}(\dots \mathbf{h}^{(1)}(\mathbf{x}))$ and these new inputs are classified with $\mathbf{h}^{(L)}(\dots \mathbf{h}^{(i)}(\mathbf{x}))$. Hence, all these transformations can be interpreted as a new representation of these inputs. In the following, we denote by $\mathbf{w} \in \mathbb{R}^D$ the vector of the network's weights and biases concatenated all together; we have thus D weights/bias in the networks. Besides, when D is large compared to the number of data m , we say in this case that the model is over-parametrized, and when L is large, we consider that the network is “deep”². This

²The term deep learning comes from the fact that the number of layers L is large.

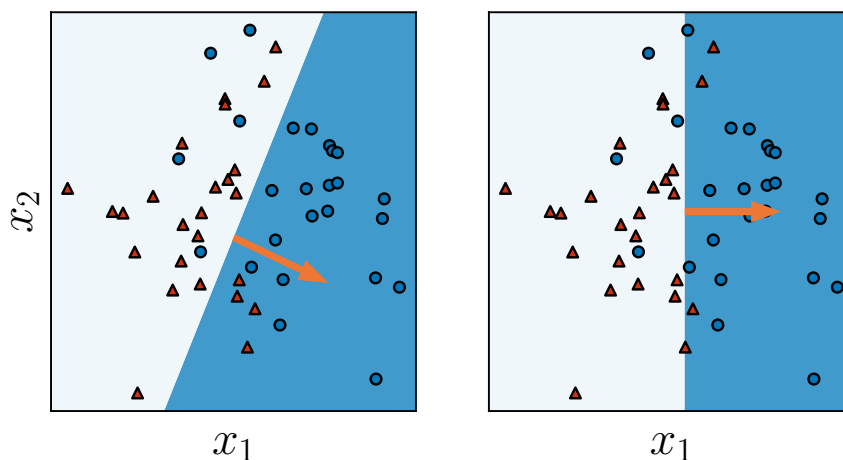


Figure 1.6. The decision of two linear classifiers when the activation function is the sign. On the left plot, the weights are $\mathbf{w} = [3.2, -1.4]^\top$ and the bias is $b = -1$. On the right plot, the weights are $\mathbf{w} = [1, 0]^\top$ and the bias is $b = -0.54$; this linear classifier is equivalent to the decision stump in Figure 1.3.

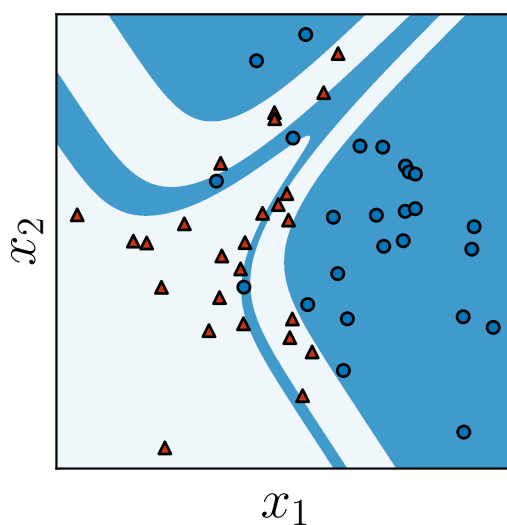


Figure 1.7. The decision boundary of a neural network classifier of a with 5 layers and the dimensions $d^{(1)}=d^{(2)}=d^{(3)}=d^{(4)}=2$ and $d^{(5)}=1$. The activation functions $\sigma^{(1)}(), \sigma^{(2)}(), \sigma^{(3)}(), \sigma^{(4)}()$ are tanh applied element-wise and $\sigma^{(5)}()$ is the threshold function.

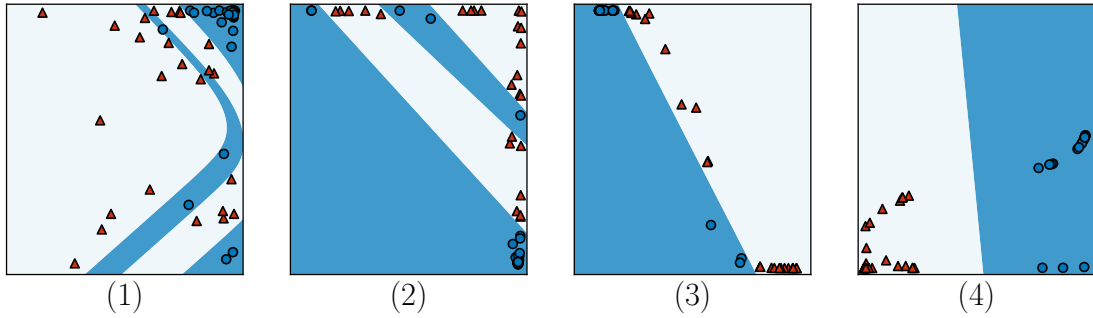


Figure 1.8. Each plot corresponds to the representation of the i -layer (where $i = 1$ on the left plot and $i = 4$ on the right plot). On each plot, the original examples $(\mathbf{x}, y) \in \mathcal{S}$ (with the red triangles and the blue dots) are transformed and plotted by the i -th layer $(\mathbf{h}^{(1)} \circ \dots \circ \mathbf{h}^{(i)})(\mathbf{x})$. Moreover, the decision boundary of the classifier $\mathbf{h}^{(i)} \circ \dots \circ \mathbf{h}^{(L)}$ is given.

model is studied in Part III when the neural network is over-parametrized and “deep”.

1.1.2.2 Loss and Risk

We need to assess to which extent the learned hypothesis $h \in \mathbb{H}$ predicts correctly an example. This kind of measure is often called a *loss function*, defined as follows.

Definition 1.1.2 (Loss function). A *loss function* is a function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ that, given a hypothesis $h \in \mathbb{H}$, evaluates the quality of the prediction $h(\mathbf{x})$ compared to the true label $y \in \mathcal{Y}$. The lower the loss, the better the quality of the hypothesis.

For a classification task, the most natural loss function is the 01-loss defined by

$$\forall h \in \mathbb{H}, \quad \forall (\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}, \quad \ell^{01}(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y] \triangleq \begin{cases} 1 & \text{if } h(\mathbf{x}) \neq y \\ 0 & \text{otherwise} \end{cases}.$$

The 01-loss returns 1 when the hypothesis $h \in \mathbb{H}$ misclassifies the example $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ and 0 otherwise. However, since many learning methods are based on the gradient $\frac{\partial \ell}{\partial h}(h, (\mathbf{x}, y))$ to choose a hypothesis h that lowers the loss function, the 01-loss $\ell^{01}()$ is not practical. Its gradient $\frac{\partial \ell^{01}}{\partial h}(h, (\mathbf{x}, y)) = 0$ for all examples $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$, which does not indicate a descent direction (that lowers the loss). To address this issue, in practice, one uses relaxation of the 01-loss. In particular, some relaxed losses rely on a notion of *margin*: in binary classification the margin is defined as $m_h(\mathbf{x}, y) = yh(\mathbf{x})$ for all $h : \mathcal{X} \rightarrow [-1, +1]$. When the margin is positive, the hypothesis makes a correct

prediction. In other words, the margin is defined such that we have $\ell^{01}(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y] = \mathbb{I}[m_h(\mathbf{x}, y) \leq 0]$. Thanks to the margin and MARKOV's inequality (Theorems A.2.1 and A.3.1), one can prove two upper bounds of the 01-loss. These two upper bounds are defined by

$$\ell^{1st}(h, (\mathbf{x}, y)) = 1 - m_h(\mathbf{x}, y) \quad (\text{LANGFORD and SHAWE-TAYLOR, 2002}),$$

and $\ell^{2nd}(h, (\mathbf{x}, y)) = \left[1 - m_h(\mathbf{x}, y)\right]^2 \quad (\text{MASEGOSA et al., 2020}).$

We plot in Figure 1.9 the 01-loss and the two upper bounds that we use in Chapters 3 and 4 to derive learning algorithms. We introduce these upper bounds with more details in Chapter 2.

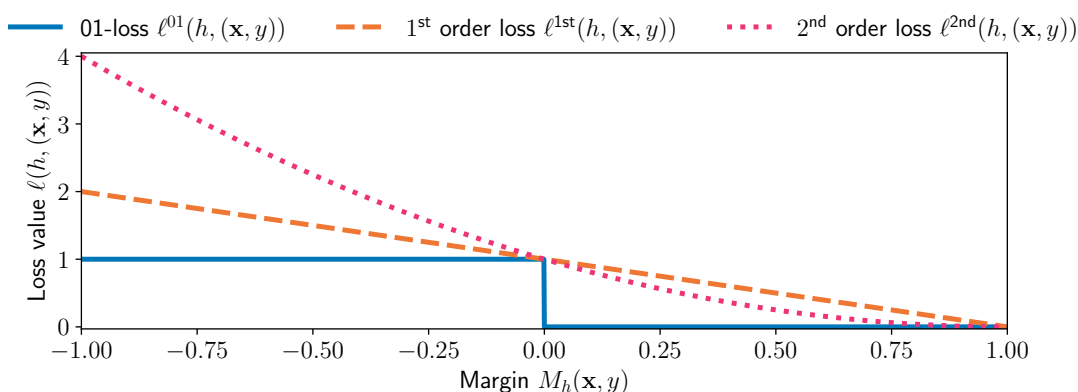


Figure 1.9. Plot of the losses ($\ell^{01}()$, $\ell^{1st}()$, and $\ell^{2nd}()$) where the x-axis represents the value of the margin $m_h(\mathbf{x}, y)$ and the y-axis represents the value of the losses $\ell(h, (\mathbf{x}, y))$ for a given example (\mathbf{x}, y) and hypothesis h .

While the loss $\ell(h, (\mathbf{x}, y))$ is computed for an example (\mathbf{x}, y) and a hypothesis h , we are usually interested in the loss computed over all the examples $(\mathbf{x}, y) \in \mathcal{S}$. The value of the loss averaged over the learning sample \mathcal{S} is called *empirical risk* and represents to what extent the hypothesis h predicts correctly all the examples $(\mathbf{x}, y) \in \mathcal{S}$. The empirical risk is defined as follows.

Definition 1.1.3 (Empirical Risk). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any learning sample $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m \sim \mathcal{D}^m$, we define the empirical risk as

$$\mathbf{R}_{\mathcal{S}}^{\ell}(h) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \ell(h, (\mathbf{x}, y)) = \frac{1}{m} \sum_{i=1}^m \ell(h, (\mathbf{x}_i, y_i)),$$

where its associated uniform distribution \mathcal{S} on \mathbb{S} can be defined as $\mathcal{S}((\mathbf{x}_i, y_i)) \triangleq \frac{1}{m}$ for all $i \in \{1, \dots, m\}$. Moreover, we define the empirical risk with the 01-loss by

$$R_{\mathcal{S}}(h) \triangleq R_{\mathcal{S}}^{\ell^{01}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(\mathbf{x}_i) \neq y_i].$$

Since the distribution \mathcal{D} is unknown, the performance of the hypothesis $h \in \mathbb{H}$ can be computed by means of the learning sample \mathbb{S} (through the empirical risk). However, we are interested in the performance of h on the task, *i.e.*, over all the examples. The performance is thus defined through the true risk $R_{\mathcal{D}}^{\ell}(h)$ defined as follows.

Definition 1.1.4 (True Risk). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any loss function $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, we define the true risk as

$$R_{\mathcal{D}}^{\ell}(h) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \ell(h, (\mathbf{x}, y)).$$

Moreover, we define the true risk with the 01-loss by

$$R_{\mathcal{D}}(h) \triangleq R_{\mathcal{D}}^{\ell^{01}}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[h(\mathbf{x}) \neq y] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [h(\mathbf{x}) \neq y].$$

The true risk can be interpreted as the expected loss over all examples $(\mathbf{x}, y) \sim \mathcal{D}$. Actually, when the hypothesis h does not depend on \mathbb{S} , the true risk of h is the expected empirical risk. In other words, the empirical risk can be seen as an unbiased estimator of the true risk, *i.e.*, we have

$$\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} R_{\mathcal{S}}^{\ell}(h) = \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \frac{1}{m} \sum_{i=1}^m \ell(h, (\mathbf{x}_i, y_i)) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{(\mathbf{x}_i, y_i) \sim \mathcal{D}} \ell(h, (\mathbf{x}_i, y_i)) = R_{\mathcal{D}}^{\ell}(h).$$

Besides, selecting a hypothesis with a low true risk boils down to finding a hypothesis with low empirical risk on average. However, since we only have access to one learning sample \mathbb{S} , the hypothesis selection cannot be performed on the distribution \mathcal{D} but rather on the learning sample \mathbb{S} .

1.2 Hypothesis Selection

In order to obtain a hypothesis $h \in \mathbb{H}$ that solves the task, the hypothesis can be selected such that the empirical risk $R_{\mathcal{S}}^{\ell}(h)$ is low. We recall two approaches from statistical learning: Empirical Risk Minimization (ERM) and Structural Risk Minimization (SRM).

1.2.1 Empirical Risk Minimization

Given the hypothesis set \mathbb{H} , a common approach in statistical learning is to minimize the empirical risk $R_S^\ell(h)$ with respect to the hypothesis $h \in \mathbb{H}$. This approach – known as *Empirical Risk Minimization* (ERM) – has been pioneered by VAPNIK and CHERVONENKIS (1968, 1971, 1974) in the 70s (see VAPNIK (1998) for an introduction). Formally, given a learning sample $\mathcal{S} \sim \mathcal{D}^m$, we select the hypothesis such that the empirical risk $R_S^\ell(h)$ is minimum; this approach is summarized in Algorithm 1.3.

Algorithm 1.3 Empirical Risk Minimization

Given: Learning sample \mathcal{S} , Hypothesis set \mathbb{H} , Loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$

$$h = \operatorname{argmin}_{h' \in \mathbb{H}} R_S^\ell(h')$$

return hypothesis h

However, if the empirical risk $R_S^\ell(h)$ of a hypothesis $h \in \mathbb{H}$ is approximately 0, then h potentially overfits the data; see Figure 1.10. Roughly speaking, in the case of overfitting, we interpret that $h \in \mathbb{H}$ has learned by heart the examples in the learning sample. Such a phenomenon can arise when the hypothesis h is complex; the complexity is manifested differently for each type of hypothesis. To measure such complexity, a real-valued function can be defined $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$. Given a learning sample \mathcal{S} and a complexity function, the hypothesis $h \in \mathbb{H}$ is considered complex when its associated complexity measure $\mu(h, \mathcal{S})$ is large. In Chapter 7, we derive generalization bounds that integrate a user-specified complexity function.

1.2.2 Structural Risk Minimization

To avoid overfitting, we can select a hypothesis h with a small complexity measure $\mu(h, \mathcal{S})$. The *Structural Risk Minimization* approach – introduced by VAPNIK and CHERVONENKIS (1974) – minimizes a trade-off between the empirical risk $R_S^\ell(h)$ and the complexity measure $\mu(h, \mathcal{S})$. The complexity measures $\mu(\cdot)$ used in this approach is constant over the hypothesis $h \in \mathbb{H}$ and the learning sample $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$. In other words, if we denote by abuse of notation $\mu(\mathbb{H})$ the complexity measure associated to the hypothesis set \mathbb{H} , we have $\forall h \in \mathbb{H}, \forall \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m, \mu(h, \mathcal{S}) = \mu(\mathbb{H})$. However, the complexity measure $\mu(\mathbb{H})$ may not be representative of overfitting for a single hypothesis since $\mu(\mathbb{H})$ is constant. Hence, the hypothesis set \mathbb{H} is structured into countable and nested hypothesis subsets $\mathbb{H}_1 \subseteq \mathbb{H}_2 \subseteq \dots \subseteq \mathbb{H}$. By doing so, the complexity measure may be increasing when the hypothesis set grows: we can have $\mu(\mathbb{H}_i) \leq \mu(\mathbb{H}_{i+1})$ when $\mathbb{H}_i \subseteq \mathbb{H}_{i+1}$. Then, thanks to these nested complexity measure, one can find a hypothesis h belonging to a set \mathbb{H}_i with $i \in \mathbb{N}_*$ that minimizes the trade-off between the empirical risk $R_S^\ell(h)$ and the complexity measure $\mu(\mathbb{H}_i)$ of the set

\mathbb{H}_i . The minimization of the trade-off performed by the *Structural Risk Minimization* approach is summarized in Algorithm 1.4.

Algorithm 1.4 Structural Risk Minimization

Given: Learning sample \mathcal{S} , Hypothesis set \mathbb{H}

$$h = \underset{i \in \mathbb{N}_*, h' \in \mathbb{H}_i}{\operatorname{argmin}} \left[R_{\mathcal{S}}^{\ell}(h') + \mu(\mathbb{H}_i) \right]$$

return hypothesis h

If the complexity measure $\mu(\mathbb{H}')$ of the subset $\mathbb{H}' \subseteq \mathbb{H}$ approximates well the difference $R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$, SRM can be used to obtain a hypothesis $h \in \mathbb{H}'$ that *generalizes* well. In other words, a hypothesis $h \in \mathbb{H}'$ that *generalizes* has a small true risk $R_{\mathcal{D}}^{\ell}(h)$ and empirical risk $R_{\mathcal{S}}^{\ell}(h)$. This situation is represented in Figure 1.10. In other words, when h generalizes, the *generalization gap*³ $R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$ is close to 0 (that is $R_{\mathcal{D}}^{\ell}(h) \approx R_{\mathcal{S}}^{\ell}(h)$). However, since the difference $R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$ is not computable because of the true risk $R_{\mathcal{D}}^{\ell}(h)$, we have to bound it.

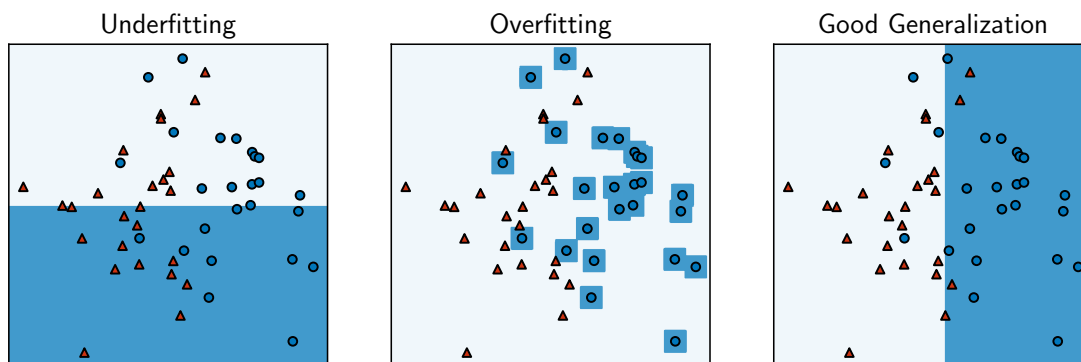


Figure 1.10. Representation of the three situations (one per scatter plot) that can arise in statistical learning. The left figure corresponds to a hypothesis h that underfits the data: both the empirical risk $R_{\mathcal{S}}(h)$ and the true risk $R_{\mathcal{D}}(h)$ are high. In the middle figure, the hypothesis h overfits the data: the empirical risk $R_{\mathcal{S}}(h) = 0$ but the true risk $R_{\mathcal{D}}(h)$ is high. The right figure represents the perfect case: the difference between the true risk $R_{\mathcal{D}}(h)$ and the empirical risk $R_{\mathcal{S}}(h)$ is small as well as $R_{\mathcal{S}}(h)$.

³We consider sometimes the generalization gap $\left| R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h) \right|$ instead, but, we are mostly interested in upper bounding the true risk $R_{\mathcal{D}}^{\ell}(h)$.

1.3 Generalization Bounds

To upper-bound the generalization gap $R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$, we can provide a bound on the true risk $R_{\mathcal{D}}^{\ell}(h)$ with high probability over the random choice of $\mathcal{S} \sim \mathcal{D}^m$. Such a bound is a *Probably Approximately Correct* (PAC) generalization bound (VALIANT, 1984) and can be defined as follows.

Definition 1.3.1 (PAC Generalization Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any data-dependent hypothesis set $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m}$ where $h_{\mathcal{S}}$ is a hypothesis dependent on the learning sample $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$, a PAC generalization bound is defined by

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[R_{\mathcal{D}}^{\ell}(h_{\mathcal{S}}) \leq \Phi \right] \geq 1 - \delta.$$

With high probability (with probability at least $1 - \delta$), the true risk of the hypothesis $h_{\mathcal{S}}$ is *Approximately Correct*, i.e., upper-bounded by Φ . Based on this definition, VALIANT (1984) has brought a computational framework: the PAC-Learnability.⁴ A data-dependent hypothesis set $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m}$ is PAC-Learnable if this high probability bound holds when the number of examples m is polynomial in $1/\delta$ and $1/\Phi$. In practice, to obtain such an upper bound Φ , we make use of *concentration inequalities* (see BOUCHERON *et al.* (2013) for an extensive introduction on the subject). These inequalities allow us to bound an expectation (the true risk in our case) with its empirical counterpart (the empirical risk). Depending on the concentration inequalities, one can obtain different frameworks, and so, different upper bounds Φ .

1.3.1 Uniform Convergence Bounds

The first framework introduced in the literature to obtain PAC generalization bounds is referred to as *uniform convergence* bounds (VAPNIK and CHERVONENKIS, 1968, 1971). Given a hypothesis set \mathbb{H} (not necessarily data-dependent), a uniform convergence bound holds for all hypotheses of \mathbb{H} . By doing so, the true risk $R_{\mathcal{D}}^{\ell}(h_{\mathcal{S}})$ of a data-dependent hypothesis $h_{\mathcal{S}} \in \mathbb{H}$ is upper-bounded to obtain a PAC generalization bound. This type of bounds takes the following form.

Definition 1.3.2 (Uniform Convergence Bound). Let $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ a generalization gap. A uniform convergence bound is defined such that if for any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis

⁴VALIANT won the Turing prize in 2010 for the definition of PAC-Learnability.

set \mathbb{H} , there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$, such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\sup_{h \in \mathbb{H}} \phi(\mathbf{R}_{\mathcal{D}}^\ell(h), \mathbf{R}_{\mathcal{S}}^\ell(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta, \quad (1.1)$$

where usually $\phi(\mathbf{R}_{\mathcal{D}}^\ell(h), \mathbf{R}_{\mathcal{S}}^\ell(h)) = \mathbf{R}_{\mathcal{D}}^\ell(h) - \mathbf{R}_{\mathcal{S}}^\ell(h)$.

A uniform convergence bound consists in bounding the supremum of the generalization gap over the hypotheses $h \in \mathbb{H}$ with the upper bound $\Phi_u(\delta)$ (with probability at least $1 - \delta$ over the learning sample $\mathcal{S} \sim \mathcal{D}^m$). In other words, the inequality boils down to bound the generalization gap of all hypotheses $h \in \mathbb{H}$ by an upper bound $\Phi_u(\delta)$ on the largest generalization gap; the hypothesis associated to the largest generalization gap is considered as the “worst” hypothesis in \mathbb{H} . As we will see in the following, the upper bound $\Phi_u(\delta)$ depends also on the number of examples m in \mathcal{S} ; we have simplified $\Phi_u(\delta)$ for readability. Hence, if $\lim_{m \rightarrow +\infty} \Phi_u(\delta) = 0$, with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, the empirical risk *converges uniformly* on \mathbb{H} to the true risk $\mathbf{R}_{\mathcal{D}}^\ell(h)$. Because of the uniform convergence, as we see in the two examples below, the upper bound $\Phi_u(\delta)$ depends on a *complexity* $\mu(\mathbb{H})$ that measures, in some sense, the performance of the worst hypothesis in \mathbb{H} .

1.3.1.1 Uniform Convergence Bounds for Finite Hypothesis Sets

The simplest complexity measure appearing in generalization bounds arises when considering a finite set \mathbb{H} of hypotheses. We recall an instantiation of uniform convergence bounds (Definition 1.3.2) in the following theorem (see, e.g., MOHRI *et al.* (2012) for more details).

Theorem 1.3.1 (Generalization Bound for Finite \mathbb{H}). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any finite hypothesis set \mathbb{H} ($\text{card}(\mathbb{H}) < +\infty$), for any loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\sup_{h \in \mathbb{H}} \left[\mathbf{R}_{\mathcal{D}}^\ell(h) - \mathbf{R}_{\mathcal{S}}^\ell(h) \right] \leq \underbrace{\sqrt{\frac{\ln \text{card}(\mathbb{H}) + \ln \frac{1}{\delta}}{2m}}}_{\Phi_u(\delta)} \right] \geq 1 - \delta,$$

where $\text{card}(\mathbb{H})$ is the cardinal of the set \mathbb{H} . Equivalently, with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, we have

$$\forall h \in \mathbb{H}, \quad \mathbf{R}_{\mathcal{D}}^\ell(h) \leq \mathbf{R}_{\mathcal{S}}^\ell(h) + \sqrt{\frac{\ln \text{card}(\mathbb{H}) + \ln \frac{1}{\delta}}{2m}}.$$

1.3. Generalization Bounds

In this case, the complexity measure for a finite hypothesis set is defined by $\mu(\mathbb{H}) = \text{card}(\mathbb{H})$. According to this complexity measure, the more hypotheses in \mathbb{H} , the more complex the hypothesis set \mathbb{H} is. The computation of this complexity can be fast, making this generalization bound easy to evaluate and to obtain an upper bound on the generalization gap $R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$. However, the bound does not converge when the cardinality of the set is large, e.g., when $\text{card}(\mathbb{H}) \geq e^m$ for $m \in \mathbb{N}$, the upper bound of Theorem 1.3.1 is $\sqrt{\frac{1}{2}} \leq \sqrt{\frac{1}{2m}(\ln \text{card}(\mathbb{H}) + \ln \frac{1}{\delta})}$. Hence, there is a need to develop generalization bounds for hypothesis sets \mathbb{H} with large $\text{card}(\mathbb{H})$ or infinite hypothesis sets. Indeed, these types of set are usually considered in practice, e.g., the set of linear classifiers.

1.3.1.2 VC-Dimension-based Generalization Bounds

When the loss is the 01-loss, a bound for infinite hypothesis sets is proposed based on a complexity measure $\mu(\cdot)$ called the VAPNIK-CHEVONENKIS (VC)-Dimension (VAPNIK and CHEVONENKIS, 1968, 1971).

Definition 1.3.3 (VAPNIK-CHEVONENKIS (VC) Dimension). Given a hypothesis set \mathbb{H} with hypotheses $h : \mathcal{X} \rightarrow \{-1, +1\}$ (for binary classification), the VC-dimension of the set \mathbb{H} is defined as

$$\text{vc}(\mathbb{H}) \triangleq \max \{m : \forall \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m, \exists h \in \mathbb{H} \text{ s.t. } R_{\mathcal{S}}(h) = 0\},$$

where $R_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^m I[h(\mathbf{x}) \neq y]$.

Put into words, the VC-Dimension $\mu(\mathbb{H}) = \text{vc}(\mathbb{H})$ of a hypothesis set \mathbb{H} is the maximum number of examples that a hypothesis h from \mathbb{H} can perfectly fit in binary classification. Note that if the hypothesis set \mathbb{H} is infinite, its associated VC-dimension $\text{vc}(\mathbb{H})$ can be finite. For example, the VC-Dimension of the d -dimensional linear classifiers is $d + 1$ (see MOHRI *et al.* (2012, Example 3.12) for a proof) even if the set of linear classifiers \mathbb{H} is infinite; we illustrate in Figure 1.11 the 2-dimensional case. Thanks to this complexity measure, we can prove the following generalization bound (MOHRI *et al.*, 2012, Theorem 3.17, Corollaries 3.18 and 3.19).

Theorem 1.3.2 (VC-Dimension-based Generalization Bounds). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any set \mathbb{H} with hypotheses $h : \mathcal{X} \rightarrow \{-1, +1\}$ and VC-Dimension $\text{vc}(\mathbb{H})$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\sup_{h \in \mathbb{H}} [R_{\mathcal{D}}(h) - R_{\mathcal{S}}(h)] \leq \underbrace{\sqrt{\frac{2\text{vc}(\mathbb{H}) \left(1 + \ln \frac{m}{\text{vc}(\mathbb{H})}\right)}{m}}}_{\Phi_{\mathbf{u}}(\delta)} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

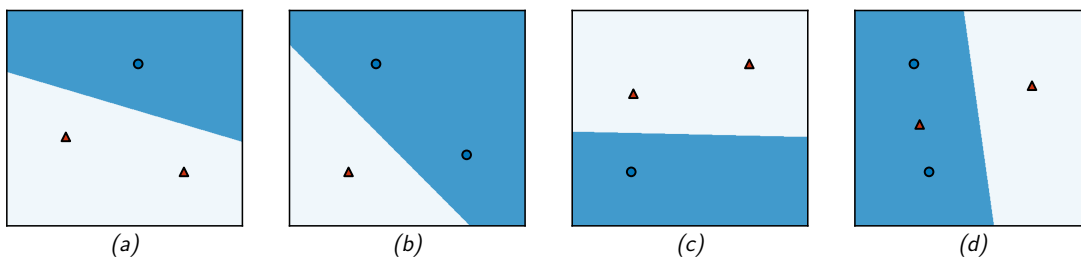


Figure 1.11. Illustration of the VC-dimension for linear classifiers in 2-dimension. When $m=3$, a linear classifier can always perfectly fit the data, i.e., $R_S(h) = 0$ (as shown in the plot (a), (b), and (c)). However, when $m=4$, the empirical risk $R_S(h) \geq 0$ for all linear classifiers (as illustrated in the plot (d)). Hence, in this case, the VC-Dimension is $\text{vc}(\mathbb{H}) = 3$.

Equivalently, with probability at least $1-\delta$ over $\mathcal{S} \sim \mathcal{D}^m$, we have

$$\forall h \in \mathbb{H}, \quad R_{\mathcal{D}}(h) \leq R_S(h) + \sqrt{\frac{2\text{vc}(\mathbb{H}) \left(1 + \ln \frac{m}{\text{vc}(\mathbb{H})}\right)}{m}} + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}.$$

Hence, if we know the VC-dimension of the hypothesis set \mathbb{H} of interest, the bound becomes easily computable to obtain an upper bound of the generalization gap. Furthermore, it is possible to prove that the ERM algorithm is consistent (VAPNIK, 1998) thanks to this bound (see Proposition 4.1 MOHRI *et al.*, 2012). When ERM is consistent, (a) the empirical risk $R_S(h)$ of the classifier h obtained by ERM converges in probability to $\inf_{h' \in \mathbb{H}} R_{\mathcal{D}}(h')$, and (b) its true risk $R_S(h)$ converges in probability to $\inf_{h' \in \mathbb{H}} R_{\mathcal{D}}(h')$ as well. Nevertheless, one drawback of the VC-Dimension is that it is only for binary classifiers and the 01-loss; there are several extensions for the multi-class setting, we introduce one of them: the Rademacher complexity (BARTLETT and MENDELSON, 2002).

1.3.1.3 Rademacher-complexity-based Generalization Bounds

The Rademacher complexity (BARTLETT and MENDELSON, 2002) of a hypothesis set \mathbb{H} can be defined as follows.

Definition 1.3.4 (Rademacher Complexity). Given a hypothesis set \mathbb{H} and a loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, the Rademacher complexity $\text{rad}(\mathbb{H})$ of the set \mathbb{H}

1.3. Generalization Bounds

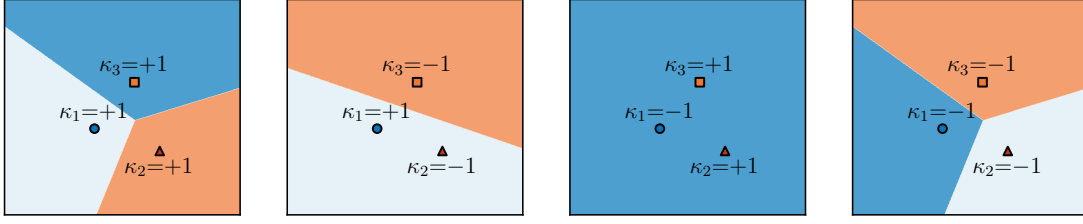


Figure 1.12. Illustration of the Rademacher complexity $\text{rad}(\mathbb{H})$ in multi-class classification with $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$ and $m = 3$. Given a learning sample $\mathbb{S} \in (\mathcal{X} \times \mathcal{Y})^2$, we show for each $(\kappa_1, \kappa_2) \in \{-1, +1\}^2$ the examples and a classifier h s.t. $R_S^\ell(h) = \sup_{h \in \mathbb{H}} \frac{1}{m} \sum_{i=1}^m \kappa_i \ell(h, (\mathbf{x}_i, y_i))$.

is defined as

$$\text{rad}(\mathbb{H}) = \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{\{\kappa_1, \dots, \kappa_m\} \sim \mathcal{K}^m} \sup_{h \in \mathbb{H}} \frac{1}{m} \sum_{i=1}^m \kappa_i \ell(h, (\mathbf{x}_i, y_i)),$$

where \mathcal{K} is the Rademacher distribution, i.e., $\mathcal{K}(+1) = \mathcal{K}(-1) = \frac{1}{2}$.

Given a learning sample $\mathbb{S} \sim \mathcal{D}^m$ and the Rademacher variables $\{\kappa_1, \dots, \kappa_m\} \sim \mathcal{K}^m$, the supremum is attained when the loss $\ell(h, (\mathbf{x}_i, y_i))$ is maximized (resp. minimized) when $\kappa_i = +1$ (resp. $\kappa_i = -1$). Then, the Rademacher complexity is the expected supremum over the learning sample and the Rademacher variables. For example, in binary classification, the Rademacher complexity measures the capacity of the hypotheses in \mathbb{H} to learn examples with random labels. As illustrated in Figure 1.12 when $m = 3$, (multi-class) linear classifiers are able to fit any data points with random labels. In this case, the Rademacher complexity $\text{rad}(\mathbb{H}) = 1$, but hopefully when $m \rightarrow +\infty$ the Rademacher complexity of the linear classifiers tends to 0 at a rate of $O\left(\sqrt{\frac{1}{m}}\right)$ (see, e.g., KAKADE *et al.* (2008), for the binary classification case).

From McDiarmid's concentration inequality (MCDIARMID, 1989), BARTLETT and MENDELSON (2002) derived the following generalization bound that depends on the Rademacher complexity.

Theorem 1.3.3 (Rademacher-complexity-based Generalization Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\sup_{h \in \mathbb{H}} \left[\underbrace{R_{\mathcal{D}}^\ell(h) - R_S^\ell(h)}_{\Phi_u(\delta)} \right] \leq 2\text{rad}(\mathbb{H}) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

Equivalently, with probability at least $1-\delta$ over $\mathbb{S} \sim \mathcal{D}^m$, we have

$$\forall h \in \mathbb{H}, \quad R_{\mathcal{D}}^{\ell}(h) \leq R_{\mathbb{S}}^{\ell}(h) + 2\text{rad}(\mathbb{H}) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}.$$

Put into words, the generalization gap is upper-bounded by the Rademacher complexity of the set \mathbb{H} along with a tight term when the number of examples m is large. As for the VC-Dimension-based bound, if we have an analytic expression of the complexity (or an upper bound), we can compute the generalization bound for all hypotheses h in \mathbb{H} . Since the bound holds for all $h \in \mathbb{H}$ (due to $\sup_{h \in \mathbb{H}}$), it can be seen as a *worst-case* analysis. Indeed, the same upper bound (*i.e.*, $\Phi_{\mathfrak{u}}(\delta)$) holds for all $h \in \mathbb{H}$, including the best, but also the worst (with the largest generalization gap $R_{\mathcal{D}}^{\ell}(h) - R_{\mathbb{S}}^{\ell}(h)$). This *worst-case* analysis makes the derivation of non-vacuous bounds hard (*i.e.*, with $\Phi_{\mathfrak{u}}(\delta) < 1$). Hence, other generalization bounds have been developed to avoid this worst-case analysis as we see in the next section.

1.3.2 Algorithmic-dependent Generalization Bounds

Let consider an algorithm that takes a learning sample $\mathbb{S} \in (\mathcal{X} \times \mathcal{Y})^m$ as input and outputs the hypothesis $h_{\mathbb{S}}$ belonging to the hypothesis set $\mathbb{H} = \{h_{\mathbb{S}}\}_{\mathbb{S} \in (\mathcal{X} \times \mathcal{Y})^m}$. Thanks to this algorithm, one can derive generalization bounds for the hypothesis $h_{\mathbb{S}}$ given $\mathbb{S} \sim \mathcal{D}^m$ that has the following form.

Definition 1.3.5 (Algorithmic-dependent Generalization Bound). Let $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ a generalization gap. An algorithmic-dependent generalization bound is defined such that if for any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, there exists a function $\Phi_{\mathfrak{a}} : (0, 1] \rightarrow \mathbb{R}$, such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\phi(R_{\mathcal{D}}^{\ell}(h_{\mathbb{S}}), R_{\mathbb{S}}^{\ell}(h_{\mathbb{S}})) \leq \Phi_{\mathfrak{a}}(\delta) \right] \geq 1 - \delta, \quad (1.2)$$

where usually $\phi(R_{\mathcal{D}}^{\ell}(h), R_{\mathbb{S}}^{\ell}(h)) = R_{\mathcal{D}}^{\ell}(h) - R_{\mathbb{S}}^{\ell}(h)$ and $h_{\mathbb{S}}$ is the hypothesis learned from an algorithm with $\mathbb{S} \sim \mathcal{D}^m$.

According to Definition 1.3.5, one can derive an upper bound $\Phi_{\mathfrak{a}}(\delta)$ of the generalization gap $R_{\mathcal{D}}^{\ell}(h_{\mathbb{S}}) - R_{\mathbb{S}}^{\ell}(h_{\mathbb{S}})$ for the algorithmic-dependent hypothesis $h_{\mathbb{S}}$. Actually, to derive such an upper bound $\Phi_{\mathfrak{a}}(\delta)$, a property of the considered learning algorithm must be considered; we recall in the following the algorithmic stability (BOUSQUET and ELISSEEFF, 2002) and the algorithmic robustness of XU and MANNOR (2010, 2012).

1.3. Generalization Bounds

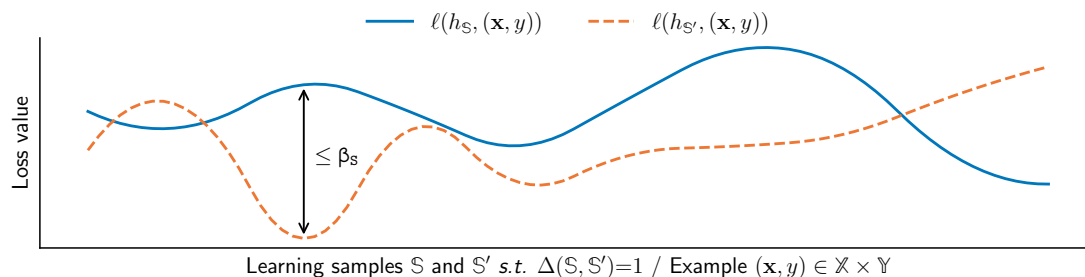


Figure 1.13. Schematic representation of the notion of stability β_S . The two curves represents the losses $\ell(h_S, (\mathbf{x}, y))$ and $\ell(h_{S'}, (\mathbf{x}, y))$ (the y-axis) for different combinations of examples $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ and learning samples \mathcal{S} and \mathcal{S}' that differs from one example (the x-axis).

1.3.2.1 Stability-based Generalization Bounds

The first algorithmic property we recall is the uniform stability (BOUSQUET and ELISSEEFF, 2002) and is defined as follows.

Definition 1.3.6 (Uniform Stability). Given the hypothesis set $\mathbb{H} = \{h_S\}_{S \in (\mathcal{X} \times \mathcal{Y})^m}$ and a loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, an algorithm is β_S -uniformly stable if

$$\sup_{\substack{S, S' \in (\mathcal{X} \times \mathcal{Y})^m \\ \text{s.t. } \Delta(S, S')=1}} \sup_{(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}} \left| \ell(h_S, (\mathbf{x}, y)) - \ell(h_{S'}, (\mathbf{x}, y)) \right| \leq \beta_S,$$

where $\Delta(S, S') = \sum_{i=1}^m \mathbb{I}[(\mathbf{x}_i, y_i) \neq (\mathbf{x}'_i, y'_i)]$ is the Hamming distance between the learning samples \mathcal{S} and \mathcal{S}' .

For two learning samples \mathcal{S} and \mathcal{S}' that differ from only one example (*i.e.*, that are very similar), the uniform stability β_S measures how much stable the algorithm is under small changes in the learning sample. To be stable, the difference of losses between the hypotheses h_S and $h_{S'}$ must be small (and upper bounded by β_S) for all examples $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$. Typically, β_S can be seen as a complexity measure that depends on the number of examples m . When $\beta_S = O(\frac{1}{\sqrt{m}})$ or $\beta_S = O(\frac{1}{m})$, the algorithm becomes more stable as the number of examples in \mathcal{S} increases. The notion of algorithmic stability is illustrated and summarized in Figure 1.13. From this notion of uniform stability and McDiarmid's inequality (MCDIARMID, 1989), the following generalization bounds can be derived (BOUSQUET and ELISSEEFF, 2002).

Theorem 1.3.4 (Stability-based Bounds). Given the hypothesis set $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m}$, for any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, for any $\beta_{\mathcal{S}}$ -uniformly stable algorithm, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[R_{\mathcal{D}}^{\ell}(h_{\mathcal{S}}) - R_{\mathcal{S}}^{\ell}(h_{\mathcal{S}}) \leq 2\beta_{\mathcal{S}} + (4m\beta_{\mathcal{S}} + 1) \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

Equivalently, with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, we have

$$R_{\mathcal{D}}^{\ell}(h_{\mathcal{S}}) \leq R_{\mathcal{S}}^{\ell}(h_{\mathcal{S}}) + 2\beta_{\mathcal{S}} + (4m\beta_{\mathcal{S}} + 1) \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}.$$

Note that, for this particular bound, if the algorithm is $O(1)$ -uniformly stable, the above bound does not converge and is vacuous. In contrast, an $O(\frac{1}{\sqrt{m}})$ -uniformly stable algorithm gives an $O(\frac{1}{\sqrt{m}})$ convergence rate. Recently, tighter generalization bounds have been improved: the bound converges with a $O(\frac{1}{m})$ -uniformly stable algorithm (FELDMAN and VONDRÁK, 2018, 2019; BOUSQUET *et al.*, 2020). As for the uniform-convergence-based bounds, an analytical expression of the stability $\beta_{\mathcal{S}}$ has to be derived to compute the bound.

1.3.2.2 Robustness-based Generalization Bounds

Another learning algorithm property used to derive generalization bounds is the algorithmic robustness. This notion can be defined as follows.

Definition 1.3.7 (Robustness). Given the hypothesis set $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m}$, a loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ and N disjoint sets *s.t.* $(\mathcal{X} \times \mathcal{Y}) = \bigcup_{i=1}^N \mathbb{Z}_i$, the algorithm is $(\{\mathbb{Z}_i\}_{i=1}^N, \beta_{\mathcal{R}})$ -robust⁵ if

$$\sup_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m} \sup_{i \in \{1, \dots, N\}} \sup_{(\mathbf{x}, y), (\mathbf{x}', y') \in \mathbb{Z}_i} \left| \ell(h_{\mathcal{S}}, (\mathbf{x}, y)) - \ell(h_{\mathcal{S}}, (\mathbf{x}', y')) \right| \leq \beta_{\mathcal{R}}.$$

For each subset $\mathbb{Z}_i \subseteq \mathcal{X} \times \mathcal{Y}$, the difference of losses between two examples $(\mathbf{x}, y) \in \mathbb{Z}_i$ and $(\mathbf{x}', y') \in \mathbb{Z}_i$ has to be upper-bounded by $\beta_{\mathcal{R}}$; for pedagogical purposes, we summarize in Figure 1.14 the algorithmic robustness. The following bound can be derived

⁵In the original definition, the robust parameter $\beta_{\mathcal{R}}$ can depend on the learning sample $\mathcal{S} \sim \mathcal{D}^m$. However, in the examples of XU and MANNOR (2012) for the classification setting, $\beta_{\mathcal{R}}$ depends only on the number of examples m .

1.3. Generalization Bounds

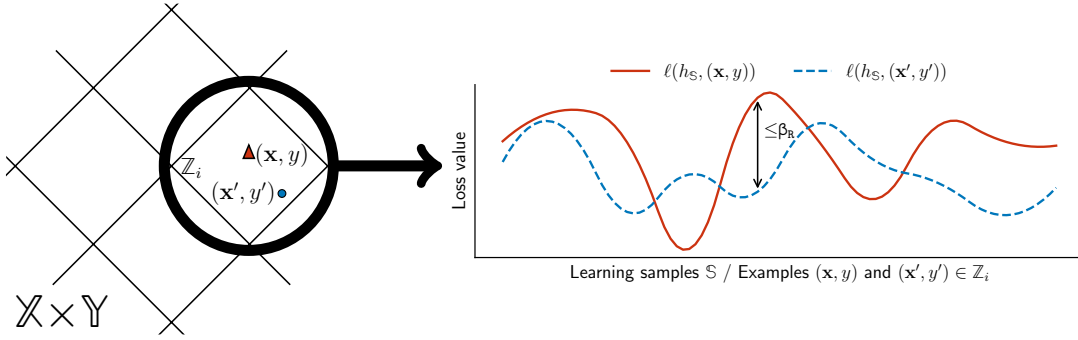


Figure 1.14. Schematic representation of the notion of algorithmic robustness. On the left plot, an example of partition $\{\mathbb{Z}_i\}_{i=1}^N$ for the set $\mathbb{X} \times \mathbb{Y}$ is shown. Then, for a given subset \mathbb{Z}_i , we represent on the right plot the losses $\ell(h_S, (\mathbf{x}, y))$ and $\ell(h_S, (\mathbf{x}', y'))$ for any learning sample S , and examples $(\mathbf{x}, y) \in \mathbb{Z}_i$ and $(\mathbf{x}', y') \in \mathbb{Z}_i$.

based on this property and thanks to the Bretaganolle-Huber-Carol concentration inequality (VAART and WELLNER, 1996, Proposition A.6.6).

Theorem 1.3.5 (Robustness-based Bounds). Given the hypothesis set $\mathbb{H} = \{h_S\}_{S \in (\mathbb{X} \times \mathbb{Y})^m}$, for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any loss function $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $(\{\mathbb{Z}_i\}_{i=1}^N, \beta_R)$ -robust algorithm, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[R_{\mathcal{D}}^{\ell}(h_S) \leq R_S^{\ell}(h_S) + \beta_R + \sqrt{\frac{2N \ln 2 + 2 \ln \frac{1}{\delta}}{2m}} \right] \geq 1 - \delta.$$

Equivalently, with probability at least $1 - \delta$ over $S \sim \mathcal{D}^m$, we have

$$R_{\mathcal{D}}^{\ell}(h_S) \leq R_S^{\ell}(h_S) + \beta_R + \sqrt{\frac{2N \ln 2 + 2 \ln \frac{1}{\delta}}{2m}}.$$

The bound is computable if we have an analytical expression of the algorithmic robustness parameter β_R . Ideally, the robustness parameter β_R must depend on m for the bound to converge. Furthermore, we can remark that there is a trade-off to find between the number of disjoint subsets N and the robust upper bound β_R . Indeed, the larger N is, the smaller we expect the parameter β_R to be. However, if $N \geq m$, the bound cannot be tighter than $\sqrt{\ln(2)}$, and hence, does not converge towards 0 when $m \rightarrow +\infty$.

The major drawback of the algorithmic-based bounds is that we have to derive the parameter β_S or β_R for each algorithm. Hence, such derivation can be tedious, and the

upper-bound $\Phi_a(\delta)$ is constant over the learning sample $\mathcal{S} \sim \mathcal{D}^m$. Another kind of bounds – called *PAC-Bayesian Bounds* – does not have such drawbacks and is appealing because of its facility to derive learning algorithms from it, as we will see in Part II. An in-depth introduction of these bounds is done in Chapter 2, but we give a quick overview in the rest of the chapter.

1.3.3 PAC-Bayesian Generalization Bounds

PAC-Bayesian bounds were introduced notably by SHAWE-TAYLOR and WILLIAMSON (1997) and MCALLESTER (1999), but it has been improved since then (SEEGER, 2002; MAURER, 2004; CATONI, 2007). The PAC-Bayesian bounds differ significantly from the ones of Sections 1.3.1 and 1.3.2: they require a probability distribution (denoted by ρ) on the hypothesis set \mathbb{H} . This distribution is used to assign a weight $\rho(h)$ to each hypothesis $h \in \mathbb{H}$. Hopefully, when the hypothesis h generalizes well, its weight $\rho(h)$ should be high. Thanks to this assumption, we can consider a *stochastic* hypothesis: given an input $\mathbf{x} \in \mathcal{X}$, its output is obtained by (1) sampling a hypothesis h from ρ , and by (2) computing the prediction $h(\mathbf{x})$. Actually, the PAC-Bayesian bounds allow us to upper-bound the true risk of *stochastic* hypotheses which is the *expected* true risk of a hypothesis sampled from a distribution ρ . More precisely, the PAC-Bayesian bounds study the *expected* generalization gap $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)$. We present one bound derived by MAURER (2004) in the following theorem.

Theorem 1.3.6 (PAC-Bayesian Bound of MAURER (2004)). For any distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, for any distribution π on \mathbb{H} (defined *a priori*), for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\text{For all distributions } \rho \text{ on } \mathbb{H}, \right. \\ \left. \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \leq \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]} \right] \geq 1 - \delta,$$

where the Kullback–Leibler (KL) divergence is defined as $\text{KL}(\rho \parallel \pi) = \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$. Equivalently, with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, we have

$$\forall \text{ distributions } \rho \text{ on } \mathbb{H}, \quad \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \leq \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]}.$$

The bound of this theorem depends on the KL divergence between ρ and π , which measures the complexity of the hypothesis sampled from ρ . Indeed, given a distribution π selected *a priori* before having the learning sample $\mathcal{S} \sim \mathcal{D}^m$, the higher $\text{KL}(\rho \parallel \pi)$,

the more different the two distributions ρ and π are. Hence, in some sense, the KL divergence captures how much the distribution ρ depends on the data \mathcal{S} .

1.4 Conclusion and Summary

In this chapter, we have seen an introduction to statistical learning. It includes an overview of algorithms to perform hypothesis selection: Empirical Risk Minimization and Structural Risk Minimization. Additionally, we recall some generalization bounds based on the uniform convergence (*e.g.*, with the VC-Dimension or the Rademacher complexity) or an algorithmic property (*e.g.*, the stability or the robustness) Moreover, we recall a generalization bound from a key theory in our contributions: the PAC-Bayesian theory. This is why we recall, with more details, in the next chapter the PAC-Bayesian theory.

THE PAC-BAYESIAN THEORY AND THE MAJORITY VOTE

2

Contents

2.1	Introduction	56
2.2	PAC-Bayesian Majority Votes	56
2.2.1	Definition	56
2.2.2	Upper Bounds on the Majority Vote's Risk	60
2.3	PAC-Bayesian Bounds	64
2.3.1	General PAC-Bayesian Bound of GERMAIN <i>et al.</i> (2009)	66
2.3.2	General PAC-Bayesian Bound of BÉGIN <i>et al.</i> (2016)	73
2.4	Disintegrated PAC-Bayesian Bounds	75
2.4.1	General Disintegrated Bound of RIVASPLATA <i>et al.</i> (2020)	76
2.4.2	Disintegrated Bound of CATONI (2007)	76
2.4.3	Disintegrated Bound of BLANCHARD and FLEURET (2007)	77
2.5	Conclusion and Summary	78

Abstract

In this chapter, we introduce, with more details, the PAC-Bayes theory that we outlined in Chapter 1. This theory allows us to upper-bound the risk of the *stochastic* classifier which samples, for each input, a hypothesis to predict the output. Moreover, the risk of the *stochastic* classifier can be linked to the *majority vote's* risk; we remind in this chapter the *majority vote classifier* which can be seen as a weighted combination of hypotheses. However, when we want to consider only one hypothesis, the *disintegrated* PAC-Bayesian theory becomes more adapted. Indeed, it upper-bounds the true risk of a *single* hypothesis associated with a high weight. Such generalization bounds are recalled as well.

2.1 Introduction

In this chapter, we details one statistical learning theory that is key in this thesis: the PAC-Bayesian theory. It is introduced by SHAWE-TAYLOR and WILLIAMSON (1997) and MCALLESTER (1998) for which we recall some bounds in Section 2.3. This theory assumes that each hypothesis $h \in \mathbb{H}$ is associated with a (positive) weight $\rho(h)$ that forms a probability distribution. Thanks to this assumption, the expected generalization gap is upper-bounded over $h \sim \rho$. The expected generalization gap allows to bound the true risk of the *stochastic* classifier; for each input \mathbf{x} , the model (i) samples a hypothesis $h \sim \rho$ and (ii) predicts the label with $h(\mathbf{x})$.

The stochastic classifier is related to a model in which we are particularly interested: the majority vote (see Section 2.2.2). For instance, majority votes are considered in boosting (FREUND and SCHAPIRE, 1996) or bagging (BREIMAN, 1996). In this context, a hypothesis is called *voter*, and a weight (modeled by the probability distribution ρ) is used to define the importance of each voter. Then, the majority vote is defined as a weighted combination of all the hypotheses from \mathbb{H} .

A majority vote might bring no significant improvements when the voters are strong (*i.e.*, when their individual risks are small). Hence, considering a *single* hypothesis associated with a high weight $\rho(h)$ can be a better option. Such a classifier is obtained by sampling from the probability distribution ρ . After sampling, generalization guarantees on the classifier can be derived through the *disintegrated* PAC-Bayesian bounds introduced independently by BLANCHARD and FLEURET (2007) and CATONI (2007). This type of bounds is recalled in Section 2.4.

For the sake of completeness, we defer in Appendix B the proofs.

2.2 PAC-Bayesian Majority Votes

2.2.1 Definition

Given a set of hypotheses \mathbb{H} , which is called set of voters in this context, the goal of a majority vote learning algorithm is to find a probability distribution ρ on \mathbb{H} . The distribution ρ defines the weights of the voters, *i.e.*, the importance of each voter in the majority vote. The majority vote is defined in the following way.

Definition 2.2.1 (Majority Vote). For any hypothesis set \mathbb{H} with voters $h : \mathcal{X} \rightarrow [-1, +1]$, for a distribution ρ on \mathbb{H} , the majority vote in the binary

classification setting ($\mathbb{Y} = \{-1, +1\}$) is defined as

$$\forall \mathbf{x} \in \mathbb{X}, \quad \text{MV}_\rho(\mathbf{x}) \triangleq \text{sign} \left(\mathbb{E}_{h \sim \rho} h(\mathbf{x}) \right).$$

In the multi-class setting ($\mathbb{Y} = \{1, 2, \dots, l\}$), for any hypothesis set \mathbb{H} with voters $h : \mathbb{X} \rightarrow \mathbb{Y}$, the ρ -weighted majority vote is defined as

$$\forall \mathbf{x} \in \mathbb{X}, \quad \text{MV}_\rho(\mathbf{x}) \triangleq \underset{y' \in \mathbb{Y}}{\text{argmax}} \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y'] = \underset{y' \in \mathbb{Y}}{\text{argmax}} \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}) = y'].$$

In the binary setting, the majority vote predicts as label the sign associated with the ρ -weighted average of the voters' outputs. In the multi-class setting, the majority vote predicts the label $y' \in \mathbb{Y}$ with the highest associated score $\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y']$. When the hypothesis set \mathbb{H} is composed of voters $h : \mathbb{X} \rightarrow \{-1, +1\}$ in the binary setting, the majority vote for the multi-class setting can be seen as a generalization. Indeed, we have $\text{sign}(\mathbb{E}_{h \sim \rho} h(\mathbf{x})) = +1$ if $\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = +1] \geq \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = -1]$ and $\text{sign}(\mathbb{E}_{h \sim \rho} h(\mathbf{x})) = -1$ otherwise.

While this definition of the majority vote classifier might appear a bit restrictive, it encompasses multiple widespread classifiers. For example, the Support Vector Machine (CORTES and VAPNIK, 1995) can be implicitly expressed as a majority vote where each voter depends on one example (GRAEPEL *et al.*, 2005). The well-known k -Nearest Neighbors (COVER and HART, 1967) is a majority vote (BELLET *et al.*, 2014). Additionally, when the voters depend on the whole learning sample \mathbb{S} , neural networks can be seen as a majority vote (KAWAGUCHI *et al.*, 2017; VIALARD *et al.*, 2019).

There are many approaches to learn a majority vote classifier based on ensemble methods. For example, the bagging (BREIMAN, 1996) method splits the learning sample and learns one voter with each subset. Then, the majority vote classifier averages the decision of the voters to take the final decision; random forest (BREIMAN, 2001) is, for instance, a particular bagging algorithm for decision trees. Moreover, the weights ρ can be learned greedily: this is the purpose of boosting algorithms such as Adaboost (FREUND and SCHAPIRE, 1996). This algorithm has been improved in various directions, *e.g.*, for the multi-class setting (SCHAPIRE and SINGER, 1998, 1999, 2000; ZHU *et al.*, 2009) or the ranking setting with RankBoost (FREUND *et al.*, 1998, 2003). Boosting algorithms have also been generalized for differentiable loss functions in a method called gradient boosting (FRIEDMAN, 2001).

Given a distribution \mathcal{D}' on $\mathcal{X} \times \mathcal{Y}$ (that encompasses the distributions \mathcal{D} and \mathcal{S}), the learner wants to learn MV_ρ that commits as few errors as possible on \mathcal{D} . To reduce the number of errors of the majority vote MV_ρ on \mathcal{D}' , the learner aims to minimize the risk $R_{\mathcal{D}'}(MV_\rho)$ under the 01-loss in the following definition.

Definition 2.2.2 (Risk of the Majority Vote). For any distribution \mathcal{D}' on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the risk of the majority vote is defined as

$$R_{\mathcal{D}'}(MV_\rho) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \mathbb{I}[MV_\rho(\mathbf{x}) \neq y] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'} [MV_\rho(\mathbf{x}) \neq y].$$

Specifically, when $\mathcal{D}' = \mathcal{D}$, the risk $R_{\mathcal{D}'}(MV_\rho) = R_{\mathcal{D}}(MV_\rho)$ is the true risk of the majority vote and when $\mathcal{D}' = \mathcal{S}$ the risk $R_{\mathcal{D}'}(MV_\rho) = R_{\mathcal{S}}(MV_\rho)$ is the empirical risk of the majority vote. In order to gain insight into the majority vote's decision, the margin captures how much the classifier makes errors. Indeed, the majority vote's risk can be expressed in terms of the margin defined in the following way.

Definition 2.2.3 (Margin of the Majority Vote). For any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the margin of the majority vote is defined as

$$m_\rho(\mathbf{x}, y) \triangleq \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] - \max_{y' \in \mathcal{Y}, y' \neq y} \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y'].$$

The margin is positive if the score $\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y]$ is higher than the score associated with the other labels $y' \neq y$ and is negative otherwise. It captures if an example $(\mathbf{x}, y) \sim \mathcal{D}'$ is misclassified. Indeed, thanks to the margin, the majority vote's risk $R_{\mathcal{D}'}(MV_\rho)$ can be rewritten as follows:

$$\begin{aligned} R_{\mathcal{D}'}(MV_\rho) &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'} \left[\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] \leq \max_{y' \in \mathcal{Y}, y' \neq y} \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y'] \right] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'} [m_\rho(\mathbf{x}, y) \leq 0]. \end{aligned}$$

Hence, the risk is the probability that one of the scores is higher than the one of the correct label. Unfortunately, deriving a learning algorithm to optimize the margin can be challenging since it is non-convex *w.r.t.* the posterior ρ because of the \max . To overcome this issue, LAVIOLETTE *et al.* (2017) propose to consider a convex lower bound of the true margin called $\frac{1}{2}$ -margin¹.

¹We multiply the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017) by two.

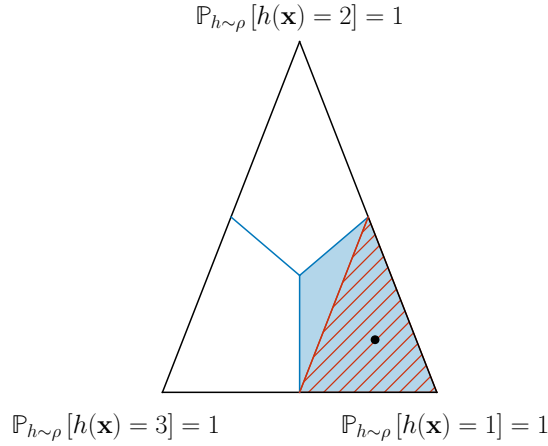


Figure 2.1. Illustration of the majority vote's margin in the multi-class setting with 3 classes, i.e., $\mathbb{Y} = \{1, 2, 3\}$. The triangle represents the ternary plot of the three scores where each triangle's vertex is the three possible maximum scores with $\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = i] = 1$ for all $i \in \mathbb{Y}$. The black dot represents the prediction of an example $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$ by the majority vote: the scores are $\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = 1] = 0.7$ and $\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = 2] = \mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = 3] = 0.15$. The blue area is where the majority vote predicts the label y for the input \mathbf{x} and where the margin $m_\rho(\mathbf{x}, y)$ is positive. Whereas the red area represents the predictions where $\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = 1] \geq \frac{1}{2}$ (i.e., where the $\frac{1}{2}$ -margin $\widehat{m}_\rho(\mathbf{x}, y)$ is positive).

Definition 2.2.4 ($\frac{1}{2}$ -Margin of the Majority Vote). For any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the $\frac{1}{2}$ -margin is defined as

$$\widehat{m}_\rho(\mathbf{x}, y) = 2 \left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y] - \frac{1}{2} \right].$$

When the score $\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y]$ exceeds $\frac{1}{2}$, the majority vote surely classifies correctly the example $(\mathbf{x}, y) \sim \mathcal{D}'$. Hence, the idea of this margin is to compute the difference between the score and $\frac{1}{2}$: when the margin is positive, the example is correctly classified. We illustrate the difference between the two margins in Figure 2.1. In binary classification, the $\frac{1}{2}$ -margin boils down to $\widehat{m}_\rho(\mathbf{x}, y) = y \mathbb{E}_{h \sim \rho} h(\mathbf{x}) = \mathbb{E}_{h \sim \rho} m_h(\mathbf{x}, y)$ (which is the margin introduced in Section 1.1.2.2). From Definition 2.2.4, we can deduce the following upper bound on the majority vote's risk, i.e., we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'} \left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y] \leq \frac{1}{2} \right] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y) \leq 0].$$

In general, the majority vote's risk $R_{\mathcal{D}'}(\text{MV}_\rho)$ is not practical in learning algorithms because the gradient is zero (see Section 1.1.2.2). Hence, in the next section, we recall different upper bounds on the majority vote's risk considered as surrogate losses.

2.2.2 Upper Bounds on the Majority Vote's Risk

We recall three surrogates on the majority vote's risk introduced in the literature, namely, the Gibbs Risk (LANGFORD and SHAWE-TAYLOR, 2002; MCALLESTER, 2003), the joint error (LACASSE *et al.*, 2006; GERMAIN *et al.*, 2015; MASEGOSA *et al.*, 2020), and the C-Bound (BREIMAN, 2001; LACASSE *et al.*, 2006; ROY *et al.*, 2011). We make use of these surrogate in Part II to derive self-bounding learning algorithms.² More precisely, in Chapter 3, we leverage the Gibbs risk for the adversarially robust setting and, in Chapter 4, we learn majority vote classifiers through the C-Bound.

2.2.2.1 The Gibbs Risk

The first surrogate on the risk $R_{\mathcal{D}'}(\text{MV}_\rho)$ introduced in the PAC-Bayesian literature is the *Gibbs risk* (LANGFORD and SHAWE-TAYLOR, 2002; MCALLESTER, 2003). This risk gives the average performance of individual voters in the majority vote. It is defined in the following way.

Definition 2.2.5 (Gibbs Risk). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the *Gibbs risk* is defined as

$$r_{\mathcal{D}'}(\rho) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}', h \sim \rho} [h(\mathbf{x}) \neq y].$$

Put into words, the Gibbs risk is the ρ -weighted average of the voters' risk. Moreover, we can write the Gibbs Risk with respect to the $\frac{1}{2}$ -margin. Indeed, we have

$$r_{\mathcal{D}'}(\rho) = \frac{1}{2} \left[1 - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \right]. \quad (2.1)$$

LANGFORD and SHAWE-TAYLOR (2002) show that the majority vote's risk is upper-bounded by twice the Gibbs risk. The inequality is recalled in the following theorem.

²A self-bounding algorithm, coined by FREUND (1998), is a learning algorithm that comes with a generalization guarantee.

Theorem 2.2.1 (Risk Upper Bound Based on the Gibbs Risk). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq 2 r_{\mathcal{D}'}(\rho). \quad (2.2)$$

Proof. Deferred to Appendix B.1. ■

However, in ensemble methods where one wants to combine voters efficiently, the Gibbs risk appears to be an imprecise surrogate since the combination of voters might compensate for individual errors. Hence, other surrogates need to be taken into account.

2.2.2.2 Joint Error

The joint error, introduced by LACASSE *et al.* (2006), takes better the voters' correlation into account; this quantity is defined in the following way.

Definition 2.2.6 (Joint Error). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the *joint error* is defined as

$$\begin{aligned} e_{\mathcal{D}'}(\rho) &\triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) \neq y] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}', h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq y, h'(\mathbf{x}) \neq y]. \end{aligned}$$

Similarly to Equation (2.1) for the Gibbs risk, we can reinterpret this surrogate with the $\frac{1}{2}$ -margin. Indeed, from GERMAIN *et al.* (2015), we have

$$\begin{aligned} e_{\mathcal{D}'}(\rho) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \left(\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) \neq y] \right)^2 \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \left(\frac{1}{2} [1 - \widehat{m}_\rho(\mathbf{x}, y)] \right)^2 \\ &= \frac{1}{4} \left(1 - 2 \mathbb{E}_{h \sim \rho} \widehat{m}_\rho(\mathbf{x}, y) + \mathbb{E}_{h \sim \rho} \widehat{m}_\rho(\mathbf{x}, y)^2 \right). \end{aligned} \quad (2.3)$$

Recently, MASEGOSA *et al.* (2020) proposed to deal directly with the joint error to bound the majority vote's risk; Their bound is presented in the following theorem.

Theorem 2.2.2 (Risk Upper Bound Based on the Joint Error). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , we have

$$R_{\mathcal{D}'}(\text{MV}_{\rho}) \leq 4e_{\mathcal{D}'}(\rho). \quad (2.4)$$

Proof. Deferred to Appendix B.2. ■

This inequality captures the fact that the voters need to be sufficiently diverse and commit errors on different points. However, when the joint error $e_{\mathcal{D}'}(\rho)$ exceeds $\frac{1}{4}$, the bound exceeds 1 and is uninformative. In some cases, these two upper bounds appear to be less tight than another one: the C-Bound (BREIMAN, 2001; LACASSE *et al.*, 2006).

2.2.2.3 The C-Bound

The C-Bound³ (BREIMAN, 2001; LACASSE *et al.*, 2006) is another surrogate (and upper bound) of the majority vote's risk. It is derived from the CHEBYSHEV-CANTELLI Inequality (see Theorem A.4.1). It depends on the *Gibbs risk*, the *joint error*, and the *disagreement*. This latter is defined in the following way.

Definition 2.2.7 (Disagreement). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , the *disagreement* is defined as

$$\begin{aligned} d_{\mathcal{D}'}(\rho) &\triangleq 2 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) = y] \\ &= 2 \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}', h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq y, h'(\mathbf{x}) = y]. \end{aligned}$$

The higher the disagreement, the more the voters do not perform the same prediction for a given $(\mathbf{x}, y) \sim \mathcal{D}'$. Similarly to the Gibbs risk and the joint error, we can express the disagreement with the margin. Indeed, we have

$$d_{\mathcal{D}'}(\rho) = \frac{1}{2} \left[1 - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \widehat{m}_{\rho}(\mathbf{x}, y)^2 \right]. \quad (2.5)$$

Interestingly, by developing Equation (2.3), we can relate the Gibbs risk and the joint error to the disagreement with

$$d_{\mathcal{D}'}(\rho) = 2 [r_{\mathcal{D}'}(\rho) - e_{\mathcal{D}'}(\rho)]. \quad (2.6)$$

³The term ‘‘C-Bound’’ was introduced in the PAC-Bayesian literature by (LACASSE *et al.*, 2006).

2.2. PAC-Bayesian Majority Votes

This inequality tells us two facts: (i) the lower the Gibbs risk, the lower the disagreement, and (ii) the disagreement increases as the joint error decreases. The expression of the disagreement in binary classification can be simplified (with voters $h : \mathcal{X} \rightarrow \{-1, +1\}$). Indeed, given an example $(\mathbf{x}, y) \sim \mathcal{D}'$ with $y \in \{-1, +1\}$, we have

$$\begin{aligned} \mathbb{P}_{h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq h'(\mathbf{x})] &= \mathbb{P}_{h \sim \rho, h' \sim \rho} [h(\mathbf{x}) = y, h'(\mathbf{x}) \neq y] + \mathbb{P}_{h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq y, h'(\mathbf{x}) = y] \\ &= 2 \mathbb{P}_{h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq y, h'(\mathbf{x}) = y], \end{aligned}$$

which gives $d_{\mathcal{D}'}(\rho) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}', h \sim \rho, h' \sim \rho} [h(\mathbf{x}) \neq h'(\mathbf{x})]$.

Thanks to this quantity, we are now able to recall the C-Bound. Note that it was first introduced by LACASSE *et al.* (2006) for the PAC-Bayesian majority vote in binary classification. The generalization to the multi-class setting has been introduced by LAVIOLETTE *et al.* (2017, Theorem 2 and Corollary 1); we recall the C-Bound in the following theorem.

Theorem 2.2.3 (The C-Bound). For any distribution \mathcal{D}' on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , if

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) > 0 \iff r_{\mathcal{D}'}(\rho) < \frac{1}{2} \iff 2e_{\mathcal{D}'}(\rho) + d_{\mathcal{D}'}(\rho) < 1,$$

we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq 1 - \frac{\left(\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)]\right)^2}{\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \quad (2.7)$$

$$= 1 - \frac{(1 - 2r_{\mathcal{D}'}(\rho))^2}{1 - 2d_{\mathcal{D}'}(\rho)} \quad (2.8)$$

$$= 1 - \frac{\left(1 - [2e_{\mathcal{D}'}(\rho) + d_{\mathcal{D}'}(\rho)]\right)^2}{1 - 2d_{\mathcal{D}'}(\rho)} \quad (2.9)$$

$$= C_{\mathcal{D}'}(\rho).$$

Proof. Deferred to Appendix B.3. ■

This surrogate is expressed as a trade-off between the first and the second statistical moment of the $\frac{1}{2}$ -margin $\widehat{m}_\rho(\mathbf{x}, y)$ (Equation (2.7)). Based on Equations (2.1) and (2.5), the trade-off can be expressed through the disagreement and the Gibbs risk (see Equation (2.8)). Moreover, from Equation (2.6), we can also see the C-Bound

as a trade-off between the joint error and the disagreement in Equation (2.9).

The three surrogates of Theorems 2.2.1 to 2.2.3 can be easily compared. For instance, when $r_{\mathcal{D}'}(\rho) \leq d_{\mathcal{D}'}(\rho)$, the C-Bound $C_{\mathcal{D}'}(\rho)$ (Equation (2.7)) is tighter than $2r_{\mathcal{D}'}(\rho)$ (Equation (2.2)) and $4e_{\mathcal{D}'}(\rho)$ (Equation (2.4)). Hence, when the disagreement increases, the C-Bound appears to be a good trade-off between the Gibbs risk and the disagreement. More precisely, the main interest of the C-bound compared to Equation (2.4) is that when $e_{\mathcal{D}'}(\rho)$ is close to $\frac{1}{4}$, the C-Bound can be close to 0 depending on the value of the disagreement $d_{\mathcal{D}'}(\rho)$: the C-bound is then more precise. Moreover, it is important to notice that the C-Bound is tighter than $4e_{\mathcal{D}'}(\rho)$ for all cases. We summarize the relationships between the different surrogates in the next theorem and illustrate it in Figure 2.2; this relation is notably given by GERMAIN *et al.* (2015) and MASEGOSA *et al.* (2020).

Theorem 2.2.4 (Relationship between Theorems 2.2.1 to 2.2.3). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any voters set \mathbb{H} , for any distribution ρ on \mathbb{H} , if $r_{\mathcal{D}'}(\rho) < \frac{1}{2}$ (i.e., $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \widehat{m}_{\rho}(\mathbf{x}, y) > 0$), we have

- (i) $R_{\mathcal{D}'}(\text{MV}_{\rho}) \leq C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho)$, if $r_{\mathcal{D}'}(\rho) \leq d_{\mathcal{D}'}(\rho)$,
- (ii) $R_{\mathcal{D}'}(\text{MV}_{\rho}) \leq 2r_{\mathcal{D}'}(\rho) \leq C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho)$, otherwise.

Proof. Deferred to Appendix B.4. ■

Since the distribution \mathcal{D} is unknown, the true risk of the majority vote $R_{\mathcal{D}}(\text{MV}_{\rho})$ is not computable. Hence, one can minimize the empirical risk of the majority vote $R_S(\text{MV}_{\rho})$ through the Empirical Risk Minimization approach (see Algorithm 1.3). However, this minimization does not necessarily lead to a low true risk $R_{\mathcal{D}}(\text{MV}_{\rho})$ since overfitting can occur (see Section 1.2). To tackle this issue, one solution is to deal with the minimization of a generalization bound to get a self-bounding algorithm (FREUND, 1998), i.e., the minimization of the risk with generalization guarantees. We will see in the next section PAC-Bayesian generalization bounds that further allow us to upper-bound the true risk $R_{\mathcal{D}}(\text{MV}_{\rho})$ in Part II.

2.3 PAC-Bayesian Bounds

The PAC-Bayesian theory⁴, introduced by SHAWE-TAYLOR and WILLIAMSON (1997) and MCALLESTER (1999), aims to provide PAC generalization bounds for Bayesian-like algorithms. Such Bayesian algorithms assume a probability distribution defined

⁴See (GUEDJ, 2019; ALQUIER, 2021) for recent on the PAC-Bayesian theory.

2.3. PAC-Bayesian Bounds

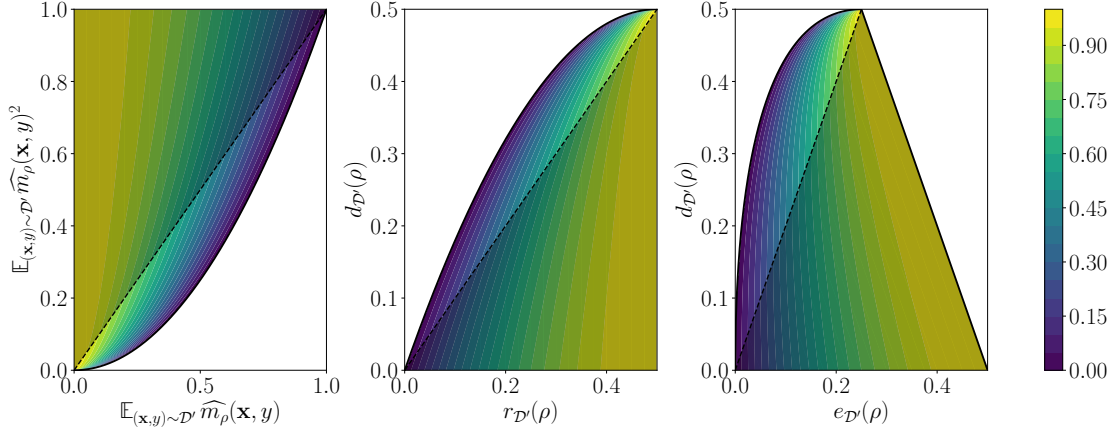


Figure 2.2. Plots (from left to right) of the different C-Bounds in Equations (2.7) to (2.9). For each plot, the darker area represents the cases where $2r_{\mathcal{D}'}(\rho)$ is tighter than the C-Bound $C_{\mathcal{D}'}(\rho)$. The dashed line represents the cases where $4e_{\mathcal{D}'}(\rho)$ matches the C-Bound $C_{\mathcal{D}'}(\rho)$.

a priori on the hypothesis set \mathbb{H} , and thanks to Bayes' theorem, obtain an *a posteriori* probability distribution on \mathbb{H} thanks to the learning sample; see BISHOP (2007) for more details on this Bayesian inference procedure. Contrary to the classical Bayesian inference where the *posterior* distribution is proportional to the product of the *prior* distribution and the likelihood of the data, in PAC-Bayesian theory, an arbitrary *prior* distribution can be considered. Actually, the term ‘‘Bayesian’’ in the PAC-Bayesian theory comes from the fact that we usually upper-bound the *expected* generalization gap $|\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h)|$ where $h \in \mathbb{H}$ is sampled from a data-dependent distribution $\rho_{\mathcal{S}}$ called *posterior*. This theory has been extended to various settings such as transductive learning (DERBEKO *et al.*, 2004; BÉGIN *et al.*, 2014), regression (GERMAIN *et al.*, 2016; SHALAEVA *et al.*, 2020), structured prediction (LAVIOLETTE *et al.*, 2017), domain adaptation (GERMAIN *et al.*, 2020), or randomized learning (LONDON, 2017).

In order to define a PAC-Bayesian bound more formally, we denote by $\rho_{\mathcal{S}}$ (or ρ) the *posterior* distribution on \mathbb{H} and π the *prior* distribution. Each probability distribution ρ is defined through its probability density function $h \mapsto \rho(h)$ with respect to a reference measure⁵ on \mathbb{H} ; we denote by $\mathbb{M}(\mathbb{H})$ the set of probability density function on \mathbb{H} . Hence, the distribution $\rho_{\mathcal{S}} \in \mathbb{M}(\mathbb{H})$ is the Radon–Nikodym derivative of a probability measure *w.r.t.* the reference measure. We also denote by $\mathbb{M}^*(\mathbb{H}) \subseteq \mathbb{M}(\mathbb{H})$ the set of strictly positive probability densities on \mathbb{H} . Moreover, for convenience, we assume that the support of the posterior $\rho_{\mathcal{S}}$ is in the support of π , *i.e.*, if $\pi(h) = 0$ then $\rho_{\mathcal{S}}(h) = 0$

⁵For instance, if $\mathbb{H} = \mathbb{R}^d$, then the reference measure is the Lebesgue one.

(the absolute continuity); hence $\pi \in \mathbb{M}^*(\mathbb{H})$.

With the setting in place, we can now define the general form of a PAC-Bayesian bound in the following.

Definition 2.3.1 (PAC-Bayesian Generalization Bound). Let $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ a generalization gap. A PAC-Bayesian bound is defined such that if for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , there exists a function $\Phi : \mathbb{M}(\mathbb{H}) \times \mathbb{M}^*(\mathbb{H}) \times (0, 1] \rightarrow \mathbb{R}$, such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\phi(\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h), \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h)) \leq \Phi(\rho_{\mathcal{S}}, \pi, \delta) \right] \geq 1 - \delta,$$

where e.g. $\phi(\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h), \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h)) = \left| \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h) \right|$.

Definition 2.3.1 is a general definition in the sense that the expected generalization gap $\left| \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h) \right|$ is not the only usable deviation between the true risk $\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h)$ and the empirical risk $\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h)$. For example, one can consider the one-sided difference $\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h)$ or the squared difference $\left[\mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho_{\mathcal{S}}} R_{\mathcal{S}}^{\ell}(h) \right]^2$. As we will see in this chapter, several gaps $\phi(\cdot)$ have been considered in the literature. These gaps are upper-bounded with probability at least $1 - \delta$ with a function $\Phi(\cdot)$ that depends on δ and two probability distributions on \mathbb{H} . Usually, the lower the parameter $\delta \in (0, 1]$, the higher the upper bound $\Phi(\cdot)$, *i.e.*, the function $\Phi(\cdot)$ is decreasing *w.r.t.* δ . Moreover, this upper bound $\Phi(\cdot)$ depends on a data-dependent *posterior* distribution $\rho_{\mathcal{S}} \in \mathbb{M}(\mathbb{H})$ and a *prior* distribution $\pi \in \mathbb{M}^*(\mathbb{H})$. The prior distribution is data-free and can incorporate prior knowledge, *e.g.*, coming from an expert knowledge or an additional learning sample (PARRADO-HERNÁNDEZ *et al.*, 2012; DZIUGAITE *et al.*, 2021). In the rest of the section, we present different instantiations of Definition 2.3.1, making an overview of the PAC-Bayesian bounds in the literature.

2.3.1 General PAC-Bayesian Bound of Germain *et al.* (2009)

There are different kinds of bounds in the PAC-Bayesian literature (*e.g.*, SEEGER, 2002; MCALLESTER, 2003; CATONI, 2007). Several existing bounds can be proved from a general theorem by GERMAIN *et al.* (2009); it is recalled in the following theorem.

Theorem 2.3.1 (General PAC-Bayesian Bound of GERMAIN *et al.* (2009)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \leq \text{KL}(\rho \| \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S}')} \right) \right] \geq 1 - \delta,$$

where $\text{KL}(\rho \| \pi) = \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$ is the Kullback-Leibler (KL) divergence between the distributions ρ and π .

Proof. Deferred to Appendix B.6. ■

Note that this bound holds for all posterior distributions $\rho \in \mathcal{M}(\mathbb{H})$, which includes notably the prior distribution π , or any data-dependent posterior $\rho_{\mathcal{S}}$. Furthermore, this bound is penalized by the KL divergence⁶ between ρ and π . The closer the posterior ρ is to the prior π , the smaller the divergence and the bound. Moreover, the bound holds for a function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$ that further captures a deviation between the true risk $R_{\mathcal{D}}^{\ell}(h)$ and the empirical risk $R_{\mathcal{S}}^{\ell}(h)$. For example, with $\varphi(h, \mathcal{S}) = m\phi(R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{S}}^{\ell}(h))$ where $\phi(\cdot)$ is convex, we are able to retrieve Definition 2.3.1. By setting this function $\phi(\cdot)$ accordingly, one can retrieve some classical bounds that have been previously introduced in the literature (SEEGER, 2002; MCALLESTER, 2003; CATONI, 2007) as we show further.

2.3.1.1 McAllester-like bound

First of all, we can retrieve Theorem 1.3.6 presented in Chapter 1 which is a tighter version of the bound derived by MCALLESTER (2003, Theorem 1). Indeed, by setting $\varphi(h, \mathcal{S}) = 2m \left[R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h) \right]^2$ in Theorem 2.3.1, one can deduce the following result.

Theorem 2.3.2 (PAC-Bayesian Bound of MCALLESTER (2003)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathcal{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \left| \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right| \leq \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]} \right] \geq 1 - \delta. \quad (2.10)$$

⁶The principal properties of the KL divergence are given in Appendix B.5

Proof. Deferred to Appendix B.7. ■

According to Theorem 2.3.2, the gap $|\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)|$ tends to zero when the number of examples increase. Hence, the more examples we have, the closer the expected empirical risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)$ is from expected true empirical risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)$ for all $\rho \in \mathbb{M}(\mathbb{H})$. Actually, bounding this gap gives an upper bound on the expected true risk. Indeed, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$, we have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \leq \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta}]} \quad (2.11)$$

$$\text{and} \quad \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \geq \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) - \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta}]} \quad (2.12)$$

Put into words, the expected true risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)$ can be lower and upper bounded by the expected empirical risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)$ and the bound $\sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta}]}$. When our objective is to obtain a low expected true risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)$, one can obtain a distribution $\rho \in \mathbb{M}(\mathbb{H})$ minimizing the bound with the minimization problem

$$\min_{\rho \in \mathbb{M}(\mathbb{H})} \left\{ \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta}]} \right\}.$$

Since the bound holds for all $\rho \in \mathbb{M}(\mathbb{H})$ (with high probability), it also holds for the optimal solution. The contributions in Part II rely on such a minimization problem to obtain models with a certified expected test risk. Nevertheless, as we will see further, the bound of Theorem 2.3.2 is not the tightest.

2.3.1.2 Catoni-like bound

CATONI (2007, Theorem 1.2.1) proposed a bound that can be tighter than the one of Theorem 2.3.2 by considering a parameter $c > 0$. His proposed bound can be retrieved from Theorem 2.3.1 by defining $\varphi(h, \mathcal{S}) = m\phi(R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{S}}^{\ell}(h))$ where the deviation is $\phi(R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{S}}^{\ell}(h)) = -\ln(1 - (1 - e^{-c}) R_{\mathcal{D}}^{\ell}(h)) - c R_{\mathcal{S}}^{\ell}(h)$. We state the bound in the following theorem.

Theorem 2.3.3 (PAC-Bayesian Bound of CATONI (2007)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss ℓ :

2.3. PAC-Bayesian Bounds

$\mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $c > 0$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), -\ln \left(1 - [1 - e^{-c}] \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right. \\ \left. \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta. \quad (2.13) \end{aligned}$$

Proof. Deferred to Appendix B.8. ■

The result is difficult to interpret because of the gap $-\ln \left(1 - [1 - e^{-c}] \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)$ is not easy to analyze. However, rewriting it as an upper bound of the expected true risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)$ makes its interpretation easier. Indeed, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$, we have

$$\forall \rho \in \mathcal{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \leq \frac{1}{1 - e^{-c}} \left[1 - \exp \left(-c \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) - \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} \right] \right) \right].$$

Put into words, the expected true risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)$ is upper-bounded by a trade-off, controlled by the parameter c , between the expected empirical risk $\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)$ and the term $\frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} \right]$. This is in contrast with Equations (2.10) to (2.12) since this bound depends on a parameter $c > 0$, which allows one to tune the tightness of the bound.

In practice, it is hard to set the parameter c since the bound holds with high probability on $\mathcal{S} \sim \mathcal{D}^m$ for all parameters $c > 0$. Hence, it is not possible to condition c on $\mathcal{S} \sim \mathcal{D}^m$. To tackle this issue, one usually applies a union bound to get a bound holding for any c belonging to a countable set. Hopefully, as shown further, one can derive a bound that avoids this parameter and is potentially tighter.

2.3.1.3 Seeger-like bound

One of the tightest PAC-Bayesian bound (that avoids the parameter c) is the one proven by SEEGER (2002). This bound depends on the KL divergence between two Bernoulli distributions. This function is defined in the following way.

Definition 2.3.2 (KL Divergence Between Two Bernoulli Distributions). For any distribution $q \in [0, 1]$ and $p \in [0, 1]$, the small kl is defined as

$$\text{kl}(q \parallel p) \triangleq \text{KL}(\mathcal{B}(q) \parallel \mathcal{B}(p)) = q \ln \frac{q}{p} + (1 - q) \ln \frac{1 - q}{1 - p},$$

where $\mathcal{B}(q)$ and $\mathcal{B}(p)$ are two Bernoulli distributions with bias q and p .

The first PAC-Bayesian theorem based on the divergence $\text{kl}()$ was proved by SEEGER (2002). Few years later, by setting $\varphi(h, \mathbb{S}) = m \text{kl}(\mathbb{R}_S^\ell(h) \parallel \mathbb{R}_D^\ell(h))$, GERMAIN *et al.* (2009) retrieved an improved PAC-Bayesian generalization bound proved by MAURER (2004). MAURER's bound and stated in the following theorem.

Theorem 2.3.4 (PAC-Bayesian Bound of SEEGER (2002)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \text{kl} \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_D^\ell(h) \parallel \mathbb{E}_{h \sim \rho} \mathbb{R}_S^\ell(h) \right) \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right] \geq 1 - \delta. \quad (2.14)$$

Proof. Deferred to Appendix B.9. ■

The bound of Theorem 2.3.4 consists in upper-bounding a deviation between $\mathbb{E}_{h \sim \rho} \mathbb{R}_D^\ell(h)$ and $\mathbb{E}_{h \sim \rho} \mathbb{R}_S^\ell(h)$. This deviation is hard to interpret. Thanks to the PINSKER's inequality (Theorem B.5.1), *i.e.*, $\forall (p, q) \in [0, 1]^2, 2(q - p)^2 \leq \text{kl}(q \parallel p)$, it can be shown that the bound of Theorem 2.3.2 is tighter than the one of Theorem 2.3.4. Interestingly, GERMAIN *et al.* (2009, Proposition 2.1) and LACASSE (2010, Proposition 6.2.2) related the bound of CATONI (Theorem 2.3.3) and Theorem 2.3.4 with the equality

$$\max_{c > 0} \left\{ -\ln(1 - [1 - e^{-c}]p) - cq \right\} = \text{kl}(q \parallel p).$$

Put into words, given $p \in (0, 1]$ and $q \in [0, 1]$, the $\text{kl}(q \parallel p)$ matches the function $-\ln(1 - [1 - e^{-c}]p) - cq$ with the optimal parameter $c \geq 0$. Expressed as it is, Equation (2.14) does not permit to upper or lower bound the expected true risk $\mathbb{R}_D^\ell(h)$ contrary to Theorems 2.3.2 and 2.3.3. In order to rewrite Equation (2.14) (of Theorem 2.3.4), we define the inverting functions of $\text{kl}()$ in the following way.

Definition 2.3.3 (Inverting Functions of $\text{kl}()$). Given $\tau \geq 0$, for any $q \in [0, 1]$, the inverting functions of the $\text{kl}()$ are defined as

$$\begin{aligned} \overline{\text{kl}}(q \mid \tau) &\triangleq \max \left\{ p \in (0, 1) \mid \text{kl}(q \parallel p) \leq \tau \right\}, \\ \text{and } \underline{\text{kl}}(q \mid \tau) &\triangleq \min \left\{ p \in (0, 1) \mid \text{kl}(q \parallel p) \leq \tau \right\}. \end{aligned}$$

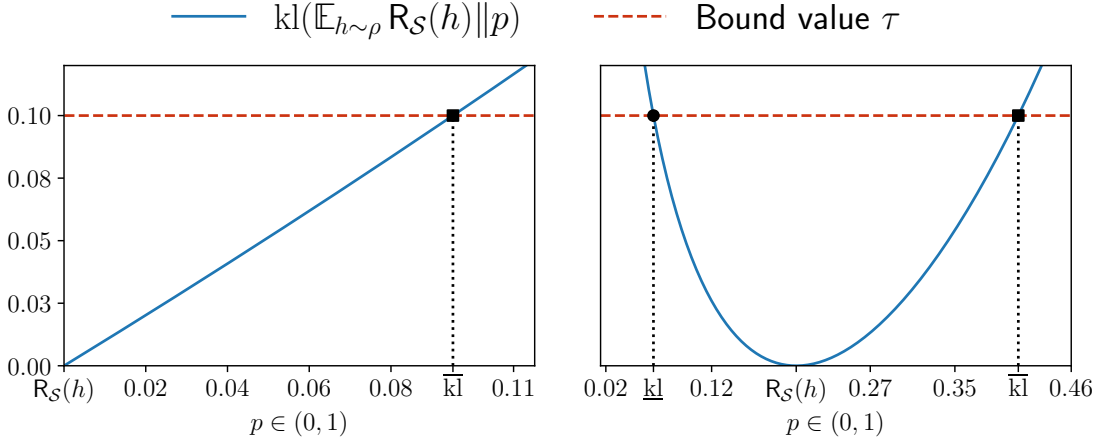


Figure 2.3. Illustration of the inverting functions of $\text{kl}(\cdot)$ for $q = R_S^\ell(h)$ s.t. $R_S^\ell(h) \in \{0.0, 0.2\}$ in two different plots. We plot the curve of the function $\text{kl}(R_S^\ell(h)||\cdot)$, the bound value $\tau = 0.1$ and the two inverting functions $\underline{\text{kl}}(R_S^\ell(h)|\tau)$ and $\overline{\text{kl}}(R_S^\ell(h)|\tau)$ (abbreviated $\underline{\text{kl}}$ and $\overline{\text{kl}}$). The dot (resp. the square) corresponds to the solution of the minimization (resp. maximization) problem associated with $\underline{\text{kl}}$ (resp. $\overline{\text{kl}}$).

The function $\overline{\text{kl}}(\cdot)$ (resp. $\underline{\text{kl}}(\cdot)$) denotes the maximum (resp. minimum) value $p \in (0, 1)$ such that the inequality $\text{kl}(q||p) \leq \tau$ holds. Figure 2.3 gives a graphical illustration of these inverting functions. The values associated with the inverting functions $\underline{\text{kl}}(\cdot)$ and $\overline{\text{kl}}(\cdot)$ can be approximated and easily computed from PINSKER's inequality (Theorem B.5.1). Indeed, we have

$$\overline{\text{kl}}(q|\tau) \leq q + \sqrt{\frac{1}{2}\tau} \quad \text{and} \quad q - \sqrt{\frac{1}{2}\tau} \leq \underline{\text{kl}}(q|\tau). \quad (2.15)$$

We present in Figure 2.4 an illustration of the tightness of this approximation.

To calculate $\underline{\text{kl}}(\cdot)$ and $\overline{\text{kl}}(\cdot)$ exactly, two optimization problems need to be solved; REEB *et al.* (2018) proposed an algorithm based on the bisection method. We recall its pseudo-code in Algorithm 2.1. The principle of this algorithm is to iteratively refine the interval $[p_{\min}, p_{\max}]$ to which the bias $p \in (0, 1]$ belongs. When the equality $\text{kl}(q||p) = \tau$ is attained or when the interval $[p_{\min}, p_{\max}]$ is small enough, the bias p is found.

Moreover, REEB *et al.* (2018) found the expression of the derivatives with respect to q and τ , allowing them to derive bound minimization algorithms based on gradient descent; these derivatives are defined by

$$\frac{\partial \text{kl}(q|\psi)}{\partial q} = \frac{\ln \frac{1-q}{1-\text{k}(q|\psi)} - \ln \frac{q}{\text{k}(q|\psi)}}{\frac{1-q}{1-\text{k}(q|\psi)} - \frac{q}{\text{k}(q|\psi)}}, \quad \text{and} \quad \frac{\partial \text{kl}(q|\psi)}{\partial \psi} = \frac{1}{\frac{1-q}{1-\text{k}(q|\psi)} - \frac{q}{\text{k}(q|\psi)}}, \quad (2.16)$$

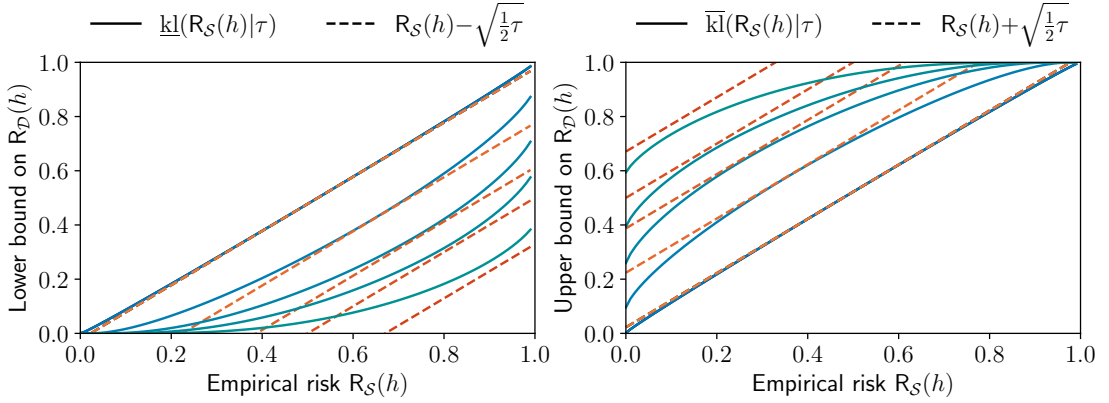


Figure 2.4. Illustration of the tightness of $\underline{\text{kl}}(\cdot)$ and $\overline{\text{kl}}(\cdot)$. The line represents the functions $\underline{\text{kl}}(\cdot|\tau)$ and $\overline{\text{kl}}(\cdot|\tau)$ for different values of $\tau \in \{0.001, 0.1, 0.3, 0.5, 0.9\}$. On the left plot, we represent the function $\underline{\text{kl}}(\cdot|\tau)$ while the right plot represents the function $\overline{\text{kl}}(\cdot|\tau)$. In the left (resp. right) plot, the lower (resp. higher) the inverting function, the smaller (resp. larger) the value of τ . Moreover, for each inverting function, we plot with the dotted lines its approximation through PINSKER'S inequality.

Algorithm 2.1 Compute $\overline{\text{kl}}(q|\tau)$ resp. $\underline{\text{kl}}(q|\tau)$ through the bisection method

Given: Bias $q \in [0, 1]$ (the empirical risk), the bound value $\tau \geq 0$

Hyperparameters: tolerance ϵ , maximal number of iterations T_{\max}

$p_{\max} \leftarrow 1$ and $p_{\min} \leftarrow q$ (resp. $p_{\max} \leftarrow q$ and $p_{\min} \leftarrow 0$)

for $t \leftarrow 1$ to T_{\max} **do**

$p = \frac{1}{2} [p_{\min} + p_{\max}]$

if $\text{kl}(q||p) = \tau$ or $(p_{\min} - p_{\max}) < \epsilon$ **then return** p

if $\text{kl}(q||p) > \tau$ **then** $p_{\max} = p$ (resp. $p_{\min} = p$)

if $\text{kl}(q||p) < \tau$ **then** $p_{\min} = p$ (resp. $p_{\max} = p$)

return p

where $\text{k}(\cdot)$ is either $\underline{\text{kl}}(\cdot)$ or $\overline{\text{kl}}(\cdot)$.

Thanks to Definition 2.3.3, we can rewrite the bound of Theorem 2.3.4 to upper-bound the expected true risk $\mathbb{E}_{h \sim \rho} R_D^\ell(h)$. With probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$, we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$\mathbb{E}_{h \sim \rho} R_D^\ell(h) \leq \overline{\text{kl}} \left(\mathbb{E}_{h \sim \rho} R_S^\ell(h) \mid \frac{1}{m} [\text{KL}(\rho||\pi) + \ln \frac{2\sqrt{m}}{\delta}] \right) \quad (2.17)$$

$$\text{and } \mathbb{E}_{h \sim \rho} R_D^\ell(h) \geq \underline{\text{kl}} \left(\mathbb{E}_{h \sim \rho} R_S^\ell(h) \mid \frac{1}{m} [\text{KL}(\rho||\pi) + \ln \frac{2\sqrt{m}}{\delta}] \right). \quad (2.18)$$

Thanks to the approximation in Equation (2.15), one can prove that the bound of

Theorem 2.3.4 is tighter than the one of Theorem 2.3.2. More precisely, if we apply Equation (2.15) to Equations (2.17) and (2.18), we retrieve Equations (2.11) and (2.12).

2.3.2 General PAC-Bayesian Bound of Bégin *et al.* (2016)

As shown previously, the KL divergence between the posterior distribution ρ and the prior distribution π is ubiquitous in PAC-Bayesian bounds. By considering this divergence, as illustrated in Theorem 2.3.1, the term $\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp[\varphi(h', S')]$ appears inevitably. Indeed, the KL divergence $\text{KL}(\rho \parallel \pi)$ between ρ and π can be expressed as a difference between two terms: $\mathbb{E}_{h \sim \rho} \varphi(h, S)$ and $\mathbb{E}_{h \sim \pi} \exp[\varphi(h, S)]$ thanks to DONSKER and VARADHAN (1976). This representation is actually used to obtain a *change of measure inequality* which quantifies how much two expectations (coming from two densities ρ and π) differ. This change of measure inequality and the expression of the KL divergence are in the following.

Proposition 2.3.1 (DONSKER-VARADHAN Variational Representation). For any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$ s.t. $\mathbb{E}_{h' \sim \pi} e^{\varphi(h', S)} < +\infty$ for all $S \in (\mathcal{X} \times \mathcal{Y})^m$, we have

$$\begin{aligned} \forall S \in (\mathcal{X} \times \mathcal{Y})^m, \quad \forall \rho \in \mathcal{M}(\mathbb{H}), \quad \mathbb{E}_{h \sim \rho} \varphi(h, S) - \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, S)} \right) &\leq \text{KL}(\rho \parallel \pi) \\ \iff \mathbb{E}_{h \sim \rho} \varphi(h, S) &\leq \text{KL}(\rho \parallel \pi) + \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, S)} \right). \end{aligned}$$

When the distribution ρ is defined as $\rho(h) = \pi(h) \frac{e^{\varphi(h, S)}}{\mathbb{E}_{h' \sim \pi} e^{\varphi(h', S)}}$, we have

$$\begin{aligned} \forall S \in (\mathcal{X} \times \mathcal{Y})^m, \quad \mathbb{E}_{h \sim \rho} \varphi(h, S) - \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, S)} \right) &= \text{KL}(\rho \parallel \pi), \\ \iff \mathbb{E}_{h \sim \rho} \varphi(h, S) &= \text{KL}(\rho \parallel \pi) + \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, S)} \right). \end{aligned}$$

Proof. Deferred to Appendix B.10. ■

As we can remark, this inequality resembles to the general bound of GERMAIN *et al.* (2009) (in Theorem 2.3.1). Hence, the change of measure inequality appears indirectly in the PAC-Bayesian bounds making the constant term $\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', S')}$ arising in the bound. Other divergences can be considered to obtain different constant terms (that are further upper-bounded). For example, OHNISHI and HONORIO (2021) prove

other change of measure inequalities for several divergences. Prior to this work, BÉGIN *et al.* (2016) derive a new general PAC-Bayesian bound with the RÉNYI divergence defined as $D_\lambda(\rho\|\pi) = \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \left[\frac{\rho(h)}{\pi(h)} \right]^\lambda \right)$ for any $\lambda > 1$. Their bound is recalled in the following theorem.

Theorem 2.3.5 (General PAC-Bayesian Bound of BÉGIN *et al.* (2016)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_*^+$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \right] \right. \\ \left. \leq D_\lambda(\rho\|\pi) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right] \right] \geq 1 - \delta. \end{aligned}$$

Proof. Deferred to Appendix B.11. ■

Unlike Theorem 2.3.1, the function $(h, \mathcal{S}) \mapsto \frac{\lambda}{\lambda-1} \ln [\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S})]$ represents the generalization gap and is upper-bounded by the RÉNYI divergence $D_\lambda(\rho\|\pi)$ and a constant term $\ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right]$. At first sight, Theorem 2.3.5 seems very different from Theorem 2.3.1, however, they are actually related. Indeed, if we replace $\varphi(h, \mathcal{S})$ by $\exp(\frac{\lambda-1}{\lambda} \varphi(h, \mathcal{S}))$ and we apply JENSEN's inequality (Theorem A.1.1) on the left-hand side, we obtain with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \leq D_\lambda(\rho\|\pi) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S}')} \right].$$

This bound is slightly looser than the one of Theorem 2.3.1: for any $\lambda > 1$ and for any distributions ρ and π , we have $\text{KL}(\rho\|\pi) \leq D_\lambda(\rho\|\pi)$ and $\lim_{\lambda \rightarrow 1^+} D_\lambda(\rho\|\pi) = \text{KL}(\rho\|\pi)$ (ERVEN and HARREMOËS, 2014). As for the general bound in Theorem 2.3.1, setting the function $\varphi(\cdot)$ and upper-bounding the term $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \varphi(h, \mathcal{S})^{\frac{\lambda}{\lambda-1}}$ gives a computable bound. Namely, Theorem 2.3.5 allows one to obtain a MCALLESTER, CATONI or a SEEGER-like PAC-Bayesian bound based on the RÉNYI divergence. For the sake of completeness, we show in Appendix B.12 the proof of the three types of bounds based on Theorem 2.3.5. Compared to Theorems 2.3.2 to 2.3.4, only the KL divergence is replaced by the looser RÉNYI divergence.

Generally speaking, all the PAC-Bayesian bounds share a common property: they bound the *expectation* of $\varphi()$ *w.r.t.* the posterior ρ . However, it might be relevant to upper-bound $\varphi()$ for a *unique* hypothesis $h \sim \rho$: this is the purpose of the *disintegrated* bounds.

2.4 Disintegrated PAC-Bayesian Bounds

Deriving PAC-Bayesian guarantees for a unique hypothesis is tedious. Indeed, the PAC-Bayesian theory is tailored for bounding the *expected* true risk. Hence, additional derivations are needed to derive a PAC-Bayesian guarantee for a unique hypothesis. See LANGFORD and SHAWE-TAYLOR (2002), LANGFORD (2005), and GERMAIN *et al.* (2009) for some examples. To avoid this issue and to get a bound for a single hypothesis, another possible solution is to sample the hypothesis $h \in \mathbb{H}$ from the posterior distribution $\rho_{\mathbb{S}} \in \mathbb{M}(\mathbb{H})$. By doing so, the gap $|\mathbb{R}_{\mathcal{D}}^{\ell}(h) - \mathbb{R}_{\mathbb{S}}^{\ell}(h)|$ can be upper-bounded with a generalization bound. This bound is defined in the following way.

Definition 2.4.1 (Disintegrated PAC-Bayesian Generalization Bound). Let $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow [0, 1]$ a generalization gap. A *disintegrated* PAC-Bayesian bound is defined such that if for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, there exists a function $\Phi : \mathbb{M}(\mathbb{H}) \times \mathbb{M}^*(\mathbb{H}) \times (0, 1] \rightarrow \mathbb{R}$ such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\phi(\mathbb{R}_{\mathcal{D}}^{\ell}(h), \mathbb{R}_{\mathbb{S}}^{\ell}(h)) \leq \Phi(\rho_{\mathbb{S}}, \pi, \delta) \right] \geq 1 - \delta,$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A and $\phi()$ is, for example, $\phi(\mathbb{R}_{\mathcal{D}}^{\ell}(h), \mathbb{R}_{\mathbb{S}}^{\ell}(h)) = |\mathbb{R}_{\mathcal{D}}^{\ell}(h) - \mathbb{R}_{\mathbb{S}}^{\ell}(h)|$.

Compared to the PAC-Bayesian bounds (see Definition 2.3.1), the expectation $\mathbb{E}_{h \sim \rho_{\mathbb{S}}}[\cdot]$ is moved outside the indicator function: this is the disintegration. Moreover, unlike the PAC-Bayesian bounds, the posterior $\rho_{\mathbb{S}}$ is obtained from an algorithm that depends on the prior $\pi \in \mathbb{M}^*(\mathbb{H})$ and the learning sample \mathbb{S} . This type of bounds has been introduced in two concurrent works, *i.e.*, CATONI (2007, Theorem 1.2.7) and BLANCHARD and FLEURET (2007). Moreover, there exists a general disintegrated bound as for the PAC-Bayesian bounds. We now present these three bounds by starting with the most general one.

2.4.1 General Disintegrated Bound of Rivasplata *et al.* (2020)

A general disintegrated PAC-Bayesian bound has been actually proposed very recently by RIVASPLATA *et al.* (2020, Theorem 1-(i)). This bound is presented below.

Theorem 2.4.1 (General Disintegrated Bound of RIVASPLATA *et al.* (2020)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \underbrace{\ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathcal{S}')) \right]}_{\Phi(\rho_{\mathcal{S}}, \pi, \delta)} \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Proof. Deferred to Appendix B.13. ■

Compared to the classical PAC-Bayesian bounds, this bound holds with high probability over the random choice of the learning sample $\mathcal{S} \sim \mathcal{D}^m$ and the hypothesis $h \sim \rho_{\mathcal{S}}$. Moreover, instead of depending on the KL divergence, the bound depends on the *disintegrated* KL divergence⁷ $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$. This term is the log ratio of the density of the prior $\pi(h)$ and the posterior distribution $\rho_{\mathcal{S}}(h)$ for the sampled hypothesis $h \sim \rho_{\mathcal{S}}$. Intuitively, the closer the posterior density $\rho_{\mathcal{S}}(h)$ to the prior density $\pi(h)$ for $h \sim \rho_{\mathcal{S}}$, the lower the disintegrated KL. As for GERMAIN *et al.* (2009)'s general bound, we need to define $\varphi(\cdot)$ and upper-bound the term $\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathcal{S}'))$ to obtain a bound that can be computed. Similarly to Theorems 2.3.1 and 2.3.5, this disintegrated bound generalizes other bounds such as the one derived by CATONI (2007, Theorem 1.2.7).

2.4.2 Disintegrated Bound of Catoni (2007)

The bound of CATONI (2007, Theorem 1.2.7), which is one of the first disintegrated bound in the literature, is recalled in the following theorem.

Theorem 2.4.2 (Disintegrated Bound of CATONI (2007)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell :$

⁷The disintegration of the KL divergence has a slightly different meaning than the disintegration of the PAC-Bayesian bound: in the former case, the integration/expectation “is removed” from the divergence.

2.4. Disintegrated PAC-Bayesian Bounds

$\mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [0, 1]$, for any $c > 0$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), -\ln \left(1 - \left[1 - e^{-c} \right] \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \right. \\ \left. \leq \frac{1}{m} \left[\ln \frac{\rho_{\mathbb{S}}(h)}{\pi(h)} + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta, \end{aligned}$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A .

Proof. Deferred to Appendix B.14. ■

After sampling the learning sample $\mathbb{S} \sim \mathcal{D}^m$ and the hypothesis $h \sim \rho_{\mathbb{S}}$, the obtained guarantee is similar to the one of Theorem 2.3.3: only the KL divergence is replaced by its disintegrated counterpart. Similar to the non-disintegrated PAC-Bayesian bound, other deviations between the true risk $R_{\mathcal{D}}^{\ell}(h)$ and the empirical risk $R_{\mathbb{S}}^{\ell}(\rho)$ can be considered. For instance, BLANCHARD and FLEURET (2007) proposed a bound on $\text{kl}(R_{\mathbb{S}}^{\ell}(h) \| R_{\mathcal{D}}^{\ell}(h))$.

2.4.3 Disintegrated Bound of Blanchard and Fleuret (2007)

The bound of BLANCHARD and FLEURET (2007) is based on another proof technique called *Occam's hammer*. For instance, they prove a SEEGER (2002)-like disintegrated PAC-Bayesian bound from this framework; it is recalled in the following theorem.

Theorem 2.4.3 (Disintegrated Bound of BLANCHARD and FLEURET (2007)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow \{0, 1\}$, for any $k > 1$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\text{kl}_+(R_{\mathbb{S}}^{\ell}(h) \| R_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[\ln \frac{k+1}{\delta} + \left(1 + \frac{1}{k} \right) \ln_+ \frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right] \right] \geq 1 - \delta,$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A , the $\ln_+(x) = \max(\ln(x), 0)$ and $\text{kl}_+(R_{\mathbb{S}}^{\ell}(h) \| R_{\mathcal{D}}^{\ell}(h)) = \text{kl}(R_{\mathbb{S}}^{\ell}(h) \| R_{\mathcal{D}}^{\ell}(h))$ if $R_{\mathbb{S}}^{\ell}(h) < R_{\mathcal{D}}^{\ell}(h)$ and 0 otherwise.

Proof. Deferred to Appendix B.15. ■

Similarly to CATONI (2007), this bound is parametrized: the tightness depends on the parameter $k > 1$. However, the optimal parameter depends on the log ratio $\ln_+ \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$ and cannot be set in advance since we don't know the sampled $\mathcal{S} \sim \mathcal{D}^m$.

2.5 Conclusion and Summary

This chapter introduces generalization bounds from the PAC-Bayesian literature. These bounds allow deriving theoretical guarantees for some machine learning models, *e.g.*, the majority vote. They are further useful to derive practical learning algorithms guaranteeing that the model is not too sensitive to overfitting; see Part II. Indeed, in Chapter 3, we derive a new self-bounding learning algorithm that minimizes a PAC-Bayesian generalization bound for the adversarially robust setting. Roughly speaking, we derive surrogates similar to Equation (2.2) to obtain guarantees and derive learning algorithms that “robustify” the majority vote. Then, Chapter 4 introduces our contributions to minimizing of the PAC-Bayesian C-Bound, which aims to minimize the true risk of the majority vote. Lastly in Part II, we introduce the stochastic majority vote (where the distribution ρ follows a Dirichlet distribution) in Chapter 5. This majority vote allows using easily a PAC-Bayesian bound on the expected true risk to learn such a classifier.

However, the main drawback of the PAC-Bayesian generalization bounds is that we bound $\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S})$ instead of $\varphi(h, \mathcal{S})$. In contrast, the disintegrated bounds allow us to bound the term $\varphi(h, \mathcal{S})$, which makes more sense if we want to deal with a unique hypothesis $h \sim \rho$. Part III shows the potential of such bounds for the analysis of the generalization of over-parametrized models. Indeed, Chapter 6 shows the first application of such bounds in practice, notably with over-parametrized models; this also leads to the derivation of new disintegrated bounds that are more appealing to optimization. Lastly, Chapter 7 introduces new perspectives based on these bounds since we can derive generalization bounds that do not depend on classical complexity measures such as the VC-Dimension or the Rademacher complexity (see Section 1.3).

PART II

**PAC-Bayesian Majority Vote:
Theory and Self-bounding Algorithms**

PAC-BAYESIAN THEORY FOR THE ROBUST MAJORITY VOTE

This chapter is based on the following paper

PAUL VIALLARD, GUILLAUME VIDOT, AMAURY HABRARD, and EMILIE MORVANT. A PAC-Bayes Analysis of Adversarial Robustness. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021b)

Contents

3.1	Introduction	82
3.2	Adversarially Robust Majority Vote	83
3.2.1	Setting	83
3.2.2	Related Works	84
3.2.3	PAC-Bayesian Adversarial Risks	86
3.3	Adversarially Robust PAC-Bayes	87
3.3.1	Relations Between the Adversarial Risks	87
3.3.2	PAC-Bayesian Bounds on the Adversarially Robust Majority Vote	89
3.3.3	From the Bounds to an Algorithm	92
3.4	Experimental Evaluation on Differentiable Decision Trees	95
3.4.1	Experiments	95
3.5	Conclusion and Summary	98

Abstract

In this chapter, we derive the first general PAC-Bayesian generalization bounds for adversarial robustness, that estimate, how much the majority vote will be robust to imperceptible perturbations in the input. Instead of deriving a worst-case analysis of the risk of the majority vote over all the possible perturbations, we leverage the PAC-Bayesian framework recalled in Chapter 2 to bound the averaged risk on the perturbations. Our theoretically founded analysis has the advantage to provide general bounds *(i)* that are valid for any kind of adversarial attacks, *(ii)* that are tight, *(iii)* that can be directly minimized in a self-bounding algorithm to obtain a robust majority vote. We empirically show this robustness on different attacks.

3.1 Introduction

In this chapter, we first formalize in Section 3.2 the notion of majority vote for the adversarial robustness setting. To do so, we adapt the majority vote (recalled in Chapter 2 for supervised learning) by assuming that the inputs can be slightly modified/perturbed to fool the prediction of the majority vote, often in a malicious way; this setting is called *adversarial robustness*. The existence of such modified inputs, known as *adversarial examples* (BIGGIO *et al.*, 2013; SZEGEDY *et al.*, 2014) and illustrated in Figure 3.1.

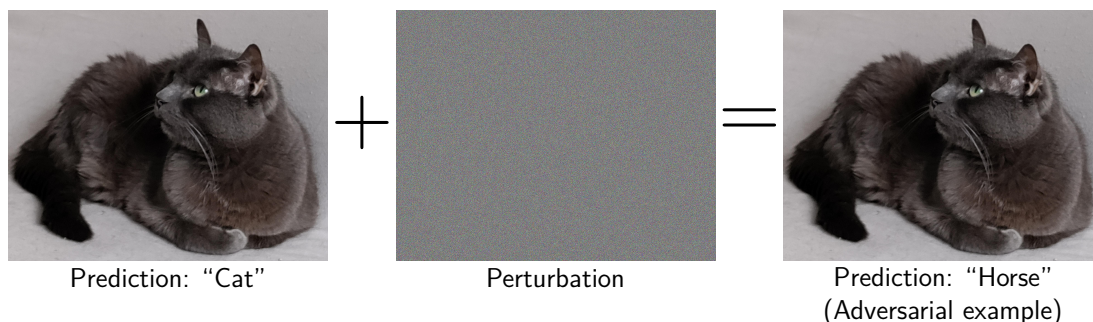


Figure 3.1. On the left, the original image is predicted correctly by the classifier as “cat”. The middle image corresponds to a perturbation/noise added to this original image. On the right, by applying the perturbation on the original image, the result is an image that looks identical to the original to the human eye, however, the prediction changes radically. This result image with the imperceptible perturbation is called an adversarial example.

The models (*i.e.*, the majority vote in our case) must be robust to these small inputs’ perturbations to better guarantees the user safety. Indeed, when machine learning models are applied to real problems, such as autonomous vehicles, the perturbations must not compromise the safety of the users. The perturbed examples is obtained from an *adversarial attack* that fools the considered model while the *adversarial defense* techniques enhance the adversarial robustness to make the attacks useless (see *e.g.*, GOODFELLOW *et al.*, 2015; PAPERNOT *et al.*, 2016; CARLINI and WAGNER, 2017; KURAKIN *et al.*, 2017; ZANTEDESCHI *et al.*, 2017; MADRY *et al.*, 2018). However, the majority votes and many other models lack guarantees on the robustness. To tackle this issue, we propose to formulate the adversarial robustness through the lens of the PAC-Bayesian theory recalled in Chapter 2; we call here our setting the *adversarially robust PAC-Bayes*.

The idea consists in considering an *averaged adversarial robustness risk* corresponding to the probability that the model misclassifies a perturbed example (this can be interpreted as an averaged risk over the perturbations). We also define an *averaged-max*

adversarial risk as the probability that there exists at least one perturbation that leads to a misclassification. These definitions, based on averaged quantities, have the advantage (i) of being suitable for the PAC-Bayesian framework and majority vote classifiers, and (ii) of being related to the classical adversarial robustness risk. Then, for each of our adversarial risks, we derive a PAC-Bayesian generalization bound that is valid to any kind of attack. From an algorithmic point of view, these bounds are directly minimizable to learn a majority vote robust in average to attacks. Since we directly minimize a generalization bound, our algorithms stand in the class of *self-bounding algorithms* (FREUND, 1998). We empirically illustrate that our framework is able to provide generalization guarantees with non-vacuous bounds for the adversarial risk while ensuring efficient protection to adversarial attacks.

Note that all the proofs of this chapter are deferred in Appendix C.

3.2 Adversarially Robust Majority Vote

3.2.1 Setting

We mainly adopt the setting of Chapter 2. We tackle *binary* classification tasks with the input space $\mathbb{X}=\mathbb{R}^d$ and the output/label space $\mathbb{Y} = \{-1, +1\}$. We assume that \mathcal{D} is a fixed but unknown distribution on $\mathbb{X}\times\mathbb{Y}$. An example is denoted by $(\mathbf{x}, y) \in \mathbb{X}\times\mathbb{Y}$. Let $\mathbb{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ be the learning sample of m examples *i.i.d.* sampled from \mathcal{D} ; We denote the distribution of such m -sample by \mathcal{D}^m . Let \mathbb{H} be a set of real-valued voters from \mathbb{X} to $[-1, +1]$. Assuming the voters set \mathbb{H} and a learning sample \mathbb{S} , our goal is to learn a well-performing ρ -weighted majority vote defined in Definition 2.2.1 by

$$\forall \mathbf{x} \in \mathbb{X}, \quad \text{MV}_\rho(\mathbf{x}) = \text{sign} \left[\mathbb{E}_{h \sim \rho} h(\mathbf{x}) \right].$$

One wants to find a ρ -weighted majority vote that minimizes the true risk $R_{\mathcal{D}}(\text{MV}_\rho)$ on \mathcal{D} defined in Definition 2.2.2 as

$$R_{\mathcal{D}}(\text{MV}_\rho) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\text{MV}_\rho(\mathbf{x}) \neq y].$$

However, in real-life applications, an imperceptible perturbation of the input can have a bad influence on the classification performance on unseen data (SZEGEDY *et al.*, 2014): the usual generalization guarantees do not stand anymore. Such an imperceptible perturbation can be modeled by a (relatively small) additive noise ϵ applied an input \mathbf{x} leading to a perturbed input $\mathbf{x} + \epsilon$. Let $b > 0$ and $\|\cdot\|$ be an arbitrary norm, the set of possible noises \mathbb{B} is defined by¹

$$\mathbb{B} = \left\{ \epsilon \in \mathbb{R}^d \mid \|\epsilon\| \leq b \right\}.$$

¹The most used norms in the set of possible noises are the ℓ_1 , ℓ_2 and ℓ_∞ -norms.

The learner aims now to find an *adversarial robust* classifier that is robust in average to all noises in \mathbb{B} over $(\mathbf{x}, y) \sim \mathcal{D}$. More formally, one wants to minimize the true adversarial risk $A_{\mathcal{D}}(\text{MV}_{\rho})$ defined in the following definition.

Definition 3.2.1 (True/Empirical Adversarial Risk). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any distribution ρ on \mathbb{H} , the *true adversarial risk* is defined as

$$A_{\mathcal{D}}(\text{MV}_{\rho}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \max_{\epsilon \in \mathbb{B}} \mathbb{I}[\text{MV}_{\rho}(\mathbf{x} + \epsilon) \neq y].$$

Since \mathcal{D} is unknown, $A_{\mathcal{D}}(\text{MV}_{\rho})$ cannot be directly computed, then one usually deals with the empirical adversarial risk defined as

$$A_{\mathcal{S}}(\text{MV}_{\rho}) = \frac{1}{m} \sum_{i=1}^m \max_{\epsilon \in \mathbb{B}} \mathbb{I}[\text{MV}_{\rho}(\mathbf{x}_i + \epsilon) \neq y_i].$$

In this chapter, our objective is to try to make the majority vote classifier MV_{ρ} robust to *adversarial attacks* that aim at finding an *adversarial example* $\mathbf{x} + \epsilon^*(\mathbf{x}, y)$ to fool $\text{MV}_{\rho}(\cdot)$ for given example (\mathbf{x}, y) , where $\epsilon^*(\mathbf{x}, y)$ is defined as

$$\epsilon^*(\mathbf{x}, y) \in \operatorname{argmax}_{\epsilon \in \mathbb{B}} \mathbb{I}[\text{MV}_{\rho}(\mathbf{x} + \epsilon) \neq y]. \quad (3.1)$$

In consequence, *adversarial defense* mechanisms often rely on the adversarial attacks by replacing the original examples (\mathbf{x}, y) with the adversarial ones $(\mathbf{x} + \epsilon^*(\mathbf{x}, y), y)$ during the learning phase; this procedure is called adversarial training. Even if there are other defenses, as we will see later, adversarial training appears to be one of the most efficient defense mechanisms (REN *et al.*, 2020). Optimizing Equation (3.1) is however intractable due to the non-convexity of MV_{ρ} induced by the sign function. The adversarial attacks of the existing frameworks in the literature (that we discuss in Section 3.2.2) aim at finding the optimal perturbation $\epsilon^*(\mathbf{x}, y)$, but, in practice, one considers an approximation of this perturbation.

3.2.2 Related Works

Adversarial Attacks/Defenses. Numerous methods² exist to solve— or approximate —the optimization of Equation (3.1). Among them, the Fast Gradient Sign Method (FGSM) of GOODFELLOW *et al.* (2015) is an attack consisting in generating a noise ϵ in the direction of the gradient of the loss function with respect to the input \mathbf{x} . KURAKIN *et al.* (2017) introduced IFGSM, an iterative version of FGSM: at

²The reader can refer to REN *et al.* (2020) for a survey on adversarial attacks and defenses.

3.2. Adversarially Robust Majority Vote

each iteration, one repeats FGSM and adds to \mathbf{x} a noise, that is the sign of the gradient of the loss with respect to \mathbf{x} . Following the same principle as IFGSM, MADRY *et al.* (2018) proposed a method based on Projected Gradient Descent (PGD) that includes a random initialization of \mathbf{x} before the optimization. Another technique known as the *Carlini and Wagner Attack* (CARLINI and WAGNER, 2017) aims at finding adversarial examples $\mathbf{x} + \epsilon^*(\mathbf{x}, y)$ that are as close as possible to the original \mathbf{x} , *i.e.*, they want an attack being the most imperceptible as possible. However, producing such imperceptible perturbation leads to a high-running time in practice. Contrary to the most popular techniques that look for a model with a low adversarial robust risk, our work stands in another line of research where the idea is to relax this worst-case risk measure by considering an *averaged* adversarial robust risk over the noises instead of a *max*-based formulation (see, *e.g.*, ZANTEDESCHI *et al.*, 2017; HENDRYCKS and DIETTERICH, 2019). Our averaged formulation is introduced in the Section 3.2.

Generalization Bounds for Adversarial Robustness. Recently, few generalization bounds for adversarial robustness have been introduced (*e.g.*, KHIM and LOH, 2018; COHEN *et al.*, 2019; MONTASSER *et al.*, 2019; PINOT *et al.*, 2019; SALMAN *et al.*, 2019; YIN *et al.*, 2019; MONTASSER *et al.*, 2020; PINOT *et al.*, 2022). KHIM and LOH, and YIN *et al.*'s results are Rademacher complexity-based bounds. The former makes use of a surrogate of the adversarial risk; the latter provides bounds in the specific case of neural networks and linear classifiers and involves an unavoidable polynomial dependence on the dimension of the input. MONTASSER *et al.* study robust PAC-learning for PAC-learnable classes with finite VC-dimension for unweighted majority votes that have been "robustified" with a boosting algorithm. However, their algorithm requires to consider all possible adversarial perturbations for each example, which is intractable in practice, and their bound also suffers from a large constant as indicated at the end of the MONTASSER *et al.* (Theorem 3.1 2019)'s proof. COHEN *et al.* provide bounds that estimate what is the minimum noise to get an adversarial example (in the case of perturbations expressed as Gaussian noise) while our results give the probability of being fooled by an adversarial example. SALMAN *et al.* leverage COHEN *et al.*'s method and adversarial training in order to get tighter bounds. Moreover, FARNIA *et al.* present margin-based bounds on the adversarial robust risk for specific neural networks and attacks (such as FGSM or PGD). While they made use of a classical PAC-Bayes bound, their result is not a PAC-Bayesian analysis and stands in the family of uniform-convergence bounds (see NAGARAJAN and KOLTER, 2019b, Ap. J for details). In this thesis, we provide PAC-Bayesian bounds for general models expressed as majority votes, their bounds are thus not directly comparable to ours.

3.2.3 PAC-Bayesian Adversarial Risks

Instead of looking for the noise from Equation (3.1) that maximizes the chance of fooling the algorithm, we propose to model the perturbation according to an example-dependent distribution. This example-dependent distribution is further used to define our new risks. First let us define $\mathcal{B}_{(x,y)}$ a distribution, on the set of possible noises \mathbb{B} , that is dependent on an example $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$. Then, we denote as \mathcal{E} the distribution on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$ defined as

$$\mathcal{E}((\mathbf{x}, y), \epsilon) = \mathcal{D}(\mathbf{x}, y) \cdot \mathcal{B}_{(x,y)}(\epsilon),$$

which further permits to generate *perturbed examples*. For a given example $(\mathbf{x}_i, y_i) \sim \mathcal{D}$, we consider a set of n perturbations sampled from $\mathcal{B}_{(x_i, y_i)}$ denoted by $\mathfrak{C}_i = \{\epsilon_j^i\}_{j=1}^n$. Then we consider as a learning set the $m \times n$ -sample $\widehat{\mathcal{S}} = \{((\mathbf{x}_i, y_i), \mathfrak{C}_i)\}_{i=1}^m \in (\mathbb{X} \times \mathbb{Y} \times \mathbb{B}^n)^m$. In other words, each $((\mathbf{x}_i, y_i), \mathfrak{C}_i) \in \widehat{\mathcal{S}}$ is sampled from a distribution that we denote by \mathcal{E}^n such that

$$\mathcal{E}^n((\mathbf{x}_i, y_i), \mathfrak{C}_i) = \mathcal{D}(\mathbf{x}_i, y_i) \cdot \prod_{j=1}^n \mathcal{B}_{(x_i, y_i)}(\epsilon_j^i).$$

Furthermore, we denote as $(\mathcal{E}^n)^m$ the empirical distribution on the perturbed learning sample consisted of m examples and n perturbations for each example. Then, inspired by the works of ZANTEDESCHI *et al.* (2017) and HENDRYCKS and DIETTERICH (2019), we define our *robustness averaged adversarial risk* as follows.

Definition 3.2.2 (Averaged Adversarial Risk). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , the averaged adversarial risk of MV_ρ is defined as

$$\begin{aligned} R_{\mathcal{E}}(\text{MV}_\rho) &= \mathbb{P}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}}(\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y) \\ &= \mathbb{E}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y]. \end{aligned}$$

The empirical averaged adversarial risk is computed on a $m \times n$ -sample $\widehat{\mathcal{S}} = \{((\mathbf{x}_i, y_i), \mathfrak{C}_i)\}_{i=1}^m$ is

$$R_{\widehat{\mathcal{S}}}(\text{MV}_\rho) = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \mathbb{I}[\text{MV}_\rho(\mathbf{x}_i + \epsilon_j^i) \neq y_i].$$

As we will show in Proposition 3.3.1, the risk $R_{\mathcal{E}}(\text{MV}_\rho)$ is considered optimistic regarding $\epsilon^*(\mathbf{x}, y)$ of Equation (3.1). Indeed, instead of taking the ϵ maximizing the loss, the ϵ is drawn from a distribution. Hence, it can lead to a non-informative risk if

3.3. Adversarially Robust PAC-Bayes

the ϵ are not informative enough to fool the classifier. To overcome this, we propose an extension that we refer as *averaged-max adversarial risk*. Note that we abuse the notation of the adversarial true risk defined in Definition 3.2.1.

Definition 3.2.3 (Averaged-Max Adversarial Risk). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , the averaged-max adversarial risk of MV_ρ is defined as

$$A_{\mathcal{E}^n}(MV_\rho) = \mathbb{P}_{((\mathbf{x}, y), \mathbb{C}) \sim \mathcal{E}^n} \left(\exists \epsilon \in \mathbb{C}, MV_\rho(\mathbf{x} + \epsilon) \neq y \right).$$

The empirical averaged-max adversarial risk computed on a $m \times n$ -sample $\widehat{\mathcal{S}} = \{((\mathbf{x}_i, y_i), \mathbb{C}_i)\}_{i=1}^m$ is

$$A_{\widehat{\mathcal{S}}}(MV_\rho) = \frac{1}{m} \sum_{i=1}^m \max_{\epsilon \in \mathbb{C}_i} \mathbb{I} [MV_\rho(\mathbf{x}_i + \epsilon) \neq y_i].$$

For an example $(\mathbf{x}, y) \sim \mathcal{D}$, instead of checking if one perturbed example $\mathbf{x} + \epsilon$ is adversarial, we sample n perturbed examples $\mathbf{x} + \epsilon_1, \dots, \mathbf{x} + \epsilon_n$ and we check if at least one example is adversarial.

3.3 Adversarially Robust PAC-Bayes

We show in Section 3.3.1 the relations between the different risks. Section 3.3.2 introduces the PAC-Bayesian bounds to assess the robustness of the majority vote.

3.3.1 Relations Between the Adversarial Risks

Proposition 3.3.1 below shows the intrinsic relationships between the classical adversarial risk $A_{\mathcal{D}}(MV_\rho)$ and our two relaxations $R_{\mathcal{E}}(MV_\rho)$ and $A_{\mathcal{E}^n}(MV_\rho)$. In particular, Proposition 3.3.1 shows that the larger number of perturbed examples n , the higher is the chance to get an adversarial example and then to be close to the adversarial risk $A_{\mathcal{D}}(MV_\rho)$.

Proposition 3.3.1 (Relations Between the Averaged Adversarial Risks). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , for any $(n, n') \in \mathbb{N}^2$, with $1 \leq n' \leq n$, we have

$$R_{\mathcal{E}}(MV_\rho) \leq A_{\mathcal{E}^{n'}}(MV_\rho) \leq A_{\mathcal{E}^n}(MV_\rho) \leq A_{\mathcal{D}}(MV_\rho). \quad (3.2)$$

Proof. Deferred to Appendix C.1. ■

The left-hand side of Equation (3.2) confirms that the averaged adversarial risk $R_{\mathcal{E}}(\text{MV}_{\rho})$ is optimistic regarding the classical $A_{\mathcal{D}}(\text{MV}_{\rho})$. Proposition 3.3.2 estimates how close $R_{\mathcal{E}}(\text{MV}_{\rho})$ can be to $A_{\mathcal{D}}(\text{MV}_{\rho})$.

Proposition 3.3.2 (Classical and Averaged Adversarial Risks). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , we have

$$A_{\mathcal{D}}(\text{MV}_{\rho}) - \text{TV}(\gamma \parallel \Gamma) \leq R_{\mathcal{E}}(\text{MV}_{\rho}),$$

where Γ and γ are distributions on $\mathbb{X} \times \mathbb{Y}$ and $\text{TV}(\gamma \parallel \Gamma) = \mathbb{E}_{(\mathbf{x}', y') \sim \Gamma} \frac{1}{2} \left| \frac{\gamma(\mathbf{x}', y')}{\Gamma(\mathbf{x}', y')} - 1 \right|$, is the Total Variation (TV) distance between γ and Γ .

The density $\Gamma(\mathbf{x}', y')$ corresponds to the probability of drawing a perturbed example $(\mathbf{x}', y') = (\mathbf{x} + \epsilon, y)$ with $((\mathbf{x}, y), \epsilon) \sim \mathcal{E}$, i.e., we have

$$\Gamma(\mathbf{x}', y') = \Pr_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} [\mathbf{x} + \epsilon = \mathbf{x}', y = y'].$$

The density $\gamma(\mathbf{x}', y')$ is the probability to draw an adversarial example $(\mathbf{x}', y') = (\mathbf{x} + \epsilon^*(\mathbf{x}, y), y)$ with $(\mathbf{x}, y) \sim \mathcal{D}$, i.e., we have

$$\gamma(\mathbf{x}', y') = \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathbf{x} + \epsilon^*(\mathbf{x}, y) = \mathbf{x}', y = y'].$$

Proof. Deferred to Appendix C.2. ■

Note that $\epsilon^*(\mathbf{x}, y)$ depends on ρ , hence γ depends on ρ . From Proposition 3.3.2 and the distributions Γ and γ , the risks $A_{\mathcal{D}}(\text{MV}_{\rho})$ and $R_{\mathcal{E}}(\text{MV}_{\rho})$ can be rewritten as

$$\begin{aligned} R_{\mathcal{E}}(\text{MV}_{\rho}) &= \Pr_{(\mathbf{x}', y') \sim \Gamma} [\text{MV}_{\rho}(\mathbf{x}') \neq y'], \\ \text{and } A_{\mathcal{D}}(\text{MV}_{\rho}) &= \Pr_{(\mathbf{x}', y') \sim \gamma} [\text{MV}_{\rho}(\mathbf{x}') \neq y']. \end{aligned}$$

Finally, Propositions 3.3.1 and 3.3.2 relate the adversarial risk $R_{\mathcal{E}}(\text{MV}_{\rho})$ to the “standard” adversarial risk $A_{\mathcal{D}}(\text{MV}_{\rho})$. Indeed, from the two propositions we obtain

$$A_{\mathcal{D}}(\text{MV}_{\rho}) - \text{TV}(\gamma \parallel \Gamma) \leq R_{\mathcal{E}}(\text{MV}_{\rho}) \leq A_{\mathcal{E}^n}(\text{MV}_{\rho}) \leq A_{\mathcal{D}}(\text{MV}_{\rho}). \quad (3.3)$$

Hence, the smaller the TV distance $\text{TV}(\gamma \parallel \Gamma)$, the closer the averaged adversarial risk $R_{\mathcal{E}}(\text{MV}_{\rho})$ is from $A_{\mathcal{D}}(\text{MV}_{\rho})$ and the more probable an example $((\mathbf{x}, y), \epsilon)$ sampled

from \mathcal{E} would be adversarial, *i.e.*, when our “averaged” adversarial example looks like a “specific” adversarial example. Moreover, Equation (3.3) justifies that the PAC-Bayesian point of view makes sense for adversarial learning with theoretical guarantees: the PAC-Bayesian guarantees we derive in the next section for our adversarial risks implies guarantees on the adversarial risk $A_{\mathcal{D}}(MV_{\rho})$.

3.3.2 PAC-Bayesian Bounds on the Adversarially Robust Majority Vote

To derive PAC-Bayesian generalization bounds on the risk $R_{\mathcal{E}}(MV_{\rho})$, respectively on $A_{\mathcal{E}^n}(MV_{\rho})$, we consider one of the classical surrogates, *i.e.*, the Gibbs risk (Definition 2.2.5), defined below in Equation (3.4), respectively Equation (3.5).

Definition 3.3.1 (Surrogates on the Averaged Adversarial Risks). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{C}$, for any hypothesis set \mathbb{H} , for any ρ on \mathbb{H} ,

$$r_{\mathcal{E}}(\rho) = \mathbb{E}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} \frac{1}{2} \left[1 - \mathbb{E}_{h \sim \rho} y h(\mathbf{x} + \epsilon) \right], \quad (3.4)$$

$$\text{and } a_{\mathcal{E}^n}(\rho) = \mathbb{E}_{((\mathbf{x}, y), \mathbb{C}) \sim \mathcal{E}^n} \frac{1}{2} \left[1 - \min_{\epsilon \in \mathbb{C}} \left(y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \right) \right]. \quad (3.5)$$

Put into words, these surrogates are expressed as the expectation over ρ of the individual risks of the voters involved in \mathbb{H} . From an algorithmic perspective, $r_{\mathcal{E}}(\rho)$ and $a_{\mathcal{E}^n}(\rho)$ have the advantages (i) of being differentiable contrary to $R_{\mathcal{E}}(MV_{\rho})$ and $A_{\mathcal{E}^n}(MV_{\rho})$, and (ii) to upper-bound to $R_{\mathcal{E}}(MV_{\rho})$ and $A_{\mathcal{E}^n}(MV_{\rho})$ as follows.

Theorem 3.3.1 (Upper Bounds on the Surrogates). For any distributions \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$ and ρ on \mathbb{H} , for any $n > 1$, we have

$$R_{\mathcal{E}}(MV_{\rho}) \leq 2r_{\mathcal{E}}(\rho), \quad \text{and} \quad A_{\mathcal{E}^n}(MV_{\rho}) \leq 2a_{\mathcal{E}^n}(\rho).$$

Proof. Deferred to Appendix C.3. ■

This theorem implies that a generalization bound on $r_{\mathcal{E}}(\rho)$, resp $a_{\mathcal{E}^n}(\rho)$ leads to a generalization bound on $R_{\mathcal{E}}(MV_{\rho})$, resp., $A_{\mathcal{E}^n}(MV_{\rho})$. Theorem 3.3.2 resp. Theorem 3.3.3 below presents our PAC-Bayesian generalization bounds for $r_{\mathcal{E}}(\rho)$ resp. $a_{\mathcal{E}^n}(\rho)$.

Theorem 3.3.2 (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathcal{B}$, for any set of voters \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathcal{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , we have

$$\text{kl}(r_{\widehat{\mathcal{S}}}(\rho) \| r_{\mathcal{E}}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{m+1}{\delta} \right], \quad (3.6)$$

$$\text{and } r_{\mathcal{E}}(\rho) \leq r_{\widehat{\mathcal{S}}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{m+1}{\delta} \right]}, \quad (3.7)$$

$$\text{where } r_{\widehat{\mathcal{S}}}(\rho) = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \frac{1}{2} \left[1 - y_i \mathbb{E}_{h \sim \rho} h(\mathbf{x}_i + \epsilon_j^i) \right].$$

Proof. Deferred to Appendix C.4. ■

It is important to mention that the empirical counterpart of $r_{\mathcal{E}}(\rho)$ is computed on $\widehat{\mathcal{S}}$ which is composed of non identically independently distributed samples, meaning that a “classical” proof technique is not applicable. The trick here is to make use of a result of RALAIVOLA *et al.* (2010) that provides a *chromatic PAC-Bayesian bound*, *i.e.*, a bound which supports non-independent data. Surprisingly, this theorem states bounds that do not depend on the number of perturbed examples n but only on the number of original examples m . The reason is that the n perturbed examples are inter-dependent (see the proof in Appendix). Note that Equation (3.6) is expressed as a SEEGER (2002)’s bound and is tighter but less interpretable than Equation (3.7) expressed as a MCALLESTER (1998)’s bound; these bounds involve the usual trade-off between the empirical risk $r_{\widehat{\mathcal{S}}}(\rho)$ and $\text{KL}(\rho \| \pi)$.

We now state a generalization bound for $a_{\mathcal{E}^n}(\rho)$. Since this value involves a minimum term, we cannot use the same trick as for Theorem 3.3.2. To bypass this issue, we use the TV distance between two “artificial” distributions on \mathbb{C}_i . Given $((\mathbf{x}_i, y_i), \mathbb{C}_i) \in \widehat{\mathcal{S}}$, let Θ_i be an arbitrary distribution on \mathbb{C}_i , and given $h \in \mathbb{H}$, let θ_i^h be a Dirac distribution on \mathbb{C}_i such that $\theta_i^h(\epsilon) = 1$ if $\epsilon = \text{argmax}_{\epsilon \in \mathbb{C}_i} \frac{1}{2} [1 - y_i h(\mathbf{x}_i + \epsilon)]$ (*i.e.*, if ϵ is maximizing the linear loss), and 0 otherwise.

Theorem 3.3.3 (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathcal{B}$, for any set of voters \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathcal{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for all $i \in \{1, \dots, m\}$, for all distributions Θ_i on \mathbb{C}_i independent from a voter $h \in \mathbb{H}$,

we have

$$a_{\mathcal{E}^n}(\rho) \leq \frac{1}{m} \mathbb{E}_{h \sim \rho} \sum_{i=1}^m \max_{\epsilon \in \mathbb{C}_i} \frac{1}{2} (1 - y_i h(\mathbf{x}_i + \epsilon)) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m}}{\delta}]} \quad (3.8)$$

$$\leq a_{\widehat{\mathcal{S}}}(\rho) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \|\Theta_i) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m}}{\delta}]}, \quad (3.9)$$

where the empirical risk $a_{\widehat{\mathcal{S}}}(\rho) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \left[1 - \min_{\epsilon \in \mathbb{C}_i} \left(y_i \mathbb{E}_{h \sim \rho} h(\mathbf{x}_i + \epsilon) \right) \right]$, and the TV distance $\text{TV}(\theta \|\Theta) = \mathbb{E}_{\epsilon \sim \Theta} \frac{1}{2} \left| \left[\frac{\theta(\epsilon)}{\Theta(\epsilon)} \right] - 1 \right|$.

Proof. Deferred to Appendix C.5. ■

To minimize the true averaged-max risk $a_{\mathcal{E}^n}(\rho)$ from Equation (3.8), we have to minimize a trade-off between $\text{KL}(\rho \|\pi)$ (*i.e.*, how much the posterior weights are close to the prior ones) and the empirical risk $\frac{1}{m} \mathbb{E}_{h \sim \rho} \sum_{i=1}^m \max_{\epsilon \in \mathbb{C}_i} \frac{1}{2} (1 - y_i h(\mathbf{x}_i + \epsilon))$. However, to compute the empirical risk, the loss for each voter and each perturbation has to be calculated and can be time-consuming. With Equation (3.9), we propose an alternative, which can be efficiently optimized using $\frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \|\Theta_i)$ and the empirical averaged-max risk $a_{\widehat{\mathcal{S}}}(\rho)$. Intuitively, Equation (3.9) can be seen as a trade-off between the empirical risk, which reflects the robustness of the majority vote, and two penalization terms: the KL term and the TV term. The KL-divergence $\text{KL}(\rho \|\pi)$ controls how much the posterior ρ can differ from the prior ones π . While the TV term $\mathbb{E}_h \text{TV}(\theta_i^h \|\Theta_i)$ controls the diversity of the voters, *i.e.*, the ability of the voters to be fooled on the same adversarial example. From an algorithmic view, an interesting behavior is that the bound of Equation (3.9) stands for all distributions Θ_i on \mathbb{C}_i . This suggests that given (\mathbf{x}_i, y_i) , we want to find Θ_i minimizing $\mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \|\Theta_i)$. Ideally, this term tends to 0 when Θ_i is close³ to θ_i^h and all voters have their loss maximized by the same perturbation $\epsilon \in \mathbb{C}_i$.

To learn a well-performing majority vote, one solution is to minimize the right-hand side of the bounds, meaning that we would like to find a good trade-off between a low empirical risk $r_{\widehat{\mathcal{S}}}(\rho)$ or $a_{\widehat{\mathcal{S}}}(\rho)$ and a low divergence between the prior weights and the learned posterior ones $\text{KL}(\rho \|\pi)$. However, the bounds of Equation (3.8) and Equation (3.6) are, in their form, not appealing for optimization. Firstly, Equation (3.6) is not directly optimizable since we upper-bound the $\text{kl}()$ function between the empirical

³Since θ_i^h is a Dirac distribution, we have $\mathbb{E}_h \text{TV}(\theta_i^h \|\Theta_i) = \frac{1}{2} \left[1 - \mathbb{E}_h \Theta_i(\epsilon_h^*) + \mathbb{E}_h \sum_{\epsilon \neq \epsilon_h^*} \Theta_i(\epsilon) \right]$, with $\epsilon_h^* = \text{argmax}_{\epsilon \in \mathbb{C}_i} \frac{1}{2} [1 - y_i h(\mathbf{x}_i + \epsilon)]$.

and true risk. To obtain an optimizable bound, we can use the $\overline{\text{kl}}()$ function introduced in Definition 2.3.3. Secondly, from an algorithmic perspective, the prior π is fixed and cannot depend on the learning sample \mathbb{S} . To overcome this issue, we propose to use the union bound by considering T priors that can be selected *a posteriori* with \mathbb{S} ; the two new bounds are presented in the following corollaries.

Corollary 3.3.1 (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any set of voters \mathbb{H} , for any $T \in \mathbb{N}_*$, for any priors' set $\{\pi_1, \dots, \pi_T\} \in \mathbb{M}^*(\mathbb{H})^T$, for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathbb{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for any $\pi \in \{\pi_1, \dots, \pi_T\} \in \mathbb{M}^*(\mathbb{H})^T$ we have

$$r_{\mathcal{E}}(\rho) \leq \overline{\text{kl}} \left(r_{\widehat{\mathbb{S}}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{T(m+1)}{\delta} \right] \right. \right). \quad (3.10)$$

Proof. Deferred to Appendix C.6. ■

Corollary 3.3.2 (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any set of voters \mathbb{H} , for any prior π on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathbb{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for all $i \in \{1, \dots, m\}$, for all distributions Θ_i on \mathbb{C}_i independent from a voter $h \in \mathbb{H}$, we have

$$a_{\mathcal{E}^n}(\rho) \leq a_{\widehat{\mathbb{S}}}(\rho) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \parallel \Theta_i) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2T\sqrt{m}}{\delta} \right]}. \quad (3.11)$$

Proof. Deferred to Appendix C.6. ■

Thanks to Corollaries 3.3.1 and 3.3.2, we now derive an algorithm that minimizes such bounds.

3.3.3 From the Bounds to an Algorithm

We are now able to derive a learning algorithm that minimizes either the bound in Equation (3.10) or Equation (3.11): it is a *self-bounding* algorithm (FREUND, 1998). We consider a finite set of voters \mathbb{H} that are differentiable and where each $h \in \mathbb{H}$ is parametrized by a weight vector \mathbf{w}^h . Inspired by MASEGOSA *et al.*, 2020, the voters of \mathbb{H} and the data-dependent prior distribution π is learned from a first learning

3.3. Adversarially Robust PAC-Bayes

set \mathcal{S}' (independent from \mathcal{S}); this is a common approach in PAC-Bayes (PARRADO-HERNÁNDEZ *et al.*, 2012; LEVER *et al.*, 2013; DZIUGAITE and ROY, 2018; DZIUGAITE *et al.*, 2021). Then, the posterior distribution is learned from the second learning set \mathcal{S} by minimizing the bounds of Corollaries 3.3.1 and 3.3.2. Concretely, we minimize an objective function that is approximated with a mini-batch $\mathcal{U} \subseteq \mathcal{S}$. The objective function to optimize Equation (3.10) *resp.* Equation (3.11) is defined as

$$G_{\mathcal{U}}(\rho) = \overline{\text{kl}} \left(r_{\mathcal{U}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{T(m+1)}{\delta} \right] \right. \right),$$

resp. $G_{\mathcal{U}}(\rho) = a_{\mathcal{U}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2T\sqrt{m}}{\delta} \right]}.$

The TV distance does not appear in the objective function, since we make the choice to set $n=1$, *i.e.*, we sample one noise per example. Indeed, when $n=1$, the value of the TV distance is 0. Note that, if we had $n > 1$ we would have to minimize it.

We propose now an adversarial training algorithm which is based on a two-step learning procedure presented in Algorithm 3.1. The first step of the algorithm aims at building the set of voters \mathbb{H} and the associated prior π , the second is dedicated to the learning of the majority vote parameter ρ by minimizing the objective function associated to our bound. These steps are presented below.

Attacking the examples. The attacks in Algorithm 3.1 differ from the attack that generates the perturbed set $\hat{\mathcal{S}}$ (to compute the bound). Indeed, at each iteration (in both steps), we attack an example with the current model while $\hat{\mathcal{S}}$ is generated with the prior majority vote MV_{π} (the output of Step 1).

Step 1. Starting from an initial prior π_0 (*e.g.*, the uniform distribution) and an initial set of voters \mathbb{H}_0 , where each voter h is parametrized by a weight vector \mathbf{w}_0^h , the objective of this step is to construct the hypothesis set \mathbb{H} and the prior distribution π to give as input to Step 2 for minimizing the bound. To do so, at each epoch t of Step 1, we learn from \mathcal{S}' an “intermediate” prior π_t on an “intermediate” hypothesis set \mathbb{H}_t consisting of voters h parametrized by the weights \mathbf{w}_t^h ; note that the optimization is done with respect to $\mathbf{w}_t = \{\mathbf{w}_t^h\}_{h \in \mathbb{H}_t}$. At each iteration of the optimizer, for each (\mathbf{x}, y) of the current mini-batch \mathcal{U} , we attack the majority vote MV_{π_t} to obtain a perturbed example $\mathbf{x} + \epsilon$. Then, we perform a forward pass in the majority vote with the perturbed examples and update the weights \mathbf{w}_t and the prior π_t according to the linear loss. To sum up, at the end of Step 1, the prior π and the hypothesis set \mathbb{H} constructed for Step 2 are the ones associated to the best epoch $t^* \in \{1, \dots, T'\}$ that permits to minimize $r_{\mathcal{S}_t}(\text{MV}_{\pi_t})$, where \mathcal{S}_t is the perturbed set obtained by attacking the majority vote MV_{π_t} with the examples of \mathcal{S} . Our selection of the prior π with \mathcal{S} may seem like “cheating”, but this remains a valid strategy since Equations (3.10) and (3.11) hold

for all prior $\pi \in \{\pi_1, \dots, \pi_T\}$.

Step 2. Starting from the prior π on \mathbb{H} and the learning set \mathbb{S} , we perform the same process as in Step 1 except that the considered objective function corresponds to the desired bound to optimize (denoted $G(\cdot)$). Note that the “intermediate” priors do not depend on \mathbb{S} , since they are learned from \mathbb{S}' : the bounds are then valid.

Algorithm 3.1 Average Adversarial Training with Guarantee

Given: disjoint learning samples \mathbb{S} and \mathbb{S}' , initial prior π_0 on \mathbb{H}_0 (with \mathbf{w}_0), the objective function $G(\cdot)$

Hyperparameters: number of epochs T, T' , the attack function

Step 1 – Prior and Voters’ Set Construction

for $t \leftarrow 1$ to T' **do**

$\pi_t \leftarrow \pi_{t-1}$ and $\mathbb{H}_t \leftarrow \mathbb{H}_{t-1}$ ($\mathbf{w}_t \leftarrow \mathbf{w}_{t-1}$)

for all mini-batch $\mathbb{U} \subseteq \mathbb{S}'$ **do**

$\mathbb{U} \leftarrow$ Attack MV_{π_t} with the examples (\mathbf{x}, y) in the mini-batch \mathbb{U}

Update π_t with $\nabla_{\pi_t} r_{\mathbb{U}}(\pi_t)$

Update \mathbf{w}_t with $\nabla_{\mathbf{w}_t} r_{\mathbb{U}}(\pi_t)$

$\mathbb{S}_t \leftarrow$ Attack MV_{π_t} with the examples of \mathbb{S}

$(\pi, \mathbb{H}) \leftarrow (\pi_{t^*}, \mathbb{H}_{t^*})$ with $t^* \leftarrow \operatorname{argmin}_{t' \in \{1, \dots, t\}} r_{\mathbb{S}_{t'}}(\pi_{t'})$

Step 2 – Bound Minimization

$\rho_0 \leftarrow \pi$

for $t \leftarrow 1$ to T **do**

$\rho_t \leftarrow \rho_{t-1}$

for all mini-batch $\mathbb{U} \subseteq \mathbb{S}$ **do**

$\mathbb{U} \leftarrow$ Attack MV_{π} with the examples (\mathbf{x}, y) in the mini-batch \mathbb{U}

Update ρ_t with $\nabla_{\rho_t} G_{\mathbb{U}}(\rho_t)$

$\mathbb{S}_t \leftarrow$ Attack MV_{π} with the examples of \mathbb{S}

$\rho \leftarrow \rho_{t^*}$ with $t^* \leftarrow \operatorname{argmin}_{t' \in \{1, \dots, t\}} G_{\mathbb{S}_{t'}}(\rho_{t'})$

3.4 Experimental Evaluation on Differentiable Decision Trees

3.4.1 Experiments

In this section, we empirically illustrate that our PAC-Bayesian framework for adversarial robustness is able to provide generalization guarantees with non-vacuous bounds for the adversarial risk.

Setting. We stand in a white-box setting meaning that the attacker knows the voters set \mathbb{H} , the prior distribution π , and the posterior one ρ . The set of voters is composed of 25 differentiable decision trees (KONTSCHIEDER *et al.*, 2016); see Appendix C.7 for more details. We empirically study two attacks with the ℓ_2 -norm and ℓ_∞ -norm: the Projected Gradient Descent (PGD, MADRY *et al.* (2018)) and the iterative version of FGSM (IFGSM, KURAKIN *et al.* (2017)). We fix the number of iterations at $k=20$ and the step size at $\frac{b}{k}$ for PGD and IFGSM (where $b=1$ for ℓ_2 -norm and $b=0.1$ for ℓ_∞ -norm). One specificity of our setting is that we deal with the perturbation distribution $\mathcal{B}_{(x,y)}$. However, in order to obtain valid bounds, $\mathcal{B}_{(x,y)}$ must be defined *a priori*. Since the prior π is defined *a priori* as well, $\mathcal{B}_{(x,y)}$ can depend on π . Hence, $\mathcal{B}_{(x,y)}$ boils down to generating a perturbed example $(\mathbf{x}+\epsilon, y)$ by attacking the prior majority vote MV_π . Based on this fact, we propose PGD_U and IFGSM_U , two variants of PGD and IFGSM. To attack an example with PGD_U or IFGSM_U we proceed with the following steps.

- (i) We attack the prior majority vote MV_π with the attack PGD or IFGSM: we will obtain a first perturbation ϵ' ;
- (ii) We sample n uniform noises ζ_1, \dots, ζ_n between -10^{-2} and $+10^{-2}$;
- (iii) We set the i -th perturbation as $\epsilon_i = \epsilon' + \zeta_i$.

Note that, for PGD_U and IFGSM_U , after one attack we end up with $n=100$ perturbed examples. For Algorithm 3.1, when these attacks are used as a defense mechanism, we set $n=1$. This makes sound since our adversarial training is iterative, we do not need to sample numerous perturbations for each example: we sample a new perturbation each time the example is forwarded through the decision trees. We additionally consider a naive defense referred to as UNIF that only adds a noise uniformly such that the ℓ_p -norm of the added noise is lower than b .

We study the following scenarios of defense/attack: they correspond to all the pairs (Defense, Attack) belonging to the set $\{\text{---}, \text{UNIF}, \text{PGD}, \text{IFGSM}\} \times \{\text{---}, \text{PGD}, \text{IFGSM}\}$ for the baseline, and $\{\text{---}, \text{UNIF}, \text{PGD}_U, \text{IFGSM}_U\} \times \{\text{---}, \text{PGD}_U, \text{IFGSM}_U\}$, where “---” means that we do not defend, *i.e.*, the attack returns the original example (note that PGD_U

and IFGSM_U when “Attack without U” refers to PGD and IFGSM for computing the classical adversarial risk).

Datasets and algorithm description. We perform our experiment on six binary classification tasks from MNIST (LECUN *et al.*, 1998) (1vs7, 4vs9, 5vs6) and Fashion MNIST (XIAO *et al.*, 2017) (Coat vs Shirt, Sandal vs Ankle Boot, Top vs Pullover). We decompose the learning set into two disjoint subsets \mathcal{S}' of around 7,000 examples (to learn the prior and the voters) and \mathcal{S} of exactly 5,000 examples (to learn the posterior). We keep as test set \mathbb{T} the original test set that contains around 2,000 examples. Moreover, we need a perturbed test set, denoted by $\hat{\mathbb{T}}$, to compute our averaged(-max) adversarial risks. Depending on the scenario, $\hat{\mathbb{T}}$ is constructed from \mathbb{T} by attacking the prior model MV_π with PGD_U or IFGSM_U with $n=100$. We run our Algorithm 3.1 for Equation (3.6) (Theorem 3.3.2), respectively Equation (3.9) (Theorem 3.3.3), and we compute our risk $R_{\hat{\mathcal{T}}}(MV_\rho)$, respectively $A_{\hat{\mathcal{T}}}(MV_\rho)$, the bound value and the usual adversarial risk associated to the model learned $A_{\mathcal{T}}(MV_\rho)$. Note that, during the evaluation of the bounds, we have to compute our relaxed adversarial risks $R_{\hat{\mathcal{S}}}(MV_\rho)$ and $A_{\hat{\mathcal{S}}}(MV_\rho)$ on \mathcal{S} . For Step 1, the initial prior P_0 is fixed to the uniform distribution, the initial set of voters \mathbb{H}_0 is constructed with weights initialized with Xavier Initializer (GLOROT and BENGIO, 2010) and bias initialized at 0. During Step 2, to optimize the bound, we fix the confidence parameter $\delta=0.05$, and we consider two settings for \mathbb{H} : the set \mathbb{H} as it is output by Step 1, and the set $\mathbb{H}^{\text{SIGN}} = \{h'(\cdot) = \text{sign}(h(\cdot)) \mid h \in \mathbb{H}\}$ for which the theoretical results are still valid. Note that for all attacks on the majority votes with \mathbb{H}^{SIGN} , in order to be differentiable with respect to the input, we remove the $\text{sign}()$ function on the voters’ outputs during the attacks. For the two steps, we use Adam optimizer (KINGMA and BA, 2015) for $T=T'=20$ epochs with a learning rate at 10^{-2} and a batch size at 64.

Analysis of the results. For the sake of readability, we exhibit the detailed results for one task (MNIST:1vs7) and all the pairs (Defense,Attack) with ℓ_2 -norm in Table 3.1, and we report in Figure 3.2 the influence of the TV term in the bound of Theorem 3.3.3 (Equation (3.9)). The detailed results on the other tasks are reported in Appendix C.8. We provide in Figure 3.3 an overview of the results we obtained on all the tasks for the pairs (Defense,Attack) where “Defense=Attack” and with \mathbb{H}^{SIGN} .

First of all, from Table 3.1 the bounds of Theorem 3.3.2 are tighter than the ones of Theorem 3.3.3: this is an expected result since we showed that the averaged-max adversarial risk $A_{\mathcal{E}^n}(MV_\rho)$ is more pessimistic than its averaged counterpart $R_{\mathcal{E}}(MV_\rho)$. Note that the bound values of Equation (3.8) are tighter than the ones of Equation (3.9). This is expected since Equation (3.8) is a lower bound on Equation (3.9).

3.4. Experimental Evaluation on Differentiable Decision Trees

Table 3.1. Test risks and bounds for MNIST:1vs7 with $n=100$ perturbations for all pairs (Defense, Attack) with the two voters’ set \mathbb{H} and \mathbb{H}^{SIGN} . The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} . To quantify the gap between our risks and the classical definition we put in italic the risk of our models against the classical attacks: we replace PGD_U and IFGSM_U by PGD or IFGSM (i.e., we did not sample from the uniform distribution). Since Eq. (3.9) upper-bounds Eq. (3.8) thanks to the TV term, we compute the two bound values of Theorem 3.3.3.

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_{\rho})$		$R_{\hat{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.2		$A_{\mathcal{T}}(MV_{\rho})$		$A_{\hat{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.005	.005	.005	.005	.017	.019	.005	.005	.005	.005	.099	0.100	.099	.100
—	PGD_U	.245	.255	.263	.276	.577	.448	.315	.313	.325	.326	.801	1.667	.684	.515
—	IFGSM_U	.084	.086	.066	.080	.170	.185	.117	.113	.106	.110	.356	1.431	.286	.251
UNIF	—	.005	.005	.005	.005	.018	.019	.005	.005	.005	.005	.099	0.100	.099	.100
UNIF	PGD_U	.151	.146	.151	.158	.355	.292	.183	.178	.190	.189	.531	1.620	.454	.355
UNIF	IFGSM_U	.063	.061	.031	.035	.088	.114	.071	.070	.056	.054	.248	1.405	.200	.186
PGD_U	—	.006	.007	.006	.007	.023	.024	.006	.007	.006	.007	.102	0.103	.102	.103
PGD_U	PGD_U	.028	.030	.021	.025	.065	.064	.028	.029	.025	.028	.143	1.389	.137	.136
PGD_U	IFGSM_U	.021	.022	.013	.016	.043	.045	.022	.022	.018	.019	.125	1.362	.121	.119
IFGSM_U	—	.006	.007	.006	.007	.019	.021	.006	.007	.006	.007	.100	0.102	.100	.102
IFGSM_U	PGD_U	.040	.041	.033	.035	.086	.094	.040	.039	.040	.038	.184	1.368	.166	.163
IFGSM_U	IFGSM_U	.021	.022	.013	.014	.039	.049	.021	.022	.018	.021	.131	1.329	.122	.123

Second, the bounds with \mathbb{H}^{SIGN} are all informative (lower than 1) and give insightful guarantees for our models. For Theorem 3.3.3 (Equation (3.9)) with \mathbb{H} , while the risks are comparable to the risks obtained with \mathbb{H}^{SIGN} , the bound values are greater than 1, meaning that we have no more guarantee on the model learned. As we can observe in Figure 3.2, this is due to the TV term involved in the bound. Considering \mathbb{H}^{SIGN} when optimizing $A(\cdot)$ helps to control the TV term. Even if the bounds are non-vacuous for Theorem 3.3.2 with \mathbb{H} , the best models with the best guarantees are obtained with \mathbb{H}^{SIGN} . This is confirmed by the columns $A_{\mathcal{T}}(MV_{\rho})$ that are always worse than $R_{\hat{\mathcal{T}}}(MV_{\rho})$ and mostly worse than $A_{\hat{\mathcal{T}}}(MV_{\rho})$ with \mathbb{H}^{SIGN} . The performance obtained with \mathbb{H}^{SIGN} can be explained by the fact that the sign “saturates” the output of the voters which makes the majority vote more robust to noises. Thus, we focus the rest of the analysis on results obtained with \mathbb{H}^{SIGN} .

Third, we observe that the naive defense UNIF is able to improve the risks $R_{\hat{\mathcal{T}}}(MV_{\rho})$ and $A_{\hat{\mathcal{T}}}(MV_{\rho})$, but the improvement with the defenses based on PGD_U and IFGSM_U is much more significant specifically against a PGD_U attack (up to 13 times better). We observe the same phenomenon for both bounds (Theorems 3.3.2 and 3.3.3). This is an interesting fact because this behavior confirms that we are able to learn models

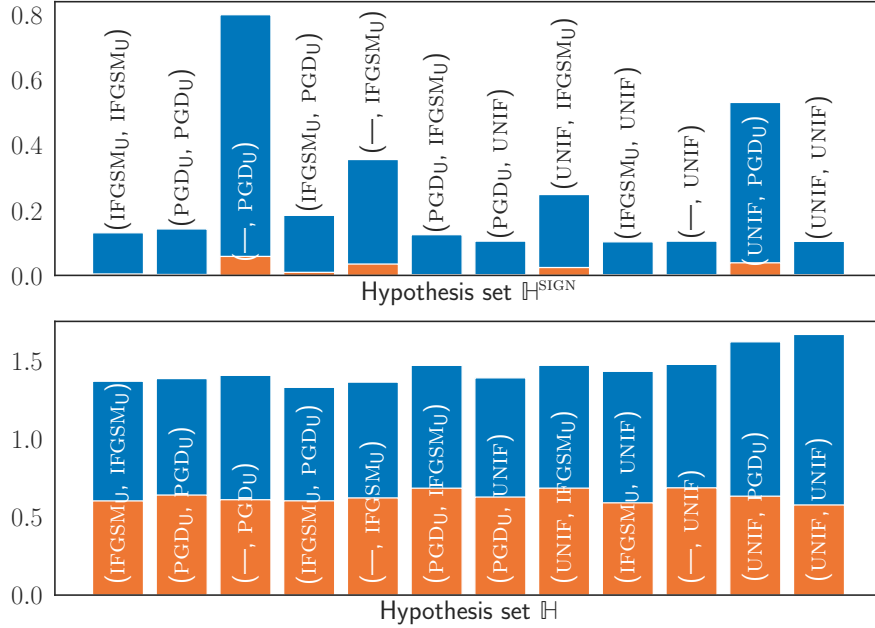


Figure 3.2. Visualization of the impact of the TV term in Equation (3.9). The top, respectively the bottom, bar plot show the bounds for the set of voters \mathbb{H}^{SIGN} , respectively \mathbb{H} . We plot the bounds for all the scenarios of Table 3.1 that use the TV distance, *i.e.*, all except the pairs $(\cdot, \text{—})$. In orange we represent the value of the TV term while in blue we represent all the remaining terms of the bound.

that are robust against the attacks tested with theoretical guarantees.

Lastly, from Figure 3.3 and Table 3.1, it is important to notice that the gap between the classical risk and our relaxed risks is small, meaning that our relaxation are not too optimistic. Despite the pessimism of the classical risk $A_{\mathcal{T}}(MV_{\rho})$, it remains consistent with our bounds, *i.e.*, it is lower than the bounds. In other words, in addition to giving upper bounds for our risks $R_{\hat{\mathcal{T}}}(MV_{\rho})$ and $A_{\hat{\mathcal{T}}}(MV_{\rho})$, our bounds give non-vacuous guarantees on the classical risks $A_{\mathcal{T}}(MV_{\rho})$.

3.5 Conclusion and Summary

To the best of our knowledge, our work is the first one that studies adversarial robustness through the PAC-Bayesian theory for the ρ -weighted majority vote. We have started by formalizing a new adversarial robustness setting (for binary classification) with our new averaged risks. This formulation allowed us to derive PAC-Bayesian generalization bounds on the majority vote’s adversarial risk. We illustrated the usefulness of this setting on the training of (differentiable) decision trees. The main objective

3.5. Conclusion and Summary

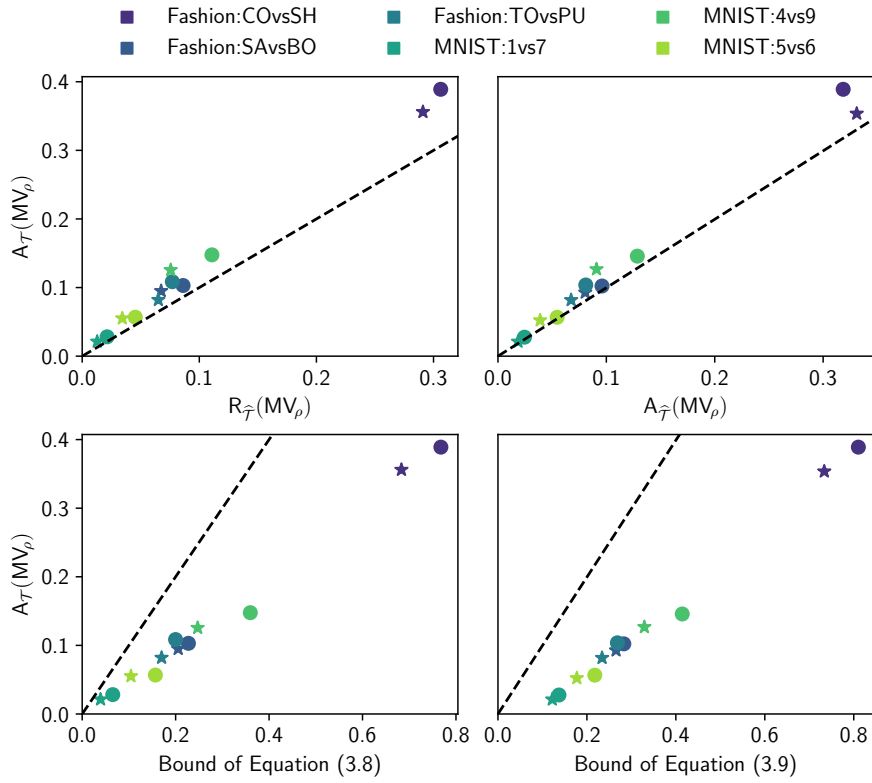


Figure 3.3. Visualization of the risk and bound values for “Defense=Attack” when the set of voters is \mathbb{H}^{SIGN} . Results obtained with the PGD_U , respectively IFGSM_U , defense are represented by a star \star , respectively a circle \bullet (reminder: $A_{\mathcal{T}}(MV_{\rho})$ is computed with a PGD , respectively IFGSM , attack). The dashed line corresponds to bisecting line $y=x$. For $R_{\hat{\mathcal{T}}}(MV_{\rho})$ and $A_{\hat{\mathcal{T}}}(MV_{\rho})$, the closer the datasets are to the bisecting line, the more accurate our relaxed risk is compared to the classical adversarial risk $A_{\mathcal{T}}(MV_{\rho})$. For the bounds, the closer the datasets are to the bisecting line, the tighter the bound.

of this work was to provide some theoretical guarantees for adversarial training. Our aim was not to improve directly the performance of the state of the art which would require a dedicated work.

One of the limitation of this work is that this PAC-Bayesian analysis holds for a majority vote only. One perspective is then to extend this work for other classifiers such as neural networks. To do so, we can leverage the disintegrated PAC-Bayesian generalization bounds (introduced in Section 2.4) to bound the adversarial risk of a *single classifier* belonging to \mathbb{H} and sampled from ρ . Another perspective of this work is to continue analyzing the adversarial true risk of the majority vote. Indeed, the C-Bound (LACASSE *et al.*, 2006) or the joint error (MASEGOSA *et al.*, 2020) (introduced in Section 2.2.2) adapted to our averaged risks might be a better choice to obtain a self-bounding algorithm since the it is more precise than twice the Gibbs risk (see Theorem 2.2.4). Moreover, thanks to the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017) (recalled in Definition 2.2.4), the multi-class case can be considered.

In the next chapters, we consider the classical supervised setting, *i.e.*, where no inputs' perturbations are added. For instance, Chapters 4 and 5 focus on deriving new learning algorithms for the majority vote in the classical supervised setting. More precisely, we derive, in Chapter 4, the first self-bounding algorithms based on the minimization of the PAC-Bayesian C-Bound (*i.e.* PAC-Bayesian bounds on the C-Bound). Chapter 5 presents a algorithm to minimize the risk of a PAC-Bayesian stochastic majority vote (where the distribution ρ are sampled from another distribution).

SELF-BOUNDING ALGORITHMS FOR THE MAJORITY VOTE

This chapter is based on the following paper

PAUL VIALARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. Self-bounding Majority Vote Learning Algorithms by the Direct Minimization of a Tight PAC-Bayesian C-Bound. *Machine Learning and Knowledge Discovery in Databases (ECML PKDD)*. (2021a)

Contents

4.1	Introduction	102
4.2	Setting	103
4.3	State of the Art: PAC-Bayesian Bounds for the Majority Vote	104
4.3.1	PAC-Bayesian Bound on the Gibbs Risk	104
4.3.2	PAC-Bayesian Bound on the Joint Error	105
4.3.3	PAC-Bayesian C-Bound of ROY <i>et al.</i>	105
4.3.4	PAC-Bayesian C-Bound of GERMAIN <i>et al.</i>	106
4.3.5	PAC-Bayesian C-Bound of LACASSE <i>et al.</i>	107
4.4	Contribution: Algorithms based on the PAC-Bayesian C-Bounds	108
4.4.1	Algorithm based on Equation (4.2)	108
4.4.2	Algorithm based on Equation (4.3)	110
4.4.3	Algorithm based on Theorem 4.3.5	111
4.5	Experiments	114
4.5.1	Setting	114
4.5.2	Analysis of the Results	122
4.6	Conclusion and Summary	123

Abstract

As we have seen in Chapter 2, the C-Bound is an insightful upper bound on the risk of a majority vote classifier. Learning algorithms in the literature minimize the empirical version of the C-Bound, instead of explicit PAC-Bayesian generalization bounds. In this chapter, we derive self-bounding majority vote learning algorithms to directly optimize PAC-Bayesian guarantees on the C-Bound. Our algorithms based on gradient descent are scalable and lead to accurate predictors paired with non-vacuous guarantees.

4.1 Introduction

In this chapter, we introduce new learning algorithms for the majority vote in the context of supervised classification. The goal of this algorithm is to minimize the true risk of the majority vote. To do so, one way to minimize such a risk is to minimize the empirical C-Bound (BREIMAN, 2001; LACASSE *et al.*, 2006) introduced in Section 2.2.2 and estimated on the learning sample \mathcal{S} . This bound has the advantage of involving the performance of the individual voters and the diversity in the voters' set. Indeed, these elements are important when one learns a combination (DIETTERICH, 2000; KUNCHEVA, 2014). A good majority vote is made up of voters that are “sufficiently diverse”.

Previous algorithms have been developed to minimize the *empirical* C-Bound such as MinCq (ROY *et al.*, 2011), P-MinCq (BELLET *et al.*, 2014), CqBoost (ROY *et al.*, 2016), or CB-Boost (BAUVIN *et al.*, 2020). ROY *et al.* (2011) first proposed MinCq which consist in minimizing a quadratic problem to learn a majority vote. MinCq considers a specific voters' set to regularize the minimization process; the algorithm P-MinCq generalizes MinCq by allowing prior distributions different from the uniform one. One drawback of MinCq and P-MinCq is that the optimization problem is not scalable to large datasets. Lately, BAUVIN *et al.* (2020) proposed CB-Boost that minimizes in a boosting-based procedure with the advantage to be more scalable while obtaining sparser majority vote. However, since both MinCq and CB-Boost minimize the empirical C-Bound, the PAC-Bayesian generalization bound associated with their learned majority vote predictors can be vacuous. Note that CB-Boost has been proposed to improve another algorithm called CqBoost (ROY *et al.*, 2016). Despite being empirically efficient and justified by theoretical analyses based on the C-Bound, all these methods minimize the empirical C-Bound and not directly a PAC-Bayesian generalization bound on the C-Bound. This can lead to vacuous generalization bound values and, thus, to poor risk certificates. When it comes to deriving a learning algorithm that directly minimizes a PAC-Bayesian bound, it is mentioned in the literature that optimizing a PAC-Bayesian bound on the C-bound is not trivial (LORENZEN *et al.*, 2019; MASEGOSA *et al.*, 2020). This underlines the need for other majority vote learning algorithms based on the C-Bound, which motivates our contributions of Section 4.4.

We cover in this chapter three different PAC-Bayesian viewpoints on generalization bounds for the C-Bound (SEEGER, 2002; MCALLESTER, 2003; LACASSE *et al.*, 2006). We derive three algorithms from these three views to optimize generalization bounds on the C-Bound. By doing so, we achieve *self-bounding algorithms* (FREUND, 1998): the predictor returned by the learner comes with a statistically valid risk upper bound. Importantly, our algorithms rely on fast gradient descent procedures. As far as we know, this is the first work that proposes both efficient algorithms for C-Bound

optimization and non-trivial risk-bound values.

We provide all the proofs in Appendix D for completeness.

4.2 Setting

We stand in supervised classification by following Chapter 2. In this context, let $\mathcal{X} \subseteq \mathbb{R}^d$ be a d -dimensional input space, and \mathbb{Y} the label space defined by $\mathbb{Y} = \{-1, +1\}$ (in binary classification) or $\mathbb{Y} = \{1, 2, \dots, l\}$ (in multi-class classification). We assume an unknown data distribution \mathcal{D} on $\mathcal{X} \times \mathbb{Y}$ and a learning sample $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where each example (\mathbf{x}_i, y_i) is drawn *i.i.d.* from \mathcal{D} ; we denote by $\mathcal{S} \sim \mathcal{D}^m$ the random draw of such a sample. Given \mathbb{H} a hypothesis set constituted by voters $h : \mathcal{X} \rightarrow \mathbb{Y}$, and \mathcal{S} , the learner aims to find a weighted combination of the voters from \mathbb{H} ; a distribution models the weights on \mathbb{H} . To learn such a combination in the PAC-Bayesian framework, we assume a *prior* distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , and—after the observation of \mathcal{S} —we learn a *posterior* distribution $\rho \in \mathcal{M}(\mathbb{H})$ on \mathbb{H} . More precisely, we aim to learn a well-performing classifier that is expressed as a ρ -*weighted majority vote* MV_ρ defined as

$$\forall \mathbf{x} \in \mathcal{X}, \quad MV_\rho(\mathbf{x}) \triangleq \operatorname{argmax}_{y' \in \mathbb{Y}} \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y'] = \operatorname{argmax}_{y' \in \mathbb{Y}} \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}) = y'].$$

We thus want to learn MV_ρ that commits as few errors as possible on unseen data from \mathcal{D} , *i.e.*, that leads to a low true risk $R_{\mathcal{D}}(MV_\rho)$ under the 01-loss defined as

$$R_{\mathcal{D}}(MV_\rho) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[MV_\rho(\mathbf{x}) \neq y] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [MV_\rho(\mathbf{x}) \neq y].$$

Since the majority vote's risk is not appealing for optimization (because its gradient is zero everywhere), some surrogates have been introduced (see Section 2.2.2). For instance, the Gibbs risk (Definition 2.2.5) is the average risk of the voters and is defined by

$$r_{\mathcal{D}}(\rho) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}, h \sim \rho} [h(\mathbf{x}) \neq y].$$

Unlike the Gibbs risk, the disagreement (Definition 2.2.7) defined as

$$d_{\mathcal{D}}(\rho) \triangleq 2 \cdot \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) = y],$$

takes the diversity of the voters into account. Moreover, the joint error (Definition 2.2.6) can be seen as a trade-off between these two quantities. It is defined as

$$\begin{aligned} e_{\mathcal{D}}(\rho) &\triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) \neq y] \\ &= r_{\mathcal{D}}(\rho) - \frac{1}{2} d_{\mathcal{D}}(\rho). \end{aligned}$$

By combining these surrogates, one can prove an upper-bound on the majority vote true risk called the C-Bound (Theorem 2.2.3) and defined as

$$\begin{aligned} R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq 1 - \frac{(1 - 2r_{\mathcal{D}}(\rho))^2}{1 - 2d_{\mathcal{D}}(\rho)} \\ &= 1 - \frac{(1 - [2e_{\mathcal{D}}(\rho) + d_{\mathcal{D}}(\rho)])^2}{1 - 2d_{\mathcal{D}}(\rho)}. \end{aligned}$$

However, these surrogates and the C-Bound are not computable because the distribution \mathcal{D} is considered *unknown*. Hence, we need to use PAC-Bayesian generalization bounds in order to upper-bound the majority vote's true risk with a C-Bound based on the empirical counterparts of these surrogates. Combined with the C-Bound, the PAC-Bayesian theory offers a natural way to analyze the risk of the majority vote. The principal PAC-Bayesian bounds for the majority vote are recalled in the next section.

4.3 State of the Art: PAC-Bayesian Bounds for the Majority Vote

This section recalls different PAC-Bayesian bounds upper-bounding the majority vote's true risk. In particular, we remind two PAC-Bayesian bounds used based on two surrogates from Section 2.2.2: Gibbs risk $r_{\mathcal{D}}(\rho)$ and the joint error $e_{\mathcal{D}}(\rho)$. Additionally, we recall three PAC-Bayesian bound on the C-Bound, that we call *PAC-Bayesian C-Bound*, which is key in our contribution of this chapter. Note that the PAC-Bayesian C-Bounds were initially developed for the binary setting but the extension for the multi-class is direct with the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017).

4.3.1 PAC-Bayesian Bound on the Gibbs Risk

One PAC-Bayesian bound, originally derived by GERMAIN *et al.* (2015), is based on the Gibbs risk $r_{\mathcal{D}}(\rho)$. It is recalled in the following theorem.

Theorem 4.3.1 (PAC-Bayesian Bound Based on the Gibbs Risk). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have

$$\begin{aligned} \forall \rho \in \mathbb{M}(\mathbb{H}), \quad r_{\mathcal{D}}(\rho) &\leq \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right), \\ \text{and } \forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq 2 \left[\overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right) \right], \end{aligned} \quad (4.1)$$

where $\bar{\text{kl}}(q|\tau) \triangleq \max \left\{ p \in (0, 1) \mid \text{kl}(q||p) \leq \tau \right\}$ (see Section 2.3.1.3).

Proof. Deferred to Appendix D.1. ■

However, since the Gibbs risk does not consider the voters' correlation, the majority votes obtained by minimizing this bound do not perform well in practice. Hence, other PAC-Bayesian bounds have therefore been derived to address this issue.

4.3.2 PAC-Bayesian Bound on the Joint Error

Another PAC-Bayesian bound based on the joint error $e_{\mathcal{D}}(\rho)$ can be derived (GERMAIN *et al.*, 2015, Theorem 25). Compared to Theorem 4.3.1, the bound of Theorem 4.3.2 takes better into account the voters' diversity (see Section 2.2.2).

Theorem 4.3.2 (PAC-Bayesian Bound Based on the Joint Error). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have

$$\begin{aligned} \forall \rho \in \mathbb{M}(\mathbb{H}), \quad e_{\mathcal{D}}(\rho) &\leq \bar{\text{kl}} \left(e_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[2 \text{KL}(\rho||\pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right), \\ \text{and } \forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq 4 \left[\bar{\text{kl}} \left(e_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[2 \text{KL}(\rho||\pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right]. \end{aligned}$$

Proof. Deferred to Appendix D.2. ■

Note that a looser bound based on the $\text{kl}()$ relaxation of THIEMANN *et al.* (2017) is presented by MASEGOSA *et al.* (2020). However, from Theorem 2.2.4, we know that there is a tighter bound on the majority vote's risk: the C-Bound (BREIMAN, 2001; LACASSE *et al.*, 2006).

4.3.3 PAC-Bayesian C-Bound of Roy *et al.*

PAC-Bayesian bounds can be used jointly with the C-Bound to obtain a computable bound on the majority vote's true risk; we call such a bound a *PAC-Bayesian C-Bound*. The most intuitive and interpretable PAC-Bayesian C-Bound has been derived by ROY *et al.* (2016) and LAVIOLETTE *et al.* (2017). The first proof of this PAC-Bayesian bound has been developed by (ROY *et al.*, 2016) in the binary setting; it has

been extended to the multi-class setting by LAVIOLETTE *et al.* (2017, Theorem 3). It consists in upper-bounding separately the Gibbs risk $r_{\mathcal{D}}(\rho)$ and the disagreement $d_{\mathcal{D}}(\rho)$ with the MCALLESTER's PAC-Bayesian bound (Theorem 2.3.2). This intuitive PAC-Bayesian bound is recalled in the following theorem.

Theorem 4.3.3 (PAC-Bayesian C-Bound of ROY *et al.* (2016)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, r_{\mathcal{S}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \right]\right)^2}{\underbrace{1 - 2 \max \left[0, d_{\mathcal{S}}(\rho) - \sqrt{\frac{1}{2m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \right]}_{\triangleq C_{\mathcal{S}}^{\text{M}}(\rho)}}. \quad (4.2)$$

Proof. Deferred to Appendix D.3. ■

While there is no algorithm that directly minimizes Theorem 4.3.3, this kind of interpretable bound can be seen as a justification of the optimization of $r_{\mathcal{S}}(\rho)$ and $d_{\mathcal{S}}(\rho)$ in the empirical C-Bound such as for MinCq (ROY *et al.*, 2011) or CB-Boost (BAUVIN *et al.*, 2020). In Section 4.4.1, we derive the first algorithm to directly minimize it. However, this PAC-Bayesian C-Bound can have a severe disadvantage with a small m and a Gibbs risk close to $\frac{1}{2}$: even for a $\text{KL}(\rho \parallel \pi)$ close to 0, the value of the PAC-Bayesian C-Bound will be close to 1. To overcome this drawback, one solution is to follow another tighter PAC-Bayesian bound, the one proposed by SEEGER (2002) (Theorem 2.3.4). Actually, we further recall two bounds based on this approach: the first one in Theorem 4.3.4 involves the Gibbs risk $r_{\mathcal{S}}(\rho)$ and the disagreement $d_{\mathcal{S}}(\rho)$ (as Theorem 4.3.3) and the second one in Theorem 4.3.5 involves the joint error $e_{\mathcal{S}}(\rho)$ and the disagreement $d_{\mathcal{S}}(\rho)$.

4.3.4 PAC-Bayesian C-Bound of Germain *et al.*

The PAC-Bayesian generalization bounds based on the SEEGER's approach are known to produce tighter bounds. As for Theorem 4.3.3, the result below bounds independently the Gibbs risk $r_{\mathcal{D}}(\rho)$ and the disagreement $d_{\mathcal{D}}(\rho)$; see the PAC-Bound 1 of GERMAIN *et al.* (2015). Note that GERMAIN *et al.* (2015) proved the bound for the binary setting but the multi-class setting is also handled. It is recalled in the following theorem.

Theorem 4.3.4 (PAC-Bayesian C-Bound of GERMAIN *et al.* (2015)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathcal{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right] \right)^2}{1 - 2 \max \left[0, \underline{\text{kl}} \left(d_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right]}. \quad (4.3)$$

$\underbrace{\hspace{15em}}_{\triangleq C_{\mathcal{S}}^{\mathcal{S}}(\rho)}$

Proof. Deferred to Appendix D.4. ■

Note that this PAC-Bayesian C-Bound is tighter than Equation (4.2) because the SEEGER's PAC-Bayesian bound is tighter than the one of MCALLESTER's one (see Section 2.3). However, one drawback of this PAC-Bayesian C-Bound is that the Gibbs risk $r_{\mathcal{D}}(\rho)$ and the disagreement $d_{\mathcal{D}}(\rho)$ are upper-bounded independently.

4.3.5 PAC-Bayesian C-Bound of Lacasse *et al.*

LACASSE *et al.* (2006) proposed to bound simultaneously the joint error $e_{\mathcal{D}}(\rho)$ and the disagreement $d_{\mathcal{D}}(\rho)$. Here, to compute the bound, we need to find the worst C-Bound value that can be obtained with a couple of joint error and disagreement denoted by (e, d) belonging to the set $\mathbb{A}_{\mathcal{S}}(\rho)$ that is defined by

$$\mathbb{A}_{\mathcal{S}}(\rho) = \left\{ (e, d) \mid \text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \parallel e, d) \leq \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m+m}}{\delta} \right], \right. \\ \left. d \leq 2\sqrt{e} - 2e, 2e + d < 1 \right\},$$

$$\text{where } \text{kl}(q_1, q_2 \parallel p_1, p_2) = q_1 \ln \frac{q_1}{p_1} + q_2 \ln \frac{q_2}{p_2} + (1 - q_1 - q_2) \ln \frac{1 - q_1 - q_2}{1 - p_1 - p_2}.$$

Based on $\mathbb{A}_{\mathcal{S}}(\rho)$, LACASSE *et al.* (2006) derive the following PAC-Bayesian C-Bound.

Theorem 4.3.5 (PAC-Bayesian C-Bound of LACASSE *et al.* (2006)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$,

for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq \sup_{(e,d) \in \mathbb{A}_{\mathbb{S}}(\rho)} \left[1 - \frac{(1 - (2e + d))^2}{1 - 2d} \right],$$

where $\mathbb{A}_{\mathbb{S}}(\rho) = \left\{ (e, d) \mid \text{kl}(e_{\mathbb{S}}(\rho), d_{\mathbb{S}}(\rho)) \mid e, d \leq \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m+m}}{\delta} \right], \right.$

$$\left. d \leq 2\sqrt{e} - 2e, 2e + d < 1 \right\}.$$

Proof. Deferred to Appendix D.5. ■

This PAC-Bayesian C-Bound can be more challenging to compute: it requires to solve a (convex) optimization problem to obtain a bound value.

4.4 Contribution: Algorithms based on the PAC-Bayesian C-Bounds

In this section, we present three self-bounding algorithms minimizing directly the PAC-Bayesian C-Bounds introduced previously.

4.4.1 Algorithm based on Equation (4.2)

Algorithm 4.1 Minimization of Equation (4.2) by Stochastic Gradient Descent

Given: learning sample \mathbb{S} , prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, the objective function $G_{\mathbb{S}}^{\text{M}}(\rho)$

Hyperparameters: number of iterations T

$\rho \leftarrow \pi$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathbb{U} \subseteq \mathbb{S}$ **do**

$\rho \leftarrow$ Update ρ with $G_{\mathbb{U}}^{\text{M}}(\rho)$ by gradient descent¹

return ρ

¹The update of ρ can be done with a vanilla gradient descent or with the update of another algorithm like Adam (KINGMA and BA, 2015) or COCOB (ORABONA and TOMMASI, 2017).

4.4. Contribution: Algorithms based on the PAC-Bayesian C-Bounds

We derive in Algorithm 4.1 a method to directly minimize the PAC-Bayesian C-Bound of Theorem 4.3.3 by stochastic gradient descent. An important aspect of the optimization is that if $r_S(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} \geq \frac{1}{2}$, the gradient of the numerator in $C_S^M(\rho)$ with respect to ρ is 0 which makes the optimization impossible. Hence, we aim to minimize the following constraint optimization problem:

$$\min_{\rho \in \mathcal{M}(\mathbb{H})} \underbrace{\left\{ 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, r_S(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} \right] \right)^2}{1 - 2 \max \left[0, d_S(\rho) - \sqrt{\frac{1}{2m} [2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} \right]} \right\}}_{\triangleq C_S^M(\rho)}$$

$$\text{s.t. } r_S(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} \leq \frac{1}{2}.$$

From this formulation, we deduce a non-constrained optimization problem:

$$\min_{\rho \in \mathcal{M}(\mathbb{H})} \left\{ C_S^M(\rho) + \text{B} \left[r_S(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} - \frac{1}{2} \right] \right\},$$

where $\text{B}(\cdot)$ is the barrier function defined as $\text{B}(a) = 0$ if $a \leq 0$ and $\text{B}(a) = +\infty$ otherwise. Due to the nature of $\text{B}(\cdot)$, this problem is not suitable for optimization: the objective function will be infinite when $a > 0$. To tackle this drawback, we replace $\text{B}(\cdot)$ by the approximation introduced by KERVADEC *et al.* (2019) called the log-barrier extension and defined as

$$\text{B}_\lambda(a) = \begin{cases} -\frac{1}{\lambda} \ln(-a), & \text{if } a \leq -\frac{1}{\lambda^2}, \\ \lambda a - \frac{1}{\lambda} \ln\left(\frac{1}{\lambda^2}\right) + \frac{1}{\lambda}, & \text{otherwise.} \end{cases}$$

The parameter $\lambda \in \mathbb{R}_*^+$ parameterized the log-barrier extension $\text{B}_\lambda(\cdot)$. The function $\text{B}_\lambda(\cdot)$ tends to $\text{B}(\cdot)$ when λ tends to $+\infty$; we plot in Figure 4.1 these two functions. Compared to the standard log-barrier², the function $\text{B}_\lambda(\cdot)$ is differentiable even when the constraint is not satisfied, *i.e.*, when $a > 0$. Thanks to $\text{B}_\lambda(\cdot)$, we can take the constraint $r_S(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} \leq \frac{1}{2}$ into account. Moreover, when the number of examples m is large, we estimate the PAC-Bayesian C-Bound $C_S^M(\rho)$ and the Gibbs risk $r_S(\rho)$ with a mini-batch $\mathcal{U} \subseteq \mathcal{S}$. Concretely, our objective function that is minimized by stochastic gradient descent with Algorithm 4.1 is the following:

$$G_{\mathcal{U}}^M(\rho) = C_{\mathcal{U}}^M(\rho) + \text{B}_\lambda \left[r_{\mathcal{U}}(\rho) + \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta}]} - \frac{1}{2} \right].$$

²The reader can refer to BOYD and VANDENBERGHE (2004) for an introduction of standard log-barrier and interior-point methods.

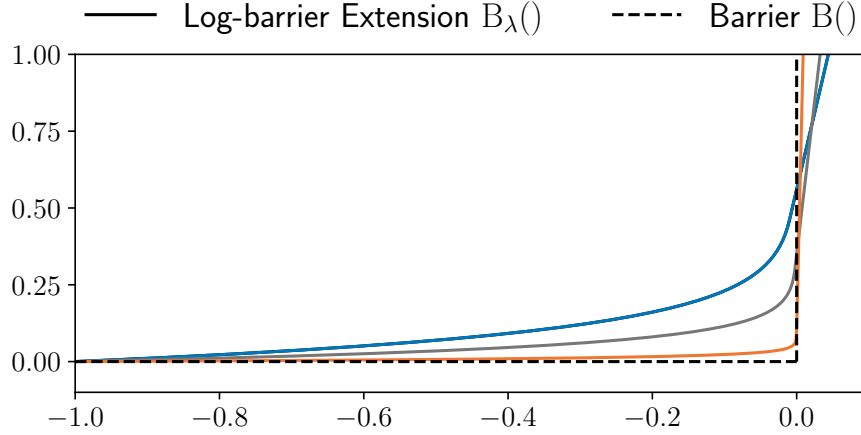


Figure 4.1. Plot of the barrier function $B()$ (with the dotted line) and the log-barrier extension $B_\lambda()$ (plain lines). We plot the function with three parameters: $\lambda \in \{10, 20, 100\}$. The higher the parameter λ , the closer the function $B_\lambda()$ to the barrier function $B()$. Thus, the blue, gray and orange curves are respectively with the parameter $\lambda = 10$, $\lambda = 20$ and $\lambda = 100$.

For a given λ , the optimizer will find a solution with a good trade-off between minimizing $C_{\mathcal{U}}^M(\rho)$ and the log-barrier extension function $B_\lambda()$. As shown in the experiments, minimizing the MCALLESTER-based bound does not lead to the tightest bound. Indeed, such a bound is looser than SEEGER-based bounds and leads to a looser PAC-Bayesian C-Bound.

4.4.2 Algorithm based on Equation (4.3)

In order to obtain better generalization guarantees, we should optimize the SEEGER-based C-bound of Theorem 4.3.4. Hence to minimize such a PAC-Bayesian C-Bound, we seek to minimize the following optimization problem:

$$\min_{\rho \in \mathcal{M}(\mathbb{H})} \underbrace{\left\{ 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right] \right)^2}{1 - 2 \max \left[0, \underline{\text{kl}} \left(d_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right]} \right\}}_{\triangleq C_{\mathcal{S}}^{\text{S}}(\rho)}$$

$$\text{s.t. } \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \leq \frac{1}{2}.$$

4.4. Contribution: Algorithms based on the PAC-Bayesian C-Bounds

For the same reasons as for deriving Algorithm 4.1, we propose to solve by stochastic gradient descent with a mini-batch $\mathcal{U} \subseteq \mathcal{S}$:

$$G_{\mathcal{U}}^{\mathcal{S}}(\rho) = C_{\mathcal{U}}^{\mathcal{S}}(\rho) + B_{\lambda} \left[\overline{\text{kl}} \left(r_{\mathcal{U}}(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) - \frac{1}{2} \right].$$

The main challenge in optimizing it is to evaluate $\overline{\text{kl}}$ or $\underline{\text{kl}}$ and to compute their derivatives. The evaluation of $\overline{\text{kl}}$ or $\underline{\text{kl}}$ is done by Algorithm 2.1 proposed by REEB *et al.* (2018). This method consists in refining iteratively an interval $[p_{\min}, p_{\max}]$ with $p \in [p_{\min}, p_{\max}]$ such that $\text{kl}(q \parallel p) = \psi$. Moreover, to compute the derivatives with respect to the posterior ρ , we use the chain rule for differentiation with a deep learning framework (such as PyTorch (PASZKE *et al.*, 2019)) and the derivatives in Equation (2.16). The global algorithm is summarized in Algorithm 4.2.

4.4.3 Algorithm based on Theorem 4.3.5

Theorem 4.3.5 jointly upper-bounds the joint error $e_{\mathcal{D}}(\rho)$ and the disagreement $d_{\mathcal{D}}(\rho)$; But as pointed out in Section 4.3.5 its optimization can be hard. To ease its manipulation, we derive below a C-Bound resulting of a reformulation of the constraints involved in the set $\mathbb{A}_{\mathcal{S}}(\rho)$.

Algorithm 4.2 Minimization of Equation (4.3) by Stochastic Gradient Descent

Given: learning sample \mathcal{S} , prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, the objective function $G_{\mathcal{S}}^{\mathcal{S}}(\rho)$
Hyperparameters: number of iterations T
 $\rho \leftarrow \pi$
for $t \leftarrow 1$ to T **do**
 for all mini-batches $\mathcal{U} \subseteq \mathcal{S}$ **do**
 Compute $G_{\mathcal{U}}^{\mathcal{S}}(\rho)$ using Algorithm 2.1
 $\rho \leftarrow$ Update ρ with $G_{\mathcal{U}}^{\mathcal{S}}(\rho)$ by gradient descent
return ρ

Theorem 4.4.1 (Reformulation of LACASSE *et al.*'s PAC-Bayesian C-Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random

choice of $S \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$\begin{aligned} R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq \sup_{(e,d) \in \mathbb{A}'_{\mathcal{S}}(\rho)} \left[1 - \frac{(1 - (2e + d))^2}{1 - 2d} \right], \quad (4.4) \\ \mathbb{A}'_{\mathcal{S}}(\rho) &= \left\{ (e, d) \mid \text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \parallel e, d) \leq \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m} + m}{\delta} \right], \right. \\ &\quad \left. d \leq 2\sqrt{\min\left(e, \frac{1}{4}\right)} - 2e, \quad d < \frac{1}{2} \right\}. \end{aligned}$$

Proof. Deferred to Appendix D.6. ■

Theorem 4.4.1 suggests then the following constrained optimization problem:

$$\min_{\rho \in \mathbb{M}(\mathbb{H})} \left\{ \sup_{(e,d) \in [0, \frac{1}{2}]^2} \left(1 - \frac{[1 - (2e + d)]^2}{1 - 2d} \right) \text{ s.t. } (e, d) \in \mathbb{A}'_{\mathcal{S}}(\rho) \right\} \text{ s.t. } 2e_{\mathcal{S}}(\rho) + d_{\mathcal{S}}(\rho) \leq 1,$$

Actually, we can rewrite this constrained optimization problem into an unconstrained one using the barrier function. We obtain

$$\begin{aligned} \min_{\rho \in \mathbb{M}(\mathbb{H})} \left\{ \max_{(e,d) \in [0, \frac{1}{2}]^2} \left(C^{\text{L}}(e, d) - \text{B} \left[d - 2\sqrt{\min\left(e, \frac{1}{4}\right)} - 2e \right] - \text{B} \left[d - \frac{1}{2} \right] \right. \right. \\ \left. \left. - \text{B} \left[\text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \parallel e, d) - \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m} + m}{\delta} \right] \right] \right) \right. \\ \left. + \text{B} \left[2e_{\mathcal{S}}(\rho) + d_{\mathcal{S}}(\rho) - 1 \right] \right\}, \quad (4.5) \end{aligned}$$

where $C^{\text{L}}(e, d) = 1 - \frac{(1 - (2e + d))^2}{1 - 2d}$ if $d < \frac{1}{2}$, and $C^{\text{L}}(e, d) = 1$ otherwise. However, this problem cannot be optimized directly by stochastic gradient descent. In this case, we have a min-max optimization problem, *i.e.*, for each descent step we need to find the couple (e, d) that maximizes the $C^{\text{L}}(e, d)$ given the three constraints that define $\mathbb{A}'_{\mathcal{S}}(\rho)$ before updating the posterior distribution ρ .

First, to derive our optimization procedure, we focus on the inner maximization problem when $e_{\mathcal{S}}(\rho)$ and $d_{\mathcal{S}}(\rho)$ are fixed in order to find the optimal (e, d) . However, the function $C^{\text{L}}(e, d)$ we aim at maximizing is not concave for all $(e, d) \in \mathbb{R}^2$, implying

4.4. Contribution: Algorithms based on the PAC-Bayesian C-Bounds

Algorithm 4.3 Minimization of Equation (4.4) by Stochastic Gradient Descent

Given: learning sample \mathcal{S} , prior $\pi \in \mathbb{M}^*(\mathbb{H})$, the objective function $G_{\mathcal{S}}^{e^*, d^*}(\rho)$

Hyperparameters: number of iterations T

$\rho \leftarrow \pi$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathcal{U} \subseteq \mathcal{S}$ **do**

$(e^*, d^*) \leftarrow \text{MAXIMIZE-}e-d(e_{\mathcal{U}}(\rho), d_{\mathcal{U}}(\rho))$

$\rho \leftarrow \text{Update } \rho \text{ with } G_{\mathcal{U}}^{e^*, d^*}(\rho) \text{ by gradient descent}$

return ρ

Given: learning sample \mathcal{S} , joint error $e_{\mathcal{S}}(\rho)$, disagreement $d_{\mathcal{S}}(\rho)$

Hyperparameters: tolerance ϵ

function $\text{MAXIMIZE-}e-d(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho))$

$\alpha_{\min} = 0$ and $\alpha_{\max} = 1$

while $\alpha_{\max} - \alpha_{\min} > \epsilon$ **do**

$\alpha = \frac{1}{2}(\alpha_{\min} + \alpha_{\max})$

$(e, d) \leftarrow \text{Solve Equation (4.6)}$

if $C^{\text{L}}(e, d) \geq 1 - \alpha$ **then** $\alpha_{\max} \leftarrow \alpha$ **else** $\alpha_{\min} \leftarrow \alpha$

return (e, d)

that the implementation of its maximization can be hard³. Fortunately, $C^{\text{L}}(e, d)$ is quasi-concave (GERMAIN *et al.*, 2015) for $(e, d) \in [0, 1] \times [0, \frac{1}{2}]$. Then by definition of quasi-concavity, we have:

$$\begin{aligned} & \forall \alpha \in [0, 1], \left\{ (e, d) \left| 1 - \frac{[1 - (2e + d)]^2}{1 - 2d} \geq 1 - \alpha \right. \right\} \\ \iff & \forall \alpha \in [0, 1], \left\{ (e, d) \left| \alpha(1 - 2d) - [1 - (2e + d)]^2 \geq 0 \right. \right\}. \end{aligned}$$

Hence, for any fixed $\alpha \in [0, 1]$ we can look for (e, d) that maximizes $C^{\text{L}}(e, d)$ and respects the constraints involved in $\mathbb{A}'_{\mathcal{S}}(\rho)$. This is equivalent to solving the following

³For example, when using CVXPY (DIAMOND and BOYD, 2016), that uses Disciplined Convex Programming (DCP (GRANT *et al.*, 2006)), the maximization of a non-concave function is not possible.

problem for a given $\alpha \in [0, 1]$:

$$\begin{aligned} \max_{(e,d) \in [0, \frac{1}{2}]^2} \quad & \alpha(1-2d) - \left[1-(2e+d)\right]^2 \\ \text{s.t.} \quad & d \leq 2\sqrt{\min\left(e, \frac{1}{4}\right)} - 2e \\ \text{and} \quad & \text{kl}\left(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \parallel e, d\right) \leq \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}+m}{\delta}\right]. \end{aligned} \quad (4.6)$$

In fact, we aim at finding $\alpha \in [0, 1]$ such that the maximization of Equation (4.6) leads to $1-\alpha$ equals to the largest value of $C^L(e, d)$ under the constraints. To do so, we make use of the ‘‘Bisection method for quasi-convex optimization’’ (BOYD and VANDENBERGHE, 2004) that is summarized in MAXIMIZE- e - d in Algorithm 4.3. We denote by (e^*, d^*) the solution of Equation (4.6). Note that, in practice, the joint error and the disagreement is approximated through the mini-batch $\mathcal{U} \subseteq \mathcal{S}$. It remains then to solve the outer minimization problem that becomes:

$$\min_{\rho \in \mathcal{M}(\mathbb{H})} \left\{ \begin{aligned} & \text{B} [2e_{\mathcal{S}}(\rho) + d_{\mathcal{S}}(\rho) - 1] \\ & - \text{B} \left[\text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \parallel e^*, d^*) - \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}+m}{\delta} \right] \right] \end{aligned} \right\}.$$

To obtain a objective function that is suitable for stochastic gradient descent, we bring two modifications to the outer minimization problem: (i) we replace $\text{B}(\cdot)$ by the log-barrier extension $\text{B}_\lambda(\cdot)$ and (ii) we approximate the disagreement and the joint error with a mini-batch $\mathcal{U} \subseteq \mathcal{S}$. Hence, we obtain the following objective function:

$$\begin{aligned} G_{\mathcal{U}}^{e^*, d^*}(\rho) = & \text{B}_\lambda [2e_{\mathcal{U}}(\rho) + d_{\mathcal{U}}(\rho) - 1] \\ & - \text{B}_\lambda \left[\text{kl}(e_{\mathcal{U}}(\rho), d_{\mathcal{U}}(\rho) \parallel e^*, d^*) - \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}+m}{\delta} \right] \right]. \end{aligned}$$

The global method is summarized in Algorithm 4.3. As a side note, we mention that the classic Danskin Theorem (DANSKIN, 1966) used in min-max optimization theory is not applicable in our case since our objective function is not differentiable for all $(e, d) \in [0, \frac{1}{2}]^2$. We discuss this point in Appendix D.8.

4.5 Experiments

4.5.1 Setting

Our experiments have a two-fold objective: (i) assessing the guarantees given by the associated PAC-Bayesian bounds, and (ii) comparing the performance of the different

4.5. Experiments

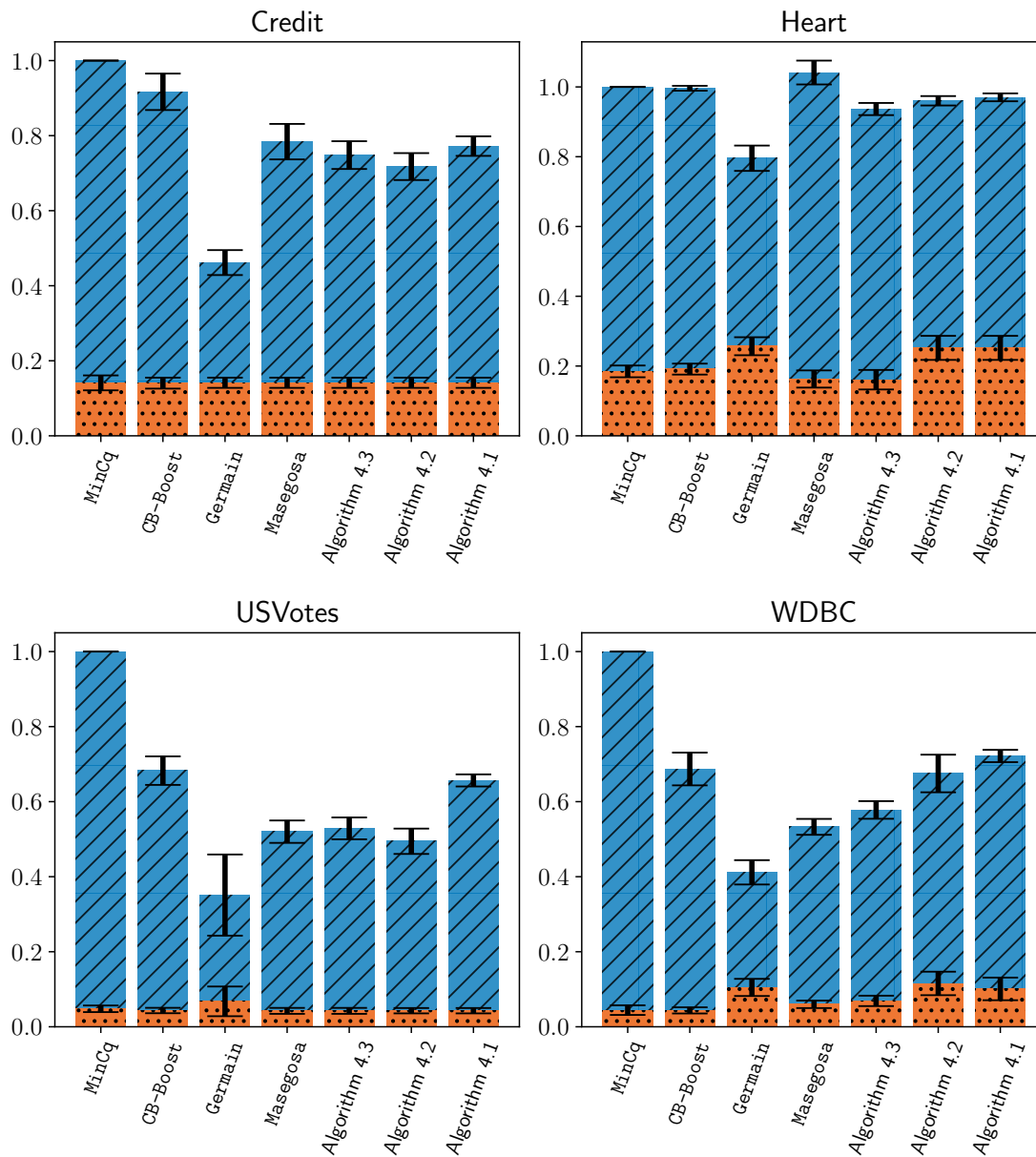


Figure 4.2. Plot of a comparison between the test risks $R_T(MV_\rho)$ and the generalization bounds in the binary setting when the voters are decision stumps. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

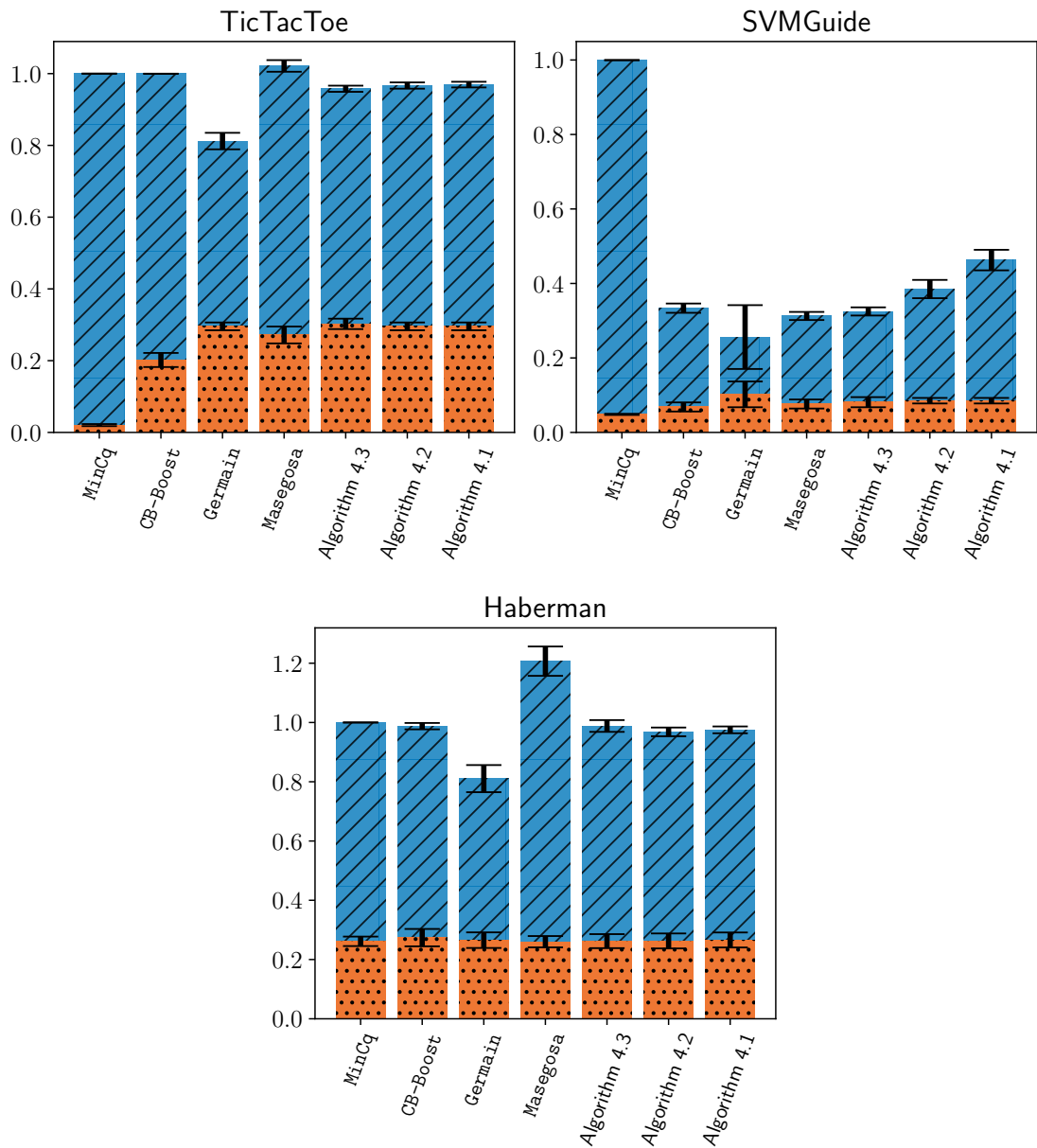


Figure 4.3. Plot of a comparison between the test risks $R_T(MV_\rho)$ and the generalization bounds in the binary setting when the voters are decision stumps. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

4.5. Experiments

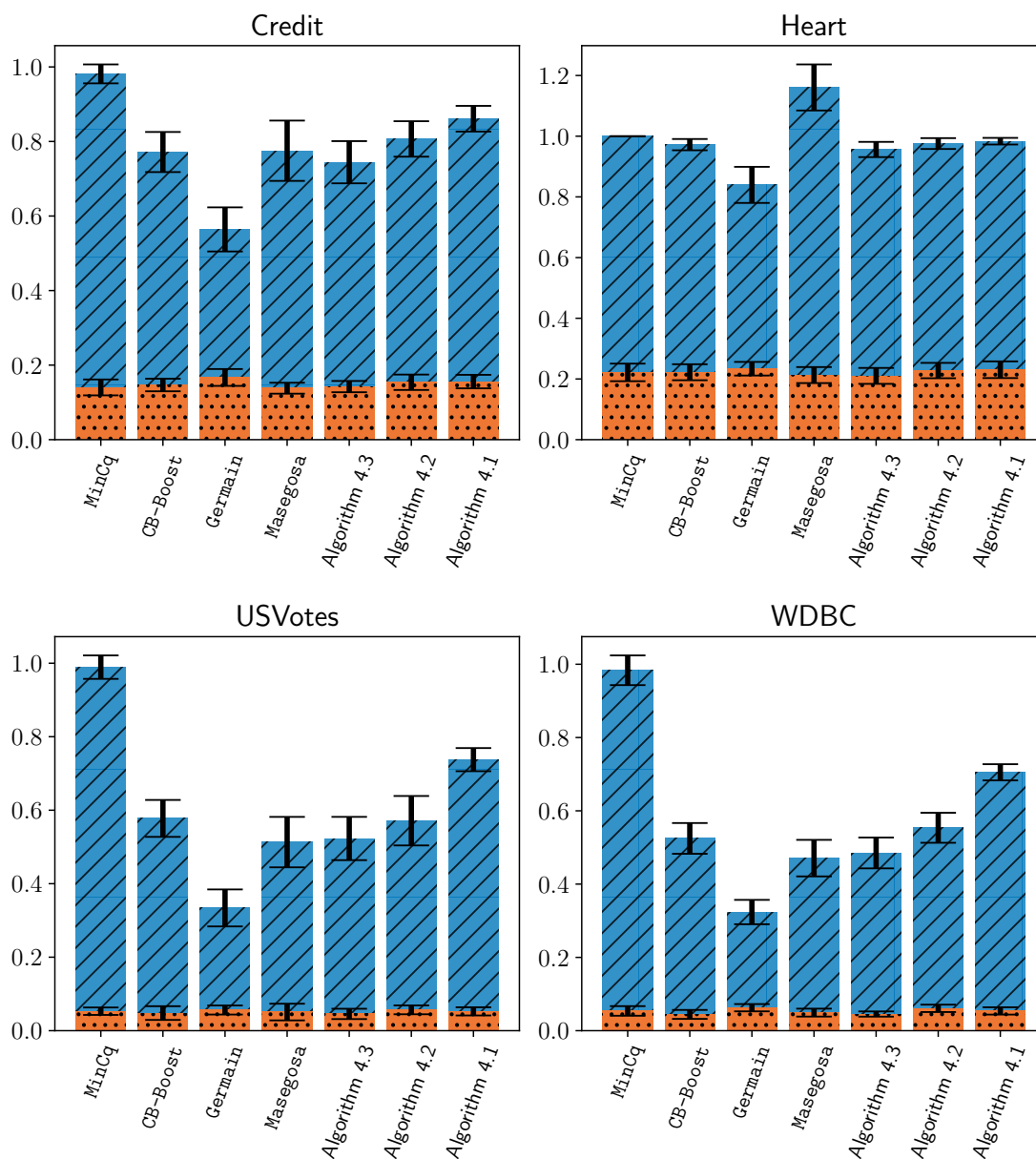


Figure 4.4. Plot of a comparison between the test risks $R_{\mathcal{T}}(MV_{\rho})$ and the generalization bounds in the binary setting when the voters are decision trees. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

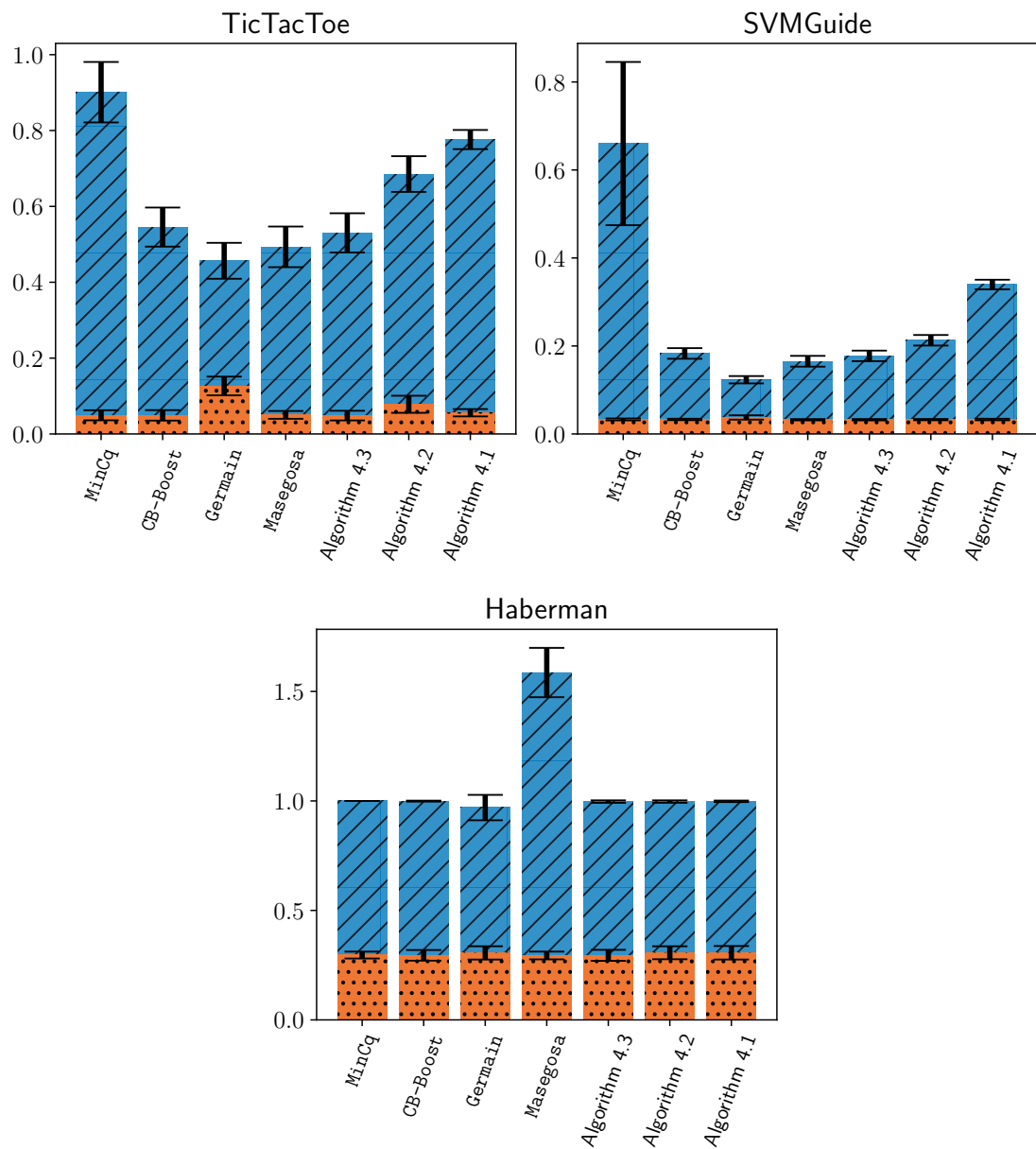


Figure 4.5. Plot of a comparison between the test risks $R_{\mathcal{T}}(MV_{\rho})$ and the generalization bounds in the binary setting when the voters are decision trees. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

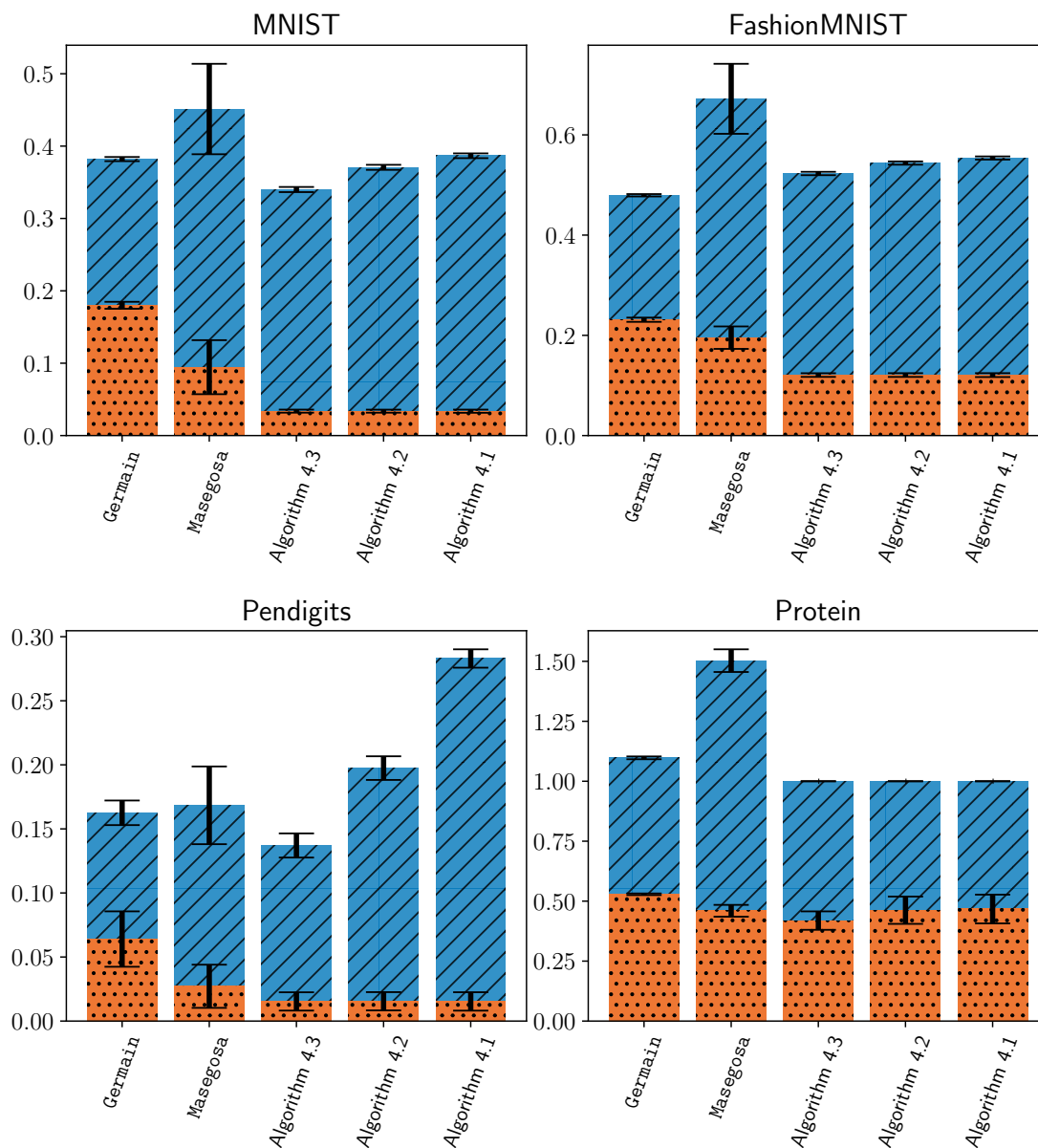


Figure 4.6. Plot of a comparison between the test risks $R_T(MV_\rho)$ and the generalization bounds in the multi-class setting when the voters are decision trees. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

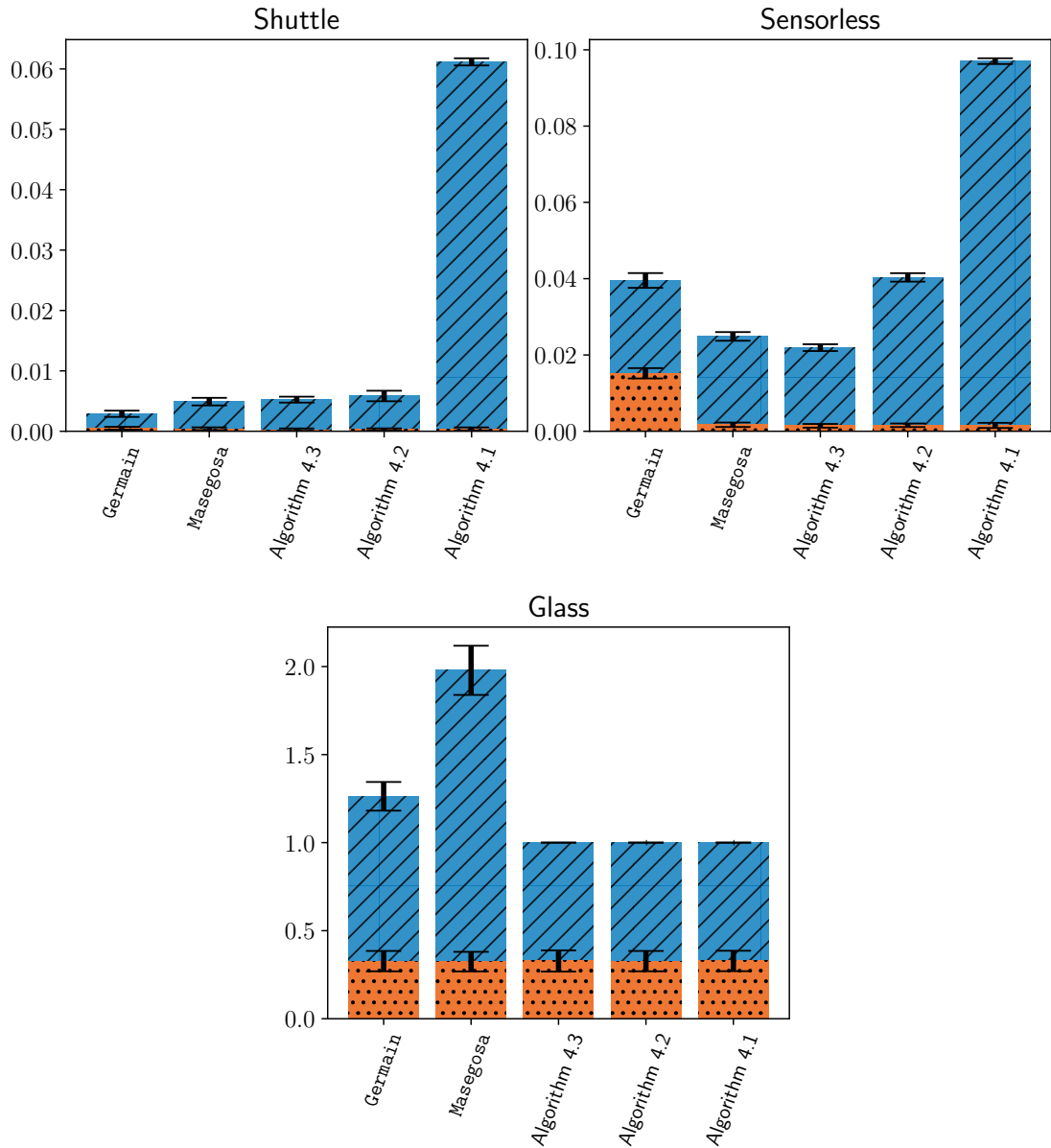


Figure 4.7. Plot of a comparison between the test risks $R_T(MV_\rho)$ and the generalization bounds in the multi-class setting when the voters are decision trees. For each algorithm, we represent the mean of the test risks in the orange bars and the bounds' mean in the blue bars. Additionally, the black lines are the standard deviations.

4.5. Experiments

C-bound-based algorithms in terms of risk optimization. We introduce the setting in the following and we report in Figures 4.2 to 4.7 the mean/standard deviation of the risks on the test set \mathbb{T} and the bound values (with $\delta = 0.05$) for 10 runs; see also Appendix D.10.1 for more details. The setting of the experiments is as follows.

Dataset. We consider several binary and multi-class datasets like FashionMNIST (XIAO *et al.*, 2017), MNIST (LECUN *et al.*, 1998) and some coming from the UCI repository (DUA and GRAFF, 2017). For each different run, we keep the same number of examples in the test or the train set as in the original split. When there is no original split, we use 50% of data in the training set \mathbb{S} and 50% in the test set (except for Sensorless where we have 15% in the test set because the original set is large).

Voters. In the binary setting, we consider either a set \mathbb{H} of decision trees or decision stumps that is complemented: if $h \in \mathbb{H}$ then there is $h' \in \mathbb{H}$ s.t. $h'(\mathbf{x}) = -h(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{X}$. Concerning the multi-class setting, we only consider decision trees. Indeed, having decision stumps would have resulted in too many voters. Following MASEGOSA *et al.* (2020), the prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ is set as the uniform distribution. For the decision stumps, following ROY *et al.* (2011) and BAUVIN *et al.* (2020), we use 10 decision stumps per feature. For the decision trees, we follow a general setting similar to the one of (MASEGOSA *et al.*, 2020). Moreover, 100 trees are learned with 50% of the training data (the remaining part serves to learn the posterior ρ). More precisely, for each tree \sqrt{d} features of the d -dimensional input space are selected, and the trees are learned by using the Gini criterion until the leaves are pure.

Algorithms' parameters. To update of the posterior ρ in Algorithms 4.1 to 4.3 is done through the COCOB-Backprop optimizer (ORABONA and TOMMASI, 2017) (its parameter remains the default one). In the binary setting, we optimize for $T = 2,000$ iterations (by batch gradient descent), and in the multi-class setting, we set 20 epochs with a batch size of 64. Lastly, we consider the parameter $\lambda=100$ for log-barrier extension $B_\lambda()$.

Comparisons. We compare the three algorithms proposed in this chapter to the following state-of-the-art PAC-Bayesian methods for majority vote learning.

- We compare with the algorithm proposed by MASEGOSA *et al.* (2020) that optimizes a PAC-Bayesian bound on $R_{\mathcal{D}}(MV_\rho) \leq 4e_{\mathcal{D}}(\rho)$ (Theorem 2.2.2); see Theorem 9 and Appendix G of (MASEGOSA *et al.*, 2020) for a description of this algorithm that we denote by Masegosa. For MASEGOSA *et al.*'s algorithm, we kept the original parameters.

- Our algorithm to optimize the PAC-Bayesian bound on $R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 2r_{\mathcal{D}}(\rho)$ (Theorem 2.2.1) recalled in Theorem 4.3.1 and derived by (GERMAIN *et al.*, 2015, PAC-Bound 0). Even though (GERMAIN *et al.*, 2015) does not optimize the bound, we denote this algorithm by `Germain`. The algorithm is similar to Algorithm 4.2, but without the numerator of the C-Bound (*i.e.*, the disagreement); more details are given in Appendix D.9.
- In the binary setting only, we compare with `MinCq` (ROY *et al.*, 2011) and `CB-Boost` (BAUVIN *et al.*, 2020) that are based on the minimization of the empirical C-Bound $C_S(\rho)$. Indeed, these algorithms are developed for this setting only. For comparison purposes and since `MinCq` and `CB-Boost` do not explicitly minimize a PAC-Bayesian bound, we report the bound values of Theorem 4.4.1 instantiated with the models learned; Moreover, for `MinCq`, we select the margin parameter among 20 values uniformly spaced between $[0, \frac{1}{2}]$ by 3-fold cross validation. For `CB-Boost`, which is based on a Boosting approach, we fix the maximal number of boosting iterations to 200.

4.5.2 Analysis of the Results

When comparing only the PAC-Bayesian C-Bounds, we observe in Figures 4.2 to 4.7 that Algorithm 4.1 provides the worst bound. Algorithm 4.3 provides usually tighter bounds than Algorithms 4.1 and 4.2 except for Harberman and USVotes. We believe that Algorithm 4.3 provides lower bounds than Algorithm 4.2 because the LACASSE *et al.*'s approach bounds simultaneously the joint error and the disagreement. Algorithm 4.3 appears then to be the best algorithm among our three self-bounding algorithms that minimize a PAC-Bayesian C-Bound. Moreover, Algorithm 4.3 gives usually the lowest true risks or it is comparable to the two other algorithms.

Compared to the baselines, `Germain` gives usually the lowest bounds among all the algorithms, but at the price of a large test risk. This clearly illustrates the limitation of considering *only* the Gibbs risk as an estimator of the majority vote risk: as discussed in Section 2.2.2, the Gibbs risk is an unfair estimator since an increase in the diversity between the voters can have a negative impact on the Gibbs risk.

Second, compared to MASEGOSA *et al.*'s approach, the results are comparable. This behavior was expected since minimizing the bound of MASEGOSA *et al.* (2020) or the PAC-Bayesian C-Bound boils down to minimize a trade-off between the risk and the disagreement. Third, in the binary setting, compared to empirical C-bound minimization algorithms, we see that Algorithm 4.3 outputs better results than `CB-Boost` and `MinCq` for which the difference is significative, and the bounds are close to 1 (*i.e.*, non-informative). Optimizing the risk bounds tends to provide better guarantees that justify that optimizing the empirical C-bound is often too optimistic (as done in

4.6. Conclusion and Summary

CB-Boost or MinCq); we provide in Appendix D.10.3 an illustration of the different solutions obtained from the algorithms.

Overall, from these experiments, our Algorithm 4.3 is the one that provides the best trade-off between having good performances in terms of risk optimization and ensuring good theoretical guarantees with informative bounds. Moreover, in Appendix D.10.2 we show that Algorithm 4.3 has a higher computation time than the others algorithms. This makes Algorithm 4.2 a good trade-off between ensuring good theoretical guarantees and having a low computation time.

4.6 Conclusion and Summary

This chapter presents learning algorithms that minimize the majority vote's risk with PAC-Bayesian generalization bounds based on the C-Bound. More precisely, we propose solving three optimization problems, each derived from an existing PAC-Bayesian bound. Our methods belong to the class of *self-bounding* learning algorithms (FREUND, 1998): the learned predictor comes with a tight and statistically valid risk upper bound. Our experimental evaluation has confirmed the quality of the learned predictor and the tightness of the bounds with respect to state-of-the-art methods minimizing the C-Bound.

As we said before, no algorithm minimizes the *empirical* C-Bound in the multi-class setting. One of the reasons is that it is not easy to find a convex program like for MinCq or P-MinCq algorithm. Hopefully, thanks to the deep learning framework, minimizing the C-Bound is possible even without convexity through a stochastic gradient descent algorithm (as we show in Chapter 5 in the binary setting). Hence, in the future, we plan to explore more the minimization of the C-Bound.

However, one drawback of the PAC-Bayesian C-Bounds and the other ones of the literature is that the majority vote true risk is not directly minimized: we need surrogates such as the Gibbs risk, the disagreement, or the joint error. Indeed, the C-Bound is already an upper bound on the true risk, which makes the PAC-Bayesian C-Bound even looser, thus, these generalization bounds cannot be tight. To avoid this issue, we bound the *expected* true risk of the majority vote. This is done in the next chapter: we (i) introduce the *stochastic* majority vote that samples a distribution ρ for each prediction, and we (ii) provide guarantee on the true risk for this classifier. As we will see, our obtained (PAC-Bayesian) guarantee is differentiable, and we can derive self-bounding learning to learn such a classifier.

TOWARD A STOCHASTIC MAJORITY VOTE

5

This chapter is based on the following paper

VALENTINA ZANTEDESCHI, PAUL VIALLARD, EMILIE MORVANT, RÉMI EMONET, AMAURY HABRARD, PASCAL GERMAIN, and BENJAMIN GUEDJ. Learning Stochastic Majority Votes by Minimizing a PAC-Bayes Generalization Bound. *Advances in Neural Information Processing Systems (NeurIPS)*. (2021)

Contents

5.1	Introduction	126
5.2	The Stochastic Majority Vote	129
5.2.1	Definitions	129
5.2.2	Approximation of the Stochastic Risk	131
5.2.3	Computing Exactly the Stochastic Risk	131
5.3	From a PAC-Bayesian Bound to an Algorithm	133
5.3.1	A PAC-Bayesian Bound for Stochastic Majority Votes	133
5.3.2	A PAC-Bayesian Bound for Data-dependent Voters	135
5.3.3	Learning Algorithms for the Stochastic Majority Vote	136
5.4	Experiments	138
5.4.1	Comparison Between the Computations of the Risk	139
5.4.2	Performance of Algorithms 5.3 and 5.4	139
5.5	Conclusion and Summary	141

Abstract

We study a stochastic counterpart of the majority vote classifier called the *stochastic majority vote*, and study its generalization properties. Unlike Chapter 4, the posterior distribution associated with the majority vote is sampled from another probability distribution. While the stochastic majority vote holds for arbitrary distributions, we instantiate it with Dirichlet distributions: this allows to derive a closed-form and differentiable expression for the expected risk. Then, we derive self-bounding algorithms for stochastic majority vote, that benefit from tight generalization bounds when compared to self-bounding algorithms studied in Chapter 4.

5.1 Introduction

In Chapter 4, we considered some self-bounding algorithms (FREUND, 1998) that minimize a PAC-Bayesian bound on the majority vote true risk. Each PAC-Bayesian bound depends on a surrogate on the majority vote risk (see Section 2.2.2). Given any distribution ρ on an hypothesis set \mathbb{H} , we have seen three surrogates of the majority vote's risk in Chapter 2:

- (i) Twice the Gibbs Risk $r_S(\rho) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \mathbb{E}_{h \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y]$ (Theorem 2.2.1),
- (ii) 4 times the joint error $e_S(\rho) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) \neq y]$ (Theorem 2.2.2),
- (iii) The C-Bound $C_S(\rho) = 1 - \frac{(1-2r_S(\rho))^2}{1-2d_S(\rho)}$, where $d_S(\rho)$ is the disagreement defined as $d_S(\rho) = 2 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) = y]$ (Theorem 2.2.3).

Figure 5.1 illustrates the models obtained by the minimization of the empirical risk with respect to the three surrogates recalled above on a simple dataset (here, moons): only the empirical C-Bound fully leverage the diversity of the voters to obtain $R_S(\text{MV}_\rho) = 0$.

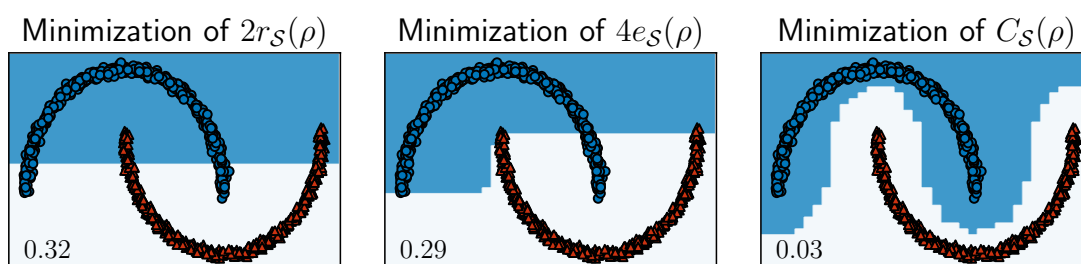


Figure 5.1. Plot of the majority vote's decision boundary obtained by minimizing the upper bounds on the majority vote's empirical risk with the surrogates recalled in Theorems 2.2.1 to 2.2.3 for moons dataset with a learning sample size $m = 1,000$. For each plot, we print the surrogate value in the bottom left corner.

As highlighted in Figure 5.1, the first two surrogates (Theorems 2.2.1 and 2.2.2) are not efficient to minimize the empirical majority vote's risk. In contrast, since the empirical C-Bound $C_S(\rho)$ (Theorem 2.2.3) is a tighter upper bound on the majority vote's empirical risk (Theorem 2.2.4), the minimization of the empirical risk through the C-Bound is more accurate. Another precise surrogate on the majority vote's empirical risk is based on the *randomized majority vote*. This model, introduced by LACASSE *et al.* (2010), consists in defining a majority vote

$$\text{MV}_\sigma(\mathbf{x}) \triangleq \operatorname{argmax}_{y' \in \mathcal{Y}} \mathbb{E}_{h \sim \sigma} \mathbb{I}[h(\mathbf{x}) = y'],$$

5.1. Introduction

where σ is constructed as follows: N voters $\mathbb{H}' = \{h_1, \dots, h_N\}$ are sampled from ρ and a uniform posterior distribution σ is defined such that $\sigma(h) = \frac{1}{N}$ for all $h \in \mathbb{H}'$. Following LACASSE *et al.* (2010), the randomized majority vote's empirical risk is defined as

$$\begin{aligned} & \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{S}, \text{MV}_\sigma \sim \rho^N} [\text{MV}_\sigma(\mathbf{x}) \neq y] \\ & \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \left[\sum_{j=\lceil \frac{N}{2} \rceil}^N \binom{N}{j} \left[\frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^j \left[1 - \frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^{(N-j)} \right] \\ & \triangleq b_{\mathcal{S}}^N(\rho), \end{aligned}$$

where $\widehat{m}_\rho(\mathbf{x}, y)$ is the $\frac{1}{2}$ -margin defined in Definition 2.2.4. Given an example $(\mathbf{x}, y) \sim \mathcal{S}$, the sum corresponds to the probability that at least $\frac{N}{2}$ voters make an error over N voters sampled from ρ . In fact, it is the complementary cumulative distribution function of the binomial distribution with parameter $\frac{1}{2}(1 - \widehat{m}_\rho(\mathbf{x}, y))$ and $\lceil \frac{N}{2} \rceil$ trials. Hence, $b_{\mathcal{S}}^N(\rho)$ is the expected complementary cumulative distribution function on the learning sample \mathcal{S} . Moreover, the *randomized* majority vote can be linked to the classical majority vote MV_ρ : the term $b_{\mathcal{S}}^N(\rho)$ is another surrogate on the majority vote empirical risks (LACASSE *et al.*, 2010). Indeed, we have

$$R_{\mathcal{S}}(\text{MV}_\rho) \leq 2b_{\mathcal{S}}^N(\rho). \quad (5.1)$$

Note that we provide a proof of this bound in Appendix E.1. Hopefully, the higher N , the better $b_{\mathcal{S}}^N(\rho)$ approximates the majority vote's empirical risk. As for the other surrogates (*i.e.*, the Gibbs risk, the joint error, and the C-Bound), when one wants to upper-bound the majority vote's true risk, PAC-Bayesian bounds are used. In the rest of the chapter, we denote by Lacasse our algorithm that minimizes a PAC-Bayesian bound depending on $b_{\mathcal{S}}^N(\rho)$ (see Appendix E.1). However, when the true risk is minimized through self-bounding algorithms, it gives even worse results on the moons dataset. This is illustrated on Figure 5.2 that plots the different self-bounding procedures.

As we can remark, except for the one labelled Lacasse and our new Algorithm 5.3, the minimization of the PAC-Bayesian bounds does not fully leverage the voters' correlations to obtain a *(i)* tight generalization bound and *(ii)* a small empirical risk. This is mainly due to the fact that they do not minimize directly the majority vote's risk (but surrogates instead). In contrast, as illustrated by LACASSE (2010)'s result, the randomized majority vote offers a way to obtain a tight generalization guarantee with a small empirical risk $R_{\mathcal{S}}(\text{MV}_\rho)$. However, the considered majority vote in LACASSE *et al.* (2010)'s approach is not the original majority vote MV_ρ , but a rather restricted form.

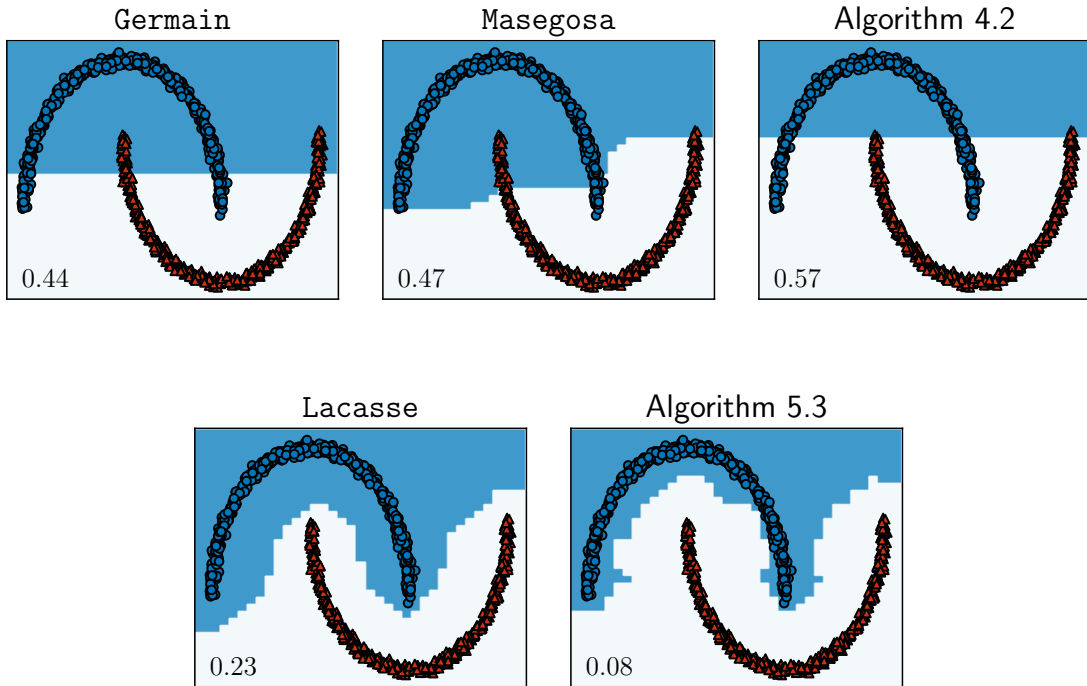


Figure 5.2. Plot of the majority vote’s decision boundary obtained by executing the self-bounding algorithms for moons dataset with $m = 1,000$ (and $N = 100$ for Lacasse). We represent in the bottom left of the plots, the value of each bound minimized with Algorithm 4.2, Algorithm 5.3 (and Algorithm 5.2), Germain (Algorithm D.1), Masegosa (MASEGOSA et al., 2020, Appendix G), and Lacasse (Algorithm E.1).

Hence, in this chapter, we introduce the *stochastic* majority vote that overcomes the drawbacks of the literature’s methods: we provide tight generalization bounds on the classical majority vote MV_ρ . The *stochastic* majority vote is defined as follows: for each input \mathbf{x} , a majority vote MV_ρ is obtained by sampling the weights from another probability distribution.

This new majority vote is presented in more details in Section 5.2; its definition is given in Section 5.2.1. Moreover, we show in Section 5.2.2 how to approximate the risk and in Section 5.2.3 the exact computation. Along with the risk computation, we derive in Sections 5.3.1 and 5.3.2 two PAC-Bayesian bounds essential to derive self-bounding algorithms in Section 5.3.3. Section 5.4 provides a study of these algorithms. The proofs are deferred in Appendix E.

5.2 The Stochastic Majority Vote

5.2.1 Definitions

For the *stochastic* majority vote, we consider that its weights $\rho \in \mathbb{M}(\mathbb{H})$ are sampled from a distribution \mathbb{P} called hyper-posterior¹; we say that \mathbb{P} is a hyper-posterior on \mathbb{H} . By considering hyper-posteriors, the majority votes become *stochastic*, *i.e.*, for each input $\mathbf{x} \in \mathbb{X}$ the weights ρ are sampled from the hyper-posterior \mathbb{P} to obtain the prediction $\text{MV}_\rho(\mathbf{x})$. The main advantage of considering a stochastic majority vote is that it allows to derive and to optimize PAC-Bayesian generalization bounds directly. The true risk and empirical risk of the proposed stochastic weighted majority vote take into account the risks of MV_ρ where ρ is sampled from the hyper-posterior \mathbb{P} ; the risk is defined in the following definition.

Definition 5.2.1 (Risks of the stochastic majority vote). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any voters' set \mathbb{H} , for any hyper-posterior distribution \mathbb{P} on \mathbb{H} , the stochastic risks are defined as

$$\begin{aligned} \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{D}}(\text{MV}_\rho) &= \mathbb{E}_{\rho \sim \mathbb{P}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[m_\rho(\mathbf{x}, y) \leq 0], \\ \text{and } \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{S}}(\text{MV}_\rho) &= \mathbb{E}_{\rho \sim \mathbb{P}} \frac{1}{m} \sum_{i=1}^m \mathbb{I}[m_\rho(\mathbf{x}_i, y_i) \leq 0]. \end{aligned}$$

We use the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017) to upper-bound the risk of the stochastic majority vote. Indeed, we have

$$\mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{D}}(\text{MV}_\rho) \leq \mathbb{E}_{\rho \sim \mathbb{P}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0] = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{\mathbb{P}}(\mathbf{x}, y), \quad (5.2)$$

$$\text{and } \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{S}}(\text{MV}_\rho) \leq \mathbb{E}_{\rho \sim \mathbb{P}} \frac{1}{m} \sum_{i=1}^m \mathbb{I}[\widehat{m}_\rho(\mathbf{x}_i, y_i) \leq 0] = \frac{1}{m} \sum_{i=1}^m s_{\mathbb{P}}(\mathbf{x}_i, y_i). \quad (5.3)$$

where $s_{\mathbb{P}}(\mathbf{x}, y) \triangleq \mathbb{E}_{\rho \sim \mathbb{P}} \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0]$ is the *stochastic* risk. Note first that the inequality is attained in the binary setting (see Section 2.2.2). Additionally, the advantage of $s_{\mathbb{P}}(\mathbf{x}, y)$ is that we are able to derive a closed form solution in Section 5.2.3. Actually, we introduce two ways to compute this risk: we can either (i) approximate it (e.g., through Monte Carlo methods) or (ii) compute its closed form. In both cases, assumptions have to be made on the distribution \mathbb{P} . When the distribution ρ is discrete, it lies in the $(\text{card}(\mathbb{H})-1)$ dimensional probability simplex: $\Delta^{(\text{card}(\mathbb{H})-1)}$. Hence,

¹The hyper-posteriors have been first used in the PAC-Bayesian theory by PENTINA and LAMPERT (2014) in the context of Life-long learning problem in order to be able to consider different majority votes adapted to the different specific tasks seen during training.

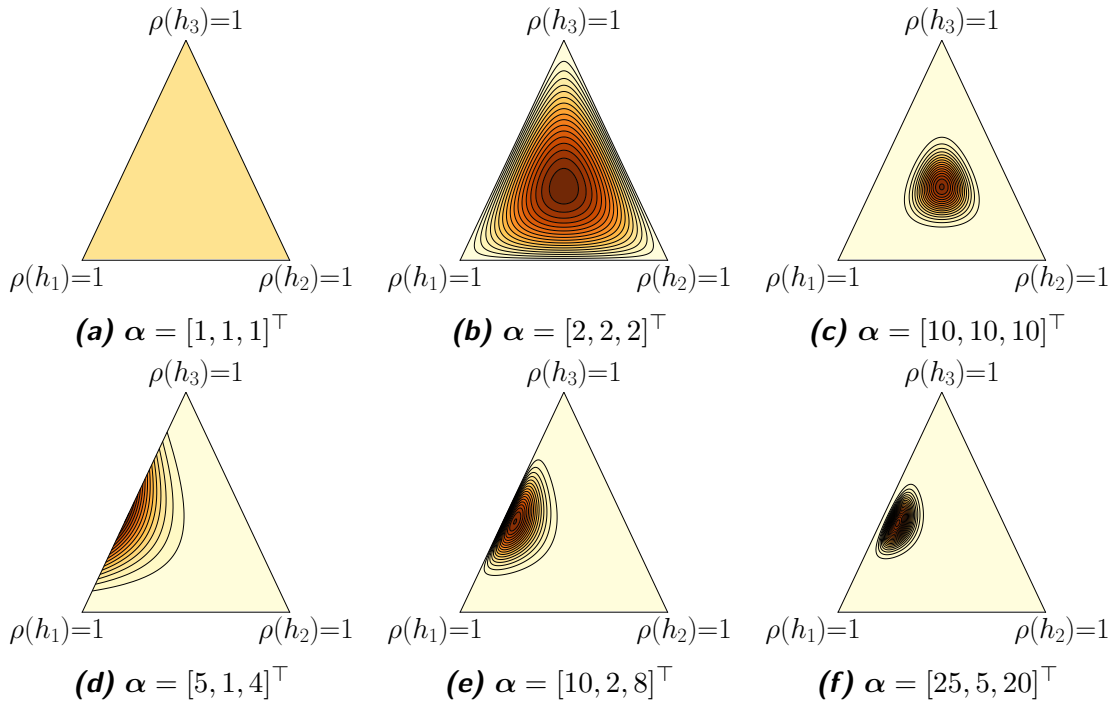


Figure 5.3. The figure shows the probability density function for different value of parameters α . More precisely, the “triangle” correspond to the 2-dimensional probability simplex Δ^2 where its extremities correspond to the extreme probability distributions. Hence, a point in the simplex is a linear combination of these extreme probability distributions.

a natural choice for the hyper-posterior is the Dirichlet distribution; its probability density function is defined as follows.

Definition 5.2.2 (Dirichlet Distribution). Let $n = \text{card}(\mathbb{H})$ be the cardinality of a finite hypothesis set \mathbb{H} . Given the concentration parameters $\alpha \in (\mathbb{R}_*^+)^n$, the Dirichlet Distribution $\text{Dir}(\alpha)$ is defined as

$$\rho \sim P \iff (\rho(h_1), \dots, \rho(h_n)) \sim \text{Dir}(\alpha),$$

$$\text{where } P(\rho) \triangleq \frac{1}{Z(\alpha)} \prod_{j=1}^n [\rho(h_j)]^{\alpha_j-1} \propto \prod_{j=1}^n [\rho(h_j)]^{\alpha_j-1}.$$

We provide in Figure 5.3 some examples of Dirichlet distributions. Notice that by taking α as the vector of all ones, the distribution corresponds to a uniform distribution over the simplex $\Delta^{(\text{card}(\mathbb{H})-1)}$.

5.2. The Stochastic Majority Vote

Under Dirichlet assumptions on the hyper-posterior distribution P , we propose in the next section an algorithm to approximate the stochastic risk $s_P(\mathbf{x}, y)$. This algorithm is actually part of our self-bounding algorithm in Section 5.3.3.

5.2.2 Approximation of the Stochastic Risk

We now propose a Monte Carlo (MC) algorithm to compute $s_P(\mathbf{x}, y)$ that is suited to speed up the optimization. For this optimization algorithm, we need to introduce a surrogate of the true risk to update α by gradient descent as the gradients of the 01-loss are always zero. We make use of a *tempered* sigmoid loss $\text{sig}_c(x) = \frac{1}{1+\exp(-cx)}$ with slope parameter $c \in \mathbb{R}^+$. Because of the surrogate, this optimization algorithm solves a relaxation of the original problem and not its exact form (NESTEROV, 2005). The MC-based optimization algorithm is described in Algorithm 5.1.

Algorithm 5.1 Approximating the Stochastic Risk

Given: Dirichlet distribution $P = \text{Dir}(\alpha)$, learning sample \mathbb{S}

Hyperparameters: number of draws K

Draw a sample $\{\rho_k\}_{k=1}^K \sim P^K = \text{Dir}(\alpha)^K$

for all example $(\mathbf{x}_i, y_i) \in \mathbb{S}$ **do**

$$s_P(\mathbf{x}_i, y_i) \approx \frac{1}{K} \sum_{k=1}^K \text{sig}_c[-\widehat{m}_{\rho_k}(\mathbf{x}_i, y_i)]$$

return $\frac{1}{m} \sum_{i=1}^m s_P(\mathbf{x}_i, y_i)$

This algorithm first samples K majority votes and computes an approximation of the stochastic risk $s_P(\mathbf{x}_i, y_i)$ for each example $(\mathbf{x}_i, y_i) \in \mathbb{S}$ by an average. A drawback of Algorithm 5.1 is that it requires to sample K majority votes and predict all the examples in the learning sample \mathbb{S} . Hence, to overcome this issue, we derive a closed-form solution of the stochastic risk $s_P(\mathbf{x}, y)$ in the next section.

5.2.3 Computing Exactly the Stochastic Risk

Under Dirichlet assumptions, a closed-form solution can be derived for the expected risk. The following lemma introduces this solution.

Lemma 5.2.1 (Computation of the Stochastic Risk). For a given $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$, let

$$\mathbb{F}(\mathbf{x}, y) = \{j : h_j(\mathbf{x}) \neq y\} \quad \text{and} \quad \mathbb{T}(\mathbf{x}, y) = \{j : h_j(\mathbf{x}) = y\}$$

be respectively the set of indices of the voters that misclassify (\mathbf{x}, y) and the set of indices of the voters that correctly classify (\mathbf{x}, y) . Then, the stochastic risk $s_P(\mathbf{x}, y)$ can be rewritten as

$$s_P(\mathbf{x}, y) = \mathbb{E}_{\rho \sim P} \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0] = I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right),$$

with $I_{0.5}(\cdot)$ the regularized incomplete beta function evaluated at 0.5. It is defined as

$$I_{0.5}(a, b) \triangleq \frac{B_{0.5}(a, b)}{B_1(a, b)}, \quad \text{where} \quad B_t(a, b) \triangleq \int_0^t x^{a-1} (1-x)^{b-1} dx$$

is the incomplete beta function.

Proof. Deferred to Appendix E.3. ■

Lemma 5.2.1 tells us that the stochastic risk given an example (\mathbf{x}, y) can be computed with a closed-form solution. In consequence, thanks to Lemma 5.2.1, we compute an upper-bound on the risk of the stochastic majority vote based on the stochastic risk. Indeed, we deduce the following corollary.

Corollary 5.2.1 (Closed-form Solution of the Stochastic Risks). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any learning sample $\mathcal{S} \sim \mathcal{D}^m$, for any finite hypothesis set \mathbb{H} , for any distribution $P = \text{Dir}(\boldsymbol{\alpha})$ with $\boldsymbol{\alpha} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, we have

$$\begin{aligned} \mathbb{E}_{\rho \sim P} R_{\mathcal{D}}(\text{MV}_\rho) &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right), \\ \text{and} \quad \mathbb{E}_{\rho \sim P} R_{\mathcal{S}}(\text{MV}_\rho) &\leq \frac{1}{m} \sum_{i=1}^m s_P(\mathbf{x}_i, y_i) = \frac{1}{m} \sum_{i=1}^m I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}_i, y_i)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}_i, y_i)} \alpha_j \right). \end{aligned}$$

Proof. Deferred to Appendix E.4. ■

From Corollary 5.2.1, we propose to compute directly the empirical stochastic risk. This is in contrast with Algorithm 5.1 that approximates the stochastic risk by Monte Carlo sampling. The computation is summarized in Algorithm 5.2. Note that we provide in 5.4.1 an empirical study showing in which regimes each algorithm is more efficient.

5.3. From a PAC-Bayesian Bound to an Algorithm

Algorithm 5.2 Computing Exactly the Stochastic Majority Vote's Risk

Given: Dirichlet distribution $P = \text{Dir}(\alpha)$, learning sample \mathbb{S}
for all example $(\mathbf{x}_i, y_i) \in \mathbb{S}$ **do**

$$s_P(\mathbf{x}_i, y_i) = I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}_i, y_i)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}_i, y_i)} \alpha_j \right)$$

return $\frac{1}{m} \sum_{i=1}^m s_P(\mathbf{x}_i, y_i)$

Thanks to Algorithm 5.1 and Algorithm 5.2, we are now able to compute the empirical stochastic risk. This is a key step to derive our self-bounding algorithm in Section 5.3.3 for the stochastic majority vote. Indeed, the PAC-Bayesian generalization bound that we derive in the next section requires the computation of the stochastic risk in order to be minimized.

5.3 From a PAC-Bayesian Bound to an Algorithm

We now derive PAC-Bayesian generalization bounds for our proposed stochastic majority vote. To do so, we upper-bound the true stochastic risk with a SEEGER-like PAC-Bayesian bound. More precisely, we propose in Section 5.3.1 a PAC-Bayesian bound for a stochastic majority vote with voters that do not depend on the learning sample \mathbb{S} and in Section 5.3.2 we derive a PAC-Bayesian for data-dependent voters.

5.3.1 A PAC-Bayesian Bound for Stochastic Majority Votes

Before presenting our PAC-Bayesian bound for the stochastic majority vote, we consider that we have an *a priori* on the majority vote weights $\rho \in \mathbb{M}(\mathbb{H})$, *i.e.*, we assume a hyper-prior distribution Π over the voters' set \mathbb{H} . By doing so, we are able to derive a bound that depends on the KL divergence $\text{KL}(P \parallel \Pi)$ between the hyper-prior Π and the hyper-posterior P . Our bound is presented in the following theorem.

Theorem 5.3.1 (PAC-Bayesian Bound for Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any finite hypothesis set \mathbb{H} , for any distribution $\Pi = \text{Dir}(\beta)$ with $\beta \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over

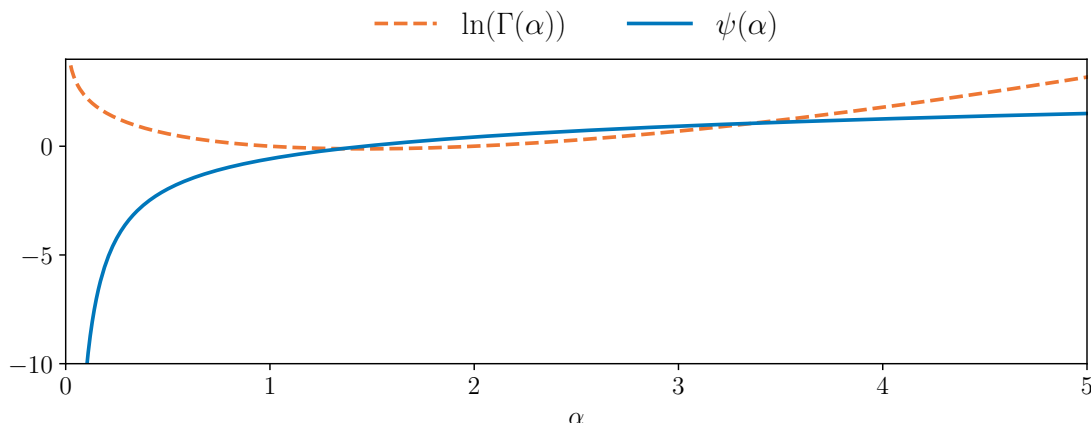


Figure 5.4. Plot of the Digamma function $\psi(\cdot)$ in the plain blue curve and its derivative $\ln(\Gamma(\cdot))$ (i.e., the logarithm of the Gamma function $\Gamma(\cdot)$) in the dotted orange curve.

the random choice of $\mathbb{S} \sim \mathcal{D}^m$, we have for all hyper-posterior P on \mathbb{H}

$$\mathbb{E}_{\rho \sim P} R_{\mathcal{D}}(\text{MV}_{\rho}) \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) \leq \bar{\text{kl}} \left(\frac{1}{m} \sum_{i=1}^m s_P(\mathbf{x}_i, y_i) \left| \frac{\text{KL}(P \parallel \Pi) + \ln \frac{2\sqrt{m}}{\delta}}{m} \right. \right),$$

$$\begin{aligned} \text{with } \text{KL}(P \parallel \Pi) &= \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\beta_j)] - \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \beta_j \right) \right] - \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\alpha_j)] \\ &\quad + \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right] + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \left[\psi(\alpha_j) - \psi \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right], \end{aligned}$$

where $\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$ is the Gamma function and the Digamma function $\Psi(\alpha)$ is defined as the derivative of $\ln[\Gamma(\alpha)]$; these two functions are plotted in Figure 5.4.

Proof. Deferred to Appendix E.5. ■

Note, while the bound shown in this theorem is based on SEEGER's form, our contribution does not restrict the choice of generalization bounds. Indeed, one can derive other versions of the bound based on MCALLESTER or CATONI's approach. As we will see in Section 5.3.3, we make use of Theorem 5.3.1 to derive a new learning algorithm for the stochastic majority vote.

5.3. From a PAC-Bayesian Bound to an Algorithm

As we have seen previously (notably in Section 2.3), the higher the KL divergence $\text{KL}(\mathbb{P} \parallel \Pi)$, the more different the two distributions \mathbb{P} and Π are. In this context, an increase of the Dirichlet parameters α can result in the increase of the KL divergence and a concentration of the Dirichlet distribution (see Figure 5.3). Indeed, when the parameters α increase, the stochastic majority vote risks tends to the risk of the majority vote with weight $\frac{\alpha_i}{\|\alpha\|_1}$ for the voter i . However, at the same time, the parameters increasing, favour a large KL divergence: the bound is a trade-off between the risk's concentration and the KL divergence. For now, Theorem 5.3.1 has a major drawback: it assumes that the voters are not dependent on the learning sample \mathcal{S} . Hence, to overcome this issue, we derive in Section 5.3.2 a PAC-Bayesian bound that allows us to have data-dependent voters.

5.3.2 A PAC-Bayesian Bound for Data-dependent Voters

Importantly, Theorem 5.3.1 holds when the hyper-prior Π and the set of voters \mathbb{H} are defined *a priori*, *i.e.*, they are independent from the data $\mathcal{S} \sim \mathcal{D}^m$. However, it is known that considering a data-dependent prior can lead to tighter PAC-Bayes bounds (PARRADO-HERNÁNDEZ *et al.*, 2012; DZIUGAITE *et al.*, 2021). Following recent works on PAC-Bayesian bounds with data-dependent priors (THIEMANN *et al.*, 2017; MHAMMEDI *et al.*, 2019), we derive a generalization bound that allows us to learn the voters from an additional set. More precisely, we consider two independent training sets \mathcal{S}_1 and \mathcal{S}_2 and we learn a set of voters on each training set (determining the set of voters \mathbb{H}_1 and \mathbb{H}_2). We refer to the hyper-prior distribution over \mathbb{H}_1 *resp.* over \mathbb{H}_2 as Π_1 *resp.* Π_2 . In the same way, we can then define a hyper-posterior distribution per voters' set: \mathbb{P}_1 and \mathbb{P}_2 . The following theorem shows that we can bound the risk of two combined stochastic majority votes, as long as their empirical risks are evaluated on the data split that was not used for learning their respective voters.

Theorem 5.3.2 (PAC-Bayesian bound with data-dependent voters). Let Π_1 and \mathbb{P}_1 be the hyper-prior and hyper-posterior distributions on \mathbb{H}_1 defined with \mathcal{S}_1 , and Π_2 and \mathbb{P}_2 the prior and posterior distributions on \mathbb{H}_2 defined with \mathcal{S}_2 . For any $\lambda \in [0, 1]$ and $\delta \in (0, 1]$ with probability at least $1 - \delta$ over samples $\mathcal{S}_1 \sim \mathcal{D}^{m_1}$ and $\mathcal{S}_2 \sim \mathcal{D}^{m_2}$, we have

$$\lambda \mathbb{E}_{\rho \sim \mathbb{P}_1} R_{\mathcal{D}}(\text{MV}_{\rho}) + (1 - \lambda) \mathbb{E}_{\rho' \sim \mathbb{P}_2} R_{\mathcal{D}}(\text{MV}_{\rho'}) \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\lambda s_{\mathbb{P}_1}(\mathbf{x}, y) + (1 - \lambda) s_{\mathbb{P}_2}(\mathbf{x}, y)] \leq \overline{\text{kl}} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_1} \frac{s_{\mathbb{P}_1}(\mathbf{x}, y)}{\frac{1}{\lambda}} + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_2} \frac{s_{\mathbb{P}_2}(\mathbf{x}, y)}{\frac{1}{1 - \lambda}} \left| \frac{\text{KL}(\mathbb{P}_1 \parallel \Pi_1) + \ln \frac{4\sqrt{m}}{\delta}}{\frac{m}{\lambda}} + \frac{\text{KL}(\mathbb{P}_2 \parallel \Pi_2) + \ln \frac{4\sqrt{m'}}{\delta}}{\frac{m'}{1 - \lambda}} \right. \right].$$

Proof. Deferred to Appendix E.6. ■

Following MHAMMEDI *et al.* (2019) we set $\lambda = 0.5$ and we applied a 50%/50% split in the training data. In the case when the number of data points is odd, we evaluate the bound with $m_2 = \text{card}(\mathcal{S}_2) - 1$ that still gives a correct bound (but simplifies the expression of the bound). In this case the bound is given in the following corollary.

Corollary 5.3.1 (PAC-Bayesian bound with data-dependent voters). Let Π_1 and P_1 be the hyper-prior and hyper-posterior distributions on \mathbb{H}_1 , and Π_2 and P_2 the prior and posterior distributions on \mathbb{H}_2 . For any $\delta \in (0, 1)$ with probability at least $1 - \delta$ over samples $\mathcal{S}_1 \sim \mathcal{D}^{m_1}$ and $\mathcal{S}_2 \sim \mathcal{D}^{m_2}$, we have

$$\frac{1}{2} \left[\mathbb{E}_{\rho \sim P_1} R_{\mathcal{D}}(\text{MV}_{\rho}) + \mathbb{E}_{\rho' \sim P_2} R_{\mathcal{D}}(\text{MV}_{\rho'}) \right] \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{1}{2} [s_{P_1}(\mathbf{x}, y) + s_{P_2}(\mathbf{x}, y)] \leq \overline{\text{kl}} \left[\frac{1}{2} \left(\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_1} s_{P_1}(\mathbf{x}, y) + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_2} s_{P_2}(\mathbf{x}, y) \right) \left| \frac{\text{KL}(P_1 \parallel \Pi_1) + \text{KL}(P_2 \parallel \Pi_2) + 2 \ln \frac{4\sqrt{m}}{\delta}}{m} \right. \right],$$

where $m = 2 \lfloor \frac{m_1 + m_2}{2} \rfloor$ and $\lfloor \cdot \rfloor$ is the floor function.

As for Theorem 5.3.1, the true risk of two combined stochastic majority votes is upper-bounded by a PAC-Bayesian bound that depends on two terms. Indeed, it depends on the empirical risk of the two stochastic majority votes and two KL divergences between the hyper-priors and the hyper-posteriors. This bound is evaluated in practice (in Section 5.4) when considering the data-dependent voters.

5.3.3 Learning Algorithms for the Stochastic Majority Vote

As in Chapters 3 and 4, we derive a self-bounding algorithm (FREUND, 1998) from Theorem 5.3.1 and Corollary 5.3.1. We based the derivation such an algorithm on the stochastic gradient descent. To do so, we consider mini-batches $\mathcal{U} \subseteq \mathcal{S}$ to optimize the generalization bounds. More precisely, the considered objective function for Theorem 5.3.1 is

$$G_{\mathcal{U}}(P) \triangleq \overline{\text{kl}} \left(\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{U}} s_P(\mathbf{x}, y) \left| \frac{\text{KL}(P \parallel \Pi) + \ln \frac{2\sqrt{m}}{\delta}}{m} \right. \right), \quad (5.4)$$

which is the upper bound applied on the mini-batch \mathcal{U} . To optimize such an objective function, we apply Algorithm 5.3 in conjunction with Algorithm 5.1 when we approximate the empirical stochastic risks or Algorithm 5.2 when we compute the risk exactly; the algorithm is summarized in the following algorithm.

5.3. From a PAC-Bayesian Bound to an Algorithm

Algorithm 5.3 Minimization of Theorem 5.3.1's Bound

Given: learning sample \mathcal{S} , hyper-prior distribution Π on \mathbb{H}

Hyperparameters: number of iterations T

$P \leftarrow \Pi$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathcal{U} \subseteq \mathcal{S}$ **do**

 Compute the empirical stochastic risk $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{U}} s_P(\mathbf{x}, y)$
 with Algorithm 5.1 or Algorithm 5.2

$P \leftarrow$ Update P with $G_{\mathcal{U}}(P)$ by gradient descent

return P

For each iteration we compute the empirical risk on the mini-batch \mathcal{U} . To do so, the risk is either approximated with Algorithm 5.1 or computed exactly with Algorithm 5.2. Then, we compute the objective function and update the hyper-posterior P with a gradient descent algorithm. For the data-dependent voters version, similarly as Corollary 5.3.1 which relies on the presence of two learning samples \mathcal{S}_1 and \mathcal{S}_2 , the objective function is defined as

$$G_{\mathcal{U}_1, \mathcal{U}_2}(P) \triangleq \overline{\text{kl}} \left[\frac{1}{2} \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{U}_1} s_{P_1}(\mathbf{x}, y) + \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{U}_2} s_{P_2}(\mathbf{x}, y) \right) \right. \\ \left. \left| \frac{\text{KL}(P_1 \parallel \Pi_1) + \text{KL}(P_2 \parallel \Pi_2) + 2 \ln \frac{4\sqrt{m}}{\delta}}{m} \right]. \quad (5.5)$$

This objective function estimated through a mini-batch \mathcal{U}_1 from \mathcal{S}_1 and \mathcal{U}_2 from \mathcal{S}_2 . Similarly as Equation (5.4), the objective function in Equation (5.5) estimates the upper-bound of Theorem 5.3.2. The algorithm considered to minimize such a bound is described in Algorithm 5.3.

Algorithm 5.4 Minimization of Corollary 5.3.1's Bound

Given: learning samples \mathcal{S}_1 and \mathcal{S}_2 , hyper-priors Π_1 on \mathbb{H}_1 and Π_2 on \mathbb{H}_2

Hyperparameters: number of iterations T

$(P_1, P_2) \leftarrow (\Pi_1, \Pi_2)$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathcal{U}_1 \subseteq \mathcal{S}_1$ and $\mathcal{U}_2 \subseteq \mathcal{S}_2$ **do**

 Compute the empirical stochastic risks $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{U}_1} s_{P_1}(\mathbf{x}, y)$ and
 $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{U}_2} s_{P_2}(\mathbf{x}, y)$ with Algorithm 5.1 or Algorithm 5.2

$(P_1, P_2) \leftarrow$ Update P_1 and P_2 with $G_{\mathcal{U}_1, \mathcal{U}_2}(P)$ by gradient descent

return (P_1, P_2)

For each iteration of the algorithm, the two stochastic risks are computed with Algorithm 5.1 or Algorithm 5.2. Then, thanks to these two values, we can compute the objective function $G_{\mathcal{U}_1, \mathcal{U}_2}(P)$ (Equation (5.5)). Finally, the hyper-posteriors P_1 and P_2 are updated through a gradient descent step.

Computing the derivatives. When the risk is computed from Algorithm 5.1, the sum is differentiable. However, since the risk is obtained by a Monte Carlo sampling, we use the implicit reparameterization trick (FIGURNOV *et al.*, 2018; JANKOWIAK and OBERMEYER, 2018) to obtain the derivatives; it is directly implemented in the automatic differentiation framework such as PyTorch (PASZKE *et al.*, 2019). Moreover, when the closed form solution (derived in Corollary 5.2.1) risk is computed through Algorithm 5.2, the risk depends on the function $I_{0.5}(\cdot)$ which is differentiable as well (see BOIK and ROBINSON-COX, 1999).

In the following section, we evaluate Algorithms 5.3 and 5.4 and compared them with the algorithms of Chapter 4.

5.4 Experiments

In this section, we compare the generalization bounds and the test risks obtained with our algorithms and the ones in Chapter 4. We show that our algorithms allow us to derive generalization bounds that are tight and non-vacuous (*i.e.*, smaller than 1) with decision stumps and decision trees.

We consider as baselines the following PAC-Bayesian self-bounding algorithms:

- (i) The algorithm `Germain` that minimizes the PAC-Bayesian of GERMAIN *et al.* (2015, PAC-Bound 0) on $2r_{\mathcal{D}}(\rho)$ (Theorem 2.2.1),
- (ii) The algorithm `Masegosa` of MASEGOSA *et al.* (2020) minimizing a PAC-Bayesian bound on $4e_{\mathcal{D}}(\rho)$ (Theorem 2.2.2),
- (iii) Algorithm 4.2 minimizing the PAC-Bayesian C-Bound based on SEEGER's approach (Theorem 4.3.4),
- (iv) `Lacasse` (LACASSE *et al.*, 2010) minimizes Equation (5.1) for a majority vote that samples N voters from ρ .

The parameters of the baseline are the same as in Section 4.5. For `Lacasse` the number of voters drawn is set to $N=100$. As in Chapter 4 the generalization bounds are evaluated with $\delta = 0.05$ and the sigmoid's slope parameter c is set to 100 for Algorithm 5.2. Moreover, the values are averaged over 10 runs.

5.4.1 Comparison Between the Computations of the Risk

For this set of experiments, we optimize Theorem 5.3.1 with Algorithm 5.3, for $T = 2,000$ iterations with COCOB-Backprop optimizer (ORABONA and TOMMASI, 2017). We study the performance of our method on the binary classification moons dataset, with 2 features, 2 classes and $\mathcal{N}(0, 0.05)$ Gaussian noise, for which we draw m points for training, and $\text{card}(\mathbb{T}) = 2,000$ points for testing.

Figure 5.5 reports a comparison of Algorithms 5.1 and 5.2 in terms of test risk $\mathbb{E}_{\rho \sim \mathcal{P}} R_{\mathcal{T}}(\text{MV}_{\rho})$, PAC-Bayesian generalization bound and training time when the number of decision stumps increases (with $m = 2,000$). We observe that the test risks and bound values can degrade for higher values of decision stumps for all methods. This is due to the KL divergence increasing with the number of voters $\text{card}(\mathbb{H})$, as highlighted in Appendix E.7, becoming a too strong regularization during training and making the bound looser. Moreover, when the number of decision stumps increases, Algorithm 5.2 can be quicker than Algorithm 5.1 especially when the Monte Carlo draws K is high compared to m .

We report in Figure 5.6 the evolution of the test risk $\mathbb{E}_{\rho \sim \mathcal{P}} R_{\mathcal{T}}(\text{MV}_{\rho})$, PAC-Bayesian bounds and training time when m increases. When m is large enough, Algorithm 5.1 achieves comparable test risk and bound values compared to Algorithm 5.2 even for $K = 1$. Increasing the number of Monte Carlo draws K unsurprisingly allows to recover Algorithm 5.2's performance, and at lower computational cost for reasonable values of m and K .

5.4.2 Performance of Algorithms 5.3 and 5.4

We now compare Algorithms 5.3 and 5.4 on different datasets namely FashionMNIST (XIAO *et al.*, 2017), MNIST (LECUN *et al.*, 1998) and datasets coming from the UCI repository (DUA and GRAFF, 2017); the processing of the dataset is the same as in Section 4.5. More precisely, the same number of examples is kept in the test or the train set as in the original split. When no original split was proposed, we use 50% of data in the training set \mathbb{S} and 50% in the test set (except for Sensorless where we have 15% in the test set). When making use of *data-independent voters*, we chose decision stumps as voters; When making use of *data-dependent voters*, we build decision trees as set of voters without bounding their maximal depth (unless stated otherwise). The voters are exactly the same as in Chapter 4.

We train the stochastic majority vote models by Stochastic Gradient Descent (SGD) using COCOB-Backprop (ORABONA and TOMMASI, 2017) with batch size equal to 64. We fix the number of epochs to 20 and for Algorithm 5.1 we fix $K = 10$ to increase randomness.

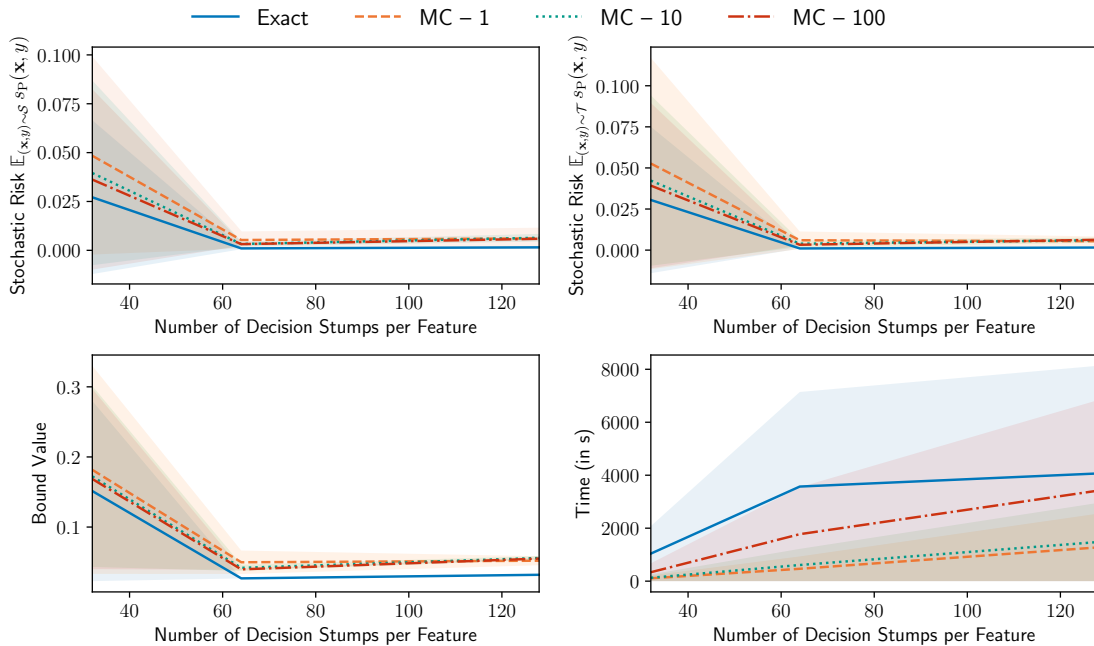


Figure 5.5. Plot of the average performance on 10 runs of Algorithm 5.1 (with $K \in \{1, 10, 100\}$) and Algorithm 5.2 as a function of the number of decision stumps per feature with learning sample size $m = 2000$.

We report the stochastic test risks $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{T}} s_P(\mathbf{x}, y)$ for our algorithms, the test risk $R_{\mathcal{T}}(\text{MV}_{\rho})$ for the others and the generalization bounds in Figures 5.7 to 5.10 (additional results are reported in Appendix E.7). More precisely, we compare the different self-bounding algorithms with Algorithm 5.3 on binary datasets in Figures 5.7 and 5.8, and on multi-class datasets with data-dependent voters and Algorithm 5.4 in Figures 5.9 and 5.10. First, we remark that Algorithms 5.3 and 5.4 have similar performance in terms of stochastic test risks and bound values. Moreover, note that we notice that the bounds obtained from Algorithms 5.3 and 5.4 are consistently non vacuous and tighter than those obtained with the other algorithms. While the risks between the methods are not comparable we remark that the stochastic test risk $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{T}} s_P(\mathbf{x}, y)$ has similar values than the test risks $R_{\mathcal{T}}(\text{MV}_{\rho})$. It means that our algorithms obtain similar test risk $R_{\mathcal{T}}(\text{MV}_{\rho})$ in expectation (where $\rho \sim P$). We believe that it is due to the fact that the stochastic risks depend on the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017). Indeed, even if our learning algorithms optimizes the 01-loss, it does not fully distinguish examples that are classified correctly or not.

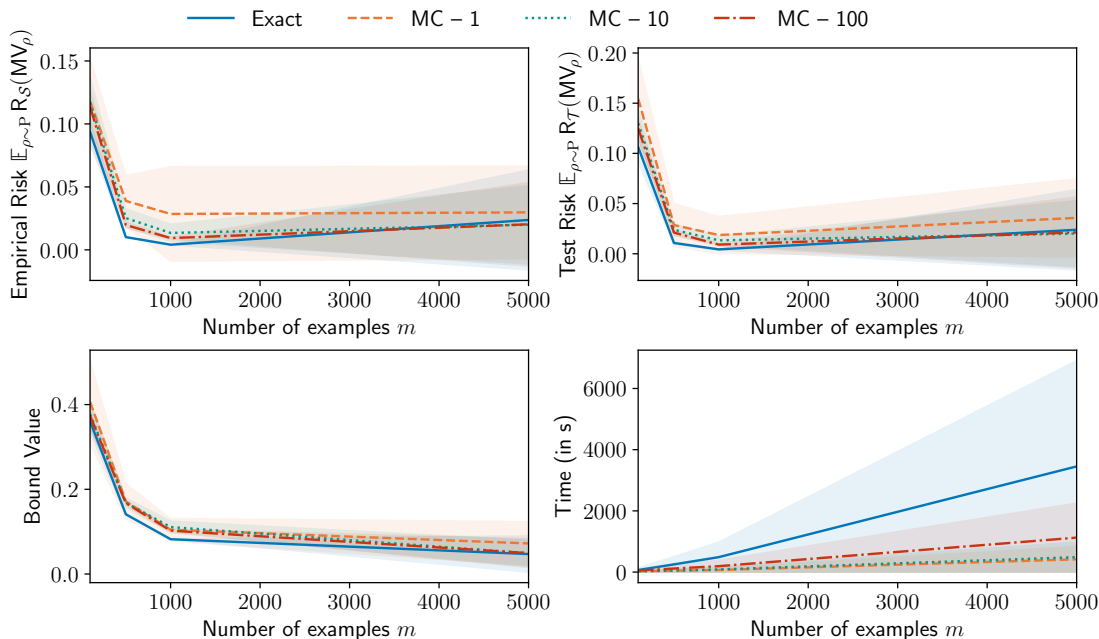


Figure 5.6. Plot of the average performance on 10 runs of Algorithm 5.1 (with $K \in \{1, 10, 100\}$) and Algorithm 5.2 as a function of the learning sample size m with 32 decision stump per feature.

5.5 Conclusion and Summary

In this chapter, we have studied a new type of majority vote: the *stochastic majority vote*. For each input $\mathbf{x} \in \mathcal{X}$, it samples a majority vote MV_ρ from a probability distribution called hyper-posterior P and output $MV_\rho(\mathbf{x})$. When the hyper-posterior is a Dirichlet distribution, the stochastic risk can be either approximated or computed exactly. This allows us to derive a self-bounding algorithm for the stochastic majority vote. The experiments show that our learning algorithm provides a tight PAC-Bayesian generalization bound along with a small empirical risk.

One of the perspectives of this work is to consider the risk of the stochastic majority vote by doing without the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017). However, we may not find the closed-form solution of the stochastic majority vote's risk, in this context. In other words, the risk might be only approximated by Monte Carlo sampling. To avoid such a drawback we can make use of a different type of bounds: the *disintegrated* PAC-Bayesian bounds. Indeed, in Part III, we study the *disintegrated* bounds in more details. As recalled in Chapter 2, these bounds have been introduced by BLANCHARD and FLEURET (2007) and CATONI (2007) and have been rediscovered lately

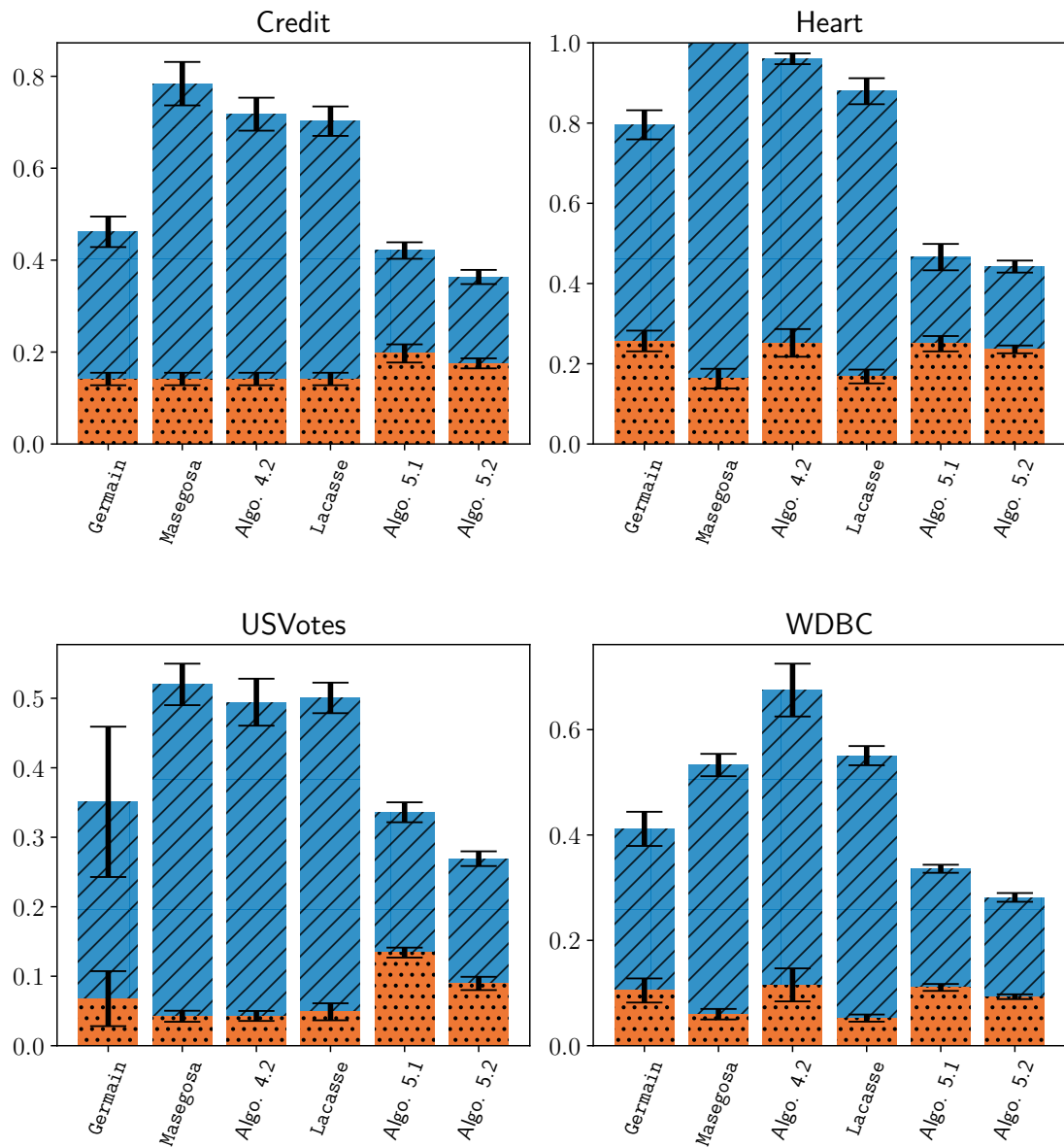


Figure 5.7. Comparison in terms of test risks, stochastic test risks and bound values. We report in the dotted orange error bars, the means and standard deviations of the test risks and stochastic risks. Moreover, the hatched blue error bars represent the mean and the standard deviations of the PAC-Bayesian bounds. The values are average over 10 different runs.

5.5. Conclusion and Summary

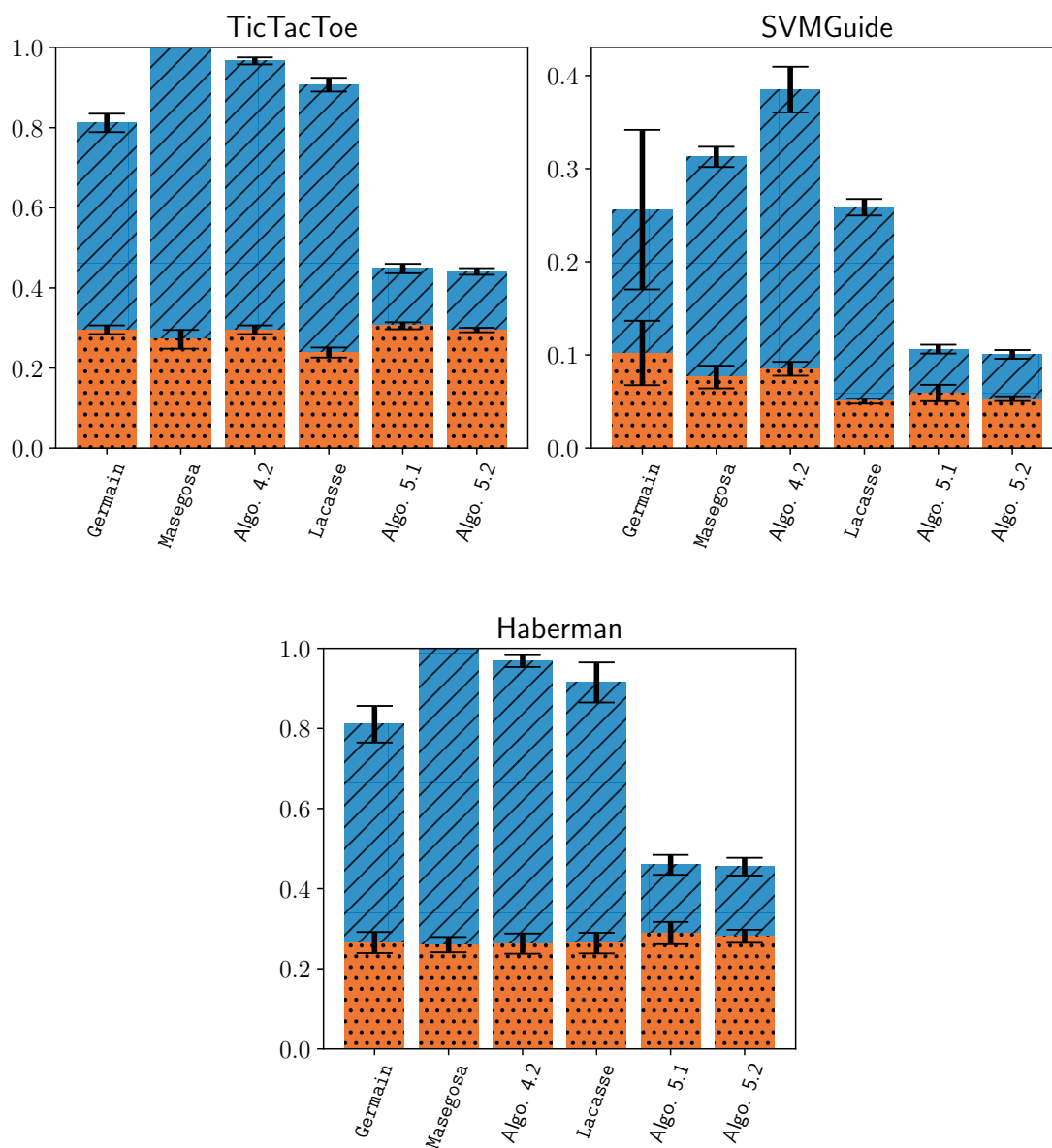


Figure 5.8. Comparison in terms of test risks, stochastic test risks and bound values. We report in the dotted orange error bars, the means and standard deviations of the test risks and stochastic risks. Moreover, the hatched blue error bars represent the mean and the standard deviations of the PAC-Bayesian bounds. The values are average over 10 different runs.

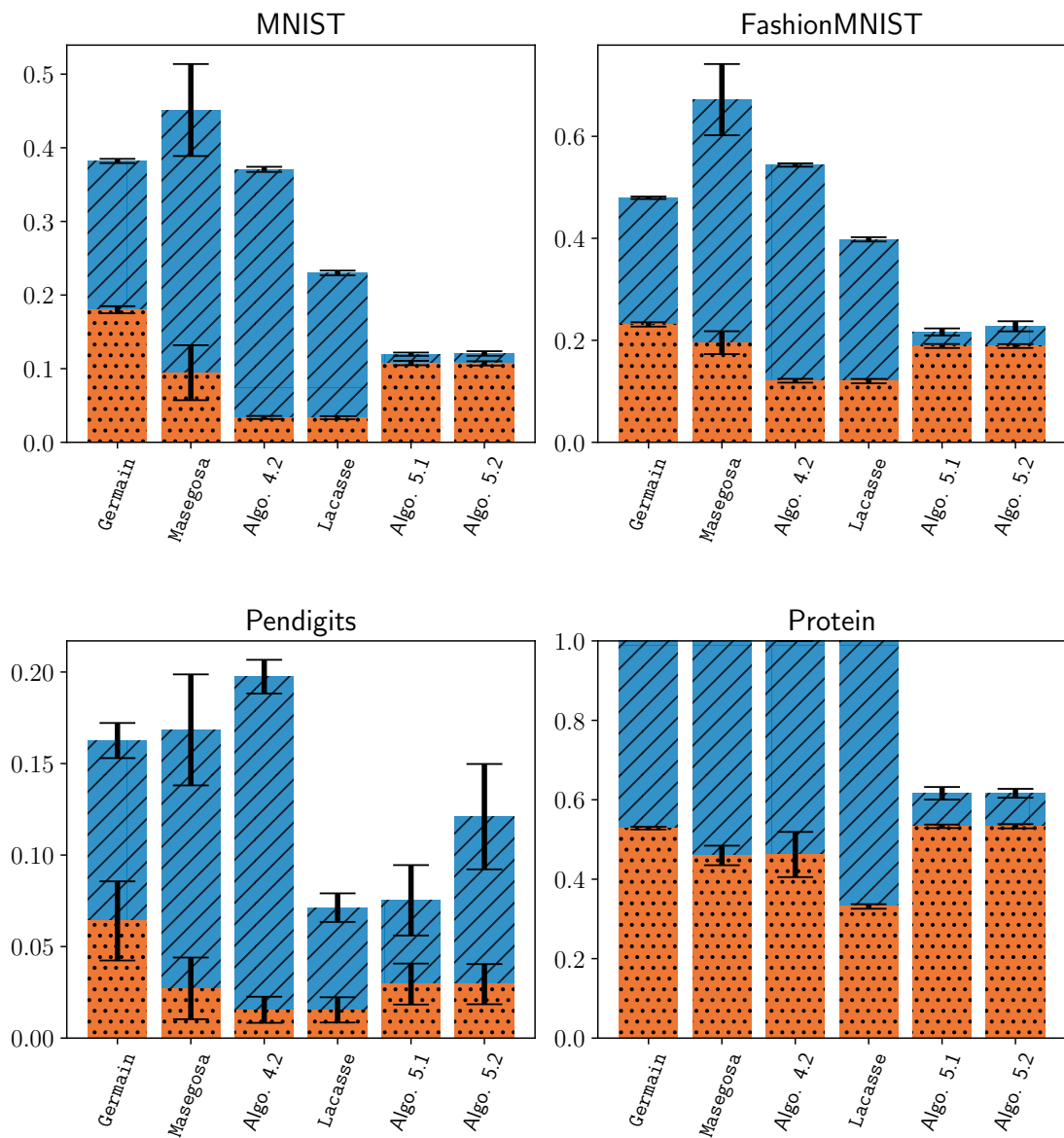


Figure 5.9. Comparison in terms of test risks, stochastic test risks and bound values. We report in the dotted orange error bars, the means and standard deviations of the test risks and stochastic risks. Moreover, the hatched blue error bars represent the mean and the standard deviations of the PAC-Bayesian bounds. The values are average over 10 different runs.

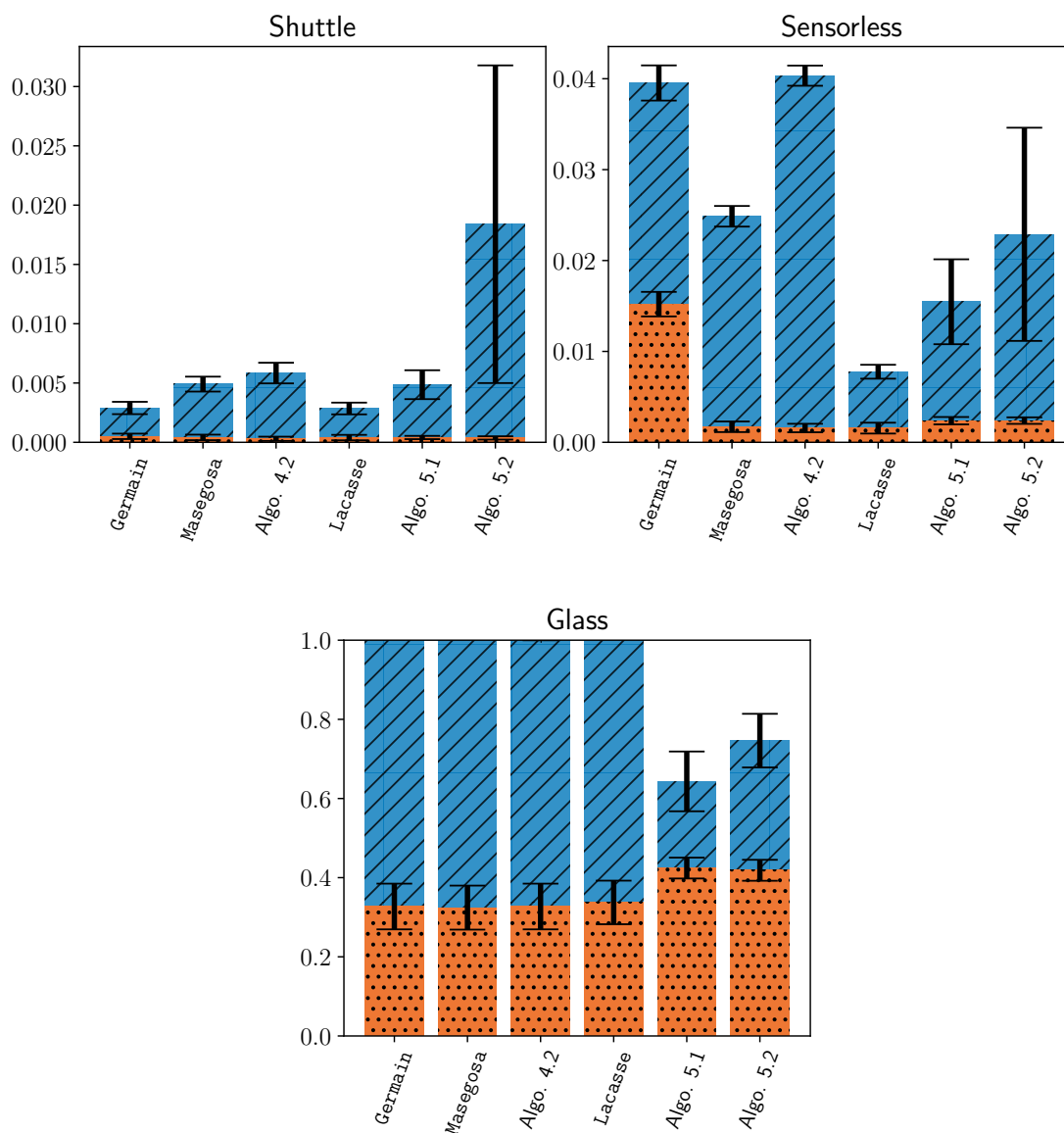


Figure 5.10. Comparison in terms of test risks, stochastic test risks and bound values. We report in the dotted orange error bars, the means and standard deviations of the test risks and stochastic risks. Moreover, the hatched blue error bars represent the mean and the standard deviations of the PAC-Bayesian bounds. The values are average over 10 different runs.

by RIVASPLATA *et al.* (2020). They allow us to derive a bound for a *single* voter or hypothesis from the hypothesis set \mathbb{H} . Chapter 6 introduces new *disintegrated* bounds based on the RÉNYI divergence (that can be used for the stochastic majority vote) and Chapter 7 presents generalization bounds that are not restrained to the KL divergence or the RÉNYI one but can depend on a complexity term defined by the users.

PART III

**Beyond PAC-Bayesian Bounds:
From Disintegration to Novel Bounds**

ON THE PRACTICAL USES OF THE DISINTEGRATED BOUNDS

This chapter is based on the following paper

PAUL VIALARD, PASCAL GERMAIN, AMAURY HABRARD, and EMILIE MORVANT. A General Framework for the Disintegration of PAC-Bayesian Bounds. *Submitted to Machine Learning Journal*. (2022b)

Contents

6.1	Introduction	150
6.2	Setting and PAC-Bayesian Bounds	151
6.3	Disintegrated PAC-Bayesian Theorems	152
6.3.1	Form of a Disintegrated PAC-Bayesian Bound	152
6.3.2	Disintegrated Bounds with the RÉNYI Divergence	153
6.4	The Disintegration in Action	158
6.4.1	Specialization to Neural Network Classifiers	158
6.4.2	A Note About Stochastic Neural Networks	161
6.5	Experiments with Neural Networks	162
6.5.1	Training Method	162
6.5.2	Optimization Procedure in Algorithms A and A_{prior}	163
6.5.3	Experimental Setting	163
6.5.4	Results	164
6.6	Perspectives for the Majority Vote	171
6.7	Summary and Conclusion	172

Abstract

PAC-Bayesian bounds are known to be tight and informative when studying the generalization ability of stochastic classifiers (see *e.g.*, Chapter 5). However, they require a loose and costly derandomization step when applied to some families of deterministic models such as neural networks. As an alternative to this step, we introduce new PAC-Bayesian generalization bounds that have the originality to provide *disintegrated* bounds, *i.e.*, they give guarantees over one *single* hypothesis instead of the usual *averaged* analysis. Our bounds are easily optimizable and can be used to design self-bounding algorithms. We illustrate this behavior on neural networks and show a significant practical improvement over the state-of-the-art framework.

6.1 Introduction

The PAC-Bayesian theory is a powerful framework for upper-bounding the true risk of stochastic models such as the stochastic majority vote (considered in Chapter 5). Remember that, in general, the stochastic model samples an hypothesis from the posterior distribution and predicts a label with the sampled hypothesis. However, the vast majority of machine learning methods nevertheless need guarantees on deterministic models (*i.e.* that are not stochastic). In this case, a *derandomization step* of the bound is required to get a bound on the deterministic model’s risk. In general, the *derandomization step* consists in obtaining a bound on the risk of a deterministic model from a bound originally valid for stochastic models. Different forms of derandomization have been introduced in the literature for specific settings. Among them, LANGFORD and SHAWE-TAYLOR (2002) proposed a derandomization for Gaussian posteriors over linear classifiers: thanks to the Gaussian symmetry, a bound on the risk of the *maximum a posteriori* (deterministic) classifier is obtainable from the bound on the average risk of the stochastic classifier. Also relying on Gaussian posteriors, LETARTE *et al.* (2019) derived a PAC-Bayesian bound for a very specific deterministic network architecture using sign functions as activations; this approach has been further extended by BIGGS and GUEDJ (2021, 2022). Another line of works derandomizes neural networks (NEYSHABUR *et al.*, 2018; NAGARAJAN and KOLTER, 2019b). While technically different, it starts from PAC-Bayesian guarantees on the stochastic classifier and uses an “output perturbation” bound to convert a guarantee from a random classifier to the mean classifier. The relative diversity and specificity of these works highlight nevertheless the lack of a general framework for the derandomization of classic PAC-Bayesian bounds.

This chapter focuses on another kind of derandomization through the *disintegration of the PAC-Bayesian bound*, proposed by CATONI (2007, Th.1.2.7) and BLANCHARD and FLEURET (2007); see Section 2.4. Despite their interest in derandomizing PAC-Bayesian bounds, these kinds of bounds have only received little study in the literature and have never been used in practice. Driven by machine learning practical purposes, our objective is thus twofold. We derive new tight and usable *disintegrated* PAC-Bayesian bounds (*i*) that directly derandomize any classifiers without any other additional step and with *almost* no impact on the guarantee, and (*ii*) that can be easily optimized to learn classifiers with strong guarantees. To achieve this objective, our contribution consists in providing a new general disintegration framework based on the RÉNYI divergence (in Theorem 6.3.1), allowing us to meet the practical goal of efficient learning. From a theoretical standpoint, due to the RÉNYI divergence term, our bound is expected to be looser than the one of RIVASPLATA *et al.* (2020, Th.1(*i*)) in which the divergence term is “disintegrated” but depends on the sampled hypothesis only. However, as we show in our experimental evaluation on neural networks, their

“disintegrated” term is, in practice, subject to high variance, making their bound harder to optimize. This variance arises because the sampled hypothesis does not influence our RÉNYI divergence term. Our bound has then the main advantage of leading to a more stable learning algorithm with better empirical results. In addition, we derive in Appendix F.9 new theoretical results based on the mutual information, giving different insights into disintegration procedures.

The rest of the chapter is organized as follows. Section 6.2 recalled the notations we follow. In Section 6.3, we derive our main contribution relying on *disintegrated* PAC-Bayesian bounds. Then, we illustrate the practical usefulness of this disintegration on deterministic neural networks in Section 6.5. Before concluding in Section 6.7, we discuss in Appendix F.9 another point of view of the disintegrated through an information-theoretic bound. For readability, we deferred the proofs of our theoretical results to Appendix F.

6.2 Setting and PAC-Bayesian Bounds

In this chapter, we consider supervised classification tasks as described in Chapter 2 with \mathcal{X} the input space, \mathcal{Y} the label set, and \mathcal{D} an unknown data distribution on $\mathcal{X} \times \mathcal{Y}$. We consider a hypothesis set \mathbb{H} of functions $h : \mathcal{X} \rightarrow \mathcal{Y}$. The learner aims to find $h \in \mathbb{H}$ that assigns a label y to an input \mathbf{x} as accurately as possible. Given an example (\mathbf{x}, y) and a hypothesis h , we assess the quality of the prediction of h with a loss function $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ evaluating to which extent the prediction is accurate. The learner wants to find the hypothesis h from \mathbb{H} that minimizes the true risk $R_{\mathcal{D}}^{\ell}(h)$. However, we cannot compute $R_{\mathcal{D}}^{\ell}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \ell(h, (\mathbf{x}, y))$ since the distribution \mathcal{D} is unknown. We only have access to a learning sample $\mathbb{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ with the empirical risk defined as $R_{\mathbb{S}}^{\ell}(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, (\mathbf{x}_i, y_i))$.

We now recall BÉGIN *et al.* (2016)’s general bound at the heart of our contribution (and introduced in Section 2.3). This bound depends on the RÉNYI divergence between ρ and π defined as $D_{\lambda}(\rho \parallel \pi) \triangleq \frac{1}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \pi} \left[\frac{\rho(h)}{\pi(h)} \right]^{\lambda} \right]$ with parameter $\lambda > 1$.

Theorem 2.3.5 (General PAC-Bayesian Bound of BÉGIN *et al.* (2016)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_*^+$, for any $\lambda > 1$, for any

$\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \right] \right. \\ \left. \leq D_\lambda(\rho \| \pi) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right] \right] \geq 1 - \delta. \end{aligned}$$

A key notion here is that the PAC-Bayesian bounds apply on the expectation over the risk of the individual classifiers in \mathbb{H} ; this randomization is the risk of the stochastic classifier. A key issue for usual machine learning tasks is then the derandomization of the PAC-Bayesian bounds to obtain a guarantee for a deterministic classifier instead of a stochastic one (by removing the expectation on \mathbb{H}). In some cases, this derandomization results from the structure of the hypotheses, such as for stochastic linear classifiers that can be directly expressed as one deterministic linear classifier (GERMAIN *et al.*, 2009). However, in other cases, the derandomization is much more complex and specific to the class of hypotheses, such as for neural networks (e.g., NEYSHABUR *et al.* (2018), NAGARAJAN and KOLTER (2019a, Ap. J), BIGGS and GUEJ (2022)).

The next section states our main contribution to this chapter: a general derandomization framework based on the RÉNYI divergence for disintegrating PAC-Bayesian bounds into a bound for a single hypothesis from \mathbb{H} .

6.3 Disintegrated PAC-Bayesian Theorems

6.3.1 Form of a Disintegrated PAC-Bayesian Bound

We recall now the main form of a disintegrated PAC-Bayesian bound used in this chapter.

Definition 2.4.1 (Disintegrated PAC-Bayesian Generalization Bound). Let $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow [0, 1]$ a generalization gap. A *disintegrated* PAC-Bayesian bound is defined such that if for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathcal{M}^*(\mathbb{H})$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$, there exists a function $\Phi : \mathcal{M}(\mathbb{H}) \times \mathcal{M}^*(\mathbb{H}) \times (0, 1] \rightarrow \mathbb{R}$ such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(R_{\mathcal{D}}^\ell(h), R_{\mathcal{S}}^\ell(h)) \leq \Phi(\rho_{\mathcal{S}}, \pi, \delta) \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A and $\phi(\cdot)$ is, for example, $\phi(\mathbb{R}_{\mathcal{D}}^{\ell}(h), \mathbb{R}_{\mathcal{S}}^{\ell}(h)) = |\mathbb{R}_{\mathcal{D}}^{\ell}(h) - \mathbb{R}_{\mathcal{S}}^{\ell}(h)|$.

More precisely, the posterior $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is defined through a given *deterministic* algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$ chosen *a priori*. The algorithm (i) takes a learning sample $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$ and a prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ as inputs, and (ii) outputs a *data-dependent* distribution $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$. Concretely, this kind of generalization bound allows one to derandomize the usual PAC-Bayes bounds as follows. Instead of considering a bound holding for all the posterior distributions on \mathbb{H} as usually done in PAC-Bayes (the “ $\forall \rho$ ” in Theorem 2.3.5), we consider only the posterior distribution $\rho_{\mathcal{S}}$ obtained through a deterministic algorithm A taking the learning sample \mathcal{S} and the prior π as inputs. Then, the above bound holds for a unique hypothesis $h \sim \rho_{\mathcal{S}}$ instead of the stochastic classifier: the individual risks are no longer averaged with respect to $\rho_{\mathcal{S}}$; this is the *PAC-Bayesian bound disintegration*. The dependence in probability on $\rho_{\mathcal{S}}$ means that the bound is valid with probability at least $1 - \delta$ over the random choice of the learning sample $\mathcal{S} \sim \mathcal{D}^m$ and the hypothesis $h \sim \rho_{\mathcal{S}}$. Under this principle, we introduce in Theorems 6.3.1 and 6.3.2 below two new general disintegrated PAC-Bayesian bounds. A key asset of our results is that the bounds are instantiable to specific settings as for the “classical” PAC-Bayesian bounds (e.g., with *i.i.d./non-i.i.d.* data, unbounded losses, etc.). By instantiating such a bound, we obtain an easily optimizable bound, leading to a self-bounding algorithm (FREUND, 1998) with theoretical guarantees. As an illustration of the usefulness of our results, we provide, in Section 6.4, such an instantiation for neural networks.

6.3.2 Disintegrated Bounds with the Rényi Divergence

6.3.2.1 Our Main Contribution: a General Disintegrated Bound

In the same spirit as Theorem 2.3.5 our main result stated in Theorem 6.3.1 is a general bound involving the RÉNYI divergence $D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi)$ of order $\lambda > 1$.

Theorem 6.3.1 (General Disintegrated PAC-Bayes Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} & \left(\frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \right. \\ & \left. \leq \frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} + D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) + \ln \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \right) \geq 1 - \delta, \end{aligned}$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A .

Proof sketch (see Appendix F.1 for details). Recall that $\rho_{\mathbb{S}}$ is obtained with the algorithm $A(\mathbb{S}, \pi)$. Applying MARKOV's inequality (Theorem A.2.1) on $\varphi(h, \mathbb{S})$ with the random variable h and using HÖLDER's inequality (Theorem A.5.1) to introduce $D_{\lambda}(\rho_{\mathbb{S}} \parallel \pi)$, we have, with probability at least $1 - \frac{\delta}{2}$ on $\mathbb{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathbb{S}}$,

$$\begin{aligned} \frac{\lambda}{\lambda-1} \ln [\varphi(h, \mathbb{S})] &\leq \frac{\lambda}{\lambda-1} \ln \left[\frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathbb{S}}} \varphi(h', \mathbb{S}) \right] \\ &\leq D_{\lambda}(\rho_{\mathbb{S}} \parallel \pi) + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \right]. \end{aligned}$$

By applying again MARKOV's inequality (Theorem A.2.1) on $\varphi(h, \mathbb{S})$ with the random variable \mathbb{S} , we have, with probability at least $1 - \frac{\delta}{2}$ on $\mathbb{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathbb{S}}$,

$$\ln \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \right] \leq \ln \left[\frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right].$$

Lastly, we combine the two bounds with a union bound argument. ■

As for the general classical PAC-Bayesian bounds (Theorem 2.3.5), the above theorem can be seen as the starting point of the derivation of generalization bounds depending on the choice of the function φ , as done in Corollary 6.4.1 in Section 6.4.1; this property makes it the main result of this chapter. In its proof, HÖLDER's inequality (Theorem A.5.1) is used differently than in the classic PAC-Bayes bound's proofs. Indeed, in the proof of BÉGIN *et al.* (2016, Th. 8), the change of measure based on HÖLDER's inequality is key for deriving a bound that holds for all posteriors ρ with high probability, while our bound holds for a unique posterior $\rho_{\mathbb{S}}$ dependent on the sample \mathbb{S} and the prior π . In fact, we use HÖLDER's inequality to introduce a prior π independent from \mathbb{S} : a crucial point for our bound instantiated in Corollary 6.4.1. Compared to Theorem 2.3.5, our bound requires an additional term $\ln 2 + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta}$. However, by setting $\varphi(h, \mathbb{S}) = m \text{kl}(\mathbb{R}_{\mathbb{S}}^{\ell}(h) \parallel \mathbb{R}_{\mathcal{D}}^{\ell}(h))$ and $\lambda=2$, the term $\ln \frac{8}{\delta^2}$ is multiplied by $\frac{1}{m}$, which turns out to be a reasonable cost to “derandomize” a bound into a disintegrated one. For instance, if $m = 5,000$ (a reasonable sample size) and $\delta = 0.05$, we have $\frac{1}{m} \ln \frac{8}{\delta^2} \approx 0.002$.

We instantiate below Theorem 6.3.1 for $\lambda \rightarrow 1^+$ and $\lambda \rightarrow +\infty$ showing that the bound converges when $\lambda \rightarrow 1^+$ and $\lambda \rightarrow +\infty$.

Corollary 6.3.1 (Extreme Cases of Theorem 6.3.1). Under the assumptions of Theorem 6.3.1, when $\lambda \rightarrow 1^+$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\ln \varphi(h, \mathcal{S}) \leq \ln \frac{2}{\delta} + \ln \left[\operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathbb{H}} \varphi(h', \mathcal{S}') \right] \right) \geq 1 - \delta,$$

when $\lambda \rightarrow +\infty$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\ln \varphi(h, \mathcal{S}) \leq \ln \operatorname{esssup}_{h' \in \mathbb{H}} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} + \ln \left[\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}') \right] \right) \geq 1 - \delta,$$

where esssup is the essential supremum defined as the supremum on a set with non-zero probability measures, *i.e.*,

$$\operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathbb{H}} \varphi(h', \mathcal{S}') = \inf \left\{ \tau \in \mathbb{R}, \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [\varphi(h, \mathcal{S}) > \tau] = 0 \right\},$$

and $\operatorname{esssup}_{h' \in \mathbb{H}} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} = \inf \left\{ \tau \in \mathbb{R}, \mathbb{P}_{h \sim \rho_{\mathcal{S}}} \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} > \tau \right] = 0 \right\}.$

Proof. Deferred to Appendix F.2. ■

This corollary illustrates that the parameter λ controls the trade-off between the RÉNYI divergence $D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi)$ and $\ln \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right]$. Indeed, when $\lambda \rightarrow 1^+$, the RÉNYI divergence vanishes ($D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) \rightarrow 0$) while the other term converges toward $\ln \left[\operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathbb{H}} \varphi(h', \mathcal{S}') \right]$, roughly speaking the maximal value possible for the second term. On the other hand, when $\lambda \rightarrow +\infty$, the RÉNYI divergence increases and converges toward $\ln \operatorname{esssup}_{h' \in \mathbb{H}} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')}$ and the other term decreases toward $\ln \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}') \right]$.

6.3.2.2 Comparison with the Bound of Rivasplata *et al.* (2020)

For the sake of comparison, we recall in Theorem 2.4.1 the bound proposed by RIVASPLATA *et al.* (2020, Th.1(i)), that is more general than the bounds of BLANCHARD and FLEURET (2007) and CATONI (2007, Th.1.2.7).

Theorem 2.4.1 (General Disintegrated Bound of RIVASPLATA *et al.* (2020)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$,

for any algorithm $A: (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \underbrace{\ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathcal{S}')) \right]}_{\Phi(\rho_{\mathcal{S}}, \pi, \delta)} \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Note that the bound can be rewritten with the logarithm, *i.e.*, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\ln(\varphi(h, \mathcal{S})) \leq \ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}') \right] \right] \geq 1 - \delta.$$

The term $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$ (also involved in BLANCHARD and FLEURET (2007) and CATONI (2007)) can be seen as a “disintegrated¹ KL divergence” depending only on the sampled $h \sim \rho_{\mathcal{S}}$. In contrast, our bound involves the RÉNYI divergence $D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi)$ between the prior π and the posterior $\rho_{\mathcal{S}}$, meaning our bound involves only one term that depends on the sample hypothesis (the risk): the divergence value is the same whatever the hypothesis. Our bound is expected to be looser because of the RÉNYI divergence (see ERVEN and HARREMOËS, 2014) and the dependence in δ (which is worse than in Theorem 2.4.1). Nevertheless, our divergence term is the main advantage of our bound. Indeed, as confirmed by our experiments (Section 6.5), our bound with $D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi)$ makes the learning procedure (in our self-bounding algorithm) more stable and efficient compared to the optimization of Theorem 2.4.1’s bound with $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$ that is subject to high variance.

6.3.2.3 A Parameterizable General Disintegrated Bound

In the PAC-Bayesian literature, parameterized bounds have been introduced (*e.g.*, CATONI (2007) and THIEMANN *et al.* (2017)) to control the trade-off between the empirical risk and the divergence along with the additional term. For the sake of completeness, we now provide a parameterized version of our bound, enlarging its practical scope. We follow a similar approach to introduce a version of a disintegrated RÉNYI divergence-based bound that has the advantage of being parameterizable.

Theorem 6.3.2 (Parameterizable Disintegrated PAC-Bayes Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$,

¹We say that the KL divergence is “disintegrated” since the log term is not averaged in contrast to the KL divergence.

6.3. Disintegrated PAC-Bayesian Theorems

for any measurable function $\varphi: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\substack{\mathbb{S} \sim \mathcal{D}^m, \\ h \sim \rho_{\mathbb{S}}}} \left(\forall \lambda > 0, \ln(\varphi(h, \mathbb{S})) \leq \ln \left[\frac{\lambda}{2} e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{8}{2\lambda\delta^3} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', \mathbb{S}')^2] \right] \right) \geq 1 - \delta,$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A .

Proof. Deferred to Appendix F.3. ■

Note that $e^{D_2(\rho_{\mathbb{S}} \parallel \pi)}$ is closely related to the χ^2 -distance. Indeed we have: $\chi^2(\rho_{\mathbb{S}} \parallel \pi) \triangleq \mathbb{E}_{h \sim \pi} \left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right]^2 - 1 = e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} - 1$. An asset of Theorem 6.3.2 is the parameter λ controlling the trade-off between the exponentiated R ENYI divergence $e^{D_2(\rho_{\mathbb{S}} \parallel \pi)}$ and $\frac{1}{\delta^3} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}')^2$. Our bound is valid for all $\lambda > 0$, thus, from a practical view, we can learn/tune the parameter λ to minimize the bound and control the possible numerical instability due to $e^{D_2(\rho_{\mathbb{S}} \parallel \pi)}$. Indeed, if $D_2(\rho_{\mathbb{S}} \parallel \pi)$ is large, the numerical computation can lead to an infinite value due to finite precision arithmetic. It is important to notice that, like other parameterized bounds (e.g., THIEMANN *et al.*, 2017), there exists a closed-form solution of the optimal parameter λ (for a fixed π and $\rho_{\mathbb{S}}$); the solution is derived in Proposition 6.3.1 and shows that the optimal bound of Theorem 6.3.2 corresponds to the bound of Theorem 6.3.1.

Proposition 6.3.1 (Optimal Bound of Theorem 6.3.2). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution π on \mathbb{H} , for any $\delta \in (0, 1]$, for any measurable function $\varphi: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, for any algorithm $A: (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, let

$$\lambda^* = \operatorname{argmin}_{\lambda > 0} \ln \left[\frac{\lambda}{2} e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [8\varphi(h', \mathbb{S}')^2]}{2\lambda\delta^3} \right],$$

$$\begin{aligned} \text{then, we have} \quad & \overbrace{2 \ln \left[\frac{\lambda^*}{2} e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\frac{8\varphi(h', \mathbb{S}')^2}{2\lambda^*\delta^3} \right) \right]}^{\text{Theorem 6.3.2}} \\ & = \underbrace{D_2(\rho_{\mathbb{S}} \parallel \pi) + \ln \left[\frac{\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\frac{8\varphi(h', \mathbb{S}')^2}{\delta^3} \right) \right]}_{\text{Theorem 6.3.1 with } \lambda = 2}, \end{aligned}$$

$$\text{where } \lambda^* = \sqrt{\frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [8\varphi(h', S')^2]}{\delta^3 \exp(D_2(\rho_S \|\pi))}}.$$

Put into words: the optimal λ^* gives the same bound for Theorem 6.3.1 and Theorem 6.3.2.

Proof. Deferred to Appendix F.4. ■

6.4 The Disintegration in Action

So far, we have introduced theoretical results to derandomize PAC-Bayesian bounds through a disintegration approach. Indeed, the disintegration allows us to obtain a bound for a unique model sampled from ρ_S instead of having a bound on the averaged risk. This section proposes to illustrate the instantiation and usefulness of Theorem 6.3.1 on neural networks compared to the classical PAC-Bayesian bounds.

6.4.1 Specialization to Neural Network Classifiers

We aim to learn the weights of the Neural Networks (NN) leading to the lowest true risk $R_{\mathcal{D}}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[h(\mathbf{x}) \neq y]$. In other words, we consider that the hypothesis set \mathbb{H} is a set of neural networks with different weights for a given architecture. Practitioners usually proceed by epochs and obtain one “intermediate” NN after each epoch. Then, they select the “intermediate” NN associated with the lowest validation risk. We propose translating this practice into our PAC-Bayesian setting by considering one prior per epoch. Given T epochs, we hence have T priors $\mathbb{P} = \{\pi_t\}_{t=1}^T$, where $\forall t \in \{1, \dots, T\}, \pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$ is a Gaussian distribution centered at \mathbf{v}_t (the weight vector associated with the t -th “intermediate” NN) with a covariance matrix of $\sigma^2 \mathbf{I}_D$ (where \mathbf{I}_D is the $D \times D$ -dimensional identity matrix). Assuming the T priors are learned from a set $\mathbb{S}_{\text{prior}}$ such that $\mathbb{S}_{\text{prior}} \cap \mathbb{S} = \emptyset$, then Corollaries 6.4.1 and 6.4.2 will guide us to learn a posterior $\rho_S = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ from the prior $\pi \in \mathbb{P}$ minimizing the empirical risk on \mathbb{S} (we give more details on the procedure after the forthcoming corollaries). Note that considering Gaussian distributions has the advantage of simplifying the expression of the KL divergence and thus is commonly used in the PAC-Bayesian literature for neural networks (e.g., DZIUGAITE and ROY, 2017; LETARTE *et al.*, 2019; ZHOU *et al.*, 2019).²

²This has been first studied in the context of linear classifiers (e.g., AMBROLADZE *et al.*, 2006; GERMAIN *et al.*, 2009, 2020). However, in this context the symmetry of the Gaussian distribution also ease the derandomization.

6.4. The Disintegration in Action

Corollary 6.4.1 below instantiates Theorem 6.3.1 to this neural networks setting. Then, for the sake of comparison, Corollary 6.4.2 instantiates other disintegrated bounds from the literature; more precisely, Equation (6.1) corresponds to RIVASPLATA *et al.* (2020)'s bound of Theorem 2.4.1, Equation (6.2) to BLANCHARD and FLEURET (2007)'s one, and Equation (6.3) to CATONI (2007)'s one.

Corollary 6.4.1 (Instantiation of Theorem 6.3.1 for Neural Networks). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\forall \pi_t \in \mathbb{P}, \text{kl}(\mathbf{R}_S^\ell(h) \| \mathbf{R}_D^\ell(h)) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{\sigma^2} + \ln \frac{16T\sqrt{m}}{\delta^3} \right] \right) \geq 1 - \delta,$$

where $\text{kl}(a \| b) = a \ln \frac{a}{b} + (1-a) \ln \frac{1-a}{1-b}$, $\rho_{\mathcal{S}} = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$, and the hypothesis $h \sim \rho_{\mathcal{S}}$ is parameterized by $\mathbf{w} + \epsilon$.

Proof. Deferred to Appendix F.5. ■

Corollary 6.4.2 (Instantiation of Known Bounds for Neural Networks). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow \{0, 1\}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the learning sample $\mathcal{S} \sim \mathcal{D}^m$ and the hypothesis $h \sim \rho_{\mathcal{S}}$ parameterized by $\mathbf{w} + \epsilon$, we have $\forall \pi_t \in \mathbb{P}$

$$\text{kl}(\mathbf{R}_S^\ell(h) \| \mathbf{R}_D^\ell(h)) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} + \epsilon - \mathbf{v}_t\|_2^2 - \|\epsilon\|_2^2}{2\sigma^2} + \ln \frac{2T\sqrt{m}}{\delta} \right], \quad (6.1)$$

$$\forall b \in \mathbb{b}, \text{kl}_+(\mathbf{R}_S^\ell(h) \| \mathbf{R}_D^\ell(h)) \leq \frac{1}{m} \left[\frac{b+1}{b} \left[\frac{\|\mathbf{w} + \epsilon - \mathbf{v}_t\|_2^2 - \|\epsilon\|_2^2}{2\sigma^2} \right]_+ + \ln \frac{(b+1)T \text{card}(\mathbb{b})}{\delta} \right], \quad (6.2)$$

$$\forall c \in \mathbb{c}, \mathbf{R}_D^\ell(h) \leq \frac{1 - \exp\left(-c\mathbf{R}_S^\ell(h) - \frac{1}{m} \left[\frac{\|\mathbf{w} + \epsilon - \mathbf{v}_t\|_2^2 - \|\epsilon\|_2^2}{2\sigma^2} + \ln \frac{T \text{card}(\mathbb{c})}{\delta} \right]\right)}{1 - e^{-c}}, \quad (6.3)$$

with $[x]_+ = \max(x, 0)$, and $\text{kl}_+(\mathbf{R}_S^\ell(h) \| \mathbf{R}_D^\ell(h)) = \text{kl}(\mathbf{R}_S^\ell(h) \| \mathbf{R}_D^\ell(h))$ if $\mathbf{R}_S^\ell(h) < \mathbf{R}_D^\ell(h)$ and 0 otherwise. Moreover, $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$ is a Gaussian noise such that $\mathbf{w} + \epsilon$ are

the weights of $h \sim \rho_S$ with $\rho_S = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$, and $\mathfrak{c}, \mathfrak{b}$ are two sets of hyperparameters fixed a priori.

Proof. Deferred to Appendix F.6. ■

Since we aim to minimize the true risk $R_{\mathcal{D}}(h)$, *i.e.*, we consider in practice the 01-loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$. As the parameter λ of the Theorem 6.3.2, $c \in \mathfrak{c}$ is a hyperparameter that controls a trade-off between the empirical risk $R_S^\ell(h)$ and $\frac{1}{m} \left[\frac{\|\mathbf{w} + \epsilon - \mathbf{v}_t\|_2^2 - \|\epsilon\|_2^2}{2\sigma^2} + \ln \frac{T \text{card}(\mathfrak{c})}{\delta} \right]$. Besides, the parameter $b \in \mathfrak{b}$ controls the tightness of the bound. These parameters can generally be tuned to minimize the bound of Equation (6.2) and Equation (6.3); however, there is no closed-form solution for the expression of the minimum of these equations. Consequently, our experimental protocol requires minimizing the bounds by gradient descent for each $b \in \mathfrak{b}$ or $c \in \mathfrak{c}$ to learn the distribution ρ_S leading to the lowest bound value. To obtain a tight bound, the divergence between one prior $\pi_t \in \mathbb{P}$ and ρ_S must be low, *i.e.*, $\|\mathbf{w} - \mathbf{v}_t\|_2^2$ (or $\|\mathbf{w} + \epsilon - \mathbf{v}_t\|_2^2 - \|\epsilon\|_2^2$) has to be small. One solution is to split the learning sample into 2 non-overlapping subsets S_{prior} and S , where S_{prior} is used to learn the prior, while S is used both to learn the posterior and compute the bound. Hence, if we “pre-learn” a good enough prior $\pi_t \in \mathbb{P}$ from S_{prior} , then we can expect to have a low $\|\mathbf{w} - \mathbf{v}_t\|_2$.

Training Method

The original training set is split in two distinct subsets: S_{prior} and S (respectively of size m_{prior} and m , that can be different).

The training has two phases.

- 1) The prior distribution π is “pre-learned” with S_{prior} and selected by early stopping, with S as validation set, using the algorithm A_{prior} (an arbitrary learning algorithm).
- 2) Given S and π , we learn the posterior ρ_S with the algorithm A (defined *a priori*).

At first sight, the selection of the prior weights with S by early stopping may appear to be “cheating”. However, this procedure can be seen as: (i) first constructing \mathbb{P} from the T “intermediate” NNs learned after each epoch from S_{prior} , then (ii) optimizing the bound with the prior that leads to the best risk on S . This gives a statistically valid result: since Corollary 6.4.1 is valid for every $\pi_t \in \mathbb{P}$, we can select the one we want, in particular the one minimizing $R_S^\ell(h)$ for a sampled $h \sim \pi_t$. This heuristic makes sense: it allows us to detect if a prior is concentrated around hypotheses that potentially overfit the learning sample S_{prior} . Usually, practitioners consider this “best” prior as the final NN. In our case, the advantage is that we refine this “best” prior with S to learn the posterior ρ_S . Note that PÉREZ-ORTIZ *et al.* (2021) have already introduced tight generalization bounds with data-dependent priors for—non-derandomized—stochastic NNs. Nevertheless, our training method to learn the prior differs greatly since (i) we

learn T NNs (*i.e.*, T priors) instead of only one, (*ii*) we fix the variance of the Gaussian in the posterior ρ_S . To the best of our knowledge, our training method for the prior is new.

6.4.2 A Note About Stochastic Neural Networks

Due to its stochastic nature, PAC-Bayesian theory has been explored to study stochastic NNs (*e.g.*, DZIUGAITE and ROY (2017, 2018), ZHOU *et al.* (2019), and PÉREZ-ORTIZ *et al.* (2021)). In Corollary 6.4.3 below, we instantiate the bound of Theorem 2.3.1 for stochastic NNs to empirically compare the stochastic and the deterministic NNs associated with the same prior and posterior distributions. We recall that, in this chapter, a deterministic NN is a *single* h sampled from the posterior distribution $\rho_S = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ output by the algorithm A . Hence, for each example, the prediction is performed by the same deterministic NN: the one parameterized by the weights $\mathbf{w} + \epsilon \in \mathbb{R}^D$. Conversely, the stochastic NN associated with a posterior distribution $\rho = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ predicts the label of a given example by (*i*) first sampling h according to ρ , and (*ii*) then returning the label predicted by h . Thus, the risk of the stochastic NN is the expected risk value $\mathbb{E}_{h \sim \rho} R_D^\ell(h)$, where the expectation is taken over *all* h sampled from ρ . We compute the empirical risk of the stochastic NN from a Monte Carlo approximation: (*i*) we sample K weight vectors, and (*ii*) we average the risk over the K associated NNs; we denote by ρ^K the distribution of such K -sample. In this context, we obtain the following PAC-Bayesian bound.

Corollary 6.4.3 (PAC-Bayesian Bound for Stochastic Neural Networks). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$, for any loss $\ell: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow \{0, 1\}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over $\mathbb{S} \sim \mathcal{D}^m$ and $\{h_1, \dots, h_K\} \sim \rho^K$, we have simultaneously $\forall \pi_t \in \mathbb{P}$,

$$\text{kl} \left(\mathbb{E}_{h \sim \rho} R_S^\ell(h) \parallel \mathbb{E}_{h \sim \rho} R_D^\ell(h) \right) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{2\sigma^2} + \ln \frac{4T\sqrt{m}}{\delta} \right], \quad (6.4)$$

$$\text{and} \quad \text{kl} \left(\frac{1}{K} \sum_{i=1}^K R_S(h_i) \parallel \mathbb{E}_{h \sim \rho} R_S^\ell(h) \right) \leq \frac{1}{n} \ln \frac{4}{\delta}, \quad (6.5)$$

where $\rho = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ and the hypothesis h sampled from ρ is parameterized by $\mathbf{w} + \epsilon$ with $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$.

Proof. Deferred to Appendix F.7. ■

This result shows two features that allow considering it as an adapted baseline for a fair comparison between disintegrated and classical PAC-Bayesian bounds, thus between deterministic and stochastic NNs. Firstly, it involves the same terms as Corollary 6.4.1.

Secondly, it is close to the bound of PÉREZ-ORTIZ *et al.* (2021, Sec. 6.2), since (i) we adapt the KL divergence to our setting (*i.e.*, $\text{KL}(\rho\|\pi)=\frac{1}{2\sigma^2}\|\mathbf{w}-\mathbf{v}_t\|_2^2$), (ii) the bound holds for T priors thanks to a union bound argument.

6.5 Experiments with Neural Networks

In this section, we do not seek state-of-the-art performance; in fact, we have a three-fold objective: (i) we check if 50%/50% is a good choice for splitting the original train set into $(\mathbb{S}_{\text{prior}}, \mathbb{S})$ (which is the most common split in the PAC-Bayesian literature (GERMAIN *et al.*, 2009; PÉREZ-ORTIZ *et al.*, 2021)); (ii) we highlight that our disintegrated bound associated with the deterministic NN is tighter than the randomized bound associated with the stochastic NN (Corollary 6.4.3); (iii) we show that our disintegrated bound (Corollary 6.4.1) is tighter and more stable than the ones based on RIVASPLATA *et al.* (2020), BLANCHARD and FLEURET (2007) and CATONI (2007) (Corollary 6.4.2).

6.5.1 Training Method

We follow our Training Method (Section 6.4.1) in which we integrate the direct minimization of all the bounds. We refer as ours the training method based on the minimization of our bound in Corollary 6.4.1, as `rivasplata` the one based on Equation (6.1), as `blanchard` the one based on Equation (6.2), and as `catoni` the one based on Equation (6.3). `stochastic` denotes the PAC-Bayesian bound with the prior and posterior distributions obtained from ours. To optimize the bound with gradient descent, we replace the non-differentiable 0-1 loss $\mathbb{I}[h(\mathbf{x}) \neq y]$ with a surrogate: the bounded cross-entropy loss (DZIUGAITE and ROY, 2018). We made this replacement since cross-entropy minimization works well in practice for neural networks (GOODFELLOW *et al.*, 2016) and because this loss is bounded between 0 and 1, which is required for the `kl()` function. The cross-entropy is defined in a multi-class setting with $y \in \mathbb{Y}$ by $\ell(h, (\mathbf{x}, y)) = -\frac{1}{Z} \ln[e^{-Z} + (1 - 2e^{-Z})h[y]] \in [0, 1]$ where $h[y]$ is the y -th output of the NN; we set $Z=4$, the default parameter of DZIUGAITE and ROY (2018). That being said, to learn a good enough prior $\pi \in \mathbb{P}$ and the posterior $\rho_{\mathbb{S}}$, we run our Training Method with two stochastic gradient descent-based algorithms A_{prior} and A . Note that the randomness in the stochastic gradient descent algorithm is fixed to have deterministic algorithms. In phase **1**) algorithm A_{prior} learns from $\mathbb{S}_{\text{prior}}$ the T priors $\pi_1, \dots, \pi_T \in \mathbb{P}$ (*i.e.*, during T epochs) by minimizing the bounded cross-entropy loss. In other words, at the end of the epoch t , the weights \mathbf{w}_t of the classifier are used to define the prior $\pi_t = \mathcal{N}(\mathbf{w}_t, \sigma^2 \mathbf{I}_D)$. Then, the best prior $\pi \in \mathbb{P}$ is selected by early stopping on \mathbb{S} . In phase **2**), given \mathbb{S} and π , algorithm A integrates the direct optimization of the bounds with the bounded cross-entropy loss.

6.5.2 Optimization³ Procedure in Algorithms A and A_{prior}

Let ω be the mean vector of a Gaussian distribution used as NN weights that we are optimizing. In algorithms A and A_{prior} , we use the Adam optimizer (KINGMA and BA, 2015), and we sample a noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$ at each iteration of the optimizer. Then, we forward the examples of the mini-batch in the NN parameterized by the weights $\omega + \epsilon$, and we update ω according to the bounded cross-entropy loss. Note that during phase **1**), at the end of each epoch t , $\pi_t = \mathcal{N}(\omega, \sigma^2 \mathbf{I}_D) = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$ and finally at the end of phase **2**) we have $\rho_S = \mathcal{N}(\omega, \sigma^2 \mathbf{I}_D) = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$.

6.5.3 Experimental Setting

6.5.3.1 Datasets

We perform our experimental study on three datasets: MNIST (LECUN *et al.*, 1998), Fashion-MNIST (XIAO *et al.*, 2017), and CIFAR-10 (KRIZHEVSKY, 2009). We divide each original train set into two independent subsets $\mathbb{S}_{\text{prior}}$ of size m_{prior} and \mathbb{S} of size m with varying split ratios defined as $\frac{m_{\text{prior}}}{m+m_{\text{prior}}} \in \{0, .1, .2, .3, .4, .5, .6, .7, .8, .9\}$. The test sets denoted by \mathbb{T} remain the original ones.

6.5.3.2 Models

For the (Fashion-)MNIST datasets, we train a variant of the All Convolutional Network (SPRINGENBERG *et al.*, 2015). The model is a 3-hidden layers convolutional network with 96 channels. We use 5×5 convolutions with a padding of size 1 and a stride of size 1 everywhere except on the second convolution where we use a stride of size 2. We adopt the Leaky ReLU activation functions after each convolution. Lastly, we use a global average pooling of size 8×8 to obtain the desired output size. Furthermore, the weights are initialized with Xavier Normal initializer (GLOROT and BENGIO, 2010) while each bias of size l is initialized uniformly between $-1/\sqrt{l}$ and $1/\sqrt{l}$. For the CIFAR-10 dataset, we train a ResNet-20 network, *i.e.*, a ResNet network from HE *et al.* (2016) with 20 layers. The weights are initialized with Kaiming Normal initializer (HE *et al.*, 2015) and each bias of size l is initialized uniformly between $-1/\sqrt{l}$ and $1/\sqrt{l}$.

6.5.3.3 Optimization

For the (Fashion-)MNIST datasets, we learn the parameters of our prior distributions π_1, \dots, π_T by using Adam optimizer for $T = 10$ epochs with a learning rate of $\eta = 10^{-3}$ and a batch size of 32 (the other parameters of Adam are left by default). Moreover, the posterior distribution's parameters are learned for one epoch with the same batch size and optimizer (except that the learning rate is either $\eta = 10^{-4}$ or $\eta = 10^{-6}$). For

³The details of the optimization and the evaluation of the bounds are described in Appendix.

the CIFAR-10 dataset, the parameters of the priors π_1, \dots, π_T are learned for $T = 100$ epochs and the posterior distribution $\rho_{\mathcal{S}}$ for 10 epochs with a batch size of 32 by using Adam optimizer as well. Additionally, the learning rate to learn the prior for CIFAR-10 is $\eta = 10^{-2}$.

6.5.3.4 Bounds

We report the bounds' values with the 0-1 loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$. Moreover, for blanchard's bounds, we define the set of hyperparameters \mathbb{b} in the following way: $\mathbb{b} = \{b \in \mathbb{N} \mid b = \sqrt{x}, (x+1) \leq 2\sqrt{m}\}$, i.e., such that blanchard's bounds can be tighter than rivasplata's ones. We fixed the set of hyperparameters for catoni as $\mathbb{c} = \{10^k \mid k \in \{-3, -2, \dots, +3\}\}$. We consider different values for the variance $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$ associated with the disintegrated KL divergence $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$ equals to $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} = \frac{1}{2\sigma^2} (\|\mathbf{w} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2)$, the "normal" RÉNYI divergence $D_2(\rho \parallel \pi) = \frac{1}{\sigma^2} \|\mathbf{w} - \mathbf{v}_t\|_2^2$ and the KL divergence $\text{KL}(\rho \parallel \pi) = \frac{1}{2\sigma^2} \|\mathbf{w} - \mathbf{v}_t\|_2^2$. For all the figures, the values are averaged over 400 deterministic NNs sampled from $\rho_{\mathcal{S}}$ (the standard deviation is small and presented in the Appendix). We additionally report as stochastic (Corollary 6.4.3) the randomized bound value and KL divergence $\text{KL}(\rho \parallel \pi)$ associated with the model learned by ours, meaning that $K=400$ and that the test risk reported for ours also corresponds to the risk of the stochastic NN approximated with these 400 NNs.

6.5.4 Results

6.5.4.1 Analysis of the Influence of the Split Ratio Between $\mathcal{S}_{\text{prior}}$ and \mathcal{S}

Figures 6.1 and 6.2 study the evolution of the bound values after optimizing the bounds with our Training Method for different parameters. Specifically, the split ratio of the original train set varies from 0.1 to 0.9 (0.1 means that $m_{\text{prior}} = 0.1(m + m_{\text{prior}})$), for four variances values σ^2 and the two learning rates ($\eta = 10^{-6}$ and $\eta = 10^{-4}$). For the sake of readability, we present detailed results when the split ratio is 0 in Table F.1. We first remark that the behavior is different for the two learning rates. On the one hand, for $\eta = 10^{-6}$, the mean bound values are close to each other, which is not surprising since the disintegrated KL divergences and the RÉNYI divergences are close to zero (see Table F.2 to Table F.10). Moreover, for MNIST and Fashion-MNIST, there is a trade-off between learning a good prior with $\mathcal{S}_{\text{prior}}$ and minimizing a generalization bound with \mathcal{S} . In this case, the split ratio 0.5 appears to be a good choice to obtain a tight (disintegrated) PAC-Bayesian bound. This ratio is widely used in the PAC-Bayesian literature (see, e.g., in the context of linear classifiers (GERMAIN *et al.*, 2009), majority votes (ZANTEDESCHI *et al.*, 2021), and neural networks (LETARTE *et al.*, 2019; PÉREZ-ORTIZ *et al.*, 2021)). On the other hand, when $\eta = 10^{-4}$, the

6.5. Experiments with Neural Networks

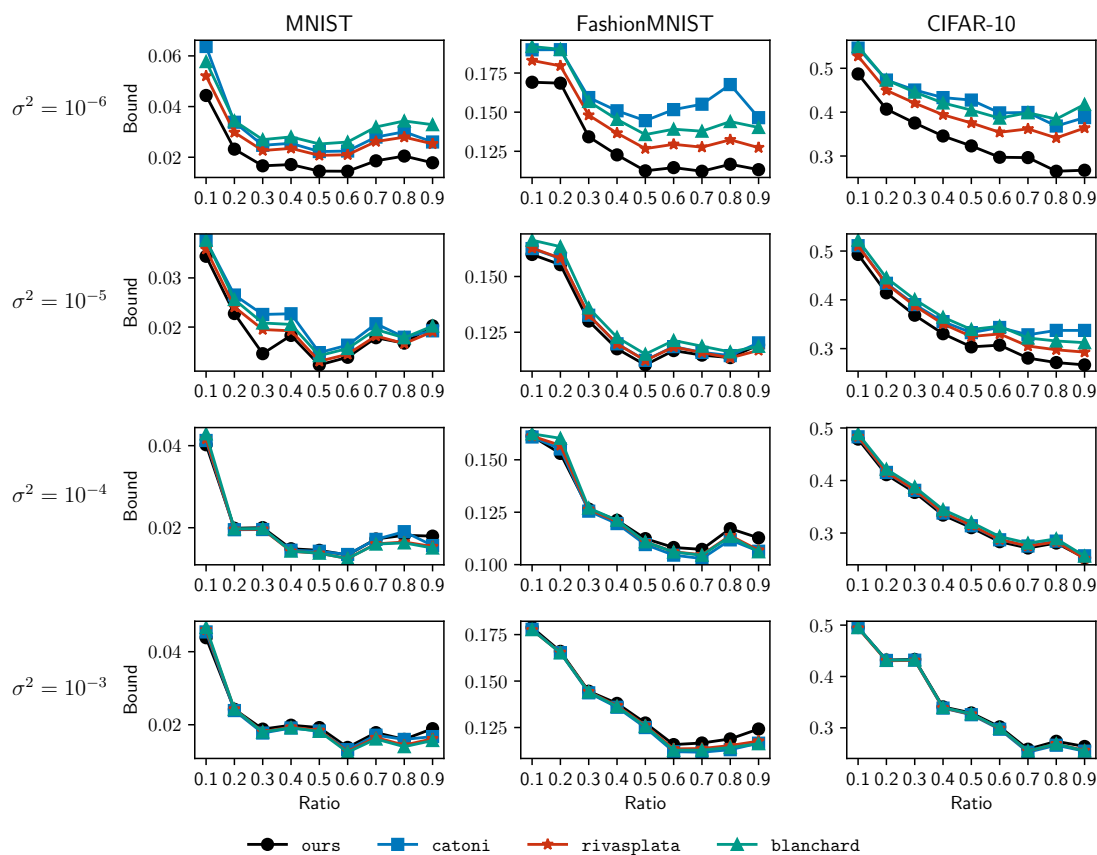


Figure 6.1. Evolution of the bound values in terms of the split ratio. The x-axis represents the different split ratios, and the y-axis represents the bound values obtained after their optimization using our Training Method. Each row corresponds to a given variance σ^2 , and each column corresponds to a dataset (MNIST, Fashion-MNIST, or CIFAR-10). In this figure, we consider a learning rate of $\eta = 10^{-6}$.

mean bound values tend to increase when the split ratio increases as well for the bounds introduced in the literature (*i.e.*, for blanchard, catoni, and rivasplata), while the mean bound values of our bound remain low. Indeed, m decreases as long as the split ratio increases, which has the effect of increasing the bound value drastically when the disintegrated KL divergence is high. We further explain why the disintegrated KL divergence can become high for the disintegrated bounds of the literature. To do so, we now restrict our study to a split ratio of 0.5 in order to (i) compare the tightness of the bounds, (ii) understand why the disintegrated bounds of the literature diverge.

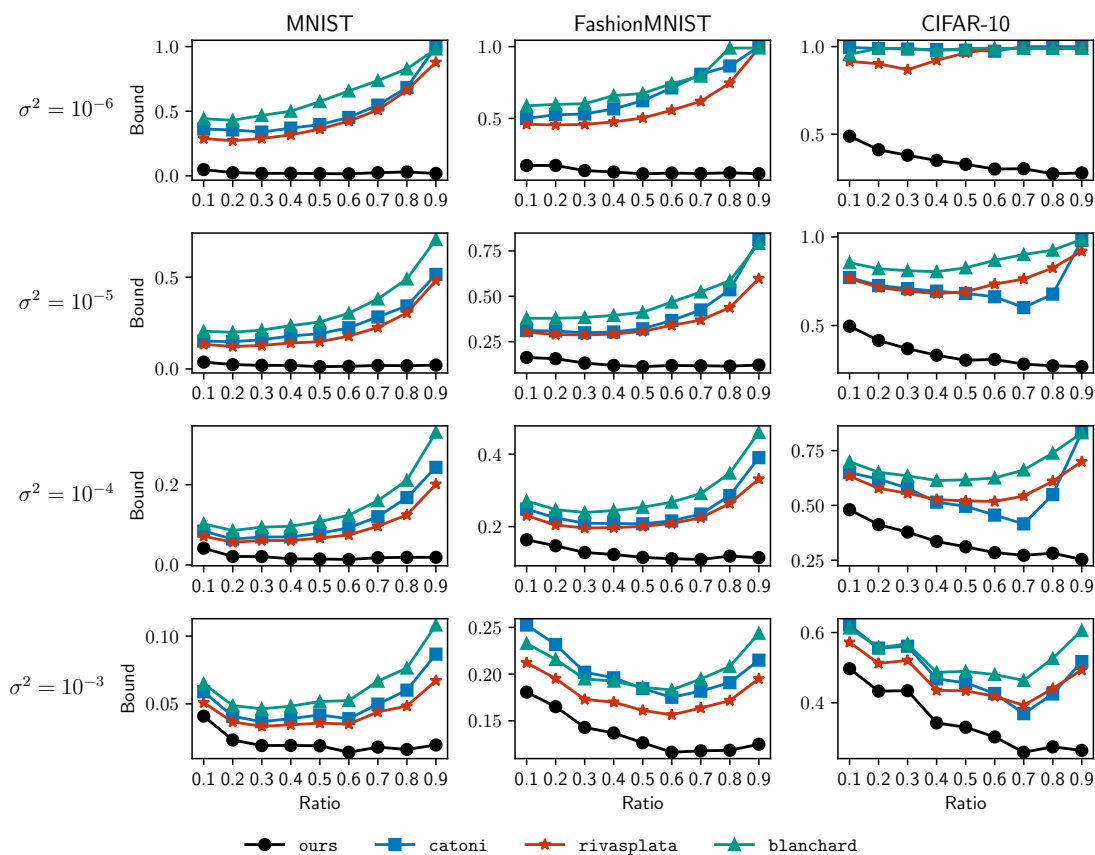


Figure 6.2. Evolution of the bound values in terms of the split ratio. The x-axis represents the different split ratios, and the y-axis represents the bound values obtained after their optimization using our Training Method. Each row corresponds to a given variance σ^2 , and each column corresponds to a dataset (MNIST, Fashion-MNIST, or CIFAR-10). In this figure, we consider a learning rate of $\eta = 10^{-4}$.

6.5.4.2 Comparison Between Disintegrated and “Classic” Bounds

We first compare the “classic” PAC-Bayesian bound (Corollary 6.4.3) and our disintegrated PAC-Bayesian bound (Corollary 6.4.1). To do so, we fix the variance $\sigma^2=10^{-3}$ (along with the split ratio equals 0.5). We report in Figure 6.3, the mean bound values associated with ours (*i.e.*, the Training Method that minimizes our bound) and stochastic (we recall that stochastic is the PAC-Bayesian bound of Corollary 6.4.3 on the model learned by ours). Actually, ours leads to more precise bounds than the randomized stochastic even if the two empirical risks are the same and the KL divergence is smaller than the RÉNYI one. This imprecision is due to the non-avoidable sampling according to ρ done in the randomized PAC-Bayesian bound

6.5. Experiments with Neural Networks

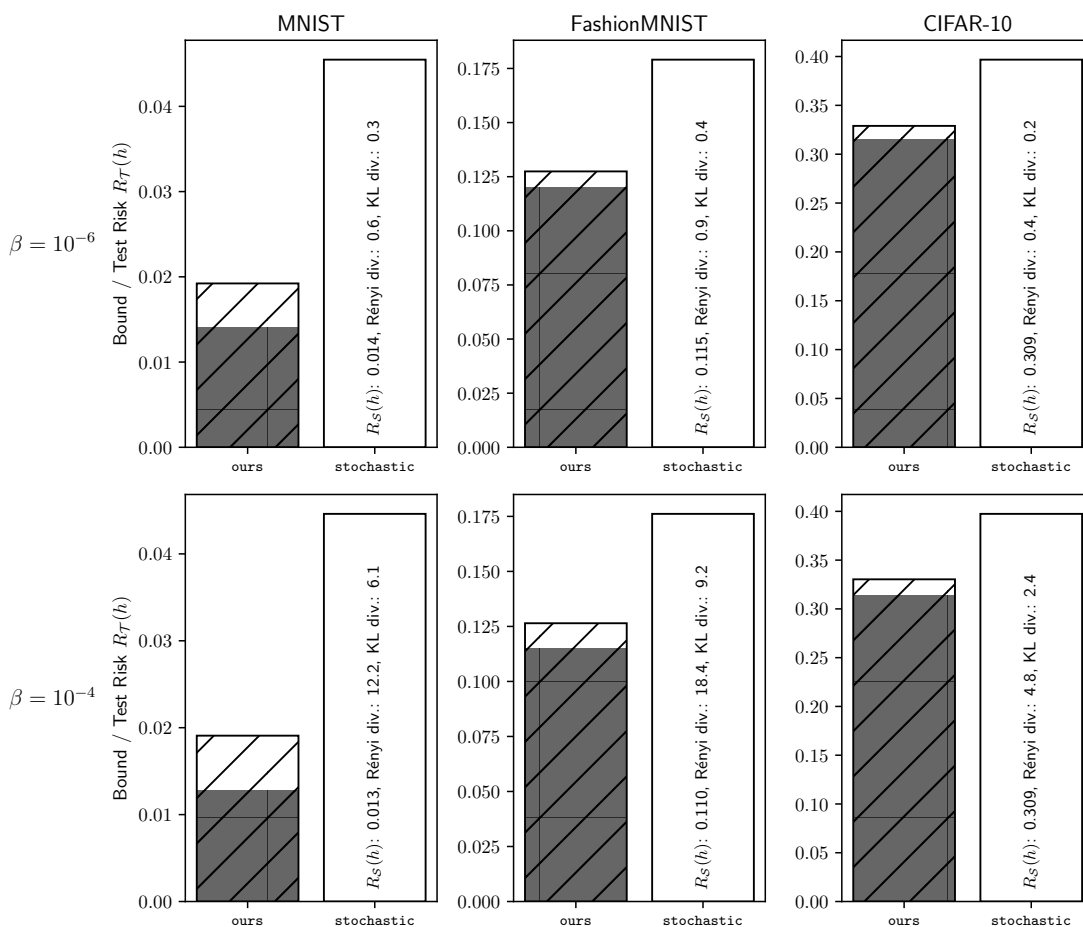


Figure 6.3. The figure illustrates the values of the PAC-Bayes bound (Corollary 6.4.3) and the values of the disintegrated bound (Corollary 6.4.1) where the learning rate is $\eta = 10^{-4}$ or $\eta = 10^{-6}$ and the split ratio is 0.5. The y-axis shows the values of the bounds (the hatched bar for ours (Corollary 6.4.1) and the white bar for stochastic (Corollary 6.4.3)) and the test risks $R_T(h)$ (grey shaded bar). We also report the values of the empirical risk $R_S^l(h)$, the RÉNYI divergence (associated with ours' bound), and the KL divergence (associated with stochastic's bound).

of Corollary 6.4.3 (the higher K , the tighter the bound). Thus, using a disintegrated PAC-Bayesian bound avoids sampling many NNs to obtain a low risk. This confirms that our framework makes sense for practical purposes and has a great advantage in terms of time complexity when computing the bounds.

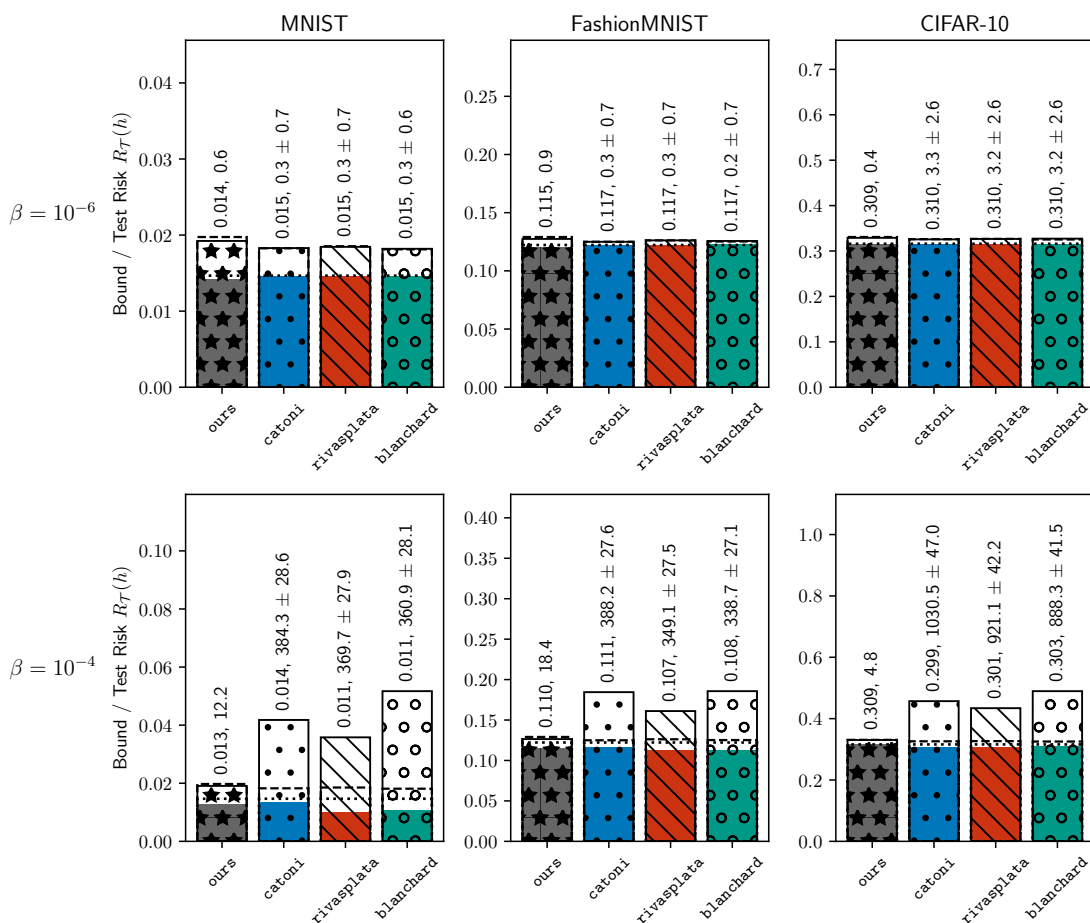


Figure 6.4. This figure shows the value of the disintegrated bounds (the colored bars) and the test risks (the hatched bars) for Corollary 6.4.1 (“ours”) and Corollary 6.4.2 (“catoni”, “rivasplata” and “blanchard”) in two different settings, i.e., with a learning rate of $\eta = 10^{-6}$ and $\eta = 10^{-4}$ and with split ratio of 0.5. We also plot the value of the bounds (the dashed lines) and the test risks (the dotted lines) before executing Step 2) of our Training Method. The y-axis shows the values of the bounds and the test risks $R_{\mathcal{T}}(h)$. The empirical risk $R_S(h)$ is presented above each bar. Moreover, the second value represents the mean value of the divergence (the standard deviations are also given for the disintegrated bounds of the literature).

6.5.4.3 Analysis of the Tightness of the Disintegrated Bounds

We now compare the tightness of the different disintegrated PAC-Bayesian bounds (i.e., our bound and the ones in the literature). We study, as before, the case where the split ratio is 0.5 and the variance $\sigma^2 = 10^{-3}$. We report in Figure 6.4

for ours, `rivasplata`, `blanchard` and `catoni`, the mean bounds values; the mean test risk $R_{\mathcal{T}}(h)$ before (*i.e.*, with the prior π) and after applying Step 2) (*i.e.*, with the posterior $\rho_{\mathcal{S}}$). Moreover, we report above the bars the mean train risks $R_{\mathcal{S}}^{\ell}(h)$ and the mean/standard deviation divergence values obtained after Step 2), *i.e.*, the RÉNYI divergence $D_2(\rho_{\mathcal{S}}\|\pi)=\frac{1}{\sigma^2}\|\mathbf{w}-\mathbf{v}_t\|_2^2$ for ours and the disintegrated KL divergence $\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}=\frac{1}{2\sigma^2} [\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}_t\|_2^2-\|\boldsymbol{\epsilon}\|_2^2]$ for the others. First of all, we can remark that we observe two different behaviors for $\eta = 10^{-4}$ and $\eta = 10^{-6}$. For $\eta = 10^{-6}$, the bound values are close to each other, as well as the empirical risks and the divergences (which are close to 0). In Figure 6.4, we observe that the bound values and the test risks are close to the one associated with the prior distribution because the divergence is close to 0. This is probably due to the fact that the learning rate is too small, implying that the bounds are not optimized. With a higher learning rate of $\eta = 10^{-4}$, we observe that our bound remains tight while the disintegrated bounds of the literature are looser. Hopefully, our bound is improved after performing Step 2) of our Training Method. However, for the bounds of the literature, the value of the disintegrated KL divergence is large, making the bounds looser after executing Step 2). We now investigate the reasons for the divergence of the bounds by looking at the influence of the variance σ^2 .

6.5.4.4 Analysis of the Influence of the Variance σ^2

Given a split ratio of 0.5 and $\eta \in \{10^{-6}, 10^{-4}\}$, we report in Figure 6.5 the evolution of the bound values associated with ours, `rivasplata`, `blanchard`, and `catoni` when the variance varies from 10^{-6} to 10^{-3} . First of all, an important point is that ours behaves differently than `rivasplata`, `blanchard`, and `catoni`. Indeed, for both learning rates, when σ^2 decreases, the value of our bound remains low, while the others increase drastically due to the explosion of the disintegrated KL divergence term (see Table F.6 in Appendix for more details). Concretely, the disintegrated KL divergence in Corollary 6.4.2 involves the noise $\boldsymbol{\epsilon}$ through $\frac{1}{2\sigma^2}\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}_t\|_2^2-\|\boldsymbol{\epsilon}\|_2^2$ compared to our divergence which is $\frac{1}{\sigma^2}\|\mathbf{w}-\mathbf{v}_t\|_2^2$ (without noise). Then, the sampled noise during the optimization procedure $\boldsymbol{\epsilon}$ influences the disintegrated KL divergence, making it prone to high variations during training (depending thus on σ^2). To illustrate the difference during the optimization, we focus on the objective function (detailed in Appendix) of Corollary 6.4.1 and Corollary 6.4.2 (Equation (6.1)). Roughly speaking, the divergence in Corollary 6.4.1 does not depend on the sampled hypothesis h (with weights $\boldsymbol{\omega} + \boldsymbol{\epsilon}$), while the divergence of Equation (6.1) does. In consequence, the derivatives are less dependent on h for Corollary 6 than for Equation (6.1). To be convinced of this, we propose to study the gradient with respect to the current mean vector $\boldsymbol{\omega}$. On the one hand, the gradient $\frac{\partial R_{\mathcal{S}}^{\ell}(h)}{\partial \boldsymbol{\omega}}$ of the risk *w.r.t.* $\boldsymbol{\omega}$ is the same for both bounds (with the loss of DZIUGAITE and ROY (2018)); hence, the phenomenon cannot come from this derivative. On the other hand, the gradients of the divergence in Equation (6.1) and

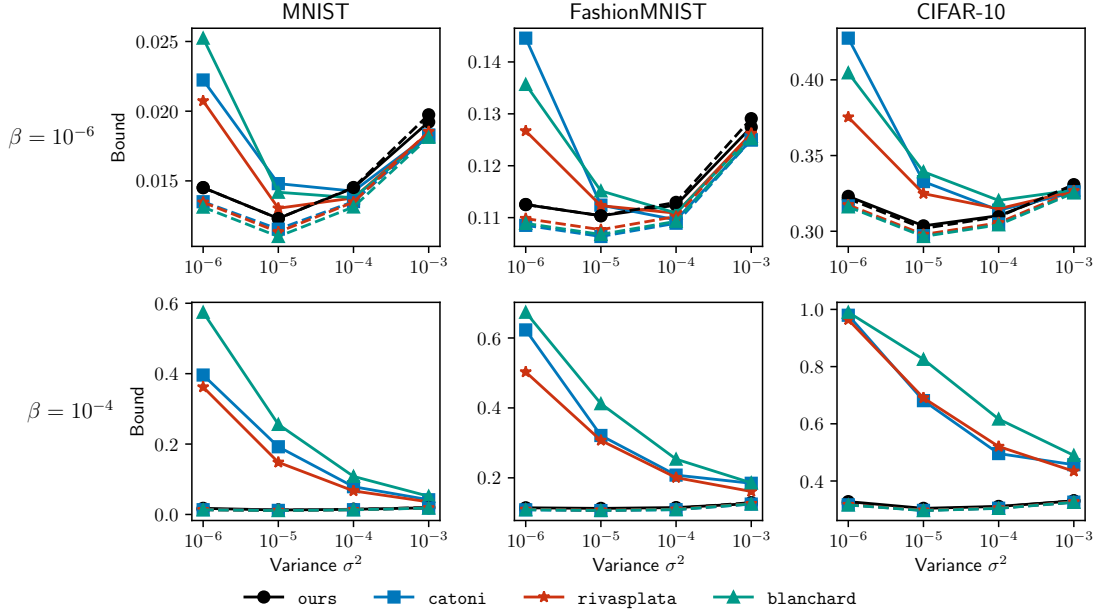


Figure 6.5. We plot the evolution of the mean bound values (the plain lines) in terms of the variance σ^2 after optimizing the bounds with our Training Method. Moreover, we plot the mean bound values (the dashed lines) obtained before executing Step 2) of our Training Method. The variance is represented on the x-axis, while the bound values are represented on the y-axis. Furthermore, each row corresponds to a given learning rate ($\eta = 10^{-6}$ or $\eta = 10^{-4}$), and each column corresponds to a dataset (either MNIST, FashionMNIST, or CIFAR-10). The split ratio considered is 0.5.

Corollary 6.4.1 are respectively

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\omega}} \left[\frac{1}{m} \left(\frac{\|\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} \right) \right] &= \frac{\partial}{\partial \boldsymbol{\omega}} \left[\frac{1}{m2\sigma^2} \|\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 \right] \\ &= \frac{1}{m\sigma^2} (\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t) = \diamond, \\ \text{and} \quad \frac{\partial}{\partial \boldsymbol{\omega}} \left[\frac{1}{m} \left(\frac{\|\boldsymbol{\omega} - \mathbf{v}_t\|_2^2}{\sigma^2} \right) \right] &= \frac{\partial}{\partial \boldsymbol{\omega}} \left[\frac{1}{m\sigma^2} \|\boldsymbol{\omega} - \mathbf{v}_t\|_2^2 \right] \\ &= \frac{2}{m\sigma^2} (\boldsymbol{\omega} - \mathbf{v}_t) = \heartsuit. \end{aligned}$$

From the two derivatives, we deduce that $\diamond = \frac{1}{2}\heartsuit + \frac{1}{m\sigma^2}\boldsymbol{\epsilon}$. Hence, each gradient step involves a noise in the gradient of the disintegrated KL divergence $\frac{1}{m\sigma^2}\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \frac{1}{m}\mathbf{I}_D)$, which is high for a small m . This randomness causes the disintegrated KL divergence $\frac{1}{2\sigma^2}\|\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2$ to be larger when σ^2 decreases since (i) the divergence is divided by $2m\sigma^2$ and (ii) the deviation between $\boldsymbol{\omega}$ and \mathbf{v}_t increases. In conclusion, this

makes the objective function (*i.e.*, the bound) subjects to high variations during the optimization, implying higher final bound values. Thus, the RÉNYI divergence has a valuable asset over the disintegrated KL divergence since it does not depend on the sampled noise ϵ .

6.5.4.5 Take Home Message from the Experiments

To summarize, our experiments show that our disintegrated bound is, in practice, tighter than the ones in the literature. This tightness allows us to precisely bound the true risk $R_{\mathcal{D}}(h)$ (or the test risk $R_{\mathcal{T}}(h)$); thus, the model selection from the disintegrated bound is effective. Moreover, we show that our bound is more easily optimizable than the others. This is mainly due to the disintegrated KL divergence, which depends on the sampled hypothesis h with weights $\omega + \epsilon$. Indeed, the gradients of the disintegrated KL divergence with respect to ω include the noise ϵ , making the gradient inaccurate (especially with a “high” learning rate and small variance σ^2).

6.6 Perspectives for the Majority Vote

Before concluding, we discuss some perspectives for the stochastic majority vote introduced in Chapter 5. Recall that the *stochastic* majority vote has its weights ρ sampled from the Dirichlet distribution P defined as $P(\rho) \triangleq \frac{1}{Z(\alpha)} \prod_{j=1}^{\text{card}(\mathbb{H})} \rho(h_j)^{\alpha_j - 1}$ (also called hyper-posterior). One drawback of the PAC-Bayesian approach for the *stochastic* majority vote in Chapter 5 is that we are not considering only one majority vote (with weights sampled from the hyper-posterior P). Moreover, based on the margin theory (initiated by SCHAPIRE *et al.*, 1998) and our work in Chapter 5, BIGGS *et al.* (2020) derived PAC-Bayesian bound for the expected majority vote. This latter work illustrates that a special care is unfortunately needed to obtain a bound for a unique majority vote through the PAC-Bayesian theory. Hopefully, thanks to the disintegrated PAC-Bayesian framework, we can bound the true risk of a single majority vote classifier. As an illustration we derive two bounds for such a classifier based on RIVASPLATA *et al.* (2020)’s bound (Theorem 2.4.1) and based on Theorem 6.3.1.

Corollary 6.6.1 (Instantiation of Theorem 2.4.1 to Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any finite set of voters \mathbb{H} , for any hyper-prior distribution $\Pi = \text{Dir}(\beta)$ on \mathbb{H} with $\beta \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, for any algorithm A that outputs a hyper-posterior given a learning sample and a hyper-prior, with probability at least $1 - \delta$ over the learning sample $\mathbb{S} \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim P_{\mathbb{S}} = \text{Dir}(\alpha)$ with $\alpha \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$

we have

$$\text{kl}(\mathbb{R}_S^\ell(\text{MV}_\rho) \parallel \mathbb{R}_D^\ell(\text{MV}_\rho)) \leq \frac{1}{m} \left[\ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \ln(\rho(h_j)) + \ln \frac{2\sqrt{m}}{\delta} \right],$$

where $\mathbb{P}_S \triangleq A(\mathbb{S}, \Pi)$ is output by the deterministic algorithm A .

Proof. Deferred to Appendix F.13 ■

Corollary 6.6.2 (Instantiation of Theorem 6.3.1 to Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any finite set of voters \mathbb{H} , for any hyper-prior distribution $\Pi = \text{Dir}(\boldsymbol{\beta})$ on \mathbb{H} with $\boldsymbol{\beta} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm A that outputs a hyper-posterior given a learning sample and a hyper-prior, with probability at least $1 - \delta$ over the learning sample $\mathbb{S} \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim \mathbb{P}_S = \text{Dir}(\boldsymbol{\alpha})$ with $\boldsymbol{\alpha} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$ we have

$$\begin{aligned} \text{kl}(\mathbb{R}_S^\ell(\text{MV}_\rho) \parallel \mathbb{R}_D^\ell(\text{MV}_\rho)) \leq \frac{1}{m} \left[\frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} + \ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} \right. \\ \left. + \frac{1}{\lambda-1} \ln \frac{Z(\lambda\boldsymbol{\alpha} + (1-\lambda)\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \ln(2\sqrt{m}) \right], \end{aligned}$$

where $\mathbb{P}_S \triangleq A(\mathbb{S}, \Pi)$ is output by the deterministic algorithm A .

Proof. Deferred to Appendix F.14 ■

These two bounds offer great perspectives to upper-bound the true risk of the majority vote. First, as we can remark, the two theorems holds for all (bounded) losses $\ell(\cdot)$ including the 01-loss. Hence, the true risk $\mathbb{R}_D(\text{MV}_\rho)$ can be upper-bounded directly without using the $\frac{1}{2}$ -margin of LAVIOLETTE *et al.* (2017) as in Chapter 5. While our bounds require to sample a single majority vote from \mathbb{P} , they might be tighter since it does not rely on margin bound (as BIGGS *et al.*, 2020) and directly deal with the 01-loss.

6.7 Summary and Conclusion

We provide a new and general disintegrated PAC-Bayesian bound (Theorem 6.3.1) in the family of disintegrated PAC-Bayesian bounds (Definition 2.4.1), *i.e.*, when the de-

6.7. Summary and Conclusion

randomization step consists in (i) learning a posterior distribution $\rho_{\mathbb{S}}$ on the classifiers set (given an algorithm, a learning sample \mathbb{S} and a prior distribution π) and (ii) sampling a hypothesis h from this posterior $\rho_{\mathbb{S}}$. While our bound can be looser than the ones of BLANCHARD and FLEURET (2007), CATONI (2007), and RIVASPLATA *et al.* (2020), it provides nice opportunities for learning deterministic classifiers. Indeed, our bound can be used not only to study the theoretical guarantees of deterministic classifiers but also to derive self-bounding algorithms (based on the bound optimization) that are more stable and efficient than the ones we obtain from the bounds of the literature. Concretely, the bounds of BLANCHARD and FLEURET (2007), CATONI (2007), and RIVASPLATA *et al.* (2020) depend on two terms related to the classifier drawn: the risk and the “disintegrated KL divergence”, while in our bound the (RÉNYI) divergence term depends on the hypothesis set, implying that the divergence remains the same whatever which classifier is drawn. In this sense, our bound is more stable as the learning algorithm seeking to minimize the bound allows, in practice, to choose a better hypothesis than with the bounds of BLANCHARD and FLEURET (2007), CATONI (2007), and RIVASPLATA *et al.* (2020). We have illustrated the interest of our bound on neural networks and provides perspectives on the the stochastic majority vote classifier introduced in Chapter 5.

One future research direction related to this work is to develop new proof techniques to convert generalization bounds holding in expectation (see *e.g.*, XU and RAGINSKY, 2017) to high-probability bounds. To do so, one can apply MARKOV’s Inequality (Theorem A.2.1) and the DONSKER-VARADHAN Variational Representation of the KL divergence to obtain a high-probability bound with the bound in expectation as upper-bound. This work might be significant since the dependence in δ is, for now, only polynomial, while a logarithmic dependence is preferable. However, the limitation would be that the bound in expectation must consider a specific distribution (related to the DONSKER-VARADHAN Variational Representation).

In the next chapter, we leverage the disintegrated KL divergence and the bound of RIVASPLATA *et al.* (2020) to obtain generalization bounds with an arbitrary complexity measure. Hence, given a complexity measure defined by the user, we are able to give a generalization bound that holds with high probability over the learning sample and a hypothesis sampled from a complexity-measure-dependent distribution. Moreover, the bound with arbitrary complexity measures encompasses the generalization bounds of the literature and is recalled in Chapter 1.

GENERALIZATION BOUNDS WITH COMPLEXITY MEASURES

7

This chapter is based on the following paper

PAUL VIALLARD, RÉMI EMONET, AMAURY HABRARD, EMILIE MORVANT, and VALENTINA ZANTEDESCHI. Generalization Bounds with Arbitrary Complexity Measures. *Submitted to ICLR 2023*. (2022a)

Contents

7.1	Introduction	176
7.2	Preliminaries	177
	7.2.1 Setting	177
	7.2.2 Reminder on Disintegrated PAC-Bayesian Bounds	178
7.3	Integrating Arbitrary Complexities in Generalization Bounds	179
	7.3.1 An Introduction to our Results	179
	7.3.2 About the Gibbs Distribution	180
	7.3.3 Our Main Results: Generalization Bound with Complexity Measures	181
7.4	Using Arbitrary Complexities in Practice	184
	7.4.1 Sampling from the Gibbs Distribution	184
	7.4.2 Experiments	186
7.5	Comparison with the Generalization Bounds of the Literature	192
	7.5.1 Bounds with Arbitrary Complexity Measures	193
	7.5.2 Uniform Convergence Bounds	193
	7.5.3 Algorithmic-Dependent Bounds	195
7.6	Conclusion and Summary	197

Abstract

In statistical learning theory, generalization bounds usually involve a complexity measure that is constrained by the considered theoretical framework limiting the scope of such analysis. Among the measured mentioned in this thesis (Chapter 1) we can cite for example the VC-dimension, the stability constant of the uniform stability framework or the KL divergence used in the PAC-Bayesian framework studied in Part II. Recently, the empirical study of JIANG *et al.* (2019), made in the context of neural networks learning, has shown that (i) common complexity measures (such as the VC-dimension) do not necessarily correlate with the generalization gap, and that (ii) there exist *arbitrary* complexity measures that are better correlated with the generalization gap, but come without generalization guarantees. In this chapter, we propose to address the second point by presenting a general framework allowing one to derive some generalization bounds able to take into account a general complexity measures.

7.1 Introduction

As shown in Chapter 1 and notably in Section 1.3, statistical learning theory offers various theoretical frameworks to assess generalization by studying whether the empirical risk is representative of the true risk thanks to an upper bounding strategy of the generalization gap. While the generalization gap represents a deviation between the true risk and the empirical risk, an upper bound on this generalization gap is generally a function of mainly two quantities: (i) the size of the training sample and (ii) a complexity measure that captures how much the model overfits the data. There are some well-known complexity measures in the literature such as the VC-Dimension (Definition 1.3.3) or the Rademacher complexity (Definition 1.3.4) that quantify how much a hypothesis from the hypothesis set can overfit the data. Following a different framework, instead of considering the whole hypothesis set, algorithmic-dependent complexity measures propose to quantify how much the hypothesis obtained by a learning algorithm overfit the data: the uniform stability (Definition 1.3.6) or the robustness (Definition 1.3.7) parameters are two examples.

In order to study generalization capabilities, a recent line of works, mainly in the context of neural networks learning, is dedicated to the empirical study of different complexity measures to find those that correlate the most with the generalization gap (JIANG *et al.*, 2019; DZIUGAITE *et al.*, 2020; JIANG *et al.*, 2021). For example, given an arbitrary complexity measure, JIANG *et al.* estimate empirically the generalization gap and the complexity measure for over-parametrized models. They are able to rank the measure and the gap to obtain the Kendall's rank coefficient (KENDALL, 1938) to evaluate how a measure reflects the generalization for a particular model. However, this correlation coefficient were criticized by DZIUGAITE *et al.* (2020) because they found that if the measures empirically correlate well *on average*, they are not robust to changes of the fixed parameters (such as the depth or the width of the model).

On the one hand, while these results are extremely important to understand generalization, notably for over-parametrized models, they remain incomplete since they are only empirical. On the other hand, as we can see in Section 1.3, the bounds in the literature are restrictive because the practitioner cannot integrate in a bound its own complexity measure. In other words, to the best of our knowledge, there is no generalization bound involving arbitrary complexity measures that are found to be good proxies for the generalization gap. In this chapter, we aim to provide generalization bounds with arbitrary complexity measures that the practitioner can define. We believe that this direction is of important interest in advancing the understanding of generalization, as the generalization gap can be *provably* upper-bounded with a term that depends on a user-specified complexity measure. To get such generalization bounds, we leverage a disintegrated PAC-Bayesian bound (Theorem 2.4.1). Such a bound further allows

to derive theoretical guarantees for arbitrary complexity measures that depend on the sampled model, and on the learning sample, which is uncommon in statistical learning theory. Hence, our novel results provide theoretical foundation to the many regularization used in practice to perform model selection (e.g., L2 regularization).

The rest of the chapter is organized as follows. In Section 7.2, we provide some preliminaries dedicated to this chapter. Then, we present the main contribution in Sections 7.3 and 7.5 before concluding in Section 7.6. As usual in this thesis, in order to improve the readability of the chapter, the proofs are deferred in Appendix G.

7.2 Preliminaries

We provide in this section some quick recaps on notions that have been introduced previously in this thesis to make the reading easier. We first remind some elements about supervised learning and the generalization gap. Then, we review quickly the disintegrated generalization bounds from which the contribution of this chapter are developed.

7.2.1 Setting

We follow the setting of Chapter 2 and Section 2.4 where we consider the supervised classification setting. We recall that \mathcal{X} denotes the input space and \mathcal{Y} the label space. Moreover, we consider that an example $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is sampled from an unknown data distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$. A learning sample $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ contains m examples drawn *i.i.d.* from \mathcal{D} ; we denote the distribution of such an m -sample by \mathcal{D}^m . Let \mathbb{H} be a potentially infinite set of hypotheses $h : \mathcal{X} \rightarrow \mathcal{Y}$ that associate a label belonging to \mathcal{Y} given an input from \mathcal{X} .

Given a learning sample \mathcal{S} , we aim to find $h \in \mathbb{H}$ that minimizes the so-called true risk $R_{\mathcal{D}}(h) = \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} [h(\mathbf{x}) \neq y]$. In practice, as the data distribution \mathcal{D} is unknown, we estimate the true risk with its empirical counterpart, *i.e.*, the empirical risk $R_{\mathcal{S}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[h(\mathbf{x}_i) \neq y_i]$. We hereafter denote the generalization gap by $\phi : [0, 1]^2 \rightarrow \mathbb{R}$, which is usually defined by $\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) = |R_{\mathcal{D}}(h) - R_{\mathcal{S}}(h)|$ that quantify how much the empirical risk is representative of the true risk. Since the generalization bound is not computable, a complexity measure can be used to capture the overfitting phenomenon.

To incorporate arbitrary complexities in the bounds, we leverage the disintegrated PAC-Bayesian bounds framework (recalled in Section 2.4), *i.e.*, we upper-bound the generalization gap for a hypothesis h sampled from $\rho_{\mathcal{S}} \in \mathcal{M}(\mathbb{H})$ with a function that depends on an *arbitrary* measure of complexity. To do so, we need to consider an *a priori* belief on the hypotheses in \mathbb{H} that is modeled by a prior distribution $\pi \in \mathcal{M}^*(\mathbb{H})$.

Hence, we aim to learn, from \mathbb{S} and π , a *posterior* distribution $\rho_{\mathbb{S}}$ to assign higher probability to the best hypotheses in \mathbb{H} ; the hypothesis h is then sampled from $\rho_{\mathbb{S}}$ to obtain the guarantee with an arbitrary complexity measure.

7.2.2 Reminder on Disintegrated PAC-Bayesian Bounds

The disintegrated PAC-Bayesian bound (recalled in Section 2.4) has been introduced by CATONI (2007, Th 1.2.7) and BLANCHARD and FLEURET (2007, Prop 3.1). To the best of our knowledge, despite their interest, they have been little used in the literature and have only recently received renewed interest for deriving tight bounds in practice (see Chapter 6). Such bounds take the following form (introduced in Chapter 2).

Definition 2.4.1 (Disintegrated PAC-Bayesian Generalization Bound). Let $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow [0, 1]$ a generalization gap. A *disintegrated* PAC-Bayesian bound is defined such that if for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, there exists a function $\Phi : \mathbb{M}(\mathbb{H}) \times \mathbb{M}^*(\mathbb{H}) \times (0, 1] \rightarrow \mathbb{R}$ such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\phi(\mathbb{R}_{\mathcal{D}}^{\ell}(h), \mathbb{R}_{\mathbb{S}}^{\ell}(h)) \leq \Phi(\rho_{\mathbb{S}}, \pi, \delta) \right] \geq 1 - \delta,$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A and $\phi(\cdot)$ is, for example, $\phi(\mathbb{R}_{\mathcal{D}}^{\ell}(h), \mathbb{R}_{\mathbb{S}}^{\ell}(h)) = |\mathbb{R}_{\mathcal{D}}^{\ell}(h) - \mathbb{R}_{\mathbb{S}}^{\ell}(h)|$.

Put into words, given a training set \mathbb{S} sampled from \mathcal{D}^m , we can learn the distribution $\rho_{\mathbb{S}}$ from \mathbb{S} , and then sample the hypothesis h from $\rho_{\mathbb{S}}$ to get a bound with high probability (at least $1 - \delta$) over the random choice of \mathbb{S} and h . In this chapter, we mainly focus on RIVASPLATA *et al.* (2020)'s bound recalled below.

Theorem 2.4.1 (General Disintegrated Bound of RIVASPLATA *et al.* (2020)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\varphi(h, \mathbb{S}) \leq \underbrace{\ln \left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}')) \right]}_{\Phi(\rho_{\mathbb{S}}, \pi, \delta)} \right] \geq 1 - \delta,$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A .

7.3. Integrating Arbitrary Complexities in Generalization Bounds

In this case, the function $\varphi(h, \mathbb{S}) = m\phi(R_{\mathcal{D}}(h), R_{\mathbb{S}}(h))$ is a deviation between the true risk $R_{\mathcal{D}}(h)$ and the empirical risk $R_{\mathbb{S}}(h)$. Moreover, the function $\Phi(\rho_{\mathbb{S}}, \pi, \delta)$ is constituted by 2 terms: (i) the *disintegrated* KL divergence $\ln \frac{\rho_{\mathbb{S}}(h)}{\pi(h)}$ defining how much the prior and posterior distributions deviate for a single h , and (ii) the term $\ln \left[\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}')) \right]$ is constant *w.r.t.* $h \in \mathbb{H}$ and $\mathbb{S} \in (\mathcal{X} \times \mathcal{Y})^m$ and usually upper-bounded to instantiate the bound. In the following we refer to the whole right-hand side of the bound, $\Phi()$, as the complexity measure. This is in slight contrast with the standard definition of complexity, where the term (ii) (related to δ and the sample size m) is not included. This additional term is, in fact, constant *w.r.t.* the hypothesis $h \sim \rho_{\mathbb{S}}$ and the learning sample $\mathbb{S} \sim \mathcal{D}^m$.

In the bound of Theorem 2.4.1, the complexity term $\Phi()$ depends on the disintegrated KL divergence and suffers from some drawbacks. Indeed, the complexity term is imposed by the framework and it can be subject to high variance in practice (see Chapter 6). However, it is important to notice that this disintegrated KL divergence has a clear advantage: it only depends on the hypothesis h and \mathbb{S} , instead of the whole hypothesis class (as it is often the case, *e.g.*, when using the KL divergence in PAC-Bayesian bounds, or the VC-dimension). This might imply a better correlation between the generalization gap and some complexity measures. In the next section, we leverage this disintegrated KL divergence to derive our main contribution: a general bound that involves arbitrary complexity measures.

7.3 Integrating Arbitrary Complexities in Generalization Bounds

We first begin with a short presentation of our result to give some preliminary intuitions and to introduce the notion of Gibbs distribution which is a key element of our contribution. We then provide formally our theoretical result in Section 7.3.3.

7.3.1 An Introduction to our Results

Let $\Phi_{\mu}(h, \mathbb{S}, \delta)$ be a real-valued function dependent on an additional function $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$ that takes a hypothesis h , a learning sample \mathbb{S} , and the parameter δ as arguments. This function $\mu()$ parametrizes the complexity measure based on the data sample \mathbb{S} and the model h , and thus allows us to introduce arbitrary complexity measures in the bound; we further denote the function $\mu()$ as *parametric function*. As an example, when \mathbb{H} is a set of hypotheses $h_{\mathbf{w}}$ parameterized by some weights $\mathbf{w} \in \mathbb{R}^d$, we can fix $\mu(h_{\mathbf{w}}, \mathbb{S}) = \|\mathbf{w}\|$, for some norm $\|\cdot\|$. This means that $\mu(h_{\mathbf{w}}, \mathbb{S})$ can be set to the regularization term of the chosen objective function, so that the complexity,

hence the bound, will depend on it. This is not really new since for example uniform stability bounds allow one to consider such norms (see *e.g.*, KAKADE *et al.*, 2008). Here we want to illustrate that we simply incorporate some norms of parameters but as the reader may guess our framework allows one to consider a broader family of complexity measures as we will see later. Given such a parametric function $\mu(\cdot)$, the bound we derive in Theorem 7.3.1 takes the following form.

Definition 7.3.1 (Generalization Bound with Complexity Measures). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap, $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$ be a parametric function. A generalization bound with arbitrary complexity measures is defined such that if for any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , there exists a real-valued function $\Phi_\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \times (0, 1] \rightarrow \mathbb{R}$ such that for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \leq \Phi_\mu(h, \mathcal{S}, \delta) \right] \geq 1 - \delta. \quad (7.1)$$

The main trick to obtain such a result is via a posterior distribution $\rho_{\mathcal{S}}$: we incorporate the function $\mu(\cdot)$ by fixing the distribution $\rho_{\mathcal{S}}$ as a Gibbs distribution defined by

$$\rho_{\mathcal{S}}(h) \propto \exp[-\alpha R_{\mathcal{S}}(h) - \mu(h, \mathcal{S})], \quad \text{where } \alpha \in \mathbb{R}^+. \quad (7.2)$$

This Gibbs distribution $\rho_{\mathcal{S}}$ is interesting from an optimization point of view: a hypothesis h is more likely to be sampled from it when the objective function $h \mapsto R_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$ is low for a given \mathcal{S} . For example, when $\mu(h, \mathcal{S}) = 0$, a hypothesis h is more likely to be sampled when its empirical risk $R_{\mathcal{S}}(h)$ is low. Conversely, when $\mu(\cdot)$ is non-constant, the function serves as a “regularizing term”, so that a hypothesis is more likely to be sampled when the trade-off $R_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$ is low. In both cases, the higher α , the more the density of $\rho_{\mathcal{S}}$ is concentrated around those hypotheses that minimize the empirical risk. Moreover, it seems that Equation (7.2) is restrictive but it can actually represent any probability density functions. Indeed, let $\rho'_{\mathcal{S}}$ be a probability distribution on \mathbb{H} , *e.g.*, a Gaussian or a Laplace distribution, by setting $\mu(h, \mathcal{S}) = -\alpha R_{\mathcal{S}}(h) - \ln \rho'_{\mathcal{S}}(h)$ we can retrieve the distribution $\rho'_{\mathcal{S}}$. Actually, the Gibbs distribution is well-known and studied in learning theory as discussed in the next section.

7.3.2 About the Gibbs Distribution

In this section, we would like to highlight two main lines of works that are related to our setting: (i) the usage of the Gibbs distribution in the “classical” PAC-Bayesian theory and (ii) the link between this distribution and optimization.

Gibbs distribution in the PAC-Bayesian theory. The Gibbs distribution has started to be studied in the PAC-Bayesian theory by CATONI (2004, 2007). Moreover, ALQUIER *et al.* (2016, Theorem 4.2 and 4.3) developed PAC-Bayesian generalization bounds with the Gibbs distribution considered in Equation (7.2) with $\mu(h, \mathbb{S}) = 0$ as posterior. However, their theorems analyze the expected true risk $\mathbb{E}_{h \sim \rho_{\mathbb{S}}} R_{\mathcal{D}}(h)$ while we are interested in a *single* hypothesis h sampled from $\rho_{\mathbb{S}}$. Moreover, their bounds involve the non-computable and hypothesis-set-dependent KL divergence between the Gibbs distribution and a prior distribution. Hence, the computation of the KL divergence must be upper-bounded to allow one to instantiate this bound in practice. As we will see further, the bounds of Theorem 7.3.1 and Corollary 7.3.1 do not have this issue since they only require to know the expression of the density (up to the normalization constant) for $h \sim \rho_{\mathbb{S}}$ and $h' \sim \pi$.

Relationship between optimization and the Gibbs distribution. Given an objective function $f : \mathbb{R}^D \rightarrow \mathbb{R}$, the Stochastic Gradient LANGEVIN Dynamics (SGLD) algorithm (WELLING and TEH, 2011) learns some parameter $\mathbf{w} \in \mathbb{R}^D$ by running several iterations of the form

$$\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \eta \nabla f(\mathbf{w}) + \sqrt{\frac{2\eta}{\alpha}} \epsilon_t, \quad \text{where } \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D), \quad (7.3)$$

where \mathbf{w}_t is the weight vector learned at the iteration $t \in \mathbb{N}$, the parameter η is the learning rate and α is the concentration parameter of the Gibbs distribution. This algorithm has an interesting feature: when the learning rate η tends to zero, the SGLD algorithm becomes a continuous-time process, called LANGEVIN diffusion, defined as the following stochastic differential equation in Equation (7.4). Indeed, Equation (7.3) can be seen as the EULER-MARUYAMA discretization (see *e.g.*, RAGINSKY *et al.*, 2017) of Equation (7.4) defined for $t \geq 0$ as

$$d\mathbf{w}(t) = -\nabla f(\mathbf{w}(t))dt + \sqrt{2\alpha}\mathbf{B}(t), \quad (7.4)$$

where $\mathbf{B}(t)$ is the Brownian motion in \mathbb{R}^D . Under some mild assumptions on the function f , CHIANG *et al.* (1987) show that the invariant distribution of the LANGEVIN diffusion is the Gibbs distribution proportional to $\exp(-\alpha f(\mathbf{w}))$.

7.3.3 Our Main Results: Generalization Bound with Complexity Measures

Thanks to the definition of $\rho_{\mathbb{S}}$, we now state our main result which consists in providing a bound on the generalization gap involving our parametric measure μ with respect to a posterior $\rho_{\mathbb{S}}(h) \propto \exp[-\alpha R_{\mathcal{S}}(h) - \mu(h, \mathbb{S})]$.

Theorem 7.3.1 (Generalization Bound with Complexity Measures). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap. For any \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h' \sim \pi, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathcal{S}}(h)) \leq \left[\alpha \mathbb{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S}) \right] - \left[\alpha \mathbb{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S}) \right] \right. \\ \left. + \ln \frac{\pi(h')}{\pi(h)} + \ln \left(\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp[\phi(\mathbb{R}_{\mathcal{D}}(h'), \mathbb{R}_{\mathcal{S}'}(h'))] \right) \right] \geq 1 - \delta. \end{aligned}$$

Proof. Deferred to Appendix G.1. ■

This theorem is general since it depends only on the functions $\phi(\cdot)$ and $\mu(\cdot)$ that must be fixed by the practitioner. Given $\phi(\cdot)$ and $\mu(\cdot)$, we can note an element that can be surprising at the first reading: it appears indeed possible to sample hypotheses with a bad objective $\mathbb{R}_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$ value and to obtain a tight generalization bound. However, by definition the Gibbs distribution $\rho_{\mathcal{S}}$, such a sampled hypothesis $h \sim \rho_{\mathcal{S}}$ is less probable since the density is higher when the objective is low. On the other hand, it is highly probable to sample a hypothesis h from $\rho_{\mathcal{S}}$ with a good objective $\mathbb{R}_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$ value. Concerning the tightness of the bound, it may appear loose. However to get a bound that converges when m increases, it is sufficient to fix $\phi(\cdot)$ as a function of m such as $\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathcal{S}}(h)) = 2m[\mathbb{R}_{\mathcal{D}}(h) - \mathbb{R}_{\mathcal{S}}(h)]^2$ or $\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathcal{S}}(h)) = m \text{kl}[\mathbb{R}_{\mathcal{S}}(h) \parallel \mathbb{R}_{\mathcal{D}}(h)]$. Then, the tightness of the bound depends on the generalization gap $\phi(\cdot)$, the parametric function $\mu(\cdot)$ and α . The remaining challenge is to upper-bound $\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp[\phi(\mathbb{R}_{\mathcal{D}}(h'), \mathbb{R}_{\mathcal{S}'}(h'))]$ and $\mathbb{E}_{h' \sim \pi} \ln \frac{\pi(h')}{\pi(h)}$ to get a practical bound. As an illustration we provide in the next corollary an instantiation of Theorem 7.3.1 with the generalization gap $\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathcal{S}}(h)) = 2m[\mathbb{R}_{\mathcal{D}}(h) - \mathbb{R}_{\mathcal{S}}(h)]^2$ and $\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathcal{S}}(h)) = m \text{kl}[\mathbb{R}_{\mathcal{S}}(h) \parallel \mathbb{R}_{\mathcal{D}}(h)]$ when π is a uniform distribution on a bounded hypothesis set \mathbb{H} .

Corollary 7.3.1 (Practical Generalization Bound with Complexity Measures). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any bounded hypothesis set \mathbb{H} , given the uniform prior distribution π on \mathbb{H} , for any $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, with

7.3. Integrating Arbitrary Complexities in Generalization Bounds

probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$, $h' \sim \pi$, $h \sim \rho_{\mathcal{S}}$ we have

$$\text{kl} [\mathcal{R}_{\mathcal{S}}(h) \parallel \mathcal{R}_{\mathcal{D}}(h)] \leq \frac{1}{m} \left[[\alpha \mathcal{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha \mathcal{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S})] + \frac{8\sqrt{m}}{\delta^2} \right]_+, \quad (7.5)$$

$$\left| \mathcal{R}_{\mathcal{D}}(h) - \mathcal{R}_{\mathcal{S}}(h) \right| \leq \sqrt{\frac{1}{2m} \left[[\alpha \mathcal{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha \mathcal{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S})] + \frac{8\sqrt{m}}{\delta^2} \right]_+}, \quad (7.6)$$

where $[a]_+ = \max(0, a)$, and $\rho_{\mathcal{S}}$ is the Gibbs distribution defined by Equation (7.2).

Proof. Deferred to Appendix G.2. ■

Interestingly, Corollary 7.3.1 gives a bound on $|\mathcal{R}_{\mathcal{D}}(h) - \mathcal{R}_{\mathcal{S}}(h)|$ and $\text{kl} [\mathcal{R}_{\mathcal{S}}(h) \parallel \mathcal{R}_{\mathcal{D}}(h)]$ where all terms except $\mathcal{R}_{\mathcal{D}}(h)$ are computable. To compute Equations (7.5) and (7.6) we can rearrange the terms to obtain a generalization bound on the true risk $\mathcal{R}_{\mathcal{D}}(h)$. Indeed, we have respectively

$$\mathcal{R}_{\mathcal{D}}(h) \leq \overline{\text{kl}} \left(\mathcal{R}_{\mathcal{S}}(h) \left| \frac{1}{m} \left[[\alpha \mathcal{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha \mathcal{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S})] + \frac{8\sqrt{m}}{\delta^2} \right]_+ \right. \right), \quad (7.7)$$

$$\text{and } \mathcal{R}_{\mathcal{D}}(h) \leq \mathcal{R}_{\mathcal{S}}(h) + \sqrt{\frac{1}{2m} \left[[\alpha \mathcal{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha \mathcal{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S})] + \frac{8\sqrt{m}}{\delta^2} \right]_+}, \quad (7.8)$$

where $\overline{\text{kl}}(q|\tau) = \max \left\{ p \in (0, 1) \mid \text{kl}(q||p) \leq \tau \right\}$ (see Definition 2.3.3). These bounds are used in Section 7.4 to illustrate the generalization guarantees with different values of the parametric function $\mu(\cdot)$ and α . Moreover, as illustrated in the experiments, Equation (7.7) is a tighter bound on the true risk $\mathcal{R}_{\mathcal{D}}(h)$ than Equation (7.8): we can prove it formally with PINSKER's inequality (see all Equation (2.15)).

Nevertheless, the complexity measures $\Phi_{\mu}(h, \mathcal{S}, \delta)$ of Equations (7.5) and (7.6) enjoys asymptotic convergence for $m \rightarrow +\infty$. However, as mentioned previously, the complexity also depends on α and $\mu(\cdot)$: the convergence rate can then be degraded by $[\alpha \mathcal{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha \mathcal{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S})]$ which may be potentially large. For example, for all $(h, \mathcal{S}) \in \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m$, if the empirical risk $\mathcal{R}_{\mathcal{S}}(h')$ is large (which is often the case when h' is sampled from a uniform prior on \mathbb{H}), $\alpha = m$, and $\mu(h, \mathcal{S}) = 0$, then

the complexity measure, simplified into $\Phi_\mu(h, \mathcal{S}, \delta) = \left[[R_{\mathcal{S}}(h') - R_{\mathcal{S}}(h)] + \frac{1}{m} \ln \frac{2\sqrt{m}}{\delta} \right]_+$ for $\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) = \text{kl}[R_{\mathcal{S}}(h) \| R_{\mathcal{D}}(h)]$, will be large. We thus have to set α and $\mu(\cdot)$ such that (i) $\rho_{\mathcal{S}}$ allows to sample a hypothesis h associated with a low objective function $h \mapsto R_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$ and (ii) the complexity measure $\Phi_\mu(h, \mathcal{S}, \delta)$ is tight, resulting in a meaningful bound. For example, with $\alpha = \sqrt{m}$ and $\mu(h, \mathcal{S}) = 0$, the distribution $\rho_{\mathcal{S}}$ will be less concentrated around the minimizers of the empirical risk, but the complexity measure will be tighter compared to the previous example: $\left[\frac{1}{\sqrt{m}} [R_{\mathcal{S}}(h') - R_{\mathcal{S}}(h)] + \frac{1}{m} \ln \frac{2\sqrt{m}}{\delta} \right]_+$.

The tightness of the bounds can be potentially improved since we choose a uniform distribution for the prior π . To obtain better bounds, data-dependent priors have been heavily used in the PAC-Bayesian literature (see e.g., PARRADO-HERNÁNDEZ *et al.*, 2012; DZIUGAITE *et al.*, 2021; PÉREZ-ORTIZ *et al.*, 2021). However, we think that the uniform distribution helps to better understand the generalization phenomenon. Indeed, the hypothesis h sampled from the uniform distribution π has a high chance of underfitting. Hence, if the hypothesis $h \sim \rho_{\mathcal{S}}$ has its associated bound value that is tight, we can easily interpret this hypothesis generalizes well. On the other hand, if a generalization bound with a data-dependent prior π is tight, it means that a posterior $\rho_{\mathcal{S}}$ (not very far from the prior) allows us to generalize well. In this case, we do not know if the generalization capability is essentially due to the choice of a good prior π or if it comes mainly from the learned posterior $\rho_{\mathcal{S}}$.

7.4 Using Arbitrary Complexities in Practice

The bound of Corollary 7.3.1 is not directly usable in its current form: the remaining challenge is to sample h from the Gibbs distribution $\rho_{\mathcal{S}}$ defined in Equation (7.2); we address the sampling issue in Section 7.4.1. Then, we make use of the proposed solution to assess our bound in practice. Section 7.4.2.1 introduces our experimental setting and Section 7.4.2.2 gives an overview of the tightness of the bound. Additionally, we give an overview of the influence of α and the number of parameters in Sections 7.4.2.3 and 7.4.2.4.

7.4.1 Sampling from the Gibbs Distribution

Sampling from the Gibbs distribution of Equation (7.2) is a hard task: naively, it requires to evaluate the function $h \mapsto -\alpha R_{\mathcal{S}}(h) - \mu(h, \mathcal{S})$ for all $h \in \mathbb{H}$, which is intractable when \mathbb{H} is infinite or even large. We tackle this issue for over-parameterized models, which we later consider in Section 7.4.2.1 in an empirical study of our bound. Let us consider a set \mathbb{H} of hypotheses $h_{\mathbf{w}}$ parameterized by $\mathbf{w} \in \mathbb{R}^D$, and a tractable

Algorithm 7.1 Stochastic MALA

```

1: Input: Learning set  $\mathcal{S}$ , weights  $\mathbf{w}$ , function  $\mu(\cdot)$ , loss function  $\ell(\cdot)$ 
2: Hyperparameters: Number of iterations  $T$ , learning rate  $\eta$ , parameter  $\alpha$ 
3: for  $t \leftarrow 1 \dots T$  do
4:    $\mathcal{U} \leftarrow$  Sample (without replacement) a mini-batch from  $\mathcal{S}$ 
5:    $\mathbf{w}' \leftarrow$  Sample from the distribution  $P_{\mathcal{U}}^{\mathbf{w}}$ 
6:    $\tau \leftarrow \min \left( 1, \frac{\rho_{\mathcal{U}}(\mathbf{w}')P_{\mathcal{U}}^{\mathbf{w}'}(\mathbf{w})}{\rho_{\mathcal{U}}(\mathbf{w})P_{\mathcal{U}}^{\mathbf{w}}(\mathbf{w}')} \right)$ 
7:    $u \leftarrow$  Sample from the distribution  $\text{Uni}(0, 1)$ 
8:   if  $u \leq \tau$  then
9:      $\mathbf{w} \leftarrow \mathbf{w}'$ 
10: return  $\mathbf{w}$ 

```

distribution denoted $P_{\mathcal{U}}^{\mathbf{w}}$ (e.g., a Gaussian distribution) such that its density approximates the density of $\rho_{\mathcal{S}}$. In this setting, to learn such an auxiliary distribution, we propose in Algorithm 7.1 a stochastic version of the Metropolis Adjusted LANGEVIN Algorithm (MALA, BESAG (1994))¹. Its objective is to generate samples from $\rho_{\mathcal{S}}$ by iteratively refining the auxiliary distribution, that we define as

$$P_{\mathcal{U}}^{\mathbf{w}} = \mathcal{N} \left(\mathbf{w} - \eta \nabla \left[R_{\mathcal{U}}^{\ell}(\mathbf{w}) + \frac{1}{\alpha} \mu(\mathbf{w}, \mathcal{U}) \right], \frac{2\eta}{\alpha} \mathbf{I} \right), \quad (7.9)$$

where $R_{\mathcal{U}}^{\ell}(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{U}} \ell(h_{\mathbf{w}}, (\mathbf{x}, y))$ is the empirical risk on the mini-batch $\mathcal{U} \subseteq \mathcal{S}$, and $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ is a loss function. Concretely, we initialize the parameters \mathbf{w} of the model as the output of an optimization algorithm (Vanilla SGD in our case). Then, we refine them as follows: at each iteration, given the current weights \mathbf{w} and a mini-batch $\mathcal{U} \subseteq \mathcal{S}$ (Line 4), we sample a candidate vector \mathbf{w}' (Line 5) according to the distribution $P_{\mathcal{U}}^{\mathbf{w}}$; then (Line 6 to 9) we decide to reject or accept the new candidate to become our current weights \mathbf{w} , depending on its ratio $\tau = \min \left(1, \frac{\rho_{\mathcal{U}}(\mathbf{w}')P_{\mathcal{U}}^{\mathbf{w}'}(\mathbf{w})}{\rho_{\mathcal{U}}(\mathbf{w})P_{\mathcal{U}}^{\mathbf{w}}(\mathbf{w}')} \right)$ being larger than a control value u sampled from the uniform distribution on $[0, 1]$. Under the mild assumption that $\rho_{\mathcal{S}}$ is absolute continuous *w.r.t.* $P_{\mathcal{S}}^{\mathbf{w}}$ (see CHIB and GREENBERG, 1995, for details), when the number of iterations tends to infinity and when $\mathcal{U} = \mathcal{S}$, the returned \mathbf{w} is sampled according to $\rho_{\mathcal{S}}$ (SMITH and ROBERTS, 1993). Note that this assumption requires that the tractable distribution $P_{\mathcal{S}}^{\mathbf{w}}$ has a strictly positive density when the density of $\rho_{\mathcal{S}}$ is strictly positive as well (see CHIB and GREENBERG, 1995).

¹See CHIB and GREENBERG (1995) for an introduction on Metropolis-Hastings Algo on which MALA is based.

7.4.2 Experiments

7.4.2.1 Experimental Setting

In this section, we investigate the tightness of our bound of Equations (7.7) and (7.8) on the MNIST (LECUN *et al.*, 1998) and FashionMNIST (XIAO *et al.*, 2017) datasets. We keep the original learning set as \mathbb{S} and the original test set \mathbb{T} to estimate the true risk, that we refer to as test risk $R_{\mathcal{T}}(h)$.

Model. We use a “Convolutional Network in Network” (LIN *et al.*, 2013) similarly to JIANG *et al.* (2019) and DZIUGAITE *et al.* (2020), that consists of several modules of 3 convolutional layers each followed by a leaky ReLU activation function (its negative slope is set to 10^{-2}). The depth of the network L is the number of convolutional layers, and the width H is the number of channels of each convolution. In addition, for each layer i , we denote its weights by \mathbf{w}_i . More precisely, the modules of this model can be described as follows. A module takes two parameters as argument: the number of input channels H_{in} and the number of output channels H_{out} and applies consecutively three convolutional layers (each followed by a leaky ReLU activation function). The first layer is composed of a 3×3 kernel (where the stride *resp.* padding is set to 2 *resp.* 1) with H_{in} channels as input and H_{out} as output. The two other layers have a 1×1 kernel with H_{out} channels as input and output. Then, the network is constructed as follows: (a) we have a module where $H_{\text{out}} = H$ and H_{in} is the number of channels in the input (b) we have $(L/3) - 1$ modules with $H_{\text{in}} = H_{\text{out}} = H$ and (c) we have a convolutional layer with a 1×1 kernel with $\text{card}(\mathbb{Y})$ channels as output followed by a leaky ReLU activation and an average pooling layer. In the experiments, we consider $L \in \{9, 12, 15\}$ and $H \in \{128, 256\}$. Furthermore, we initialize the network with the weights $\mathbf{w}^0 \in \mathbb{R}$ obtained the uniform Kaiming He initializer (HE *et al.*, 2015). The set \mathbb{H} corresponds to the hypotheses $h_{\mathbf{w}}$ that can be obtained from this initialization (and we clamp the weights during the optimization in the initialization’s interval).

Arbitrary complexity measures. We study different complexity measures parametrized by different functions $\mu(\cdot)$ from JIANG *et al.* (2019, Sec. C)². Indeed, we consider the

²Note we consider a subset of the functions studied by JIANG *et al.*: we select those that are optimizable.

7.4. Using Arbitrary Complexities in Practice

6 following parametric functions $\mu()$:

$$\begin{aligned} \text{DIST_FRO}(h_{\mathbf{w}}) &= \sum_{i=1}^L \|\mathbf{w}_i - \mathbf{w}_i^0\|_2, \quad \text{and} \quad \text{DIST_L}_2(h_{\mathbf{w}}) = \|\mathbf{w} - \mathbf{w}^0\|_2, \\ \text{and} \quad \text{PARAM_NORM}(h_{\mathbf{w}}) &= \sum_{i=1}^L \|\mathbf{w}_i\|_2^2, \quad \text{and} \quad \text{PATH_NORM}(h_{\mathbf{w}}) = \sum_{i=1}^{\text{card}(\mathbb{Y})} h_{\mathbf{w}^2}(\mathbf{1})[i], \\ \text{and} \quad \text{SUM_FRO}(h_{\mathbf{w}}) &= L \left(\prod_{i=1}^L \|\mathbf{w}_i\|_2^2 \right)^{\frac{1}{L}}, \quad \text{and} \quad \text{ZERO}(h_{\mathbf{w}}) = 0. \end{aligned}$$

We define the considered measures with α taken among 5 values uniformly spaced between $[\sqrt{m}, m]$. Note that we analyze other parametric functions that depends on the learning sample \mathcal{S} . However, since the results are similar, we decided to defer the results in Appendix G.

Bound optimization. To compute our bound in Equations (7.7) and (7.8), we aim to minimize the objective function $\mathbf{w} \mapsto R_{\mathcal{U}}^{\ell}(\mathbf{w}) + \frac{1}{\alpha} \mu(\mathbf{w}, \mathbb{U})$, via Algorithm 7.1. We set the loss function to the bounded cross entropy from DZIUGAITE and ROY (2018): $\ell(h, (\mathbf{x}, y)) = -\frac{1}{4} \ln(e^{-4} + (1 - 2e^{-4})h[y])$, where $h[y]$ is the probability assigned to label y by h . The advantage of DZIUGAITE and ROY (2018)'s cross-entropy is that it lies in $\ell(h, (\mathbf{x}, y)) \in [0, 1]$, whereas the classical cross-entropy is unbounded. Indeed, taking into account the classical cross-entropy when optimizing the objective function would lead to focus too much on the risk minimization, while we want to take into account $\frac{1}{\alpha} \mu(\mathbf{w}, \mathbb{U})$. We initialize the weights $\mathbf{w} \in \mathbb{R}^D$ to the solution found by optimizing the objective function with a Vanilla SGD (with 10 epochs, a learning rate of 10^{-1} , and a batch size of 64). Given these initial parameters \mathbf{w} , we execute Algorithm 7.1 for 1 epoch with a mini-batch of size 64, where $\eta = 10^{-4}$.

7.4.2.2 Tightness of the Bounds

For each parametric function $\mu()$, we report in Figures 7.1 and 7.2, the test risks $R_{\mathcal{T}}(h)$ and the values of the tightest bound (*w.r.t.* α) associated to Equations (7.7) and (7.8) for different parameters (depth L , width H). First of all, we can remark that certain empirical risks are high. This is due to the sampling of the hypothesis h from the distribution $\rho_{\mathcal{S}}$: the hypothesis does not necessarily minimizes the objective function $h \mapsto R_{\mathcal{S}}(h) + \frac{1}{\alpha} \mu(h, \mathcal{S})$. We can nevertheless observe that the bounds' values are higher when the empirical risk $R_{\mathcal{T}}(h)$ is low. This can be explained by the fact that $[\alpha R_{\mathcal{S}}(h') + \mu(h', \mathcal{S})] - [\alpha R_{\mathcal{S}}(h) + \mu(h, \mathcal{S})]$ is large in this case. When the empirical risks are a bit higher, the bounds become tighter for certain parametric function such as DIST_L2, SUM_FRO. This confirms that there is an interest to use a parametric function that captures information on the model during the training phase.

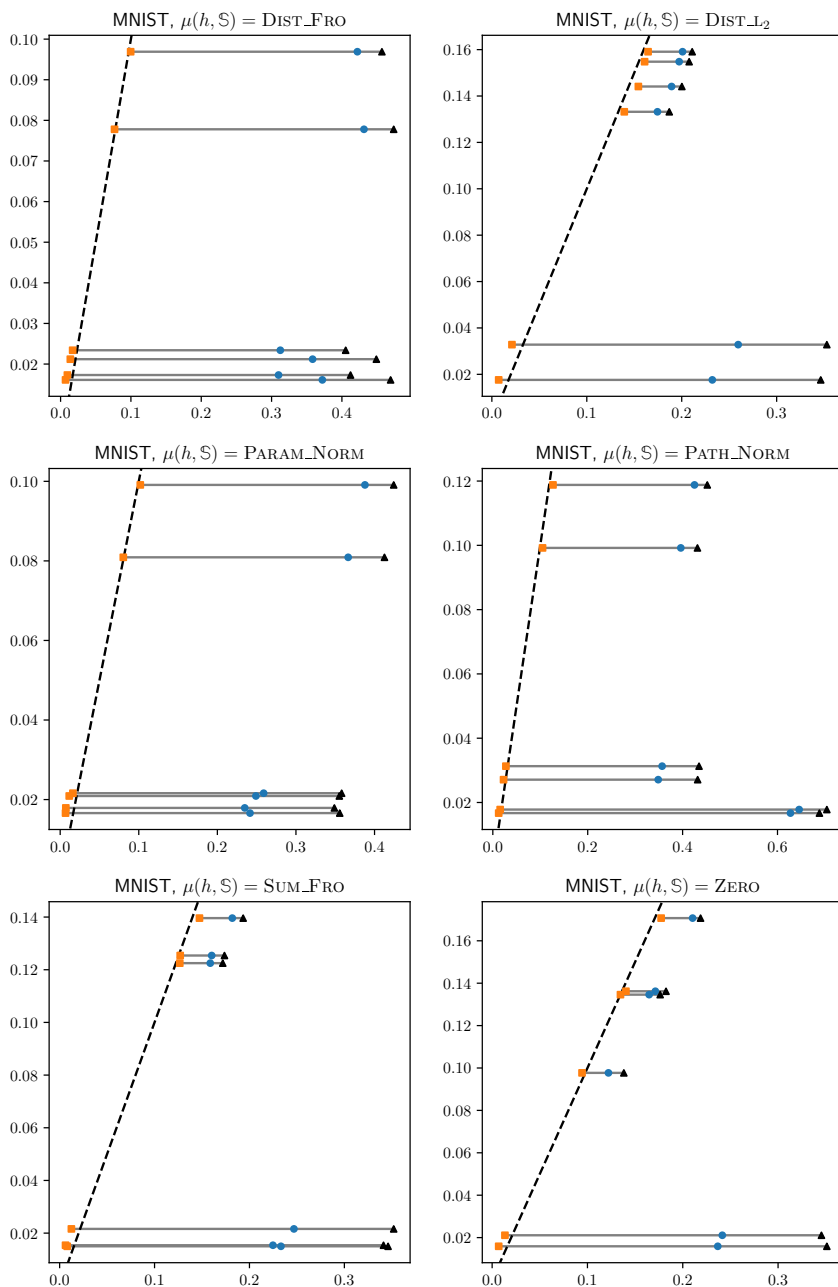


Figure 7.1. Scatter plot given a parametric function $\mu(h, \mathcal{S})$, where each segment represents a neural network $h_{\mathbf{w}}$ learned with a given α , width H and depth L . For each segment, there is a corresponding orange square and a blue circle. The orange squares corresponds to the empirical risk $R_{\mathcal{S}}(h)$ (x-axis) and test risk $R_{\mathcal{T}}(h)$ (y-axis). The blue circle resp. the black triangle represents Equation (7.7) resp. Equation (7.8) in the x-axis and the test risk $R_{\mathcal{T}}(h)$ in the y-axis. The dashed line is the identity function.

7.4. Using Arbitrary Complexities in Practice

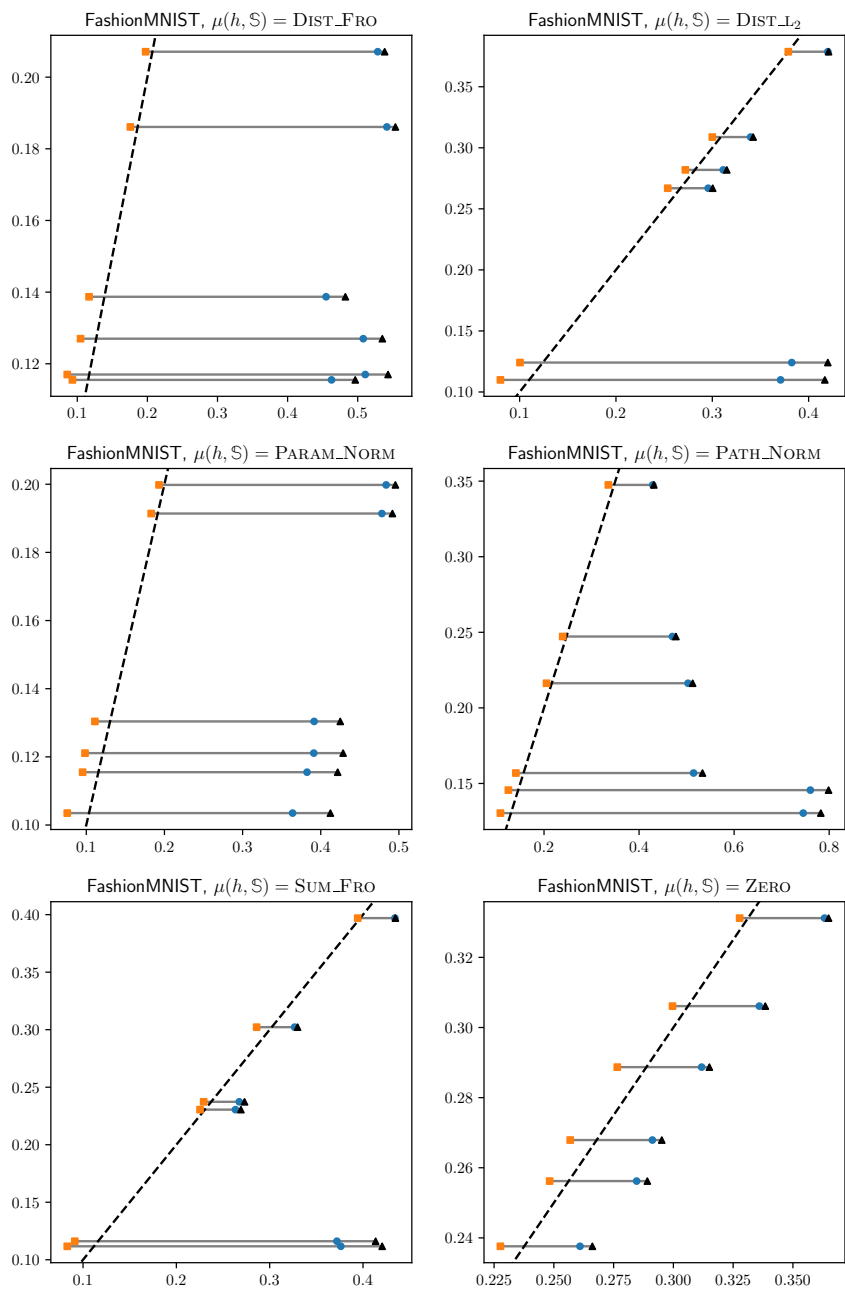


Figure 7.2. Scatter plot given a parametric function $\mu(h, \mathcal{S})$, where each segment represents a neural network h_w learned with a given α , width H and depth L . For each segment, there is a corresponding orange square and a blue circle. The orange squares corresponds to the empirical risk $R_S(h)$ (x-axis) and test risk $R_T(h)$ (y-axis). The blue circle resp. the black triangle represents Equation (7.7) resp. Equation (7.8) in the x-axis and the test risk $R_T(h)$ in the y-axis. The dashed line is the identity function.

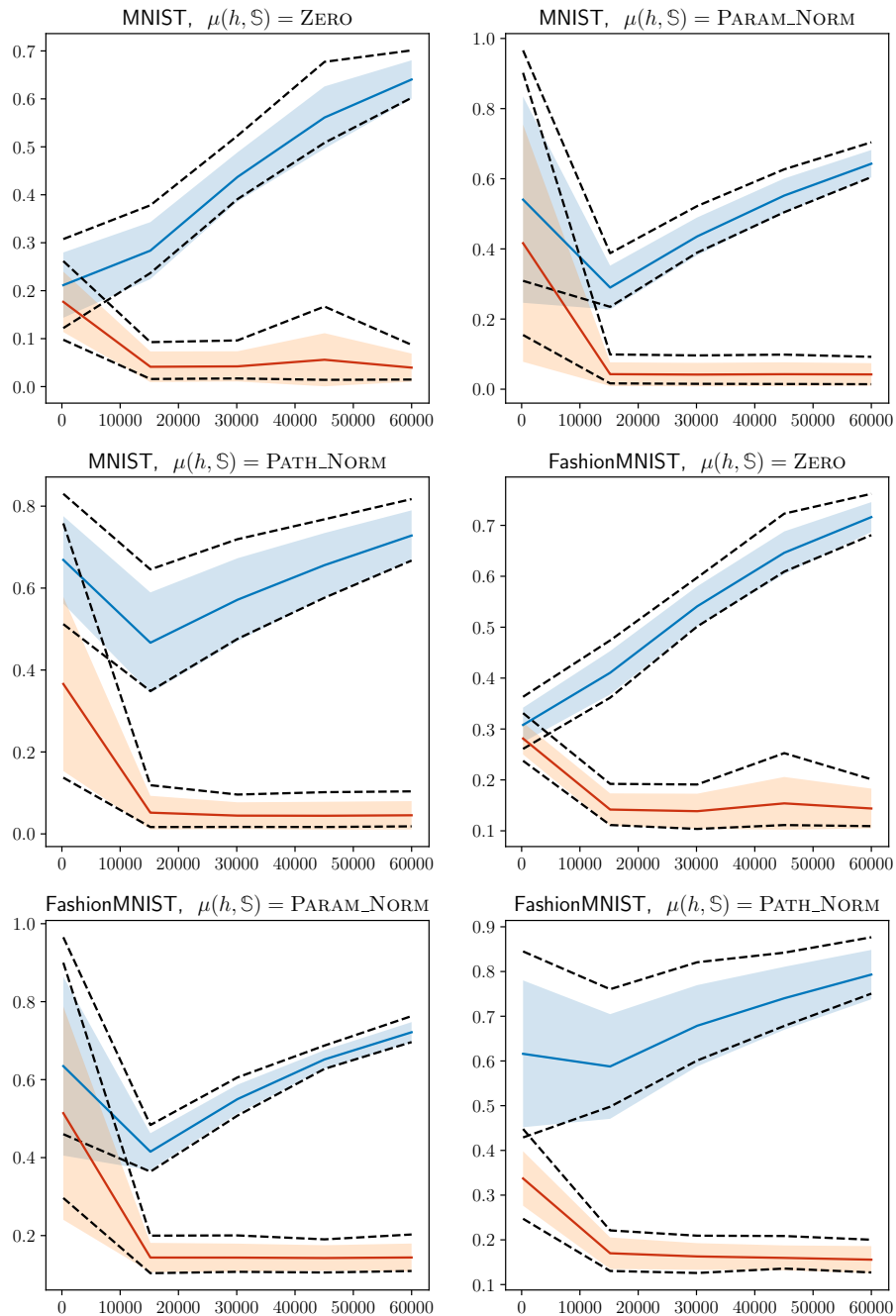


Figure 7.3. Influence of the parameter α (in the x-axis) for three parametric functions: ZERO, PARAM_NORM, and PATH_NORM for MNIST and FashionMNIST. The bound values are represented in blue and the test risk in red. The two (solid) lines are the mean values computed on the depths and widths; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

7.4. Using Arbitrary Complexities in Practice

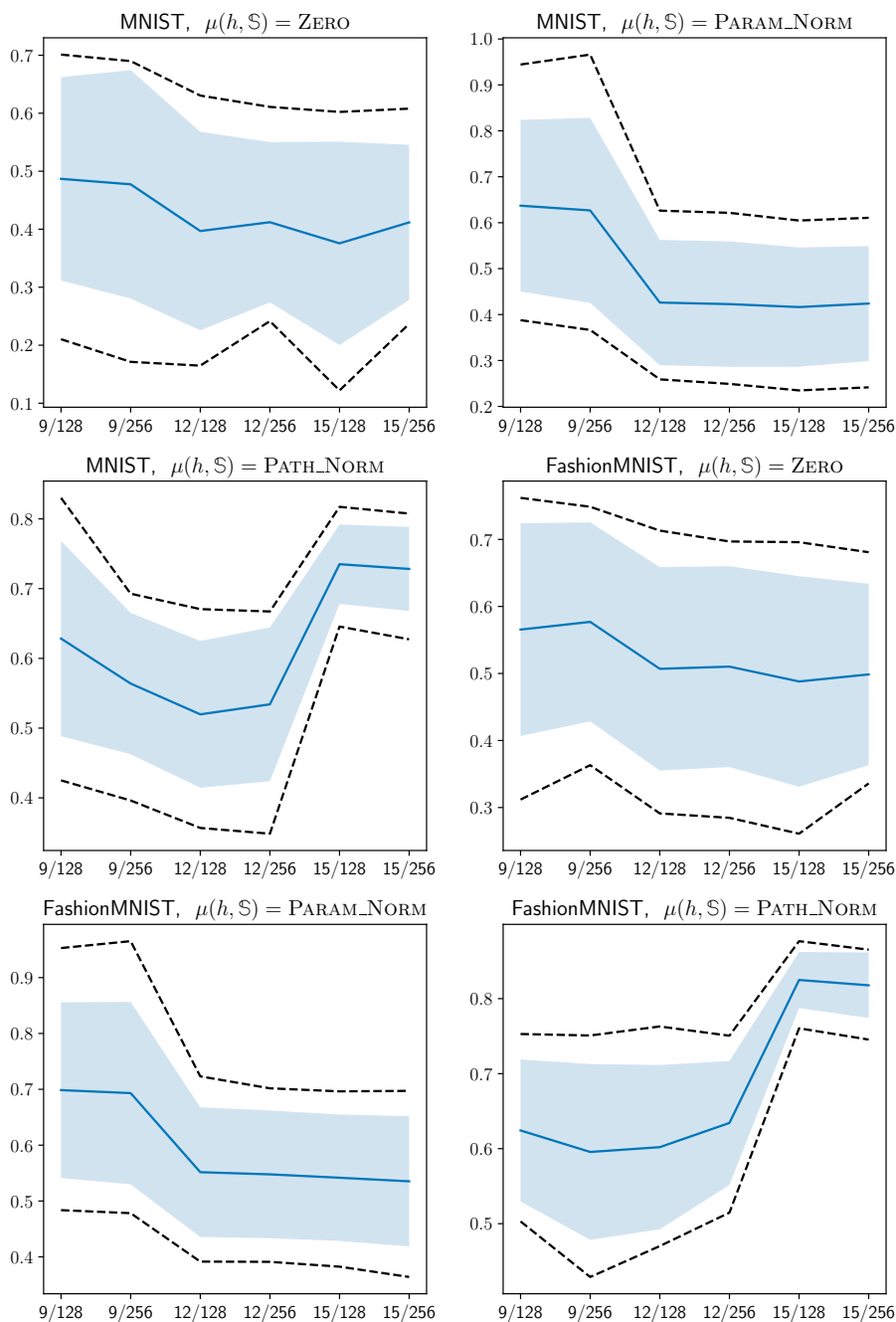


Figure 7.4. Influence of the depth and the width (in the x-axis as “depth/width”) for three parametric functions: ZERO, PARAM_NORM, and PATH_NORM for MNIST and FashionMNIST. The (solid) lines are the mean values computed on the different values of α ; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

7.4.2.3 Influence of the Parameter α

We analyze the influence of the parameter α in Equation (7.7). To do so, we plot an overview of the evolution for the bounds and the test risks $R_{\mathcal{T}}(h)$; the details are reported in Appendix G. For each parameter α , we plot in Figure 7.3 the mean, the standard deviation, the minimum and the maximum for the different parameters (depth and width). In general, the bound increases when the α tends to m but the test risks $R_{\mathcal{T}}(h)$ are less prone to variations. Indeed, the higher the parameter α , the more concentrated around the minimizers the hypothesis will be sampled. On the contrary, for a small α (e.g., $\alpha = \sqrt{m}$), the Gibbs distribution defined in Equation (7.2) is less concentrated making the test risks potentially high with a tighter generalization bound.

7.4.2.4 Influence of the Depth/Width

In Figure 7.4, we show an overview of the evolution of Equation (7.7) with respect to the depth and the width. More precisely, we report the mean, the standard deviation, the minimum and the maximum values for three parametric functions (ZERO, PARAM_NORM, and PATH_NORM).

Interestingly, the evolution of the bounds highly depends on the chosen parametric function $\mu(\cdot)$. For instance, the bound increases with PATH_NORM when the depth and the width increase. This is in contrast with PARAM_NORM that decreases when the number of parameters increases. This shows the interest of our framework: considering a user-specified complexity measure $\Phi_{\mu}(\cdot)$ can help to understand the generalization of over-parameterized models (that are sampled from $\rho_{\mathcal{S}}$).

7.5 Comparison with the Generalization Bounds of the Literature

In this section, we theoretically compare generalization bounds with arbitrary complexity measures compared to literature's bounds. We prove that our bound generalizes the uniform-convergence and the algorithmic-dependent bounds. Additionally, we show that the algorithmic-dependent bounds can be tighter than uniform-convergence bounds. To do so, we propose in Propositions 7.5.1 to 7.5.3 a reinterpretation of the high probability bounds in terms of sets. Additionally, we prove in Corollaries 7.5.1 and 7.5.2 two special cases of our bound in Theorem 7.3.1 that generalizes the two types of bounds.

7.5.1 Bounds with Arbitrary Complexity Measures

In order to compare our framework with the uniform-convergence and the algorithmic-dependent bounds, we translate Definition 7.3.1 into the following set-theoretic result.

Proposition 7.5.1 (Set-theoretic view of Definition 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \times (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 7.3.1. Under these conditions, with $\mathbb{Z}_\delta = \left\{ (h, \mathcal{S}) \in \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m : \phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \leq \Phi_\mu(h, \mathcal{S}, \delta) \right\}$, and $\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [(h, \mathcal{S}) \in \mathbb{Z}_\delta] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (7.1)} &\iff \forall (h, \mathcal{S}) \in \mathbb{Z}_\delta, \phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \leq \Phi_\mu(h, \mathcal{S}, \delta) \\ &\iff \sup_{(h, \mathcal{S}) \in \mathbb{Z}_\delta} \left\{ \phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) - \Phi_\mu(h, \mathcal{S}, \delta) \right\} \leq 0. \end{aligned}$$

Proof. Deferred to Appendix G.3. ■

For a given confidence δ , with probability at least $1 - \delta$, the bound is then valid for all (h, \mathcal{S}) belonging to a (reduced) set $\mathbb{Z}_\delta \subseteq \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m$. In other words, the bound always holds for a given hypothesis and learning sample $(h, \mathcal{S}) \in \mathbb{Z}_\delta$, and its value depends on these h and \mathcal{S} . The generality of our framework can thus generalize both uniform convergence and algorithmic dependent bounds as we see in the rest of this section.

7.5.2 Uniform Convergence Bounds

Uniform-convergence-based bounds were the first type of generalization bounds to be introduced, notably in VAPNIK and CHERVONENKIS (1971) using the VC-dimension as complexity. Other bounds were later developed based on the Gaussian/Rademacher complexity (BARTLETT and MENDELSON, 2002) instead. We recall the definition of this type of bounds encountered in Chapter 1.

Definition 1.3.2 (Uniform Convergence Bound). Let $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ a generalization gap. A uniform convergence bound is defined such that if for any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$, such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\sup_{h \in \mathbb{H}} \phi(R_{\mathcal{D}}^\ell(h), R_{\mathcal{S}}^\ell(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta, \quad (1.1)$$

where usually $\phi(R_{\mathcal{D}}^\ell(h), R_{\mathcal{S}}^\ell(h)) = R_{\mathcal{D}}^\ell(h) - R_{\mathcal{S}}^\ell(h)$.

Remember that this definition encompasses different complexity measures, such as $\Phi_u(\delta) = \text{rad}(\mathbb{H}) + \sqrt{\frac{1}{2m} \ln \frac{1}{\delta}}$ in Theorem 1.3.3, or $\Phi_u(\delta) = \sqrt{\frac{1}{m} 2\text{vc}(\mathbb{H}) \ln \frac{em}{\text{vc}(\mathbb{H})}} + \sqrt{\frac{1}{2m} \ln \frac{1}{\delta}}$ described in Theorem 1.3.2. For ease of comparison, we refine and reinterpret this type of bounds in a set-theoretic manner as follows. This result has been originally remarked by NAGARAJAN and KOLTER (2019b) (but not proved).

Proposition 7.5.2 (Set-theoretic View of Uniform Convergence Bounds). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.2. Under these conditions, with $\mathbb{Z}_u = \left\{ \mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m : \forall h \in \mathbb{H}, \phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h)) \leq \Phi_u(\delta) \right\}$, and $\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} [\mathbb{S} \in \mathbb{Z}_u] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (1.1)} &\iff \forall \mathbb{S} \in \mathbb{Z}_u, \forall h \in \mathbb{H}, \phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h)) \leq \Phi_u(\delta) \\ &\iff \sup_{\mathbb{S} \in \mathbb{Z}_u} \sup_{h \in \mathbb{H}} \left\{ \phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h)) \right\} \leq \Phi_u(\delta). \end{aligned}$$

Proof. Deferred to Appendix G.4. ■

Proposition 7.5.2 is, in fact, a reinterpretation of PAC generalization bounds by identifying the subset $\mathbb{Z}_u \subseteq (\mathbb{X} \times \mathbb{Y})^m$ for which the upper bound $\Phi_u(\delta)$ is valid. This highlights their worst-case nature: given a confidence δ , the generalization gap $\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h))$ is upper-bounded by a complexity measure $\Phi_u(\delta)$ for all $(h, \mathbb{S}) \in \mathbb{H} \times \mathbb{Z}_u$. To get a bound holding with probability at least $1 - \delta$, since the complexity $\Phi_u(\delta)$ does not depend on h or \mathbb{S} , the complexity has to upper-bound the worst hypothesis $h \in \mathbb{H}$ and the worst learning sample $\mathbb{S} \in \mathbb{Z}_u$. As a consequence, $\Phi_u(\delta)$ is lower-bounded by $\sup_{\mathbb{S} \in \mathbb{Z}_u} \sup_{h \in \mathbb{H}} \phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h))$. As we have seen in Proposition 7.5.1, our bound is more permissive than the uniform convergence bounds since the upper bound can depend on the learning sample \mathbb{S} and the hypothesis h . Hence, this dependence on \mathbb{S} and h allows us to retrieve the uniform convergence bounds with our framework. Indeed, from Theorem 7.3.1, we can obtain the following generalization bound.

Corollary 7.5.1 (Uniform Convergence Bound from Theorem 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap and assume that there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.2 such that $\Phi_u(\delta) \geq \ln \left[\frac{4}{\delta^2} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\phi(\mathbb{R}_{\mathcal{D}}(h'), \mathbb{R}_{\mathbb{S}'}(h'))) \right]$. For any \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\phi(\mathbb{R}_{\mathcal{D}}(h), \mathbb{R}_{\mathbb{S}}(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta.$$

Proof. Deferred to Appendix G.5. ■

Note that, to prove Corollary 7.5.1, we require an additional assumption: a lower-bound on $\Phi_u(\delta)$. When the generalization gap is $\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) = 2m[R_{\mathcal{D}}(h) - R_{\mathcal{S}}(h)]^2$ or $\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) = \text{kl}[R_{\mathcal{S}}(h) \| R_{\mathcal{D}}(h)]$, the lower bound is $\ln \frac{8\sqrt{m}}{\delta^2}$ (see Corollary 7.3.1) which is low enough to be a worst-case upper-bound. To sum up, our framework is general enough to retrieve classical uniform convergence bounds (under the mild assumption) such that the ones based on the Rademacher complexity (Definition 1.3.4) or the VC-Dimension (Definition 1.3.3). In practice, the sampling involved in the bound of Corollary 7.5.1 is not necessary: the bound holds for all hypothesis $h \in \mathbb{H}$ with high probability. More precisely, for $\Phi_u(h, \mathcal{S}, \delta) = \Phi_u(\delta)$, the set \mathbb{Z}_u in Proposition 7.5.2 can be seen as a subset of \mathbb{Z}_d since

$$\left\{ \mathcal{S} \in (\mathbb{X} \times \mathbb{Y})^m \mid \forall h \in \mathbb{H}, (\mathcal{S}, h) \in \mathbb{Z}_d \right\} = \mathbb{Z}_u.$$

Hence, for a well-behaved learning sample \mathcal{S} , *i.e.*, for $\mathcal{S} \in \mathbb{Z}_u$, we have that

$$\begin{aligned} \mathbb{E}_{h \in \rho_{\mathcal{S}}} \mathbb{I} \left[\phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \leq \Phi_u(\delta) \right] &= \mathbb{I} \left[\sup_{h \in \mathbb{H}} \phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \leq \Phi_u(\delta) \right] \\ &= 1, \end{aligned}$$

which gives a valid bound for all $h \in \mathbb{H}$ without sampling. This generality of this framework does not apply uniquely to these type of bounds. We can obtain a result similar for the algorithmic-dependent bounds that can be tighter than the uniform-convergence-based bounds.

7.5.3 Algorithmic-Dependent Bounds

The upper bound $\Phi_u(\delta)$ can generally be improved by considering algorithmic-dependent bounds (BOUSQUET and ELISSEEFF, 2002; XU and MANNOR, 2012). In this case, only the output $h_{\mathcal{S}}$ of a learning algorithm given \mathcal{S} is studied: we only bound the generalization gap $\phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}}))$ specific to $h_{\mathcal{S}}$ (here, $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathbb{X} \times \mathbb{Y})^m}$). The definition of such bounds encountered in Chapter 1 is recalled in the following.

Definition 1.3.5 (Algorithmic-dependent Generalization Bound). Let $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$ be a loss function and $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ a generalization gap. An algorithmic-dependent generalization bound is defined such that if for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, there exists a function $\Phi_a : (0, 1] \rightarrow \mathbb{R}$, such that for any $\delta \in (0, 1]$ we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\phi(R_{\mathcal{D}}^{\ell}(h_{\mathcal{S}}), R_{\mathcal{S}}^{\ell}(h_{\mathcal{S}})) \leq \Phi_a(\delta) \right] \geq 1 - \delta, \quad (1.2)$$

where usually $\phi(R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{S}}^{\ell}(h)) = R_{\mathcal{D}}^{\ell}(h) - R_{\mathcal{S}}^{\ell}(h)$ and $h_{\mathcal{S}}$ is the hypothesis learned from an algorithm with $\mathcal{S} \sim \mathcal{D}^m$.

Similarly to the uniform convergence bounds, these bounds can be reformulated through a similar set-theoretic lens stated in the following proposition.

Proposition 7.5.3 (Set-theoretic View of Algorithmic Dependent Bounds). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_{\mathbf{a}} : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.5. Under these conditions, with $\mathbb{Z}_{\mathbf{a}} = \left\{ \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m : \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}})) \leq \Phi_{\mathbf{a}}(\delta) \right\}$ and $\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m}[\mathcal{S} \in \mathbb{Z}_{\mathbf{a}}] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (1.2)} &\iff \forall \mathcal{S} \in \mathbb{Z}_{\mathbf{a}}, \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}})) \leq \Phi_{\mathbf{a}}(\delta) \\ &\iff \sup_{\mathcal{S} \in \mathbb{Z}_{\mathbf{a}}} \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}})) \leq \Phi_{\mathbf{a}}(\delta). \end{aligned}$$

Proof. Deferred to Appendix G.6. ■

Since the upper bound $\Phi_{\mathbf{a}}(\delta)$ is at least $\sup_{\mathcal{S} \in \mathbb{Z}_{\mathbf{a}}} \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}}))$, this result has the potential to lead to tighter guarantees than the uniform convergence ones. For example, when $\mathbb{H} = \{h_{\mathcal{S}}\}_{\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m}$ is an algorithmic-dependent hypothesis set and $\mathbb{Z}_{\mathbf{a}} \subseteq \mathbb{Z}_{\mathbf{u}}$. The complexity measure $\Phi_{\mathbf{a}}(\delta)$ can potentially be smaller than $\Phi_{\mathbf{u}}(\delta)$ since we have the inequality

$$\begin{aligned} \sup_{\mathcal{S} \in \mathbb{Z}_{\mathbf{a}}} \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}})) &\leq \sup_{\mathcal{S} \in \mathbb{Z}_{\mathbf{u}}} \phi(R_{\mathcal{D}}(h_{\mathcal{S}}), R_{\mathcal{S}}(h_{\mathcal{S}})) \\ &\leq \sup_{\mathcal{S} \in \mathbb{Z}_{\mathbf{u}}} \sup_{h \in \mathbb{H}} \phi(R_{\mathcal{D}}(h), R_{\mathcal{S}}(h)) \\ &\leq \Phi_{\mathbf{u}}(\delta). \end{aligned}$$

Even though these type of bounds can be tighter, it is still not as permissive as our framework. Indeed, the upper bound $\Phi_{\mathbf{a}}(\delta)$ is a constant *w.r.t.* the hypothesis and the learning sample (like the uniform convergence bounds). Hence, since our bound can depend on the learning sample \mathcal{S} and the hypothesis h , we retrieve the algorithmic-dependent bounds illustrating the generality of our framework (similarly to Corollary 7.5.1). The result is in the following corollary.

Corollary 7.5.2 (Algorithmic-dependent Bound from Theorem 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap and assume that there exists a function $\Phi_{\mathbf{a}} : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.5 such that $\Phi_{\mathbf{a}}(\delta) \geq$

7.6. Conclusion and Summary

$\ln \left[\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\phi(\mathbf{R}_{\mathcal{D}}(h'), \mathbf{R}_{\mathcal{S}'}(h'))) \right]$. For any \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mathbf{a}}(\delta) \right] \geq 1 - \delta.$$

Proof. Deferred to Appendix G.7. ■

Compared to the bounds of Definition 1.3.5, Corollary 7.5.2 still involves the expectation over the hypotheses. Hopefully, the bound holds with high probability for the data-dependent hypothesis $h_{\mathcal{S}}$ (see Proposition 7.5.3). Hence, when using Corollary 7.5.2's bound the sampling is not necessary since we can consider the bound only for the hypothesis of interest, *i.e.*, $h_{\mathcal{S}}$ for all $\mathcal{S} \in \mathbb{Z}_{\mathbf{a}}$ which holds with high probability. In other words, when $\Phi_{\mu}(h, \mathcal{S}, \delta) = \Phi_{\mathbf{a}}(\delta)$, the set $\mathbb{Z}_{\mathbf{a}}$ in Proposition 7.5.3 can be seen as a subset of $\mathbb{Z}_{\mathbf{d}}$ since

$$\left\{ \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m \mid (\mathcal{S}, h_{\mathcal{S}}) \in \mathbb{Z}_{\mathbf{d}} \right\} = \mathbb{Z}_{\mathbf{a}}.$$

In summary, the framework proposed in this chapter is powerful enough to cover uniform-convergence-based bound and algorithm-dependent bounds with the integration of a complexity measure. To the best of our knowledge, this has not been identified before and we think this is something novel.

7.6 Conclusion and Summary

In this chapter, we provide a novel generalization bound that involves arbitrary complexity measures, unlike classical learning theory frameworks (for which the complexity is imposed by the framework itself). These measures incorporate a data and model dependent function, that allow us to generalize the previous framework introduced in the literature (see Section 1.3). Importantly, to the best of our knowledge, our framework provides for the first time theoretical guarantees for the many arbitrary complexity measures used in practice in machine learning, *e.g.*, for regularization purposes.

The limitation of this work is clearly that the hypothesis is obtained from a distribution difficult to use, namely the Gibbs distribution. Indeed, one sampling from the Gibbs distribution is performed by an algorithm such as Algorithm 7.1. Hopefully, the generality of this framework allows to avoid the sampling if we consider uniform-convergence-type bounds as in Corollary 7.5.1. We can easily imagine continuing in

this research direction. For instance, we can try to get rid of the supremum *w.r.t.* to the hypothesis set \mathbb{H} in the Rademacher complexity (Definition 1.3.4) since it is hard to compute.

We hope that our results foster research in the topic and the development of new complexity measures for specific neural network architectures and for specific learning tasks. Indeed, we believe that this work paves the way to new research directions that try to bridge the statistical learning theory and the practice. Indeed, finding a good complexity measure becomes a practical matter since any complexity measure can be integrated in our framework.

In general, this thesis explores the disintegrated bounds in order to explain better the generalization of over-parameterized models that is largely misunderstood. We believe that this type of bounds is not the only promising type of bounds that can explain the generalization phenomenon. In Part IV, we give an idea of research direction to explore other type of generalization bounds.

PART IV

Conclusion and Perspectives

CONCLUSION AND PERSPECTIVES

Conclusion

This thesis mainly derives self-bounding algorithms that learn a model minimizing a (disintegrated) PAC-Bayesian generalization bound. This type of algorithm has received little attention in the machine learning literature and we propose some contributions in various contexts.

Indeed, Part II is dedicated to deriving self-bounding algorithms in the context of majority vote classifiers. In Chapters 3 and 4, we derived self-bounding algorithms for the majority vote classifier in two different settings: the adversarial robustness and the classical supervised setting. More precisely, Chapter 3's self-bounding algorithms robustify the majority votes against small perturbations. While Chapter 4 minimizes the majority vote's true risk through the PAC-Bayesian C-Bounds considered as challenging to optimize (LORENZEN *et al.*, 2019; MASEGOSA *et al.*, 2020). However, as shown in Chapter 5, the majority vote's self-bounding algorithms considered, *e.g.*, in Chapter 4, do not minimize tight generalization bounds on the true risk, even for simple tasks. Hence, to overcome this drawback, Chapter 5 introduces the stochastic majority vote, which samples a majority vote for each prediction. Considering such a majority vote allows us to obtain tight generalization bounds. Additionally, we derive a self-bounding algorithm that directly minimizes the risk of the stochastic majority vote in this context. However, the risk of a stochastic model is the expected risk of the hypotheses, which requires certain assumptions to be computed while we may be only interested in assessing the behavior of only one hypothesis in some situations. Hence, to overcome this drawback, we consider in Part III the *disintegrated* PAC-Bayesian bounds. Chapter 6 provides new bounds based on the Rényi divergence that are more easily optimizable (for self-bounding algorithms) than the ones of the literature (*i.e.*, BLANCHARD and FLEURET, 2007; CATONI, 2007; RIVASPLATA *et al.*, 2020). Even though RIVASPLATA *et al.* (2020)'s bound is not easily optimizable, it is a starting point to derive new generalizations bounds. Indeed, in the last contribution (Chapter 7), we leverage RIVASPLATA *et al.* (2020)'s disintegrated framework to derive generalization bounds with arbitrary complexity measures. Such work is fundamental in statistical learning theory: to the best of our knowledge, we are the first to provide generalization bounds that integrate complexity measures that can be defined by the user. This work allows the machine learning community to consider new generalization bounds by defining a new complexity measure. Hence, new works can focus on developing new complexity measures to understand better the generalization phenomenon.

Perspectives

We present several perspectives following the contributions of this thesis.

Perspectives on the Adversarial Robustness Setting

As recalled in Chapter 3, in the adversarial robustness setting, we aim to make the model robust to small perturbations in the input. Indeed, we must ensure that the model does not radically change its prediction for a slight change in the input. To do so, we consider that the model's output must not change in a ball of a given radius. This new constraint on the input actually creates a new unknown data distribution that is close, in some sense, to the original unknown data distribution.

On the other hand, the transfer learning/domain adaptation³ consider two unknown data distributions: a source (*i.e.*, the original) and a target (*i.e.*, a new) distribution. In this setting, the model learned to solve a task (represented by the source distribution) is adapted to solve a new task (represented by the target distribution). In some transfer learning scenarios, we assume that we have access to the labels and the inputs obtained from the target distribution, while in unsupervised domain adaptation, only the inputs are considered. In these two settings, the true risk on the target distribution can be upper-bounded with a generalization bound (see *e.g.*, BEN-DAVID *et al.*, 2010; GALANTI *et al.*, 2016; MCNAMARA and BALCAN, 2017; GERMAIN *et al.*, 2020).

Besides, it is known that domain adaptation and adversarial robustness are related: unlabeled examples (considered in domain adaptation) can be used to improve the adversarial robustness (ALAYRAC *et al.*, 2019; CARMON *et al.*, 2019; DENG *et al.*, 2021). As a perspective, we propose to investigate the link between these two settings from a theoretical viewpoint. First, we could explore the connection between the original distribution and the new data distribution induced by the adversarial robustness that can be respectively seen as a source and a target distribution in transfer learning. Then, this connection may help to leverage transfer learning/domain adaptation generalization bounds to obtain guarantees for the adversarial robustness setting. The new guarantees might serve to get self-bounding algorithms that (i) detect out-of-distribution examples⁴ and (ii) robustify machine learning models.

Extending the Majority Vote

In Part II, we consider that the set of voters in the PAC-Bayesian majority vote is fixed. Hence, only the weights of the majority vote are adapted to fit the examples.

³We refer the reader to REDKO *et al.* (2019, 2020) for an introduction on domain adaption.

⁴The examples that are not probable in a given distribution are called out-of-distribution examples.

Alternatively, in the (Gradient) Boosting framework (FREUND and SCHAPIRE, 1996; FRIEDMAN, 2001), the voters are greedily learned one by one. Moreover, in bagging (BREIMAN, 1996) and random forest (BREIMAN, 2001), no weights are learned while the models are learned separately. For the Support Vector Machine (GRAEPEL *et al.*, 2005) that can be interpreted as a majority vote, the voters are fixed before learning the weights by choosing a kernel. As we can remark in these approaches, the voters and the weights are not learned together. This appears as a limitation since learning the weights and the voters in an end-to-end way can offer a better accurate majority vote. Hence, one bottleneck has to be overcome: deriving differentiable voters such as differentiable decision stumps. By doing so, we may improve the voters' diversity while limiting the voters' complexity. Moreover, the disintegrated PAC-Bayesian framework (developed, *e.g.*, in Part III) may be leveraged to derive generalization guarantees for majority votes that depend on the full learning sample. Again, new generalization bounds can be further used to derive self-bounding algorithms.

Self-bounding and Optimization Algorithms

The optimization algorithms are key to obtain a good classifier in self-bounding algorithms. Specifically in Chapters 4 and 5, we use an optimization algorithm that tune automatically the learning rate, namely COCOB (ORABONA and TOMMASI, 2017). This approach, belonging to the parameter-free algorithms⁵, is interesting in machine learning because it has a clear advantage: there is no need to tune the learning rate. Hence, the parameter-free algorithms could facilitate the use of machine learning approaches for practitioners. However, we believe that more hyper-parameters can be tuned automatically in parameter-free optimization algorithms such as the batch size, which offers interesting research perspectives.

One idea to derive new parameter-free algorithms is to take inspiration from the federated learning setting.⁶ It considers different clients that learn collaboratively in a machine learning model; each client has its own learning sample and does not necessarily share it. For instance, to learn the model, each client has its own local model and runs an optimization algorithm, to obtain new weights. The new weights of each local model are aggregated to obtain a global model finally, without exchanging the data; see, *e.g.*, the FedAvg algorithm (see MCMAHAN *et al.*, 2017). A modification must be made to FedAvg to obtain a parameter-free algorithm since each client runs an algorithm with different values of hyper-parameters. The aggregation of the weights can take different forms, such as a convex combination. With this latter type of aggregation, the PAC-Bayesian theory might be helpful to obtain convergence guarantees.

⁵We refer the reader to the ICML 2020 tutorial on Parameter-free online optimization for more details on the parameter-free algorithms.

⁶Federated learning is a sub-field of machine learning; see KAIROUZ *et al.* (2021) for a survey.

Towards a New Type of Generalization Bounds

The PAC-Bayesian theory considers that the data and the models are respectively sampled from two probability distributions: the unknown distribution and the posterior distribution. While it can be convenient to derive generalization guarantees on a single model sampled from the posterior distribution, we are usually interested in a model that is not necessarily sampled.

Instead of considering the posterior distribution on the models, one could consider a distribution on the label set conditioned on the input. If this distribution is somehow learned from the learning sample, it can be seen as a machine learning model. Thanks to this distribution, we could derive generalization bounds on the expected loss when the labels are sampled from the new data-dependent distribution (associated with a classifier). We can for example hope to obtain a bound dependent on the mutual information between the predictions and the labels. Roughly speaking, mutual information measures how much information on the labels is contained in the predictions. Hence, it can be seen as a complexity measure of the data-dependent distribution (representing the classifier). For instance, this quantity has been considered in the information bottleneck framework of TISHBY *et al.* (2000). Besides, the training of neural networks has been studied through this framework (see SHWARTZ-ZIV and TISHBY, 2017).

PART V

Appendix

SOME MATHEMATICAL TOOLS



A.1 Jensen's Inequality

Theorem A.1.1 (JENSEN's Inequality). Let $X \in \mathcal{X}$ a random variable following a probability distribution \mathcal{X} with $f : \mathcal{X} \rightarrow \mathbb{R}$ a measurable convex function, we have

$$f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right) \leq \mathbb{E}_{X \sim \mathcal{X}}[f(X)].$$

Proof. Since $f()$ is a convex function, the following inequality holds, i.e., we have

$$\forall X' \in \mathcal{X}, \quad a\left(X' - \mathbb{E}_{X \sim \mathcal{X}}[X]\right) \leq f(X') - f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right),$$

where a is the tangent's slope. By taking the expectation to both sides of the inequality, we have

$$\underbrace{a\left(\mathbb{E}_{X \sim \mathcal{X}}[X] - \mathbb{E}_{X \sim \mathcal{X}}[X]\right)}_{=0} \leq \mathbb{E}_{X \sim \mathcal{X}}[f(X)] - f\left(\mathbb{E}_{X \sim \mathcal{X}}[X]\right).$$

Hence, by rearranging the terms, we prove the claimed result. ■

A.2 Markov's Inequality

Theorem A.2.1 (MARKOV's Inequality). Let $X \in \mathcal{X}$ a non-negative random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}}[X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}}[X]}{\tau}.$$

Proof. First of all, remark that we have the following inequality for any $X \in \mathbb{X}$

$$\tau \mathbb{I}[X \geq \tau] \leq X \mathbb{I}[X \geq \tau] \leq X. \quad (\text{A.1})$$

Indeed, on the one hand, if $X < \tau$, $\mathbb{I}[X \geq \tau] = 0$, the inequality holds trivially. On the other hand, if $X \geq \tau$, $\mathbb{I}[X \geq \tau] = 1$ and the inequality becomes $\tau \leq X$, which is true. By taking the expectation of Equation (A.1), we have

$$\mathbb{E}_{X \sim \mathcal{X}} \left[\tau \mathbb{I}[X \geq \tau] \right] \leq \mathbb{E}_{X \sim \mathcal{X}} \left[X \right].$$

From the fact that the expectation of a constant is the constant and by definition of the probability, we have

$$\tau \mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \mathbb{E}_{X \sim \mathcal{X}} [X] \iff \mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X]}{\tau},$$

which is the desired result. ■

A.3 2nd Order Markov's Inequality

Theorem A.3.1 (2nd Order MARKOV's Inequality). Let X a non-negative random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X^2]}{\tau^2}.$$

Proof. We apply MARKOV's inequality (Theorem A.2.1) to have

$$\mathbb{P}_{X \sim \mathcal{X}} [X^2 \geq \tau^2] \leq \frac{\mathbb{E}_{X \sim \mathcal{X}} [X^2]}{\tau^2}.$$

Moreover, since $\mathbb{I}[X \geq \tau] = \mathbb{I}[X^2 \geq \tau^2]$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} [X \geq \tau] = \mathbb{P}_{X \sim \mathcal{X}} [X^2 \geq \tau^2],$$

which proves the desired result. ■

A.4 Chebyshev-Cantelli Inequality

Theorem A.4.1 (CHEBYSHEV-CANTELLI Inequality). Let X a random variable following a probability distribution \mathcal{X} and $\tau > 0$, we have

$$\mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] \leq \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\mathbb{V}_{X' \sim \mathcal{X}} X' + \tau^2}.$$

Proof. First of all, remark that we have

$$\begin{aligned} \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] &= \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \geq \tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right] \\ &\leq \mathbb{P}_{X \sim \mathcal{X}} \left[\left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2 \geq \left[\tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2 \right], \end{aligned} \tag{A.2}$$

where $\mathbb{V}_{X \sim \mathcal{X}} X$ is the variance of the random variable $X \sim \mathcal{X}$. From Equation (A.2) and MARKOV's Inequality (Theorem A.2.1), we can deduce that

$$\begin{aligned} \mathbb{P}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' \geq \tau \right] &\leq \frac{\mathbb{E}_{X \sim \mathcal{X}} \left[X - \mathbb{E}_{X' \sim \mathcal{X}} X' + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2}{\left[\tau + \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\tau} \right]^2} \\ &= \frac{\mathbb{V}_{X' \sim \mathcal{X}} X'}{\mathbb{V}_{X' \sim \mathcal{X}} X' + \tau^2}. \end{aligned}$$

■

A.5 Hölder's Inequality

In order to prove HÖLDER's inequality, we first prove the following lemma.

Lemma A.5.1 (YOUNG's Inequality). For any $\alpha > 1$ and $\beta > 1$ such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, for any $a \geq 0$ and $b \geq 0$, we have

$$ab \leq \frac{a^\alpha}{\alpha} + \frac{b^\beta}{\beta}.$$

Proof. We first develop $\ln [ab]$ and we apply JENSEN's inequality (Theorem A.1.1) since the logarithm is concave and $\frac{1}{\alpha} + \frac{1}{\beta} = 1$. Indeed, we have

$$\ln [ab] = \ln a + \ln b = \frac{\alpha}{\alpha} \ln a + \frac{\beta}{\beta} \ln b = \frac{\ln a^\alpha}{\alpha} + \frac{\ln b^\beta}{\beta} \leq \ln \left[\frac{a^\alpha}{\alpha} + \frac{b^\beta}{\beta} \right].$$

Then, we take the exponential to both sides of the inequality and we are done. ■

We are now ready to prove HÖLDER's inequality.

Theorem A.5.1 (HÖLDER's Inequality). For any measurable function $f()$ and $g()$, for any $\alpha > 1$ and $\beta > 1$ such that $\frac{1}{\alpha} + \frac{1}{\beta} = 1$, we have

$$\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| \leq \left[\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha \right]^{\frac{1}{\alpha}} \left[\mathbb{E}_{X \sim \mathcal{X}} |g(X)|^\beta \right]^{\frac{1}{\beta}}.$$

Proof. For convenience of notation, let $\|f\|_\alpha = \left[\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha \right]^{\frac{1}{\alpha}}$ and $\|g\|_\beta = \left[\mathbb{E}_{X \sim \mathcal{X}} |g(X)|^\beta \right]^{\frac{1}{\beta}}$. If $\|f\|_\alpha = 0$ or $\|g\|_\beta = 0$, then $\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| = 0$, hence, the inequality holds in this case. Then for $\|f\|_\alpha > 0$ and $\|g\|_\beta > 0$, we upper-bound the term $\frac{|f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta}$ with YOUNG's inequality (Lemma A.5.1), i.e., we have

$$\frac{|f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} \leq \frac{|f(X)|^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{|f(X)|^\beta}{\beta \|f\|_\beta^\beta}.$$

By taking the expectation *w.r.t.* $X \sim \mathcal{X}$, we have

$$\begin{aligned} \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} &\leq \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)|^\beta}{\beta \|f\|_\beta^\beta} \\ &= \frac{\|f\|_\alpha^\alpha}{\alpha \|f\|_\alpha^\alpha} + \frac{\|f\|_\beta^\beta}{\beta \|f\|_\beta^\beta} \\ &= \frac{1}{\alpha} + \frac{1}{\beta} \\ &= 1. \end{aligned}$$

A.5. HÖLDER's Inequality

This concludes the proof since

$$\frac{\mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)|}{\|f\|_\alpha \|g\|_\beta} \leq 1 \iff \mathbb{E}_{X \sim \mathcal{X}} |f(X)g(X)| \leq \|f\|_\alpha \|g\|_\beta.$$

■

APPENDIX OF CHAPTER 2

B

B.1 Proof of Theorem 2.2.1

Theorem 2.2.1 (Risk Upper Bound Based on the Gibbs Risk). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq 2r_{\mathcal{D}'}(\rho). \quad (2.2)$$

Proof. The proof is given by GERMAIN *et al.* (2015). First of all, remark that

$$\mathbb{I}[\text{MV}_\rho(\mathbf{x}) \neq y] \leq \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0].$$

Hence, by taking the expectation, we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0].$$

From MARKOV's inequality (Theorem A.2.1), we have

$$\begin{aligned} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0] &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[1 - \widehat{m}_\rho(\mathbf{x}, y) \geq 1] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left[1 - 2\left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y] - \frac{1}{2}\right] \geq 1\right] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left[1 - \left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y]\right] \geq \frac{1}{2}\right] \\ &\leq 2 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left[1 - \mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y]\right] \\ &= 2r_{\mathcal{D}'}(\rho). \end{aligned}$$

■

B.2 Proof of Theorem 2.2.2

Theorem 2.2.2 (Risk Upper Bound Based on the Joint Error). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq 4e_{\mathcal{D}'}(\rho). \quad (2.4)$$

Proof. The proof is given by MASEGOSA *et al.* (2020). First of all, remark that

$$\mathbb{I}[\text{MV}_\rho(\mathbf{x}) \neq y] \leq \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0].$$

Hence, by taking the expectation, we have

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0].$$

From MARKOV's inequality (Theorem A.3.1), we have

$$\begin{aligned} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0] &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}[1 - \widehat{m}_\rho(\mathbf{x}, y) \geq 0] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left[1 - 2\left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y] - \frac{1}{2}\right] \geq 1\right] \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left[1 - \left[\mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y]\right] \geq \frac{1}{2}\right] \\ &\leq 4 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'}\left(\left[1 - \mathbb{P}_{h \sim \rho}[h(\mathbf{x}) = y]\right]^2\right) \\ &= 4e_{\mathcal{D}'}(\rho). \end{aligned}$$

■

B.3 Proof of Theorem 2.2.3

Theorem 2.2.3 (The C-Bound). For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , if

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'}[\widehat{m}_\rho(\mathbf{x}, y) > 0] \iff r_{\mathcal{D}'}(\rho) < \frac{1}{2} \iff 2e_{\mathcal{D}'}(\rho) + d_{\mathcal{D}'}(\rho) < 1,$$

we have

$$R_{\mathcal{D}' }(\text{MV}_\rho) \leq 1 - \frac{\left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)]\right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \quad (2.7)$$

$$= 1 - \frac{(1 - 2r_{\mathcal{D}'}(\rho))^2}{1 - 2d_{\mathcal{D}'}(\rho)} \quad (2.8)$$

$$= 1 - \frac{\left(1 - [2e_{\mathcal{D}'}(\rho) + d_{\mathcal{D}'}(\rho)]\right)^2}{1 - 2d_{\mathcal{D}'}(\rho)} \quad (2.9)$$

$$= C_{\mathcal{D}'}(\rho).$$

Proof. To prove Equation (2.7), we start from the definition of $R_{\mathcal{D}' }(\text{MV}_\rho)$ to have

$$\begin{aligned} R_{\mathcal{D}' }(\text{MV}_\rho) &\leq \mathbb{P}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y) \leq 0) \\ &= \mathbb{P}_{(\mathbf{x},y) \sim \mathcal{D}'} \left(-\widehat{m}_\rho(\mathbf{x}, y) + \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \geq \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \right) \\ &\leq \frac{\mathbb{V}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))}{\mathbb{V}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y)) + \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2} \\ &= \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2 - \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2 - \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2 + \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2} \\ &= \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2 - \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2} \\ &= 1 - \frac{\left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)\right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2}, \end{aligned}$$

where the second inequality comes from CHEBYSHEV-CANTELLI's inequality (Theorem A.4.1) and $\mathbb{V}_{A \sim \mathcal{A}}(A)$ is the variance of the random variable $A \sim \mathcal{A}$. Equation (2.9) is obtained by rewriting Equation (2.7) with Equations (2.1) and (2.5). Indeed, we have

$$\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) = 1 - 2r_{\mathcal{D}'}(\rho) \quad \text{and} \quad \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)^2 = 1 - 2d_{\mathcal{D}'}(\rho),$$

which gives Equation (2.8).

Moreover, thanks to Equation (2.6), we can rewrite the Gibbs risk as

$$r_{\mathcal{D}'}(\rho) = \left[e_{\mathcal{D}'}(\rho) + \frac{1}{2}d_{\mathcal{D}'}(\rho) \right],$$

which allows us to obtain Equation (2.9) by rewriting Equation (2.8). \blacksquare

B.4 Proof of Theorem 2.2.4

Theorem 2.2.4 (Relationship between Theorems 2.2.1 to 2.2.3). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any voters set \mathbb{H} , for any distribution ρ on \mathbb{H} , if $r_{\mathcal{D}'}(\rho) < \frac{1}{2}$ (i.e., $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) > 0$), we have

$$(i) \ R_{\mathcal{D}'}(\text{MV}_\rho) \leq C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho), \text{ if } r_{\mathcal{D}'}(\rho) \leq d_{\mathcal{D}'}(\rho),$$

$$(ii) \ R_{\mathcal{D}'}(\text{MV}_\rho) \leq 2r_{\mathcal{D}'}(\rho) \leq C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho), \text{ otherwise.}$$

Proof. We first prove that $C_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho)$ is equivalent to $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \geq \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)]^2$, i.e., we have

$$\begin{aligned} C_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho) \\ \iff 1 - \frac{\left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} &\leq 1 - \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \\ \iff \frac{\left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \right)^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} &\geq \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \\ \iff \left(\frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \right)^2 &\geq \left(\frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \right) \left(\frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \right) \\ \iff \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} &\geq \frac{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y)}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2}. \end{aligned}$$

Then, we prove that $4e_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho)$ is equivalent to $\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \geq$

B.5. About the KL Divergence

$\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)]^2$, i.e., we have

$$\begin{aligned}
 & 4e_{\mathcal{D}'}(\rho) \leq 2r_{\mathcal{D}'}(\rho) \\
 \iff & 1 - 2 \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) + \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 \leq 1 - \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \\
 \iff & 2 \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) - \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 \geq \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \\
 \iff & \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) \geq \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)]^2.
 \end{aligned}$$

Additionally, we prove that $C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho)$, i.e., we have

$$\begin{aligned}
 & C_{\mathcal{D}'}(\rho) \leq 4e_{\mathcal{D}'}(\rho) \\
 \iff & 1 - \frac{(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y))^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \leq 1 - 2 \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) + \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 \\
 \iff & \frac{(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y))^2}{\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2} \geq 2 \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) - \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 \\
 \iff & \left(\mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} \widehat{m}_\rho(\mathbf{x}, y) - \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 \right)^2 \geq 0.
 \end{aligned}$$

Finally, by merging the three equivalence, we obtain the claimed result. ■

B.5 About the KL Divergence

B.5.1 Basic Properties

The KL divergence has the following properties:

- (1) It is positive, i.e., we have $\text{KL}(\rho \parallel \pi) \geq 0$ for all $\rho \in \mathbb{M}(\mathbb{H})$ and $\pi \in \mathbb{M}^*(\mathbb{H})$.
- (2) We have $\text{KL}(\pi \parallel \pi) = 0$ for all $\pi \in \mathbb{M}(\mathbb{H})$.
- (3) In general, it is not symmetric: $\text{KL}(\rho \parallel \pi) \neq \text{KL}(\pi \parallel \rho)$ for all $\rho, \pi \in \mathbb{M}(\mathbb{H})$.

Proof. We prove the points (1), (2) and (3) separately.

Concerning (1). Since $-\ln(\cdot)$ is convex, we have from JENSEN's inequality (Theorem A.1.1)

$$\text{KL}(\rho \parallel \pi) = \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} = \mathbb{E}_{h \sim \rho} \left[-\ln \frac{\pi(h)}{\rho(h)} \right] \geq -\ln \left[\mathbb{E}_{h \sim \rho} \frac{\pi(h)}{\rho(h)} \right] = 0.$$

Concerning (2). The property follows directly by developing the KL divergence, i.e., we have

$$\text{KL}(\pi\|\pi) = \mathbb{E}_{h\sim\pi} \ln \frac{\pi(h)}{\pi(h)} = \mathbb{E}_{h\sim\pi} \ln(1) = 0.$$

Concerning (3). For example, we have $\text{KL}(\mathcal{B}(0.1)\|\mathcal{B}(0.5)) \neq \text{KL}(\mathcal{B}(0.5)\|\mathcal{B}(0.1))$, where $\mathcal{B}(p)$ is a Bernoulli distribution with bias p . ■

Moreover, since we have (from l'Hôpital's rule) $\lim_{x\rightarrow 0^+} x \ln x = 0$, we adopt several conventions. Indeed, we consider that (i) $\rho(h) \ln \frac{\rho(h)}{\pi(h)} = 0$ whenever $\rho(h) = 0$ and $\pi(h) \geq 0$ and (ii) if $\pi(h) = 0$ and $\rho(h) > 0$, $\rho(h) \ln \frac{\rho(h)}{\pi(h)} = +\infty$ (which implies that $\text{KL}(\rho\|\pi) = +\infty$).

B.5.2 Joint Convexity

The following proposition shows that the KL divergence is jointly convex.

Proposition B.5.1. For any pairs $(\rho_1, \rho_2) \in \mathbb{M}(\mathbb{H})^2$ and $(\pi_1, \pi_2) \in \mathbb{M}^*(\mathbb{H})^2$ of distributions, we have for all $\lambda \in [0, 1]$

$$\text{KL}(\lambda\rho_1 + (1-\lambda)\rho_2\|\lambda\pi_1 + (1-\lambda)\pi_2) \leq \lambda \text{KL}(\rho_1\|\pi_1) + (1-\lambda) \text{KL}(\rho_2\|\pi_2).$$

In order to provide a proof for Proposition B.5.1. We need the log-sum inequality; a proof (based on COVER and THOMAS (2006, Theorem 2.7.1)) is given bellow.

Lemma B.5.1 (Log-sum inequality). For any strictly positive reals $p_1, \dots, p_n \geq 0$ and $q_1, \dots, q_n \geq 0$, we have

$$\left[\sum_{i=1}^n q_i \right] \ln \frac{\sum_{i=1}^n q_i}{\sum_{i=1}^n p_i} \leq \sum_{i=1}^n q_i \ln \frac{q_i}{p_i}.$$

Proof. First of all, note that $f(x) = x \ln x$ is convex w.r.t. $x \in \mathbb{R}_*^+$ since $\frac{\partial f}{\partial x}(x) = \frac{1}{x}$

B.5. About the KL Divergence

for all $x \in \mathbb{R}_*^+$. Then, from JENSEN's inequality (Theorem A.1.1) we have

$$\begin{aligned} f\left(\sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} \frac{q_i}{p_i}\right) &\leq \sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} f\left(\frac{q_i}{p_i}\right) \\ \Leftrightarrow \|\mathbf{p}\|_1 f\left(\sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} \frac{q_i}{p_i}\right) &\leq \|\mathbf{p}\|_1 \sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} f\left(\frac{q_i}{p_i}\right), \end{aligned}$$

where $\|\mathbf{p}\|_1 = \sum_{i=1}^n p_i$. Then, we can develop right-hand side of the inequality, i.e., we have

$$\|\mathbf{p}\|_1 f\left(\sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} \frac{q_i}{p_i}\right) = \|\mathbf{p}\|_1 f\left(\frac{\|\mathbf{q}\|_1}{\|\mathbf{p}\|_1}\right) = \|\mathbf{q}\|_1 \ln \frac{\|\mathbf{q}\|_1}{\|\mathbf{p}\|_1} = \left[\sum_{i=1}^n q_i\right] \ln \frac{\sum_{i=1}^n q_i}{\sum_{i=1}^n p_i},$$

where $\|\mathbf{q}\|_1 = \sum_{i=1}^n q_i$. Similarly for the left-hand side, we have

$$\|\mathbf{p}\|_1 \sum_{i=1}^n \frac{p_i}{\|\mathbf{p}\|_1} f\left(\frac{q_i}{p_i}\right) = \sum_{i=1}^n p_i f\left(\frac{q_i}{p_i}\right) = \sum_{i=1}^n q_i \ln \frac{q_i}{p_i}.$$

■

We are now able to prove Proposition B.5.1 based on the proof of COVER and THOMAS (2006, Theorem 2.7.2).

Proof of Proposition B.5.1. From the log-sum inequality (Lemma B.5.1) with $q_1 = \lambda\rho_1(h)$, $q_2 = (1-\lambda)\rho_2(h)$, $p_1 = \lambda\pi_1(h)$, $p_2 = (1-\lambda)\pi_2(h)$, we have

$$\begin{aligned} [\lambda\rho_1(h) + (1-\lambda)\rho_2(h)] \ln \left[\frac{\lambda\rho_1(h) + (1-\lambda)\rho_2(h)}{\lambda\pi_1(h) + (1-\lambda)\pi_2(h)} \right] \\ \leq \lambda\rho_1(h) \ln \frac{\rho_1(h)}{\pi_1(h)} + (1-\lambda)\rho_2(h) \ln \frac{\rho_2(h)}{\pi_2(h)}. \end{aligned}$$

Hence, by integrating over all h , gives us the desired result. ■

B.5.3 Pinsker's Inequality

We present a special case of PINSKER's inequality when we deal with two Bernoulli distributions. The presented proof is due to WU (2020) (see also CANONNE (2022)).

Theorem B.5.1 (Pinsker's inequality). For any $p \in [0, 1]$ and $q \in [0, 1]$, we have

$$2(q-p)^2 \leq \text{kl}(q||p),$$

where $\text{kl}(q||p) \triangleq q \ln \frac{q}{p} + (1-q) \ln \frac{1-q}{1-p}$.

Proof. For any $p \in \{0, 1\}$ and $q \in \{0, 1\}$, we can easily verify that the inequality holds. Then, with $p \in (0, 1)$ and $q \in (0, 1)$, we have from the fundamental theorem of calculus

$$\text{kl}(q\|p) = f(q) - f(p) = \int_q^p \frac{\partial f}{\partial x}(x) dx,$$

where $f(x) = q \ln x + (1 - q) \ln(1 - x)$. Hence, we have

$$\int_q^p \frac{\partial f}{\partial x}(x) dx = \int_q^p \frac{(q - x)}{(1 - x)x} dx \leq 4 \int_q^p (q - x) dx = 4 \frac{1}{2} (q - p)^2 = 2(q - p)^2.$$

■

B.6 Proof of Theorem 2.3.1

Theorem 2.3.1 (General PAC-Bayesian Bound of GERMAIN *et al.* (2009)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \leq \text{KL}(\rho\|\pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S}')} \right) \right] \geq 1 - \delta,$$

where $\text{KL}(\rho\|\pi) = \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)}$ is the Kullback-Leibler (KL) divergence between the distributions ρ and π .

Proof. We start by developing $\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S})$, *i.e.*, we obtain

$$\begin{aligned} \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) &= \mathbb{E}_{h \sim \rho} \ln [\exp(\varphi(h, \mathcal{S}))] \\ &= \mathbb{E}_{h \sim \rho} \ln \left[\frac{\rho(h) \pi(h)}{\pi(h) \rho(h)} \exp(\varphi(h, \mathcal{S})) \right] \\ &= \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} + \mathbb{E}_{h' \sim \rho} \ln \left[\frac{\pi(h')}{\rho(h')} \exp(\varphi(h', \mathcal{S})) \right] \\ &= \text{KL}(\rho\|\pi) + \mathbb{E}_{h' \sim \rho} \ln \left[\frac{\pi(h')}{\rho(h')} \exp(\varphi(h', \mathcal{S})) \right]. \end{aligned}$$

B.7. Proof of Theorem 2.3.2

Since \ln is concave, we can apply JENSEN's inequality (Theorem A.1.1) on the right-most term to obtain

$$\begin{aligned} \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} + \mathbb{E}_{h' \sim \rho} \ln \left[\frac{\pi(h')}{\rho(h')} \exp(\varphi(h', \mathbb{S})) \right] &\leq \mathbb{E}_{h' \sim \rho} \ln \left[\frac{\pi(h')}{\rho(h')} \exp(\varphi(h', \mathbb{S})) \right] \\ &= \ln \left[\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S})) \right]. \end{aligned}$$

Hence, we can deduce the following inequality:

$$\mathbb{E}_{h \sim \rho} \varphi(h, \mathbb{S}) \leq \text{KL}(\rho \parallel \pi) + \ln \left[\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S})) \right]. \quad (\text{B.1})$$

Since the exponential function $\exp(a)$ is positive for all $a \in \mathbb{R}$, thus, the term $\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}))$ is positive for all $\mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m$ as well. We can apply MARKOV's inequality (Theorem A.2.1) to have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S})) \leq \frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}')) \right] &\geq 1 - \delta \\ \iff \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\ln \left(\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S})) \right) \leq \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}')) \right) \right] &\geq 1 - \delta \\ \iff \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \text{KL}(\rho \parallel \pi) + \ln \left(\mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S})) \right) \right. \\ &\quad \left. \leq \text{KL}(\rho \parallel \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathbb{S}')) \right) \right] \geq 1 - \delta \end{aligned} \quad (\text{B.2})$$

By combining Equations (B.1) and (B.2), we can deduce the claimed result. ■

B.7 Proof of Theorem 2.3.2

Theorem 2.3.2 (PAC-Bayesian Bound of MCALLESTER (2003)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathcal{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \left| \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \right| \leq \sqrt{\frac{1}{2m} [\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta}]} \right] \geq 1 - \delta. \quad (2.10)$$

Proof. It is a direct consequence of PINSKER's inequality (Theorem B.5.1), *i.e.*, we have

$$2 \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right)^2 \leq \text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right),$$

and from Theorem 2.3.4 by rearranging the terms. ■

B.8 Proof of Theorem 2.3.3

Theorem 2.3.3 (PAC-Bayesian Bound of CATONI (2007)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathcal{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $c > 0$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), -\ln \left(1 - [1 - e^{-c}] \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right. \\ \left. \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta. \quad (2.13) \end{aligned}$$

Proof. We apply Theorem 2.3.1 with $\varphi(h, \mathcal{S}) = m [F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathcal{S}}^{\ell}(h)]$, where $F(R_{\mathcal{D}}^{\ell}(h)) \triangleq -\ln(1 - R_{\mathcal{D}}^{\ell}(h)[1 - e^{-c}])$. We have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} [F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathcal{S}}^{\ell}(h)] \right. \\ \left. \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathcal{S}'}^{\ell}(h')]} \right) \right] \right] \geq 1 - \delta. \end{aligned}$$

Since the distribution π on \mathbb{H} does not depend on the $\mathcal{S}' \sim \mathcal{D}^m$, we can exchange the two expectations (with FUBINI's theorem), *i.e.*, we have

$$\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathcal{S}'}^{\ell}(h')]} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathcal{S}'}^{\ell}(h')]}.$$

B.9. Proof of Theorem 2.3.4

Then, from Lemma B.16.2, we have

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} e^{m[F(\mathbb{R}_{\mathcal{D}}^\ell(h)) - c\mathbb{R}_{\mathcal{S}}^\ell(h)]} \leq 1 \implies \ln \left(\frac{1}{\delta} \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m[F(\mathbb{R}_{\mathcal{D}}^\ell(h')) - c\mathbb{R}_{\mathcal{S}'}^\ell(h')] } \right) \leq \ln \frac{1}{\delta}.$$

The function $F(x)$ is convex, since its second derivative is $\frac{\partial^2 F}{\partial x^2}(x) = \frac{(e^c - 1)^2}{(x - e^c(x-1))^2} \geq 0$. In this case, easily conclude that $F(p) - cq$ is jointly convex in q and p . Hence, from JENSEN's inequality (Theorem A.1.1) we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$F \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{D}}^\ell(h) \right) - c \left[\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{S}}^\ell(h) \right] \leq \mathbb{E}_{h \sim \rho} \left[F(\mathbb{R}_{\mathcal{D}}^\ell(h)) - c\mathbb{R}_{\mathcal{S}}^\ell(h) \right].$$

Hence, this gives the bound

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), F \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{D}}^\ell(h) \right) - c \left[\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{S}}^\ell(h) \right] \right. \\ \left. \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta, \end{aligned}$$

and by rearranging the terms we obtain the desired result. \blacksquare

B.9 Proof of Theorem 2.3.4

Theorem 2.3.4 (PAC-Bayesian Bound of SEEGER (2002)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \text{kl} \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{D}}^\ell(h) \parallel \mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{S}}^\ell(h) \right) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right] \geq 1 - \delta. \quad (2.14)$$

Proof. We can apply Theorem 2.3.1 with $\varphi(h, \mathcal{S}) = m \text{kl} \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{D}}^\ell(h) \parallel \mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{S}}^\ell(h) \right)$,

i.e., we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h) \| \mathbb{R}_{\mathcal{S}}^\ell(h)) \right. \\ \left. \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h') \| \mathbb{R}_{\mathcal{S}'}^\ell(h'))} \right) \right] \right] \geq 1 - \delta. \end{aligned}$$

Since the distribution π on \mathbb{H} does not depend on the $\mathcal{S}' \sim \mathcal{D}^m$, we can exchange the two expectations (with FUBINI's theorem), *i.e.*, we have

$$\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h') \| \mathbb{R}_{\mathcal{S}'}^\ell(h'))} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h') \| \mathbb{R}_{\mathcal{S}'}^\ell(h'))}.$$

Then, from Lemma B.16.1, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h') \| \mathbb{R}_{\mathcal{S}'}^\ell(h'))} &\leq 2\sqrt{m} \\ \implies \ln \left(\frac{1}{\delta} \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h') \| \mathbb{R}_{\mathcal{S}'}^\ell(h'))} \right) &\leq \ln \frac{2\sqrt{m}}{\delta}. \end{aligned}$$

Finally, thanks to the joint convexity of the KL divergence (Proposition B.5.1), $\text{kl}(q \| p)$ is jointly convex in q and p . Hence, we have from JENSEN's inequality (Theorem A.1.1), for all $\rho \in \mathcal{M}(\mathbb{H})$

$$\text{kl} \left(\mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{D}}^\ell(h) \| \mathbb{E}_{h \sim \rho} \mathbb{R}_{\mathcal{S}}^\ell(h) \right) \leq \mathbb{E}_{h \sim \rho} \text{kl}(\mathbb{R}_{\mathcal{D}}^\ell(h) \| \mathbb{R}_{\mathcal{S}}^\ell(h)).$$

■

B.10 Proof of Proposition 2.3.1

Proposition 2.3.1 (DONSKER-VARADHAN Variational Representation). For any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$ s.t. $\mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S})} < +\infty$ for all $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$, we have

$$\begin{aligned} \forall \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m, \quad \forall \rho \in \mathcal{M}(\mathbb{H}), \quad \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) - \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, \mathcal{S})} \right) &\leq \text{KL}(\rho \| \pi) \\ \iff \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) &\leq \text{KL}(\rho \| \pi) + \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, \mathcal{S})} \right). \end{aligned}$$

B.11. Proof of Theorem 2.3.5

When the distribution ρ is defined as $\rho(h) = \pi(h) \frac{e^{\varphi(h, \mathcal{S})}}{\mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S})}}$, we have

$$\begin{aligned} \forall \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m, \quad \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) - \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, \mathcal{S})} \right) &= \text{KL}(\rho \| \pi), \\ \iff \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) &= \text{KL}(\rho \| \pi) + \ln \left(\mathbb{E}_{h \sim \pi} e^{\varphi(h, \mathcal{S})} \right). \end{aligned}$$

Proof. Let $\rho'(h) = \pi(h) \frac{e^{\varphi(h, \mathcal{S})}}{\mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S})}}$, then, we develop the term $\text{KL}(\rho \| \rho')$. We have

$$\begin{aligned} \text{KL}(\rho \| \rho') &= \mathbb{E}_{h \sim \rho} \ln \left(\frac{\rho(h)}{\rho'(h)} \right) \\ &= \mathbb{E}_{h \sim \rho} \ln \left(\frac{\rho(h)}{\pi(h) e^{\varphi(h, \mathcal{S})}} \mathbb{E}_{h' \sim \rho} e^{\varphi(h', \mathcal{S})} \right) \\ &= \mathbb{E}_{h \sim \rho} \ln \left(\frac{\rho(h)}{\pi(h)} \right) - \mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) + \ln \left(\mathbb{E}_{h' \sim \rho} e^{\varphi(h', \mathcal{S})} \right). \end{aligned}$$

Hence, by setting $\rho = \rho'$, we have $\text{KL}(\rho \| \rho') = 0$ which leads to the desired result by rearranging the terms. ■

B.11 Proof of Theorem 2.3.5

Theorem 2.3.5 (General PAC-Bayesian Bound of BÉGIN *et al.* (2016)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_*^+$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} \varphi(h, \mathcal{S}) \right] \right. \\ \left. \leq D_\lambda(\rho \| \pi) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right] \right] \geq 1 - \delta. \end{aligned}$$

Proof. We start by developing $\frac{\lambda}{\lambda-1} \ln [\mathbb{E}_{h \sim \rho} \varphi(h, \mathbb{S})]$, i.e., we have for all $\rho \in \mathcal{M}(\mathbb{H})$

$$\frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} \varphi(h, \mathbb{S}) \right] = \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} \frac{\rho(h) \pi(h)}{\pi(h) \rho(h)} \varphi(h, \mathbb{S}) \right] = \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \pi} \frac{\rho(h)}{\pi(h)} \varphi(h, \mathbb{S}) \right]. \quad (\text{B.3})$$

We apply HÖLDER's inequality (Theorem A.5.1) to have

$$\mathbb{E}_{h \sim \pi} \frac{\rho(h)}{\pi(h)} \varphi(h, \mathbb{S}) \leq \left[\mathbb{E}_{h \sim \pi} \left[\frac{\rho(h)}{\pi(h)} \right]^\lambda \right]^{\frac{1}{\lambda}} \left[\mathbb{E}_{h \sim \pi} \varphi(h, \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right]^{\frac{\lambda-1}{\lambda}}.$$

By taking the logarithm (since both sides are positive) and multiplying by $\frac{\lambda}{\lambda-1}$ both sides of the inequality, we have

$$\begin{aligned} \frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \frac{\rho(h)}{\pi(h)} \varphi(h, \mathbb{S}) \right) &\leq \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \left[\frac{\rho(h)}{\pi(h)} \right]^\lambda \right) + \ln \left(\mathbb{E}_{h \sim \pi} \varphi(h, \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \\ &= D_\lambda(\rho \| \pi) + \ln \left(\mathbb{E}_{h \sim \pi} \varphi(h, \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right). \end{aligned} \quad (\text{B.4})$$

Since the function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_*^+$ is positive and $\frac{\lambda}{\lambda-1} > 0$, we can deduce that $\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} > 0$ for all $\mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m$. Then, we can apply MARKOV's inequality (Theorem A.2.1), we have

$$\begin{aligned} &\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \leq \frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right] \geq 1 - \delta \\ \iff &\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\ln \left(\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \geq 1 - \delta \\ \iff &\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathcal{M}(\mathbb{H}), D_\lambda(\rho \| \pi) + \ln \left(\mathbb{E}_{h' \sim \pi'} \varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \right. \\ &\quad \left. \leq D_\lambda(\rho \| \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \geq 1 - \delta \end{aligned} \quad (\text{B.5})$$

By combining Equations (B.3) to (B.5), we can deduce the claimed result. \blacksquare

B.12 About the Bounds Derived From Theorem 2.3.5

In this section, we provide a SEEGER, MCALLESTER and CATONI-like PAC-Bayesian generalization bound.

B.12.1 McAllester-like Bound

Based on Corollary B.12.2, we can obtain the following MCALLESTER-like bound.

Corollary B.12.1. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \left| \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right| \leq \sqrt{\frac{1}{2m} \left[D_{\lambda}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]} \right] \geq 1 - \delta.$$

Proof. It is a direct consequence of PINSKER's inequality (Theorem B.5.1), i.e., we have

$$2 \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) - \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right)^2 \leq \text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right),$$

and from Corollary B.12.2 by rearranging the terms. ■

B.12.2 Seeger-like Bound

The SEEGER-like obtained from Theorem 2.3.5 is the following.

Corollary B.12.2. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h) \right) \leq \frac{1}{m} \left[D_{\lambda}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right] \geq 1 - \delta.$$

Proof. We apply Theorem 2.3.5 with $\varphi(h, \mathbb{S}) = \frac{\lambda-1}{\lambda} m \text{kl} \left[\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h) \right]$, i.e., we have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbf{R}_{\mathcal{S}}^{\ell}(h))} \right] \right. \\ \left. \leq D_{\lambda}(\rho \parallel \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h') \parallel \mathbf{R}_{S'}^{\ell}(h'))} \right) \right] \geq 1 - \delta. \end{aligned}$$

Since the distribution π on \mathbb{H} does not depend on the $S' \sim \mathcal{D}^m$, we can exchange the two expectations (with FUBINI's theorem), i.e., we have

$$\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h') \parallel \mathbf{R}_{S'}^{\ell}(h'))} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{S' \sim \mathcal{D}^m} e^{m \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h') \parallel \mathbf{R}_{S'}^{\ell}(h'))}.$$

Then, from Lemma B.16.1, we have

$$\begin{aligned} \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} e^{m \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbf{R}_{\mathbb{S}}^{\ell}(h))} &\leq 2\sqrt{m} \\ \implies \ln \left(\frac{1}{\delta} \mathbb{E}_{h' \sim \pi} \mathbb{E}_{S' \sim \mathcal{D}^m} e^{m \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h') \parallel \mathbf{R}_{S'}^{\ell}(h'))} \right) &\leq \ln \frac{2\sqrt{m}}{\delta}. \end{aligned}$$

Thanks to the joint convexity of the KL divergence (Proposition B.5.1), the function $q, p \mapsto \exp \left(\frac{\lambda-1}{\lambda} m \text{kl}(q \parallel p) \right)$ is jointly convex in q and p by composition (see, e.g., BOYD and VANDENBERGHE (2004)). Hence, we have from JENSEN's inequality (Theorem A.1.1), for all $\rho \in \mathbb{M}(\mathbb{H})$

$$\begin{aligned} e^{m \frac{\lambda-1}{\lambda} \text{kl}(\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h))} &\leq \mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbf{R}_{\mathcal{S}}^{\ell}(h))} \\ \iff m \text{kl} \left[\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h) \right] &\leq \frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} \text{kl}(\mathbf{R}_{\mathcal{D}}^{\ell}(h) \parallel \mathbf{R}_{\mathcal{S}}^{\ell}(h))} \right). \end{aligned}$$

Finally, by rearranging the terms, we have the stated result. ■

B.12.3 Catoni-like Bound

The derivation of a CATONI (2007)-like generalization bound based on Theorem 2.3.5 is the following.

Corollary B.12.3. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow [0, 1]$, for any $c > 0$,

for any $\lambda > 1$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), -\ln \left(1 - [1 - e^{-c}] \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h) \right. \\ \left. \leq \frac{1}{m} \left[D_{\lambda}(\rho \| \pi) + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta. \end{aligned}$$

Proof. We apply Theorem 2.3.5 with $\varphi(h, \mathcal{S}) = m \frac{\lambda-1}{\lambda} [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)]$, where $F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) \triangleq -\ln(1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h)[1 - e^{-c}])$. We have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln \left[\mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)} \right] \right. \\ \left. \leq D_{\lambda}(\rho \| \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h')) - c\mathbf{R}_{\mathcal{S}'}^{\ell}(h')]} \right) \right] \geq 1 - \delta. \end{aligned}$$

Since the distribution π on \mathbb{H} does not depend on the $\mathcal{S}' \sim \mathcal{D}^m$, we can exchange the two expectations (with FUBINI's theorem), i.e., we have

$$\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h')) - c\mathbf{R}_{\mathcal{S}'}^{\ell}(h')]} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h')) - c\mathbf{R}_{\mathcal{S}'}^{\ell}(h')]}.$$

Then, from Lemma B.16.2, we have

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} e^{m [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)]} \leq 1 \implies \ln \left(\frac{1}{\delta} \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h')) - c\mathbf{R}_{\mathcal{S}'}^{\ell}(h')]} \right) \leq \ln \frac{1}{\delta}.$$

The function $F(x)$ is convex, since its second derivative is $\frac{\partial^2 F}{\partial x^2}(x) = \frac{(e^c - 1)^2}{(x - e^c(x-1))^2} \geq 0$. Hence, we conclude that $F(p) - cq$ is jointly convex in q and p . Moreover, we can deduce that the function $q, p \mapsto \exp \left(m \frac{\lambda-1}{\lambda} [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)] \right)$ is jointly convex in q and p by composition (see, e.g., BOYD and VANDENBERGHE (2004)). Hence, from JENSEN's inequality (Theorem A.1.1) we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$\begin{aligned} e^{m \frac{\lambda-1}{\lambda} [F(\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h)]} &\leq \mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)]} \\ \iff m \left[F \left(\mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} \mathbf{R}_{\mathcal{S}}^{\ell}(h) \right] &\leq \frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \rho} e^{m \frac{\lambda-1}{\lambda} [F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c\mathbf{R}_{\mathcal{S}}^{\ell}(h)]} \right). \end{aligned}$$

Finally, by rearranging the terms, we have the stated result. ■

B.13 Proof of Theorem 2.4.1

Theorem 2.4.1 (General Disintegrated Bound of RIVASPLATA *et al.* (2020)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \underbrace{\ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\varphi(h', \mathcal{S}')) \right]}_{\Phi(\rho_{\mathcal{S}}, \pi, \delta)} \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Proof. Note that $\exp \left[\varphi(h, \mathcal{S}) - \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right]$ is a non-negative random variable. Thus, we can apply MARKOV's inequality (Theorem A.2.1) to obtain

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\exp \left(\varphi(h, \mathcal{S}) - \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right) \right. \\ \left. \leq \frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \exp \left(\varphi(h', \mathcal{S}') - \ln \frac{\rho_{\mathcal{S}'}(h')}{\pi(h')} \right) \right] \geq 1 - \delta. \end{aligned}$$

Hence, by rearranging the terms, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\exp \left(\varphi(h, \mathcal{S}) - \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right) \leq \frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \frac{\pi(h')}{\rho_{\mathcal{S}'}(h')} e^{\varphi(h', \mathcal{S}')} \right] \geq 1 - \delta \\ \iff \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\exp \left(\varphi(h, \mathcal{S}) - \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right) \leq \frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S}')} \right] \geq 1 - \delta. \end{aligned}$$

Since both sides of the inequality are strictly positive, we can apply the logarithm, *i.e.*, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\varphi(h', \mathcal{S}')} \right) \right] \geq 1 - \delta,$$

which is the desired result. ■

B.14 Proof of Theorem 2.4.2

Theorem 2.4.2 (Disintegrated Bound of CATONI (2007)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [0, 1]$, for any $c > 0$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), -\ln \left(1 - [1 - e^{-c}] \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) - c \mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \right. \\ \left. \leq \frac{1}{m} \left[\ln \frac{\rho_{\mathbb{S}}(h)}{\pi(h)} + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta, \end{aligned}$$

where $\rho_{\mathbb{S}} \triangleq A(\mathbb{S}, \pi)$ is output by the deterministic algorithm A .

Proof. We apply Theorem 2.4.1 with $\varphi(h, \mathbb{S}) = m [F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathbb{S}}^{\ell}(h)]$, where $F(R_{\mathcal{D}}^{\ell}(h)) \triangleq -\ln(1 - R_{\mathcal{D}}^{\ell}(h)[1 - e^{-c}])$. We have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), \mathbb{E}_{h \sim \rho} [F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathbb{S}}^{\ell}(h)] \right. \\ \left. \leq \frac{1}{m} \left[\ln \frac{\rho_{\mathbb{S}}(h)}{\pi(h)} + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathbb{S}'}^{\ell}(h')]} \right) \right] \right] \geq 1 - \delta. \end{aligned}$$

Since the distribution π on \mathbb{H} does not depend on the $\mathbb{S}' \sim \mathcal{D}^m$, we can exchange the two expectations (with FUBINI's theorem), *i.e.*, we have

$$\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathbb{S}'}^{\ell}(h')]} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathbb{S}'}^{\ell}(h')]}.$$

Then, from Lemma B.16.2, we have

$$\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} e^{m[F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathbb{S}}^{\ell}(h)]} \leq 1 \implies \ln \left(\frac{1}{\delta} \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} e^{m[F(R_{\mathcal{D}}^{\ell}(h')) - cR_{\mathbb{S}'}^{\ell}(h')]} \right) \leq \ln \frac{1}{\delta}.$$

The function $F(x)$ is convex, since its second derivative is $\frac{\partial^2 F}{\partial x^2}(x) = \frac{(e^c - 1)^2}{(x - e^c(x-1))^2} \geq 0$. In this case, easily conclude that $F(p) - cq$ is jointly convex in q and p . Hence,

from JENSEN's inequality (Theorem A.1.1) we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$F\left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)\right) - c\left[\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)\right] \leq \mathbb{E}_{h \sim \rho} \left[F(R_{\mathcal{D}}^{\ell}(h)) - cR_{\mathcal{S}}^{\ell}(h)\right].$$

Hence, this gives the bound

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\forall \rho \in \mathbb{M}(\mathbb{H}), F\left(\mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h)\right) - c\left[\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h)\right] \right. \\ \left. \leq \frac{1}{m} \left[\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} + \ln \frac{1}{\delta} \right] \right] \geq 1 - \delta, \end{aligned}$$

and by rearranging the terms we obtain the desired result. \blacksquare

B.15 Proof of Theorem 2.4.3

The proof of Theorem 2.4.3 relies on another theorem from BLANCHARD and FLEURET (2007, Theorem 2.4) called OCCAM's Hammer.

Lemma B.15.1 (OCCAM's Hammer). Given a measurable function measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}$, assume that for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \pi} [\varphi(h, \mathcal{S}) \geq \Phi(\delta)] \leq \delta,$$

where $\Phi : (0, 1] \rightarrow \mathbb{R}$ is a decreasing function.

It implies that for any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution π on \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [\varphi(h, \mathcal{S}) \geq \Phi'(\delta f(\Theta(h)^{-1}))] \leq \delta,$$

where $\Theta(h) \triangleq \frac{\rho_{\mathcal{S}}(h)}{\pi(h)}$ and $\Phi'(\delta) \triangleq \Phi(\max(\delta, 1))$ and $f : \mathbb{R}^+ \rightarrow \mathbb{R}_*^+$ a measurable increasing function s.t.

$$\int_{y>0} y^{-2} f(y) dy \leq 1.$$

Proof. The proof consists in upper-bounding the probability

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \geq \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right].$$

We rewrite this term as

$$\begin{aligned} & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \frac{\pi(h)}{\rho_{\mathcal{S}}(h)} \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \Theta(h) \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right]. \end{aligned}$$

We can actually express the term $\Theta(h)$ in a form of integral. Indeed, we have

$$\Theta(h) = \int_{\Theta(h)^{-1}}^{+\infty} y^{-2} dy = \int_{y>0} y^{-2} \mathbb{I} \left[y \geq \Theta(h)^{-1} \right] dy.$$

Then, thanks to FUBINI's theorem, we can rewrite the probability as

$$\begin{aligned} & \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \Theta(h) \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \left[\int_{y>0} y^{-2} \mathbb{I} \left[y \geq \Theta(h)^{-1} \right] dy \right] \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \\ &= \int_{y>0} y^{-2} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \mathbb{I} \left[y \geq \Theta(h)^{-1} \right] \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] dy. \end{aligned}$$

Since $f(\cdot)$ is increasing, *i.e.*, for all $y > 0$ s.t. $\Theta(h)^{-1} < y$, we have $f(\Theta(h)^{-1}) \leq f(y)$. Moreover, since $\Phi(\cdot)$ is decreasing, we have $\Phi'(f(y)) \leq \Phi'(f(\Theta(h)^{-1}))$.

This implies that

$$\mathbb{I} \left[y \geq \Theta(h)^{-1} \right] \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] \leq \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(y \right) \right) \right].$$

Based on this inequality, we have

$$\begin{aligned} & \int_{y>0} y^{-2} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \mathbb{I} \left[y \geq \Theta(h)^{-1} \right] \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f \left(\Theta(h)^{-1} \right) \right) \right] dy \\ & \leq \int_{y>0} y^{-2} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \mathbb{I} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f(y) \right) \right] dy \\ & = \int_{y>0} y^{-2} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \pi} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f(y) \right) \right] dy. \end{aligned}$$

By assumption, we have $\int_{y>0} y^{-2} f(y) dy \leq 1$ and $\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \pi} [\varphi(h, \mathcal{S}) \geq \Phi(\delta)] \leq \delta$, which leads to

$$\int_{y>0} y^{-2} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \pi} \left[\varphi(h, \mathcal{S}) \geq \Phi' \left(\delta f(y) \right) \right] dy \leq \int_{y>0} y^{-2} \delta f(y) dy \leq \delta. \quad \blacksquare$$

This theorem is further used in addition to the two following lemmas. We first one known as CHERNOFF's bound is originally due to CHERNOFF (1952) but the proof is from LANGFORD (2005, Lemma 3.6).

Lemma B.15.2 (CHERNOFF's bound). For any

$$\mathbb{P}_{\mathbf{X} \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq k \right] \leq e^{-m \text{kl}_+ \left(\frac{k}{m} \| p \right)},$$

where $\text{kl}_+ \left(\frac{k}{m} \| p \right) = \text{kl} \left(\frac{k}{m} \| p \right)$ if $\frac{k}{m} < p$ and 0 otherwise.

Proof. First remark that we have

$$\mathbb{P}_{\mathbf{X} \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq k \right] = \mathbb{P}_{\mathbf{X} \sim \mathcal{B}(p)^m} \left[e^{-m\lambda \frac{1}{m} \sum_{i=1}^m X_i} \geq e^{-m\lambda \frac{k}{m}} \right].$$

Then, from MARKOV's inequality (Theorem A.2.1), we have

$$\mathbb{P}_{\mathbf{X} \sim \mathcal{B}(p)^m} \left[e^{-m\lambda \frac{1}{m} \sum_{i=1}^m X_i} \geq e^{-m\lambda \frac{k}{m}} \right] \leq \frac{\mathbb{E}_{\mathbf{X} \sim \mathcal{B}(p)^m} e^{-\lambda \sum_{i=1}^m X_i}}{e^{-\lambda k}}.$$

Then, using the fact that X_1, \dots, X_m are *i.i.d.* and from the expression of the

B.15. Proof of Theorem 2.4.3

moment generating function of the Bernoulli distribution $\mathcal{B}(p)$, we have

$$\frac{\mathbb{E}_{\mathbf{X} \sim \mathcal{B}(p)^m} e^{-\lambda \sum_{i=1}^m X_i}}{e^{-\lambda k}} = e^{\lambda k} \left[\mathbb{E}_{\mathbf{X} \sim \mathcal{B}(p)} e^{-\lambda X} \right]^m = e^{\lambda k} [pe^\lambda + (1-p)]^m.$$

Actually, we can find the optimal value λ^* , which is

$$\lambda^* = \ln \left[p \left(1 - \frac{k}{m} \right) \right] - \ln \left[\frac{k}{m} (1-p) \right],$$

for all $p > \frac{k}{m}$. Finally, setting $\lambda = \lambda^*$, we obtain

$$e^{\lambda k} [pe^\lambda + (1-p)]^m = e^{-m \text{kl}_+(\frac{k}{m} \| p)}.$$

■

CHERNOFF's bound is actually used to prove the test set bound of LANGFORD (2005, Theorem 3.3 and Corollary 3.7). We prove, with more details, his theorem in the following lemma.

Lemma B.15.3. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any hypothesis $h \in \mathbb{H}$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \{0, 1\}$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\text{kl}_+(\mathbb{R}_{\mathcal{S}}^\ell(h) \| \mathbb{R}_{\mathcal{D}}^\ell(h)) \geq \frac{\ln \frac{1}{\delta}}{m} \right] \leq \delta.$$

Proof. Step 1. First of all, we prove that

$$\mathbb{P}_{Y_1, \dots, Y_m \sim \mathcal{B}(p)^m} \left[\mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq \sum_{i=1}^m Y_i \right] \leq \delta \right] \leq \delta. \quad (\text{B.6})$$

To do so, let $k^* = \max \left\{ k \in \{0, \dots, m\} \mid \mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq k \right] \leq \delta \right\}$. Then, we have

$$\begin{aligned} & \mathbb{P}_{Y_1, \dots, Y_m \sim \mathcal{B}(p)^m} \left[\mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq \sum_{i=1}^m Y_i \right] \leq \delta \right] \\ &= \mathbb{E}_{Y_1, \dots, Y_m \sim \mathcal{B}(p)^m} \mathbb{I} \left[\mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq \sum_{i=1}^m Y_i \right] \leq \delta \right] \end{aligned}$$

$$\begin{aligned}
 &= \sum_{k=0}^m \mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i = k \right] \cdot \mathbb{I} \left[\mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq \sum_{i=1}^m Y_i \right] \leq \delta \right] \\
 &= \sum_{k=0}^{k^*} \mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i = k \right] \\
 &= \mathbb{P}_{X_1, \dots, X_m \sim \mathcal{B}(p)^m} \left[\sum_{i=1}^m X_i \leq k^* \right] \\
 &\leq \delta.
 \end{aligned}$$

Step 2. From Equation (B.6), we can deduce that

$$\mathbb{P}_{R_S^\ell(h) \sim \mathcal{B}(R_D^\ell(h))^m} \left[\mathbb{P}_{\mathbf{x} \sim \mathcal{B}(R_D^\ell(h))^m} \left[\frac{1}{m} \sum_{i=1}^m X_i \leq R_S^\ell(h) \right] \leq \delta \right] \leq \delta,$$

where $R_S^\ell(h) \sim \mathcal{B}(R_D^\ell(h))^m$ is a slight abuse of notations since $mR_S^\ell(h)$ is the sum of the successes. From Chernoff's inequality (Lemma B.15.2), we can deduce that

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} \left[\exp\left(-m \text{kl}_+ \left[R_S^\ell(h) \| R_D^\ell(h) \right]\right) \leq \delta \right] \leq \delta.$$

Finally, by rearranging the terms, we obtain the desired result. \blacksquare

Finally, we are able to prove Lemma B.15.1 using Lemmas B.15.1 and B.15.2.

Theorem 2.4.3 (Disintegrated Bound of BLANCHARD and FLEURET (2007)). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \{0, 1\}$, for any $k > 1$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\text{kl}_+ \left(R_S^\ell(h) \| R_D^\ell(h) \right) \leq \frac{1}{m} \left[\ln \frac{k+1}{\delta} + \left(1 + \frac{1}{k} \right) \ln_+ \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A , the $\ln_+(x) = \max(\ln(x), 0)$ and $\text{kl}_+(R_S^\ell(h) \| R_D^\ell(h)) = \text{kl}(R_S^\ell(h) \| R_D^\ell(h))$ if $R_S^\ell(h) < R_D^\ell(h)$ and 0 otherwise.

Proof. The proof consist of applying OCCAM's Hammer (Lemma B.15.1). To do so, given $k > 1$, let $f : \mathbb{R}^+ \rightarrow \mathbb{R}_*^+$ the function defined as

$$f(y) \triangleq \frac{1}{k+1} \min \left(y^{1+\frac{1}{k}}, 1 \right)$$

Indeed, $f(\cdot)$ is increasing and we have

$$\begin{aligned}
 \int_{y>0} y^{-2} f(y) dy &= \int_0^1 y^{-2} f(y) dy + \int_1^{+\infty} y^{-2} f(y) dy \\
 &= \frac{1}{k+1} \left[\int_0^1 y^{-2} \cdot y^{1+\frac{1}{k}} dy + \int_1^{+\infty} y^{-2} dy \right] \\
 &= \frac{1}{k+1} \left[\int_0^1 y^{-2} \cdot y^{1+\frac{1}{k}} dy + 1 \right] \\
 &= \frac{1}{k+1} \left[\int_0^1 y^{\frac{1}{k}-1} dy + 1 \right] \\
 &= \frac{1}{k+1} \left[k \cdot \left[1^{\frac{1}{k}} - 0^{\frac{1}{k}} \right] + 1 \right] \\
 &= 1.
 \end{aligned}$$

Moreover, note that with Lemma B.15.3, we have $\varphi(h, \mathbb{S}) = \text{kl}_+(\mathbb{R}_{\mathbb{S}}^\ell(h) \parallel \mathbb{R}_{\mathcal{D}}^\ell(h))$ and $\Phi(\delta) = \frac{\ln \frac{1}{\delta}}{m}$. Hence, we apply OCCAM's Hammer (Lemma B.15.1) to obtain

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\varphi(h, \mathbb{S}) \leq \Phi'(\delta f(\Theta(h)^{-1})) \right] \geq 1 - \delta,$$

To obtain the final bound, we upper-bound the term $\Phi'(\delta f(\Theta(h)^{-1}))$, *i.e.*, we have

$$\begin{aligned}
 \Phi'(\delta f(\Theta(h)^{-1})) &= -\frac{1}{m} \ln \left(\min(\delta f(\Theta(h)^{-1}), 1) \right) \\
 &= \frac{1}{m} \max \left(\ln \left(\frac{1}{\delta f(\Theta(h)^{-1})} \right), 0 \right) \\
 &\leq \frac{1}{m} \ln \frac{1}{\delta} + \frac{1}{m} \max \left(-\ln \left(f(\Theta(h)^{-1}) \right), 0 \right) \\
 &= \frac{1}{m} \ln \frac{1}{\delta} + \frac{1}{m} \max \left(-\ln \left(\frac{1}{k+1} \min(\Theta(h)^{-(1+\frac{1}{k})}, 1) \right), 0 \right) \\
 &\leq \frac{1}{m} \ln \frac{k+1}{\delta} + \frac{1}{m} \max \left(-\ln \left(\min(\Theta(h)^{-(1+\frac{1}{k})}, 1) \right), 0 \right) \\
 &= \frac{1}{m} \ln \frac{k+1}{\delta} + \frac{1}{m} \ln_+ \left(\Theta(h)^{1+\frac{1}{k}} \right) \\
 &= \frac{1}{m} \left[\ln \frac{k+1}{\delta} + \left(1 + \frac{1}{k} \right) \ln_+ \frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right].
 \end{aligned}$$

■

B.16 Proof of Lemmas B.16.1 and B.16.2

In this section, we prove the lemmas necessary to prove Theorems 2.3.3, 2.3.4 and 2.4.2 and Corollary B.12.2. The proof of Lemma B.16.1 is due to MAURER (2004) and Lemma B.16.2 was proven by GERMAIN *et al.* (2009).

Lemma B.16.1. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, we have

$$\forall h \in \mathbb{H}, \quad \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \exp \left[m \text{kl}(\mathbb{R}_{\mathcal{S}}^\ell(h) \| \mathbb{R}_{\mathcal{D}}^\ell(h)) \right] \leq 2\sqrt{m}.$$

Lemma B.16.2. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, we have

$$\forall h \in \mathbb{H}, \quad \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \exp \left[m \left[F(\mathbb{R}_{\mathcal{D}}^\ell(h)) - c\mathbb{R}_{\mathcal{S}}^\ell(h) \right] \right] \leq 1,$$

where $F(\mathbb{R}_{\mathcal{D}}^\ell(h)) \triangleq -\ln(1 - \mathbb{R}_{\mathcal{D}}^\ell(h)[1 - e^{-c}])$.

However, before giving the proofs, we need to prove two lemmas.

Lemma B.16.3. For any $m \in \mathbb{N}_*$, any point $\mathbf{x} \in [0, 1]^m$ can be written as a convex combination of the extremes points $\boldsymbol{\eta} \in \{0, 1\}^m$, *i.e.*, we have

$$\forall \mathbf{x} \in [0, 1]^m, \quad \mathbf{x} = \sum_{\boldsymbol{\eta} \in \{0, 1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1 - x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] \boldsymbol{\eta},$$

where $\sum_{\boldsymbol{\eta} \in \{0, 1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1 - x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] = 1$.

Proof. We prove this fact by induction.

For $m = 1$, we can easily prove the claim, *i.e.*, we have

$$\forall \mathbf{x} \in [0, 1], \quad \mathbf{x} = x_1 \cdot 1 + (1 - x_1) \cdot 0, \quad \text{and} \quad (x_1) + (1 - x_1) = 1.$$

For $m > 1$, we assume that the claim is true for a particular m (from our induction

hypothesis) and we prove the equality for $m + 1$.

$$\begin{aligned}
 \forall \mathbf{x} \in [0, 1]^{m+1}, \quad & \sum_{\boldsymbol{\eta} \in \{0,1\}^{m+1}} \left[\prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 1}} x_i \right] \boldsymbol{\eta}, \\
 &= \sum_{\boldsymbol{\eta} \in \{0,1\}^{m+1}} \left[\prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 1}} x_i \right] [\eta_1, \dots, \eta_{m+1}]^\top \\
 &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] x_{m+1} [\eta_1, \dots, 1]^\top \\
 &+ \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] (1-x_{m+1}) [\eta_1, \dots, 0]^\top.
 \end{aligned}$$

For any $\mathbf{x} \in [0, 1]^{m+1}$, its $m+1$ -th component is

$$\underbrace{\left(\sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] \right)}_{=1 \text{ by the induction hypothesis}} (x_{m+1} \cdot 1 + (1-x_{m+1}) \cdot 0) = x_{m+1}.$$

Moreover, from the 1-st to the m -th component, we have

$$\begin{aligned}
 & \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] x_{m+1} \cdot \boldsymbol{\eta} \\
 &+ \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] (1-x_{m+1}) \cdot \boldsymbol{\eta} \\
 &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] [x_{m+1} + (1-x_{m+1})] \cdot \boldsymbol{\eta}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] \cdot \boldsymbol{\eta} \\
 &= [x_1, \dots, x_m]^\top,
 \end{aligned}$$

where the last equality holds by the induction hypothesis. Finally, we prove that it sums to $\sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] = 1$, i.e., we have

$$\begin{aligned}
 \forall \mathbf{x} \in [0, 1]^{m+1}, \quad & \sum_{\boldsymbol{\eta} \in \{0,1\}^{m+1}} \left[\prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m+1\} \\ \text{s.t. } \eta_i = 1}} x_i \right] \\
 &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] x_{m+1} \\
 &+ \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 0}} (1-x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \text{s.t. } \eta_i = 1}} x_i \right] (1-x_{m+1}) \\
 &= 1,
 \end{aligned}$$

by the induction hypothesis. ■

The second lemma that we have to prove is the following.

Lemma B.16.4. Let $X \sim \mathcal{X}$ be a random variable such that $X \in [0, 1]$ and X' be a random variable following a Bernoulli distribution of parameter $\mathbb{E}_{X \sim \mathcal{X}}[X]$, i.e., $X' \sim \mathcal{B}(\mathbb{E}_{X \sim \mathcal{X}}[X])$. We define as $\mathbf{X} \sim \mathcal{X}^m$ (resp. $\mathbf{X}' \sim \mathcal{B}(\mathbb{E}_{X \sim \mathcal{X}}[X])^m$) the m independent copies of $X \sim \mathcal{X}$ (resp. $X' \sim \mathcal{B}(\mathbb{E}_{X \sim \mathcal{X}}[X])$).

If $f : [0, 1]^m \rightarrow \mathbb{R}$ is a convex function and permutation symmetric, we have

$$\mathbb{E}_{\mathbf{X} \sim \mathcal{X}^m} [F(\mathbf{X})] \leq \mathbb{E}_{\mathbf{X}' \sim \mathcal{B}(\mathbb{E}_{X \sim \mathcal{X}}[X])^m} [F(\mathbf{X}')].$$

Proof. From Lemma B.16.3, we have

$$\mathbf{x} = \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 0}} (1 - x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 1}} x_i \right] \boldsymbol{\eta}. \quad (\text{B.7})$$

Hence from Equation (B.7) and from JENSEN's inequality (Theorem A.1.1), we have

$$\begin{aligned} F(\mathbf{x}) &= f \left(\sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 0}} (1 - x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 1}} x_i \right] \boldsymbol{\eta} \right) \\ &\leq \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 0}} (1 - x_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 1}} x_i \right] F(\boldsymbol{\eta}). \end{aligned}$$

Taking the expectation gives us

$$\begin{aligned} \mathbb{E}_{\mathbf{X} \sim \mathcal{X}^m} F(\mathbf{X}) &\leq \mathbb{E}_{\mathbf{X} \sim \mathcal{X}^m} \left(\sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 0}} (1 - X_i) \prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 1}} X_i \right] F(\boldsymbol{\eta}) \right) \\ &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 0}} \left(1 - \mathbb{E}_{X_i \sim \mathcal{X}} [X_i] \right) \prod_{\substack{i \in \{1, \dots, m\} \\ \eta_i = 1}} \mathbb{E}_{X_i \sim \mathcal{X}} [X_i] \right] F(\boldsymbol{\eta}) \\ &= \sum_{\boldsymbol{\eta} \in \{0,1\}^m} \left[\left(1 - \mathbb{E}_{X \sim \mathcal{X}} [X] \right)^{\text{card}(\{i: \eta_i = 0\})} \left(\mathbb{E}_{X \sim \mathcal{X}} [X] \right)^{\text{card}(\{i: \eta_i = 1\})} \right] F(\boldsymbol{\eta}) \\ &= \sum_{k=0}^m \binom{m}{k} \left(1 - \mathbb{E}_{X \sim \mathcal{X}} [X] \right)^{m-k} \left(\mathbb{E}_{X \sim \mathcal{X}} [X] \right)^k F(\underbrace{1, \dots, 1}_{k \text{ times}}, \underbrace{0, \dots, 0}_{m-k \text{ times}}) \\ &= \mathbb{E}_{\mathbf{X}' \sim \mathcal{B}(\mathbb{E}_{X \sim \mathcal{X}} [X])^m} F(\mathbf{X}'). \end{aligned}$$

■

We are now ready to prove Lemma B.16.1.

Proof of Lemma B.16.1. Since the KL divergence is jointly convex so is the function $\exp \left[m \text{kl}(\cdot \| \mathbf{R}_{\mathcal{D}}^{\ell}(h)) \right]$ (see, e.g., BOYD and VANDENBERGHE (2004, Section 3.2.4)). Then, we can apply Lemma B.16.4 to have

$$\begin{aligned}
 & \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \exp \left[m \text{kl}(\mathbf{R}_{\mathcal{S}}^{\ell}(h) \| \mathbf{R}_{\mathcal{D}}^{\ell}(h)) \right] \\
 & \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{B}(\mathbf{R}_{\mathcal{D}}^{\ell}(h))^m} \exp \left[m \text{kl} \left(\frac{1}{m} \sum_{i=1}^m X_i \| \mathbf{R}_{\mathcal{D}}^{\ell}(h) \right) \right] \\
 & = \mathbb{E}_{\mathbf{x} \sim \mathcal{B}(\mathbf{R}_{\mathcal{D}}^{\ell}(h))^m} \left[\frac{\frac{1}{m} \sum_{i=1}^m X_i}{\mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^{\sum_{i=1}^m X_i} \left[\frac{1 - \frac{1}{m} \sum_{i=1}^m X_i}{1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^{m - \sum_{i=1}^m X_i} \\
 & = \sum_{k=0}^m \mathbb{P}_{\mathbf{x} \sim \mathcal{B}(\mathbf{R}_{\mathcal{D}}^{\ell}(h))^m} \left[\sum_{i=1}^m X_i = k \right] \left[\frac{\frac{k}{m}}{\mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^k \left[\frac{1 - \frac{k}{m}}{1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^{m-k} \\
 & = \sum_{k=0}^m \binom{m}{k} (1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h))^{m-k} (\mathbf{R}_{\mathcal{D}}^{\ell}(h))^k \left[\frac{\frac{k}{m}}{\mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^k \left[\frac{1 - \frac{k}{m}}{1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h)} \right]^{m-k} \\
 & = \sum_{k=0}^m \binom{m}{k} \left[\frac{k}{m} \right]^k \left[1 - \frac{k}{m} \right]^{m-k}.
 \end{aligned}$$

Finally, MAURER (2004) proves that $\sum_{k=0}^m \binom{m}{k} \left[\frac{k}{m} \right]^k \left[1 - \frac{k}{m} \right]^{m-k} \leq 2\sqrt{m}$ for $m \geq 8$ and GERMAIN *et al.* (2015) verify computationally that the inequality holds also for $m \in \{1, \dots, 7\}$. ■

We can prove Lemma B.16.2.

Proof of Lemma B.16.2. We have

$$\begin{aligned}
 & \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \exp \left(m \left[F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c \mathbf{R}_{\mathcal{S}}^{\ell}(h) \right] \right) \\
 & \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{B}(\mathbf{R}_{\mathcal{D}}^{\ell}(h))^m} \exp \left(m F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - c \left(\sum_{i=1}^m X_i \right) \right) \tag{B.8}
 \end{aligned}$$

$$\begin{aligned}
 & = \sum_{k=0}^m \mathbb{P}_{\mathbf{x} \sim \mathcal{B}(\mathbf{R}_{\mathcal{D}}^{\ell}(h))^m} \left[\sum_{i=1}^m X_i = k \right] \exp \left(m F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - ck \right) \\
 & = \sum_{k=0}^m \binom{m}{k} \mathbf{R}_{\mathcal{D}}^{\ell}(h)^k (1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h))^{m-k} \exp \left(m F(\mathbf{R}_{\mathcal{D}}^{\ell}(h)) - ck \right) \\
 & = e^{m F(\mathbf{R}_{\mathcal{D}}^{\ell}(h))} \sum_{k=0}^m \binom{m}{k} (\mathbf{R}_{\mathcal{D}}^{\ell}(h) e^{-c})^k (1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h))^{m-k} \\
 & = e^{m F(\mathbf{R}_{\mathcal{D}}^{\ell}(h))} \left(1 - \mathbf{R}_{\mathcal{D}}^{\ell}(h) \left[1 - e^{-c} \right] \right)^m \tag{B.9}
 \end{aligned}$$

B.16. Proof of Lemmas B.16.1 and B.16.2

$$\begin{aligned} &= \left(1 - R_{\mathcal{D}}^{\ell}(h) [1 - e^{-c}]\right)^{-m} \left(1 - R_{\mathcal{D}}^{\ell}(h) [1 - e^{-c}]\right)^m \quad (\text{B.10}) \\ &= 1, \end{aligned}$$

where we apply Lemma B.16.4 to obtain Equation (B.8), with the Binomial theorem we have Equation (B.9), and we deduce Equation (B.10) by definition of $F(R_{\mathcal{D}}^{\ell}(h))$. ■

APPENDIX OF CHAPTER 3

C.1 Proof of Proposition 3.3.1

Proposition 3.3.1 (Relations Between the Averaged Adversarial Risks). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , for any $(n, n') \in \mathbb{N}^2$, with $1 \leq n' \leq n$, we have

$$R_{\mathcal{E}}(\text{MV}_{\rho}) \leq A_{\mathcal{E}^{n'}}(\text{MV}_{\rho}) \leq A_{\mathcal{E}^n}(\text{MV}_{\rho}) \leq A_{\mathcal{D}}(\text{MV}_{\rho}). \quad (3.2)$$

Proof. First, we prove $A_{\mathcal{E}^1}(\text{MV}_{\rho}) = R_{\mathcal{E}}(\text{MV}_{\rho})$. We have

$$\begin{aligned} A_{\mathcal{E}^1}(\text{MV}_{\rho}) &= 1 - \mathbb{P}_{((\mathbf{x}, y), \mathfrak{C}) \sim \mathcal{E}^1} (\forall \epsilon \in \mathfrak{C}, \text{MV}_{\rho}(\mathbf{x} + \epsilon) = y) \\ &= 1 - \mathbb{P}_{((\mathbf{x}, y), \mathfrak{C}) \sim \mathcal{E}^1} (\forall \epsilon \in \{\epsilon_1\}, \text{MV}_{\rho}(\mathbf{x} + \epsilon) = y) \\ &= 1 - \mathbb{P}_{((\mathbf{x}, y), \mathfrak{C}) \sim \mathcal{E}^1} (\text{MV}_{\rho}(\mathbf{x} + \epsilon_1) = y) = R_{\mathcal{E}}(\text{MV}_{\rho}). \end{aligned}$$

Then, we prove the inequality $A_{\mathcal{E}^{n'}}(\text{MV}_{\rho}) \leq A_{\mathcal{E}^n}(\text{MV}_{\rho})$ from the fact that the indicator function $I[\cdot]$ is upper-bounded by 1. Indeed, from Definition 3.2.3 we have

$$\begin{aligned} 1 - A_{\mathcal{E}^n}(\text{MV}_{\rho}) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{\mathfrak{C} \sim \mathcal{B}_{(\mathbf{x}, y)}^n} I[\forall \epsilon \in \mathfrak{C}, \text{MV}_{\rho}(\mathbf{x} + \epsilon) = y] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\prod_{i=1}^n \mathbb{E}_{\epsilon_i \sim \mathcal{B}_{(\mathbf{x}, y)}} I[\text{MV}_{\rho}(\mathbf{x} + \epsilon_i) = y] \right] \\ &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[\prod_{i=1}^{n'} \mathbb{E}_{\epsilon_i \sim \mathcal{B}_{(\mathbf{x}, y)}} I[\text{MV}_{\rho}(\mathbf{x} + \epsilon_i) = y] \right] \\ &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{\mathfrak{C}' \sim \mathcal{B}_{(\mathbf{x}, y)}^{n'}} I[\forall \epsilon \in \mathfrak{C}', \text{MV}_{\rho}(\mathbf{x} + \epsilon) = y] \\ &= 1 - A_{\mathcal{E}^{n'}}(\text{MV}_{\rho}). \end{aligned}$$

Lastly, to prove the right-most inequality, we have to use the fact that the expectation over the set \mathbb{B} is bounded by the maximum over the set \mathbb{B} . We have

$$A_{\mathcal{E}^n}(\text{MV}_{\rho}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{E}_{\epsilon_1 \sim \mathcal{B}_{(\mathbf{x}, y)}} \dots \mathbb{E}_{\epsilon_n \sim \mathcal{B}_{(\mathbf{x}, y)}} I[\exists \epsilon \in \{\epsilon_1, \dots, \epsilon_n\}, \text{MV}_{\rho}(\mathbf{x} + \epsilon) \neq y]$$

$$\begin{aligned}
 &\leq \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \max_{\epsilon_1 \in \mathbb{B}} \dots \max_{\epsilon_n \in \mathbb{B}} \mathbb{I}[\exists \epsilon \in \{\epsilon_1, \dots, \epsilon_n\}, \text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] \\
 &= \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \max_{\epsilon_1 \in \mathbb{B}} \dots \max_{\epsilon_{n-1} \in \mathbb{B}} \mathbb{I}[\exists \epsilon \in \{\epsilon_1, \dots, \epsilon^*\}, \text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] \\
 &= \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon^*) \neq y] \\
 &= \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \max_{\epsilon \in \mathbb{B}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] = A_{\mathcal{D}}(\text{MV}_\rho).
 \end{aligned}$$

Merging the three equations proves the claim. ■

C.2 Proof of Proposition 3.3.2

In this section, we provide the proof of Proposition 3.3.2 that relies on Lemmas C.2.1 and C.2.2 which are also described and proved. Lemma C.2.1 shows that $R_{\mathcal{E}}(\text{MV}_\rho)$ equals $R_{\Gamma}(\text{MV}_\rho)$.

Lemma C.2.1. For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$ and its associated distribution Γ , for any posterior ρ on \mathbb{H} , we have

$$R_{\mathcal{E}}(\text{MV}_\rho) = \Pr_{(\mathbf{x}+\epsilon, y) \sim \Gamma} [\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] = R_{\Gamma}(\text{MV}_\rho).$$

Proof. Starting from the averaged risk $R_{\mathcal{E}}(\text{MV}_\rho) = \mathbb{E}_{((\mathbf{x},y), \epsilon) \sim \mathcal{E}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y]$, we have

$$\begin{aligned}
 R_{\mathcal{E}}(\text{MV}_\rho) &= \mathbb{E}_{(\mathbf{x}'+\epsilon', y') \sim \Gamma} \frac{1}{\Gamma(\mathbf{x}'+\epsilon', y')} \left[\Pr_{((\mathbf{x},y), \epsilon) \sim \mathcal{E}} [\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y, \mathbf{x}'+\epsilon' = \mathbf{x} + \epsilon, y' = y] \right] \\
 &= \mathbb{E}_{(\mathbf{x}'+\epsilon', y') \sim \Gamma} \frac{1}{\Gamma(\mathbf{x}'+\epsilon', y')} \left[\mathbb{E}_{((\mathbf{x},y), \epsilon) \sim \mathcal{E}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] \mathbb{I}[\mathbf{x}'+\epsilon' = \mathbf{x} + \epsilon, y' = y] \right].
 \end{aligned}$$

In other words, the double expectation only rearranges the terms of the original expectation: given an example $(\mathbf{x}'+\epsilon', y')$, we gather probabilities such that $\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y$ with $(\mathbf{x} + \epsilon, y) = (\mathbf{x}'+\epsilon', y')$ in the inner expectation, while integrating over all couple $(\mathbf{x}'+\epsilon', y') \in \mathbb{X} \times \mathbb{Y}$ in the outer expectation. Then, from the fact

C.2. Proof of Proposition 3.3.2

that when $\mathbf{x}' + \epsilon' = \mathbf{x} + \epsilon$ and $y' = y$, $\mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] = \mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y']$, we have

$$\begin{aligned} R_{\mathcal{E}}(\text{MV}_\rho) &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \Gamma} \frac{1}{\Gamma(\mathbf{x}' + \epsilon', y')} \left[\mathbb{E}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} \mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] \mathbb{I}[\mathbf{x}' + \epsilon' = \mathbf{x} + \epsilon, y' = y] \right] \\ &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \Gamma} \frac{1}{\Gamma(\mathbf{x}' + \epsilon', y')} \left[\mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] \mathbb{E}_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} \mathbb{I}[\mathbf{x}' + \epsilon' = \mathbf{x} + \epsilon, y' = y] \right]. \end{aligned}$$

Finally, by definition of $\Gamma(\mathbf{x}' + \epsilon', y')$, we can deduce that

$$\begin{aligned} R_{\mathcal{E}}(\text{MV}_\rho) &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \Gamma} \frac{1}{\Gamma(\mathbf{x}' + \epsilon', y')} [\mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] \Gamma(\mathbf{x}' + \epsilon', y')] \\ &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \Gamma} \mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] = R_{\Gamma}(\text{MV}_\rho). \end{aligned}$$

■

Similarly, Lemma C.2.2 shows that $A_{\mathcal{D}}(\text{MV}_\rho)$ is equivalent to $R_{\gamma}(\text{MV}_\rho)$.

Lemma C.2.2. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$ and its associated distribution γ , for any posterior ρ on \mathbb{H} , we have

$$A_{\mathcal{D}}(\text{MV}_\rho) = \Pr_{(\mathbf{x} + \epsilon, y) \sim \gamma} [\text{MV}_\rho(\mathbf{x} + \epsilon) \neq y] = R_{\gamma}(\text{MV}_\rho).$$

Proof. The proof is similar to the one of Lemma C.2.1. Indeed, starting from the definition of $A_{\mathcal{D}}(\text{MV}_\rho) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon^*(\mathbf{x}, y)) \neq y]$, we have

$$\begin{aligned} A_{\mathcal{D}}(\text{MV}_\rho) &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \gamma} \frac{1}{\gamma(\mathbf{x}' + \epsilon', y')} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\text{MV}_\rho(\mathbf{x} + \epsilon^*(\mathbf{x}, y)) \neq y] \mathbb{I}[\mathbf{x}' + \epsilon' = \mathbf{x} + \epsilon^*(\mathbf{x}, y), y' = y] \right] \\ &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \gamma} \frac{1}{\gamma(\mathbf{x}' + \epsilon', y')} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] \mathbb{I}[\mathbf{x}' + \epsilon' = \mathbf{x} + \epsilon^*(\mathbf{x}, y), y' = y] \right]. \end{aligned}$$

Finally, by definition of $\gamma(\mathbf{x}' + \epsilon', y')$, we can deduce that

$$\begin{aligned} A_{\mathcal{D}}(\text{MV}_\rho) &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \gamma} \frac{1}{\gamma(\mathbf{x}' + \epsilon', y')} [\mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] \gamma(\mathbf{x}' + \epsilon', y')] \\ &= \mathbb{E}_{(\mathbf{x}' + \epsilon', y') \sim \gamma} \mathbb{I}[\text{MV}_\rho(\mathbf{x}' + \epsilon') \neq y'] = R_{\gamma}(\text{MV}_\rho). \end{aligned}$$

■

We can now prove Proposition 3.3.2.

Proposition 3.3.2 (Classical and Averaged Adversarial Risks). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any distribution ρ on \mathbb{H} , we have

$$A_{\mathcal{D}}(\text{MV}_{\rho}) - \text{TV}(\gamma \parallel \Gamma) \leq R_{\mathcal{E}}(\text{MV}_{\rho}),$$

where Γ and γ are distributions on $\mathbb{X} \times \mathbb{Y}$ and $\text{TV}(\gamma \parallel \Gamma) = \mathbb{E}_{(\mathbf{x}', y') \sim \Gamma} \frac{1}{2} \left| \frac{\gamma(\mathbf{x}', y')}{\Gamma(\mathbf{x}', y')} - 1 \right|$, is the Total Variation (TV) distance between γ and Γ .

The density $\Gamma(\mathbf{x}', y')$ corresponds to the probability of drawing a perturbed example $(\mathbf{x}', y') = (\mathbf{x} + \epsilon, y)$ with $((\mathbf{x}, y), \epsilon) \sim \mathcal{E}$, i.e., we have

$$\Gamma(\mathbf{x}', y') = \Pr_{((\mathbf{x}, y), \epsilon) \sim \mathcal{E}} [\mathbf{x} + \epsilon = \mathbf{x}', y = y'].$$

The density $\gamma(\mathbf{x}', y')$ is the probability to draw an adversarial example $(\mathbf{x}', y') = (\mathbf{x} + \epsilon^*(\mathbf{x}, y), y)$ with $(\mathbf{x}, y) \sim \mathcal{D}$, i.e., we have

$$\gamma(\mathbf{x}', y') = \Pr_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathbf{x} + \epsilon^*(\mathbf{x}, y) = \mathbf{x}', y = y'].$$

Proof. From Lemmas C.2.1 and C.2.2, we have

$$R_{\mathcal{E}}(\text{MV}_{\rho}) = R_{\Gamma}(\text{MV}_{\rho}), \quad \text{and} \quad A_{\mathcal{D}}(\text{MV}_{\rho}) = R_{\gamma}(\text{MV}_{\rho}).$$

Then, we apply Lemma 4 of OHNISHI and HONORIO (2021), we have

$$R_{\gamma}(\text{MV}_{\rho}) \leq \text{TV}(\gamma \parallel \Gamma) + R_{\Gamma}(\text{MV}_{\rho}) \iff A_{\mathcal{D}}(\text{MV}_{\rho}) \leq \text{TV}(\gamma \parallel \Gamma) + R_{\mathcal{E}}(\text{MV}_{\rho}).$$

■

C.3 Proof of Theorem 3.3.1

Theorem 3.3.1 (Upper Bounds on the Surrogates). For any distributions \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$ and ρ on \mathbb{H} , for any $n > 1$, we have

$$R_{\mathcal{E}}(\text{MV}_{\rho}) \leq 2r_{\mathcal{E}}(\rho), \quad \text{and} \quad A_{\mathcal{E}^n}(\text{MV}_{\rho}) \leq 2a_{\mathcal{E}^n}(\rho).$$

Proof. By the definition of the majority vote and from MARKOV's inequality (The-

orem A.2.1), we have

$$\begin{aligned}
 \frac{1}{2}R_{\mathcal{E}}(\text{MV}_{\rho}) &= \frac{1}{2} \mathbb{P}_{((\mathbf{x},y),\epsilon) \sim \mathcal{E}} \left(y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \leq 0 \right) \\
 &= \frac{1}{2} \mathbb{P}_{((\mathbf{x},y),\epsilon) \sim \mathcal{E}} \left(1 - y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \geq 1 \right) \\
 &\leq \mathbb{E}_{((\mathbf{x},y),\epsilon) \sim \mathcal{E}} \frac{1}{2} \left[1 - y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \right] \\
 &= r_{\mathcal{E}}(\rho).
 \end{aligned}$$

Similarly we have

$$\begin{aligned}
 \frac{1}{2}A_{\mathcal{E}^n}(\text{MV}_{\rho}) &= \frac{1}{2} \mathbb{P}_{((\mathbf{x},y),\mathfrak{C}) \sim \mathcal{E}^n} \left(\exists \epsilon \in \mathfrak{C}, y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \leq 0 \right) \\
 &= \frac{1}{2} \mathbb{P}_{((\mathbf{x},y),\mathfrak{C}) \sim \mathcal{E}^n} \left(\min_{\epsilon \in \mathfrak{C}} \left(y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \right) \leq 0 \right) \\
 &= \frac{1}{2} \mathbb{P}_{((\mathbf{x},y),\epsilon) \sim \mathcal{E}} \left(1 - \min_{\epsilon \in \mathfrak{C}} \left(y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \right) \geq 1 \right) \\
 &\leq \mathbb{E}_{((\mathbf{x},y),\epsilon) \sim \mathcal{E}} \frac{1}{2} \left[1 - \min_{\epsilon \in \mathfrak{C}} \left(y \mathbb{E}_{h \sim \rho} h(\mathbf{x} + \epsilon) \right) \right] \\
 &= a_{\mathcal{E}^n}(\rho).
 \end{aligned}$$

■

C.4 Proof of Theorem 3.3.2

Theorem 3.3.2 (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any set of voters \mathbb{H} , for any prior $\pi \in \mathbb{M} * (\mathbb{H})$ on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathbb{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , we have

$$\text{kl}(r_{\widehat{\mathbb{S}}}(\rho) \| r_{\mathcal{E}}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{m+1}{\delta} \right], \quad (3.6)$$

$$\text{and } r_{\mathcal{E}}(\rho) \leq r_{\widehat{\mathbb{S}}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{m+1}{\delta} \right]}, \quad (3.7)$$

$$\text{where } r_{\widehat{\mathbb{S}}}(\rho) = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n \frac{1}{2} \left[1 - y_i \mathbb{E}_{h \sim \rho} h(\mathbf{x}_i + \epsilon_j^i) \right].$$

Proof. Let $G=(V, E)$ be the graph representing the dependencies between the random variables where (i) the set of vertices is $V=\widehat{\mathcal{S}}$, (ii) the set of edges E is defined such that $((\mathbf{x}, y), \epsilon), ((\mathbf{x}', y'), \epsilon')) \notin E \Leftrightarrow x \neq \mathbf{x}'$. Then, applying Theorem 8 of RALAIVOLA *et al.* (2010) with our notations gives

$$\text{kl}(r_{\widehat{\mathcal{S}}}(\rho) \| r_{\mathcal{E}}(\rho)) \leq \frac{\chi(G)}{mn} \left[\text{KL}(\rho \| \pi) + \ln \frac{mn + \chi(G)}{\delta \chi(G)} \right],$$

where $\chi(G)$ is the fractional chromatic number of G . From a property of SCHEINERMAN and ULLMAN (2011), we have

$$c(G) \leq \chi(G) \leq \Delta(G) + 1,$$

where $c(G)$ is the order of the largest clique in G and $\Delta(G)$ is the maximum degree of a vertex in G . By construction of G , $c(G)=n$ and $\Delta(G)=n-1$. Thus, $\chi(G)=n$ and rearranging the terms proves Equation (3.6). Finally, by applying PINSKER's inequality (i.e., $|a-b| \leq \sqrt{\frac{1}{2} \text{kl}(a \| b)}$), we obtain Equation (3.7). ■

C.5 Proof of Theorem 3.3.3

Theorem 3.3.3 (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathcal{B}$, for any set of voters \mathbb{H} , for any prior $\pi \in \mathbb{M}^*(\mathbb{H})$ on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1-\delta$ over $\widehat{\mathcal{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for all $i \in \{1, \dots, m\}$, for all distributions Θ_i on \mathcal{C}_i independent from a voter $h \in \mathbb{H}$, we have

$$a_{\mathcal{E}^n}(\rho) \leq \frac{1}{m} \mathbb{E}_{h \sim \rho} \sum_{i=1}^m \max_{\epsilon \in \mathcal{C}_i} \frac{1}{2} (1 - y_i h(\mathbf{x}_i + \epsilon)) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]} \quad (3.8)$$

$$\leq a_{\widehat{\mathcal{S}}}(\rho) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \| \Theta_i) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}}{\delta} \right]}, \quad (3.9)$$

where the empirical risk $a_{\widehat{\mathcal{S}}}(\rho) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \left[1 - \min_{\epsilon \in \mathcal{C}_i} \left(y_i \mathbb{E}_{h \sim \rho} h(\mathbf{x}_i + \epsilon) \right) \right]$, and the TV distance $\text{TV}(\theta \| \Theta) = \mathbb{E}_{\epsilon \sim \Theta} \frac{1}{2} \left| \left[\frac{\theta(\epsilon)}{\Theta(\epsilon)} \right] - 1 \right|$.

Proof. Let $L_{h,(\mathbf{x},y),\epsilon} = \frac{1}{2} \left[1 - y h(\mathbf{x} + \epsilon) \right]$ for the sake of readability. Given $h \in \mathbb{H}$, the losses $\max_{\epsilon \in \mathcal{C}_1} L_{h,(\mathbf{x}_1,y_1),\epsilon}, \dots, \max_{\epsilon \in \mathcal{C}_m} L_{h,(\mathbf{x}_m,y_m),\epsilon}$ are *i.i.d.*. Hence, we can

apply Theorem 20 of GERMAIN *et al.* (2015) and PINSKER's inequality, *i.e.*, the inequality $|q-p| \leq \sqrt{\frac{1}{2} \text{kl}(q||p)}$ (Theorem B.5.1) to obtain

$$\mathbb{E}_{h \sim \rho} \mathbb{E}_{(\mathbf{x}, y), \mathfrak{C}} \max_{\epsilon \in \mathfrak{C}} L_{h, (\mathbf{x}, y), \epsilon} \leq \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \max_{\epsilon \in \mathfrak{C}_i} L_{h, (\mathbf{x}_i, y_i), \epsilon} + \sqrt{\frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{m}}{\delta}}{2m}}.$$

Then, we lower-bound the left-hand side of the inequality with $a_{\mathcal{E}^n}(\rho)$, we have

$$a_{\mathcal{E}^n}(\rho) \leq \mathbb{E}_{h \sim \rho} \mathbb{E}_{((\mathbf{x}, y), \mathfrak{C}) \sim \mathcal{E}^n} \max_{\epsilon \in \mathfrak{C}} L_{h, (\mathbf{x}, y), \epsilon}.$$

Finally, from the definition of θ_i^h , and from Lemma 4 of OHNISHI and HONORIO (2021), we have

$$\begin{aligned} \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \max_{\epsilon \in \mathfrak{C}_i} L_{h, (\mathbf{x}_i, y_i), \epsilon} &= \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\epsilon \sim \theta_i^h} L_{h, (\mathbf{x}_i, y_i), \epsilon} \\ &\leq \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \text{TV}(\theta_i^h || \Theta_i) + \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\epsilon \sim \Theta_i} L_{h, (\mathbf{x}_i, y_i), \epsilon} \\ &= \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \text{TV}(\theta_i^h || \Theta_i) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\epsilon \sim \Theta_i} \mathbb{E}_{h \sim \rho} L_{h, (\mathbf{x}_i, y_i), \epsilon} \\ &\leq \mathbb{E}_{h \sim \rho} \frac{1}{m} \sum_{i=1}^m \text{TV}(\theta_i^h || \Theta_i) + a_{\widehat{\mathcal{S}}}(\rho). \end{aligned}$$

■

C.6 Proof of Corollaries 3.3.1 and 3.3.2

We start to prove Corollary 3.3.1.

Corollary 3.3.1 (PAC-Bayesian Bound on $r_{\mathcal{E}}(\rho)$). For any distribution \mathcal{E} on $(\mathcal{X} \times \mathcal{Y}) \times \mathbb{B}$, for any set of voters \mathbb{H} , for any $T \in \mathbb{N}_*$, for any priors' set $\{\pi_1, \dots, \pi_T\} \in \mathbb{M}^*(\mathbb{H})^T$, for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathcal{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for any $\pi \in \{\pi_1, \dots, \pi_T\} \in \mathbb{M}^*(\mathbb{H})^T$ we have

$$r_{\mathcal{E}}(\rho) \leq \overline{\text{kl}} \left(r_{\widehat{\mathcal{S}}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho || \pi) + \ln \frac{T(m+1)}{\delta} \right] \right. \right). \quad (3.10)$$

Proof. Let $\mathcal{E}_1, \dots, \mathcal{E}_T$ be T distributions defined as $\mathcal{E}_1 = \mathcal{D}(\mathbf{x}, y)\mathcal{B}_{(x,y)}^1(\epsilon)$, \dots , $\mathcal{E}_T = \mathcal{D}(\mathbf{x}, y)\mathcal{B}_{(x,y)}^T(\epsilon)$ on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$ where each distribution $\mathcal{B}_{(x,y)}^t$ depends on the example (\mathbf{x}, y) and possibly on the fixed prior π_t . Then, for all distributions \mathcal{E}_t , we can derive a bound on the risk $r_{\mathcal{E}_t}(\rho)$ which holds with probability at least $1 - \frac{\delta}{T}$, we have

$$\Pr_{\widehat{\mathcal{S}}_t \sim (\mathcal{E}_t^n)^m} \left[\forall \rho, \text{kl}(r_{\widehat{\mathcal{S}}_t}(\rho) \| r_{\mathcal{E}_t}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi_t) + \ln \frac{T(m+1)}{\delta} \right] \right] \geq 1 - \frac{\delta}{T}.$$

Then, from a union bound argument, we have

$$\Pr_{\widehat{\mathcal{S}}_1 \sim (\mathcal{E}_1^n)^m, \dots, \widehat{\mathcal{S}}_T \sim (\mathcal{E}_T^n)^m} \left[\forall \rho, \text{kl}(r_{\widehat{\mathcal{S}}_1}(\rho) \| r_{\mathcal{E}_1}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi_1) + \ln \frac{T(m+1)}{\delta} \right], \right. \\ \left. \dots, \text{ and } \text{kl}(r_{\widehat{\mathcal{S}}_T}(\rho) \| r_{\mathcal{E}_T}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi_T) + \ln \frac{T(m+1)}{\delta} \right] \right] \geq 1 - \delta.$$

Hence, we have

$$\text{kl}(r_{\widehat{\mathcal{S}}}(\rho) \| r_{\mathcal{E}}(\rho)) \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \frac{T(m+1)}{\delta} \right],$$

where $\mathcal{B}_{(x,y)}$ can be dependent on the selected prior π . From Definition 2.3.3, we can obtain the claimed result. \blacksquare

We can prove Corollary 3.3.2 similarly to Corollary 3.3.1.

Corollary 3.3.2 (PAC-Bayesian Bound on $a_{\mathcal{E}^n}(\rho)$). For any distribution \mathcal{E} on $(\mathbb{X} \times \mathbb{Y}) \times \mathbb{B}$, for any set of voters \mathbb{H} , for any prior π on \mathbb{H} , for any $n \in \mathbb{N}_*$, with probability at least $1 - \delta$ over $\widehat{\mathcal{S}} \sim (\mathcal{E}^n)^m$, for all posteriors $\rho \in \mathbb{M}(\mathbb{H})$ on \mathbb{H} , for all $i \in \{1, \dots, m\}$, for all distributions Θ_i on \mathbb{E}_i independent from a voter $h \in \mathbb{H}$, we have

$$a_{\mathcal{E}^n}(\rho) \leq a_{\widehat{\mathcal{S}}}(\rho) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{h \sim \rho} \text{TV}(\theta_i^h \| \Theta_i) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \| \pi) + \ln \frac{2T\sqrt{m}}{\delta} \right]}. \quad (3.11)$$

Proof. From a union bound argument, we obtain the claimed result. \blacksquare

C.7 About the (Differentiable) Decision Trees

In this section, we introduce the differentiable decision trees, *i.e.*, the voters of our majority vote. Note that we adapt the model of KONTSCIEDER *et al.* (2016) in order to fit with our framework: a voter must output a real between -1 and $+1$. An example of such a tree is represented in Figure C.1.

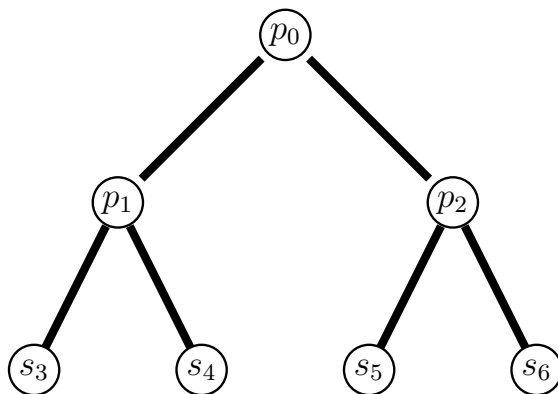


Figure C.1. Representation of a (differentiable) decision tree of depth $l = 2$; The root is the node 0 and the leafs are 4; 5; 6 and 7. The probability $p_i(\mathbf{x})$ (respectively $1-p_i(\mathbf{x})$) to go left (respectively right) at the node i is represented by p_i (we omitted the dependence on \mathbf{x} for simplicity). Similarly, the predicted label (a “score” between -1 and $+1$) at the leaf i is represented by s_i .

This differentiable decision tree is stochastic by nature: at each node i of the tree, we continue recursively to the left sub-tree with a probability of $p_i(\mathbf{x})$ and to the right sub-tree with a probability of $1-p_i(\mathbf{x})$; When we attain a leaf j , the tree predicts the label s_j . Precisely, the probability $p_i(\mathbf{x})$ is constructed by (i) selecting randomly 50% of the input features \mathbf{x} and applying a random mask $M_i \in \mathbb{R}^d$ on \mathbf{x} (where the k -th entry of the mask is 1 if the k -th feature is selected and 0 otherwise), by (ii) multiplying this quantity by a learned weight vector $v_i \in \mathbb{R}^d$, and by (iii) applying a sigmoid function to output a probability. Indeed, we have

$$p_i(\mathbf{x}) = \sigma\left(\langle v_i, M_i \odot \mathbf{x} \rangle\right),$$

where $\sigma(a) = [1 + e^{-a}]^{-1}$ is the sigmoid function; $\langle a, b \rangle$ is the dot product between the vector a and b and $a \odot b$ is the elementwise product between the vector a and b . Moreover, s_i is obtained by learning a parameter $u_i \in \mathbb{R}$ and applying a tanh function, *i.e.*, we have

$$s_i = \tanh(u_i).$$

Finally, instead of having a stochastic voter, h will output the expected label predicted by the tree (see KONTSCIEDER *et al.* (2016) for more details). It can be computed by $h(\mathbf{x}) = f(\mathbf{x}, 0, 0)$ with

$$f(\mathbf{x}, i, l') = \begin{cases} s_i & \text{if } l' = l \\ p_i(\mathbf{x})f(\mathbf{x}, 2i+1, l'+1) + (1 - p_i(\mathbf{x}))f(\mathbf{x}, 2i+2, l'+1) & \text{otherwise} \end{cases} .$$

C.8 Additional Experimental Results

In this section, we present the detailed results for the 6 tasks (3 on MNIST and 3 on Fashion MNIST) on which we perform experiments that show the test risks and the bounds for the different scenarios of (Defense, Attack). We train all the models using the same parameters as described in Section 3.4.1. Table C.1 and Appendix C.8 complement Table 3.1 to present the results for all the tasks when using the ℓ_2 -norm with $b = 1$ (the maximum noise allowed by the norm). Then, we run again the same experiment but we use the ℓ_∞ -norm with $b = 0.1$ and exhibit the results in Appendix C.8 and Table C.6. For the experiments on the 5 other tasks using the ℓ_2 -norm, we have a similar behavior than MNIST:1vs7. Indeed, using the attacks PGD_U and IFGSM_U as defense mechanism allows to obtain better risks and also tighter bounds compared to the bounds obtained with a defense based on UNIF (which is a naive defense). For the experiments on the 6 tasks using the ℓ_∞ -norm, the trend is the same as with the ℓ_2 -norm, *i.e.*, the appropriate defense leads to better risks and bounds.

We also run experiments that do not rely on the PAC-Bayesian framework. In other words, we train the models following only Step 1 of our adversarial training procedure (*i.e.*, Algorithm 3.1) using classical attacks (PGD or IFGSM): we refer to this experiment as a baseline. In our cases, it means learning a majority vote $MV_{\pi'}$ that follows a distribution π' . As a reminder, the studied scenarios for the baseline are all the pairs (Defense, Attack) belonging to the set $\{\text{---}, \text{UNIF}, \text{PGD}, \text{IFGSM}\} \times \{\text{---}, \text{PGD}, \text{IFGSM}\}$. We report the results in Table C.8 and Table C.9. With this experiment, we are now able to compare our defense based on PGD_U or IFGSM_U and a classical defense based on PGD and IFGSM. Hence, considering the test risks $A_{\mathcal{T}}(MV_{\rho})$ (columns “Attack without U” of Tables 3.1 to C.6) and $A_{\mathcal{T}}(MV_{\pi'})$ (in Tables C.8 and C.9), we observe similar results between the baseline and our framework.

C.8. Additional Experimental Results

Table C.1. Test risks and bounds for 2 tasks of MNIST with $n=100$ perturbations for all pairs (Defense,Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} . The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} . To quantify the gap between our risks and the classical definition we put in italic the risk of our models against the classical attacks: we replace PGD_U and IFGSM_U by PGD or IFGSM (i.e., we did not sample from the uniform distribution). Since Eq. (3.9) upperbounds Eq. (3.8) thanks to the TV term, we compute the two bound values of Theorem 3.3.3.

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_T(MV_\rho)$		$R_{\bar{T}}(MV_\rho)$		Th. 3.3.2		$A_T(MV_\rho)$		$A_{\bar{T}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015	.015	.015	0.060	.067	.015	.015	.015	.015	0.129	0.135	0.129	.135
—	PGD_U	.632	.628	.520	.526	1.059	.847	.672	.641	.683	.684	1.718	2.405	1.392	.962
—	IFGSM_U	.447	.443	.157	.166	0.387	.572	.461	.451	.337	.345	1.137	2.090	0.776	.669
UNIF	—	.024	.024	.024	.024	0.073	.083	.024	.024	.024	.024	0.140	0.148	0.140	.148
UNIF	PGD_U	.646	.619	.486	.500	1.016	.809	.649	.626	.648	.650	1.646	2.417	1.338	.915
UNIF	IFGSM_U	.442	.442	.128	.139	0.316	.528	.442	.442	.281	.293	0.907	2.118	0.633	.617
PGD_U	—	.024	.025	.024	.025	0.094	.101	.024	.025	.024	.025	0.158	0.163	0.158	.163
PGD_U	PGD_U	.148	.135	.111	.103	0.360	.355	.146	.136	.129	.120	0.442	2.062	0.414	.403
PGD_U	IFGSM_U	.104	.103	.072	.072	0.277	.277	.102	.102	.090	.084	0.358	1.954	0.335	.328
IFGSM_U	—	.027	.025	.027	.025	0.080	.091	.027	.025	.027	.025	0.146	0.154	0.146	.154
IFGSM_U	PGD_U	.188	.178	.111	.119	0.383	.405	.190	.178	.126	.134	0.501	2.063	0.454	.454
IFGSM_U	IFGSM_U	.126	.115	.076	.070	0.248	.290	.127	.115	.091	.085	0.371	1.918	0.329	.342

(a) MNIST 4vs9

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_T(MV_\rho)$		$R_{\bar{T}}(MV_\rho)$		Th. 3.3.2		$A_T(MV_\rho)$		$A_{\bar{T}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015	.015	.015	.043	.045	.015	.015	.015	.015	.117	0.118	.117	.118
—	PGD_U	.279	.271	.232	.234	.600	.453	.284	.274	.284	.284	.829	1.929	.724	.524
—	IFGSM_U	.143	.137	.089	.090	.204	.227	.144	.139	.125	.127	.422	1.662	.337	.293
UNIF	—	.017	.017	.017	.017	.054	.055	.017	.017	.017	.017	.124	0.125	.124	.125
UNIF	PGD_U	.219	.201	.172	.177	.433	.350	.219	.209	.217	.218	.671	1.810	.565	.419
UNIF	IFGSM_U	.122	.122	.052	.055	.119	.181	.122	.123	.077	.082	.307	1.554	.242	.248
PGD_U	—	.013	.015	.013	.015	.061	.061	.013	.015	.013	.015	.131	0.130	.131	.130
PGD_U	PGD_U	.057	.057	.045	.041	.157	.160	.057	.057	.055	.045	.227	1.536	.218	.218
PGD_U	IFGSM_U	.043	.043	.027	.031	.114	.119	.042	.043	.037	.035	.187	1.433	.179	.181
IFGSM_U	—	.014	.012	.014	.012	.057	.057	.014	.013	.014	.013	.128	0.127	.128	.127
IFGSM_U	PGD_U	.077	.072	.054	.043	.170	.174	.076	.075	.055	.052	.252	1.510	.233	.236
IFGSM_U	IFGSM_U	.055	.048	.034	.030	.105	.121	.052	.051	.039	.032	.191	1.379	.177	.185

(b) MNIST 5vs6

Table C.2. Test risks and bounds for 3 tasks Fashion MNIST with $n=100$ perturbations for all pairs (Defense,Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} . The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} . To quantify the gap between our risks and the classical definition we put in *italic* the risk of our models against the classical attacks: we replace PGD_U and IFGSM_U by PGD or IFGSM (i.e., we did not sample from the uniform distribution). Since Eq. (3.9) upperbounds Eq. (3.8) thanks to the TV term, we compute the two bound values of Theorem 3.3.3.

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_{\rho})$		$R_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.2		$A_{\mathcal{T}}(MV_{\rho})$		$A_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.021	.020	.021	.020	0.060	0.070	.019	<i>.019</i>	.019	.019	0.130	0.139	0.130	0.139
—	PGD_U	<i>.695</i>	.650	.494	.568	1.042	1.090	.677	<i>.686</i>	.588	.674	1.326	2.307	1.152	1.082
—	IFGSM_U	<i>.451</i>	<i>.451</i>	.269	.328	0.585	0.731	.405	<i>.438</i>	.295	.381	0.878	1.971	0.730	0.746
UNIF	—	<i>.071</i>	<i>.071</i>	.071	.071	0.185	0.191	<i>.071</i>	<i>.071</i>	.071	.071	0.236	0.241	0.236	0.241
UNIF	PGD_U	.423	<i>.477</i>	.418	.425	0.957	0.755	<i>.486</i>	<i>.486</i>	.513	.513	1.372	2.173	1.151	0.869
UNIF	IFGSM_U	.326	<i>.331</i>	.105	.105	0.273	0.422	<i>.333</i>	.331	.144	.142	0.496	1.642	0.397	0.504
PGD_U	—	<i>.034</i>	.032	<i>.034</i>	.032	0.094	0.114	<i>.034</i>	.032	<i>.034</i>	.032	0.158	0.174	0.158	0.174
PGD_U	PGD_U	.103	<i>.115</i>	.086	.091	0.227	0.289	.102	<i>.115</i>	.096	.101	0.299	1.985	0.283	0.338
PGD_U	IFGSM_U	.092	<i>.099</i>	.073	.076	0.195	0.248	.092	<i>.099</i>	.082	.082	0.266	1.914	0.253	0.299
IFGSM_U	—	.028	<i>.030</i>	.028	.030	0.091	0.105	.027	<i>.030</i>	.027	.030	0.155	0.166	0.155	0.166
IFGSM_U	PGD_U	<i>.115</i>	.114	.085	.085	0.254	0.287	.112	<i>.114</i>	.096	.101	0.331	2.026	0.313	0.337
IFGSM_U	IFGSM_U	.095	<i>.097</i>	.067	.068	0.206	0.232	.093	<i>.097</i>	.080	.081	0.282	1.927	0.266	0.285

(a) Fashion MNIST Sandall vs Ankle Boot

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_{\rho})$		$R_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.2		$A_{\mathcal{T}}(MV_{\rho})$		$A_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	<i>.038</i>	.037	.038	.037	0.088	.091	<i>.038</i>	.037	.038	.037	.153	0.155	.153	.155
—	PGD_U	<i>.292</i>	.248	.233	.112	.452	.363	<i>.289</i>	.272	.287	.246	.578	1.314	.525	.479
—	IFGSM_U	<i>.194</i>	.154	.132	.075	.300	.262	<i>.193</i>	.181	.176	.148	.423	1.103	.376	.359
UNIF	—	<i>.039</i>	<i>.039</i>	.039	.039	0.091	.093	<i>.041</i>	.039	.041	.039	.155	0.157	.155	.157
UNIF	PGD_U	<i>.240</i>	.220	.099	.117	.346	.332	<i>.250</i>	.231	.250	.245	.553	1.228	.490	.443
UNIF	IFGSM_U	<i>.177</i>	.171	.070	.078	.228	.247	<i>.197</i>	.185	.186	.164	.445	1.046	.371	.346
PGD_U	—	<i>.045</i>	.044	.045	.044	.108	.105	<i>.046</i>	.045	.046	.045	.172	0.167	.172	.167
PGD_U	PGD_U	<i>.108</i>	.100	.077	.082	.203	.211	<i>.104</i>	.100	.081	.087	.279	1.118	.269	.264
PGD_U	IFGSM_U	<i>.094</i>	.086	.071	.069	.184	.186	<i>.090</i>	.086	.076	.073	.257	1.015	.248	.241
IFGSM_U	—	.041	<i>.043</i>	.041	.043	.094	.101	.039	<i>.042</i>	.039	.042	.158	0.163	.158	.163
IFGSM_U	PGD_U	.106	<i>.114</i>	.078	.092	.220	.226	.109	<i>.113</i>	.084	.095	.293	1.052	.279	.275
IFGSM_U	IFGSM_U	.082	<i>.087</i>	.065	.072	.171	.176	.082	<i>.089</i>	.068	.078	.247	0.927	.234	.232

(b) Fashion MNIST Top vs Pullover

C.8. Additional Experimental Results

ℓ_2 -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_{\rho})$		$R_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.2		$A_{\mathcal{T}}(MV_{\rho})$		$A_{\bar{\mathcal{T}}}(MV_{\rho})$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.122	.122	.122	.122	0.276	0.286	.122	.122	.122	.122	0.318	0.328	0.318	0.328
—	PGD _U	.744	.738	.674	.689	1.386	1.066	.745	.740	.767	.768	1.773	2.386	1.576	1.180
—	IFGSM _U	.652	.646	.454	.474	0.947	0.887	.659	.648	.618	.632	1.597	2.214	1.276	0.992
UNIF	—	.204	.204	.204	.204	0.444	0.444	.204	.204	.204	.204	0.475	0.476	0.475	0.476
UNIF	PGD _U	.750	.714	.682	.671	1.350	1.069	.750	.719	.752	.749	1.732	2.063	1.524	1.189
UNIF	IFGSM _U	.605	.575	.423	.431	0.871	0.866	.605	.578	.530	.526	1.304	1.860	1.091	0.956
PGD _U	—	.168	.165	.168	.165	0.423	0.428	.167	.165	.167	.165	0.463	0.461	0.463	0.460
PGD _U	PGD _U	.389	.402	.306	.369	0.768	0.719	.390	.402	.319	.403	0.847	2.354	0.810	0.755
PGD _U	IFGSM _U	.361	.368	.298	.324	0.693	0.672	.362	.368	.320	.361	0.799	2.258	0.754	0.707
IFGSM _U	—	.150	.163	.150	.163	0.424	0.428	.149	.163	.149	.163	0.458	0.461	0.458	0.461
IFGSM _U	PGD _U	.391	.428	.347	.292	0.778	0.757	.390	.426	.371	.298	0.856	2.327	0.820	0.791
IFGSM _U	IFGSM _U	.356	.382	.291	.273	0.685	0.689	.354	.382	.331	.278	0.772	2.218	0.734	0.723

(a) Fashion MNIST Coat vs Shirt

Table C.4. Test risks and bounds for 3 tasks of MNIST with $n=100$ perturbations for all pairs (Defense,Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} . The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} . To quantify the gap between our risks and the classical definition we put in italic the risk of our models against the classical attacks: we replace PGD_U and IFGSM_U by PGD or IFGSM (i.e., we did not sample from the uniform distribution). Since Eq. (3.9) upperbounds Eq. (3.8) thanks to the TV term, we compute the two bound values of Theorem 3.3.3.

ℓ_∞ -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 0.1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_\rho)$		$R_{\mathcal{F}}(MV_\rho)$		Th. 3.3.2		$A_{\mathcal{T}}(MV_\rho)$		$A_{\mathcal{F}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.005	.005	.005	.005	.017	.019	.005	.005	.005	.005	0.099	0.100	.099	.100
—	PGD_U	<i>.454</i>	<i>.454</i>	.375	.384	.770	.638	.492	.484	.480	.476	1.127	2.031	.946	.716
—	IFGSM_U	<i>.428</i>	.423	.350	.361	.727	.610	.474	.465	.448	.443	1.061	2.008	.886	.686
UNIF	—	.004	.004	.004	.004	.018	.019	.004	.004	.004	.004	0.099	0.100	.099	.100
UNIF	PGD_U	.487	.491	.369	.392	.779	.667	.512	.507	.487	.484	1.179	2.083	.972	.739
UNIF	IFGSM_U	.436	<i>.442</i>	.325	.337	.664	.598	.466	.459	.417	.417	1.023	1.959	.841	.671
PGD_U	—	<i>.006</i>	<i>.006</i>	.006	.006	.024	.024	.005	<i>.006</i>	.005	.006	0.103	0.103	.103	.103
PGD_U	PGD_U	.018	<i>.020</i>	.013	.016	.046	.050	.018	<i>.020</i>	.015	.020	0.127	1.461	.122	.123
PGD_U	IFGSM_U	.020	<i>.021</i>	.012	.016	.048	.054	.019	<i>.021</i>	.015	.020	0.130	1.455	.125	.127
IFGSM_U	—	.006	<i>.007</i>	.006	.007	.023	.024	.006	<i>.007</i>	.006	.007	0.102	0.103	.102	.103
IFGSM_U	PGD_U	.018	<i>.019</i>	.016	.016	.046	.051	.018	<i>.019</i>	.018	.019	0.126	1.489	.122	.124
IFGSM_U	IFGSM_U	<i>.020</i>	<i>.020</i>	.015	.016	.050	.055	<i>.020</i>	<i>.020</i>	.020	.019	0.131	1.481	.126	.127

(a) MNIST 1 vs 7

ℓ_∞ -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 0.1$		Attack without u						Attack without u							
Defense	Attack	$A_{\mathcal{T}}(MV_\rho)$		$R_{\mathcal{F}}(MV_\rho)$		Th. 3.3.2		$A_{\mathcal{T}}(MV_\rho)$		$A_{\mathcal{F}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015	.015	.015	0.060	0.067	.015	.015	.015	.015	0.129	0.135	0.129	0.135
—	PGD_U	.929	<i>.930</i>	.651	.662	1.367	1.125	.920	<i>.925</i>	.874	.880	2.213	2.661	1.792	1.266
—	IFGSM_U	<i>.935</i>	<i>.935</i>	.601	.609	1.243	1.088	.926	<i>.928</i>	.800	.806	2.047	2.615	1.649	1.224
UNIF	—	.017	.017	.017	.017	0.062	0.072	.017	.017	.017	.017	0.131	0.139	0.131	0.139
UNIF	PGD_U	<i>.895</i>	<i>.895</i>	.615	.623	1.302	1.078	.884	<i>.888</i>	.815	.818	2.035	2.722	1.670	1.208
UNIF	IFGSM_U	<i>.898</i>	<i>.898</i>	.516	.528	1.112	1.027	.884	<i>.890</i>	.697	.706	1.875	2.658	1.497	1.153
PGD_U	—	<i>.039</i>	.037	<i>.039</i>	.037	0.093	0.094	<i>.039</i>	.037	<i>.039</i>	.037	0.156	0.157	0.156	0.157
PGD_U	PGD_U	.108	<i>.109</i>	.090	.090	0.200	0.209	.108	<i>.109</i>	.110	.112	0.337	1.874	0.290	0.271
PGD_U	IFGSM_U	.121	<i>.124</i>	.101	.103	0.229	0.235	.121	<i>.124</i>	.126	.126	0.378	1.890	0.326	0.297
IFGSM_U	—	<i>.046</i>	.044	<i>.046</i>	.044	0.102	0.119	<i>.046</i>	.044	.046	.046	0.164	0.178	0.164	0.178
IFGSM_U	PGD_U	.105	.093	.091	.078	0.203	0.214	.105	.093	.108	.108	0.321	1.810	0.286	0.269
IFGSM_U	IFGSM_U	<i>.119</i>	.095	.102	.080	0.220	0.229	<i>.119</i>	.095	.122	.122	0.357	1.821	0.309	0.283

(b) MNIST 4 vs 9

C.8. Additional Experimental Results

ℓ_∞ -norm $b = 0.1$		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
Defense	Attack	Attack without U						Attack without U							
		$A_{\mathcal{T}}(MV_\rho)$		$R_{\mathcal{T}}(MV_\rho)$		Th. 3.3.2		$A_{\mathcal{T}}(MV_\rho)$		$A_{\mathcal{T}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015	.015	.015	.043	.045	.015	.015	.015	.015	0.117	0.118	0.117	.118
—	PGD _U	.500	.499	.387	.390	.923	.744	.502	.500	.474	.475	1.361	2.275	1.146	.830
—	IFGSM _U	.519	.505	.395	.398	.915	.762	.514	.516	.481	.481	1.335	2.283	1.129	.847
UNIF	—	.015	.015	.015	.015	.052	.053	.015	.015	.015	.015	0.123	0.124	0.123	.124
UNIF	PGD _U	.529	.544	.388	.393	.925	.761	.517	.532	.481	.482	1.342	2.349	1.137	.848
UNIF	IFGSM _U	.536	.544	.372	.379	.881	.774	.523	.544	.451	.456	1.268	2.348	1.077	.857
PGD _U	—	.015	.014	.015	.014	.060	.064	.015	.014	.015	.014	0.130	0.133	0.130	.133
PGD _U	PGD _U	.055	.058	.037	.039	.131	.143	.056	.057	.046	.046	0.219	1.619	0.202	.204
PGD _U	IFGSM _U	.061	.065	.040	.043	.146	.154	.059	.062	.050	.046	0.232	1.626	0.216	.214
IFGSM _U	—	.019	.014	.019	.014	.069	.064	.018	.014	.018	.014	0.136	0.132	0.136	.132
IFGSM _U	PGD _U	.061	.061	.040	.050	.143	.142	.061	.061	.045	.061	0.218	1.694	0.208	.205
IFGSM _U	IFGSM _U	.066	.069	.044	.054	.154	.152	.065	.069	.048	.068	0.228	1.708	0.216	.214

(a) MNIST 5 vs 6

Table C.6. Test risks and bounds for 3 tasks of Fashion MNIST with $n=100$ perturbations for all pairs (Defense, Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} . The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} . To quantify the gap between our risks and the classical definition we put in *italic* the risk of our models against the classical attacks: we replace PGD_U and IFGSM_U by PGD or IFGSM (i.e., we did not sample from the uniform distribution). Since Eq. (3.9) upperbounds Eq. (3.8) thanks to the TV term, we compute the two bound values of Theorem 3.3.3.

ℓ_∞ -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 0.1$		Attack without u						Attack without u							
Defense	Attack	$A_T(MV_\rho)$		$R_{\hat{T}}(MV_\rho)$		Th. 3.3.2		$A_T(MV_\rho)$		$A_{\hat{T}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.021	.020	.021	.020	0.060	0.070	.019	<i>.019</i>	.019	.019	0.130	0.139	0.130	0.139
—	PGD_U	.951	.944	.606	.719	1.275	1.333	.935	.920	.762	.864	1.617	2.503	1.421	1.317
—	IFGSM_U	.957	.947	.588	.718	1.231	1.336	.950	<i>.950</i>	.734	.851	1.587	2.495	1.395	1.316
UNIF	—	.076	<i>.077</i>	.076	<i>.077</i>	0.178	0.184	.076	<i>.077</i>	.076	<i>.077</i>	0.230	0.235	0.230	0.235
UNIF	PGD_U	.964	.961	.714	.719	1.496	1.265	.966	.963	.853	.859	2.098	2.417	1.785	1.416
UNIF	IFGSM_U	.978	.976	.627	.632	1.306	1.259	.979	<i>.979</i>	.758	.762	1.914	2.422	1.597	1.396
PGD_U	—	<i>.041</i>	.040	.041	.040	0.114	0.111	<i>.041</i>	.040	.041	.040	0.173	0.171	0.173	0.171
PGD_U	PGD_U	.098	.097	.089	.086	0.207	0.210	.099	.097	.101	.100	0.306	1.826	0.281	0.267
PGD_U	IFGSM_U	.113	.112	.105	.101	0.244	0.246	.115	.112	.120	.113	0.353	1.853	0.321	0.302
IFGSM_U	—	.045	<i>.047</i>	.045	<i>.047</i>	0.131	0.137	.045	<i>.047</i>	.045	<i>.047</i>	0.188	0.194	0.188	0.194
IFGSM_U	PGD_U	.100	<i>.102</i>	.089	.085	0.203	0.232	.100	<i>.102</i>	.102	<i>.102</i>	0.298	1.645	0.274	0.287
IFGSM_U	IFGSM_U	.112	<i>.116</i>	.099	.096	0.232	0.260	.112	<i>.116</i>	.114	.112	0.328	1.687	0.301	0.313

(a) Fashion MNIST Sandall vs Ankle Boot

ℓ_∞ -norm		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
$b = 0.1$		Attack without u						Attack without u							
Defense	Attack	$A_T(MV_\rho)$		$R_{\hat{T}}(MV_\rho)$		Th. 3.3.2		$A_T(MV_\rho)$		$A_{\hat{T}}(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.038	.037	.038	.037	.088	.091	.038	.037	.038	.037	0.153	0.155	0.153	.155
—	PGD_U	.596	.515	.477	.218	.844	.662	.590	.576	.570	.502	1.049	1.924	0.948	.857
—	IFGSM_U	.723	.623	.573	.257	.971	.751	.716	.695	.678	.598	1.189	2.031	1.080	.980
UNIF	—	.032	<i>.032</i>	.032	<i>.032</i>	.083	.085	.032	<i>.033</i>	.032	<i>.033</i>	0.149	0.151	0.149	.151
UNIF	PGD_U	.438	<i>.439</i>	.356	.245	.813	.563	.435	<i>.435</i>	.423	.312	1.082	1.867	0.959	.688
UNIF	IFGSM_U	.546	<i>.547</i>	.453	.325	.974	.690	.544	<i>.547</i>	.530	.409	1.266	2.009	1.128	.823
PGD_U	—	.048	<i>.053</i>	.048	<i>.053</i>	.115	.130	.048	<i>.053</i>	.048	<i>.053</i>	0.177	0.188	0.177	.188
PGD_U	PGD_U	.102	<i>.116</i>	.089	.099	.205	.223	.102	<i>.116</i>	.096	.115	0.282	1.323	0.266	.278
PGD_U	IFGSM_U	.120	<i>.135</i>	.102	.115	.237	.255	.120	<i>.135</i>	.109	.133	0.318	1.380	0.299	.309
IFGSM_U	—	.051	.045	.051	.045	.120	.115	.051	.045	.051	.045	0.179	0.175	0.179	.175
IFGSM_U	PGD_U	.106	.094	.091	.085	.211	.193	.106	.094	.102	.097	0.292	1.488	0.273	.252
IFGSM_U	IFGSM_U	.120	.111	.101	.102	.239	.218	.119	.111	.113	.113	0.322	1.546	0.299	.277

(b) Fashion MNIST Top vs Pullover

C.8. Additional Experimental Results

ℓ_∞ -norm $b = 0.1$		Algo. 3.1 with Eq. (3.6)						Algo. 3.1 with Eq. (3.9)							
Defense	Attack	Attack without u						Attack without u							
		$A_T(MV_\rho)$		$R_T(MV_\rho)$		Th. 3.3.2		$A_T(MV_\rho)$		$A_T(MV_\rho)$		Th. 3.3.3 - Eq. (3.9)		Th. 3.3.3 - Eq. (3.8)	
		\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.122	.122	.122	.122	0.276	0.286	.122	.122	.122	.122	0.318	0.328	0.318	0.328
—	PGD _U	.884	.887	.781	.795	1.579	1.268	.882	.886	.864	.872	2.020	2.640	1.803	1.390
—	IFGSM _U	.901	.902	.756	.774	1.558	1.272	.901	.902	.865	.876	2.032	2.651	1.795	1.393
UNIF	—	.166	.166	.166	.166	0.352	0.357	.166	.166	.166	.166	0.389	0.394	0.389	0.394
UNIF	PGD _U	.911	.914	.796	.798	1.402	1.326	.913	.914	.896	.888	1.934	2.325	1.713	1.447
UNIF	IFGSM _U	.935	.937	.787	.798	1.392	1.350	.934	.936	.887	.882	1.905	2.378	1.693	1.469
PGD _U	—	.163	.162	.163	.162	0.386	0.395	.163	.162	.163	.162	0.419	0.430	0.419	0.430
PGD _U	PGD _U	.394	.396	.359	.329	0.764	0.673	.394	.396	.403	.394	0.954	2.321	0.865	0.726
PGD _U	IFGSM _U	.475	.480	.442	.410	0.910	0.769	.477	.480	.487	.472	1.121	2.411	1.020	0.826
IFGSM _U	—	.167	.168	.167	.168	0.411	0.395	.167	.168	.167	.168	0.445	0.429	0.445	0.429
IFGSM _U	PGD _U	.396	.373	.359	.293	0.772	0.641	.396	.373	.405	.328	0.970	2.368	0.877	0.692
IFGSM _U	IFGSM _U	.465	.428	.424	.334	0.891	0.705	.465	.429	.470	.372	1.090	2.425	0.995	0.758

(a) Fashion MNIST Coat vs Shirt

Table C.8. Test risks for 6 tasks of MNIST and Fashion MNIST datasets for all pairs (Defense, Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} using ℓ_2 -norm. The results of these tables are computed considering defenses of the literature, i.e., adversarial training using PGD or IFGSM. We also add an adversarial training using UNIF for the completeness of comparison between this baseline defense and our algorithm. The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} .

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.005	.005	—	—	.015	.015	—	—	.015	.015
—	PGD	.326	.327	—	PGD	.692	.692	—	PGD	.283	.283
—	IFGSM	.122	.121	—	IFGSM	.464	.462	—	IFGSM	.144	.144
UNIF	—	.005	.005	UNIF	—	.024	.024	UNIF	—	.017	.017
UNIF	PGD	.191	.190	UNIF	PGD	.653	.653	UNIF	PGD	.220	.219
UNIF	IFGSM	.071	.072	UNIF	IFGSM	.441	.438	UNIF	IFGSM	.122	.122
PGD	—	.007	.007	PGD	—	.024	.027	PGD	—	.014	.013
PGD	PGD	.027	.026	PGD	PGD	.136	.138	PGD	PGD	.056	.055
PGD	IFGSM	.022	.021	PGD	IFGSM	.097	.102	PGD	IFGSM	.045	.041
IFGSM	—	.005	.006	IFGSM	—	.022	.027	IFGSM	—	.013	.014
IFGSM	PGD	.041	.035	IFGSM	PGD	.166	.186	IFGSM	PGD	.077	.070
IFGSM	IFGSM	.021	.021	IFGSM	IFGSM	.113	.124	IFGSM	IFGSM	.053	.047

(a) MNIST 1 vs 7

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019	—	—	.038	.038	—	—	.122	.122
—	PGD	.709	.708	—	PGD	.286	.285	—	PGD	.768	.767
—	IFGSM	.426	.414	—	IFGSM	.188	.186	—	IFGSM	.683	.680
UNIF	—	.071	.072	UNIF	—	.041	.039	UNIF	—	.204	.204
UNIF	PGD	.531	.531	UNIF	PGD	.249	.248	UNIF	PGD	.753	.754
UNIF	IFGSM	.331	.329	UNIF	IFGSM	.197	.192	UNIF	IFGSM	.607	.606
PGD	—	.034	.036	PGD	—	.043	.045	PGD	—	.182	.178
PGD	PGD	.107	.103	PGD	PGD	.102	.117	PGD	PGD	.453	.412
PGD	IFGSM	.091	.087	PGD	IFGSM	.090	.094	PGD	IFGSM	.408	.379
IFGSM	—	.031	.029	IFGSM	—	.038	.040	IFGSM	—	.148	.146
IFGSM	PGD	.125	.108	IFGSM	PGD	.120	.106	IFGSM	PGD	.405	.411
IFGSM	IFGSM	.104	.090	IFGSM	IFGSM	.092	.080	IFGSM	IFGSM	.369	.364

(b) MNIST 4 vs 9

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019	—	—	.038	.038	—	—	.122	.122
—	PGD	.709	.708	—	PGD	.286	.285	—	PGD	.768	.767
—	IFGSM	.426	.414	—	IFGSM	.188	.186	—	IFGSM	.683	.680
UNIF	—	.071	.072	UNIF	—	.041	.039	UNIF	—	.204	.204
UNIF	PGD	.531	.531	UNIF	PGD	.249	.248	UNIF	PGD	.753	.754
UNIF	IFGSM	.331	.329	UNIF	IFGSM	.197	.192	UNIF	IFGSM	.607	.606
PGD	—	.034	.036	PGD	—	.043	.045	PGD	—	.182	.178
PGD	PGD	.107	.103	PGD	PGD	.102	.117	PGD	PGD	.453	.412
PGD	IFGSM	.091	.087	PGD	IFGSM	.090	.094	PGD	IFGSM	.408	.379
IFGSM	—	.031	.029	IFGSM	—	.038	.040	IFGSM	—	.148	.146
IFGSM	PGD	.125	.108	IFGSM	PGD	.120	.106	IFGSM	PGD	.405	.411
IFGSM	IFGSM	.104	.090	IFGSM	IFGSM	.092	.080	IFGSM	IFGSM	.369	.364

(c) MNIST 5 vs 6

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019	—	—	.038	.038	—	—	.122	.122
—	PGD	.709	.708	—	PGD	.286	.285	—	PGD	.768	.767
—	IFGSM	.426	.414	—	IFGSM	.188	.186	—	IFGSM	.683	.680
UNIF	—	.071	.072	UNIF	—	.041	.039	UNIF	—	.204	.204
UNIF	PGD	.531	.531	UNIF	PGD	.249	.248	UNIF	PGD	.753	.754
UNIF	IFGSM	.331	.329	UNIF	IFGSM	.197	.192	UNIF	IFGSM	.607	.606
PGD	—	.034	.036	PGD	—	.043	.045	PGD	—	.182	.178
PGD	PGD	.107	.103	PGD	PGD	.102	.117	PGD	PGD	.453	.412
PGD	IFGSM	.091	.087	PGD	IFGSM	.090	.094	PGD	IFGSM	.408	.379
IFGSM	—	.031	.029	IFGSM	—	.038	.040	IFGSM	—	.148	.146
IFGSM	PGD	.125	.108	IFGSM	PGD	.120	.106	IFGSM	PGD	.405	.411
IFGSM	IFGSM	.104	.090	IFGSM	IFGSM	.092	.080	IFGSM	IFGSM	.369	.364

(d) Fashion MNIST Sandall vs Ankle Boot

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019	—	—	.038	.038	—	—	.122	.122
—	PGD	.709	.708	—	PGD	.286	.285	—	PGD	.768	.767
—	IFGSM	.426	.414	—	IFGSM	.188	.186	—	IFGSM	.683	.680
UNIF	—	.071	.072	UNIF	—	.041	.039	UNIF	—	.204	.204
UNIF	PGD	.531	.531	UNIF	PGD	.249	.248	UNIF	PGD	.753	.754
UNIF	IFGSM	.331	.329	UNIF	IFGSM	.197	.192	UNIF	IFGSM	.607	.606
PGD	—	.034	.036	PGD	—	.043	.045	PGD	—	.182	.178
PGD	PGD	.107	.103	PGD	PGD	.102	.117	PGD	PGD	.453	.412
PGD	IFGSM	.091	.087	PGD	IFGSM	.090	.094	PGD	IFGSM	.408	.379
IFGSM	—	.031	.029	IFGSM	—	.038	.040	IFGSM	—	.148	.146
IFGSM	PGD	.125	.108	IFGSM	PGD	.120	.106	IFGSM	PGD	.405	.411
IFGSM	IFGSM	.104	.090	IFGSM	IFGSM	.092	.080	IFGSM	IFGSM	.369	.364

(e) Fashion MNIST Top vs Pullover

ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$		ℓ_2 -norm, $b = 1$		$A_{\mathcal{T}}(\text{MV}_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}	Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019	—	—	.038	.038	—	—	.122	.122
—	PGD	.709	.708	—	PGD	.286	.285	—	PGD	.768	.767
—	IFGSM	.426	.414	—	IFGSM	.188	.186	—	IFGSM	.683	.680
UNIF	—	.071	.072	UNIF	—	.041	.039	UNIF	—	.204	.204
UNIF	PGD	.531	.531	UNIF	PGD	.249	.248	UNIF	PGD	.753	.754
UNIF	IFGSM	.331	.329	UNIF	IFGSM	.197	.192	UNIF	IFGSM	.607	.606
PGD	—	.034	.036	PGD	—	.043	.045	PGD	—	.182	.178
PGD	PGD	.107	.103	PGD	PGD	.102	.117	PGD	PGD	.453	.412
PGD	IFGSM	.091	.087	PGD	IFGSM	.090	.094	PGD	IFGSM	.408	.379
IFGSM	—	.031	.029	IFGSM	—	.038	.040	IFGSM	—	.148	.146
IFGSM	PGD	.125	.108	IFGSM	PGD	.120	.106	IFGSM	PGD	.405	.411
IFGSM	IFGSM	.104	.090	IFGSM	IFGSM	.092	.080	IFGSM	IFGSM	.369	.364

(f) Fashion MNIST Coat vs Shirt

C.8. Additional Experimental Results

Table C.9. Test risks for 6 tasks of MNIST and Fashion MNIST datasets for all pairs (Defense, Attack) with the two voters' set \mathbb{H} and \mathbb{H}^{SIGN} using ℓ_∞ -norm. The results of these tables are computed considering defenses of the literature, i.e., adversarial training using PGD or IFGSM. We also add an adversarial training using UNIF for the completeness of comparison between this baseline defense and our algorithm. The results in **bold** correspond to the best values between results for \mathbb{H} and \mathbb{H}^{SIGN} .

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.005	.005
—	PGD	.499	.498
—	IFGSM	.479	.480
UNIF	—	.004	.004
UNIF	PGD	.516	.515
UNIF	IFGSM	.467	.467
PGD	—	.006	.007
PGD	PGD	.019	.019
PGD	IFGSM	.021	.021
IFGSM	—	.007	.007
IFGSM	PGD	.017	.018
IFGSM	IFGSM	.019	.020

(a) MNIST 1 vs 7

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015
—	PGD	.921	.921
—	IFGSM	.923	.923
UNIF	—	.017	.017
UNIF	PGD	.877	.876
UNIF	IFGSM	.877	.877
PGD	—	.041	.040
PGD	PGD	.108	.109
PGD	IFGSM	.122	.123
IFGSM	—	.057	.044
IFGSM	PGD	.109	.101
IFGSM	IFGSM	.119	.108

(b) MNIST 4 vs 9

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.015	.015
—	PGD	.498	.498
—	IFGSM	.511	.510
UNIF	—	.015	.015
UNIF	PGD	.512	.511
UNIF	IFGSM	.511	.511
PGD	—	.014	.014
PGD	PGD	.065	.058
PGD	IFGSM	.068	.065
IFGSM	—	.018	.017
IFGSM	PGD	.061	.063
IFGSM	IFGSM	.069	.071

(c) MNIST 5 vs 6

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.019	.019
—	PGD	.938	.938
—	IFGSM	.948	.949
UNIF	—	.076	.077
UNIF	PGD	.970	.969
UNIF	IFGSM	.981	.981
PGD	—	.041	.040
PGD	PGD	.098	.097
PGD	IFGSM	.115	.111
IFGSM	—	.112	.047
IFGSM	PGD	.045	.100
IFGSM	IFGSM	.101	.114

(d) Fashion MNIST Sandall vs Ankell Boot

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.038	.038
—	PGD	.574	.577
—	IFGSM	.700	.696
UNIF	—	.032	.033
UNIF	PGD	.428	.435
UNIF	IFGSM	.540	.550
PGD	—	.047	.049
PGD	PGD	.101	.097
PGD	IFGSM	.118	.112
IFGSM	—	.049	.048
IFGSM	PGD	.100	.090
IFGSM	IFGSM	.112	.108

(e) Fashion MNIST Top vs Pullover

ℓ_∞ -norm, $b = 0.1$		$A_{\mathcal{T}}(MV_{\pi'})$	
Defense	Attack	\mathbb{H}^{SIGN}	\mathbb{H}
—	—	.122	.122
—	PGD	.879	.879
—	IFGSM	.898	.898
UNIF	—	.166	.166
UNIF	PGD	.913	.911
UNIF	IFGSM	.934	.933
PGD	—	.164	.167
PGD	PGD	.398	.395
PGD	IFGSM	.479	.481
IFGSM	—	.163	.169
IFGSM	PGD	.356	.391
IFGSM	IFGSM	.422	.461

(f) Fashion MNIST Coat vs Shirt

APPENDIX OF CHAPTER 4

D

D.1 Proof of Theorem 4.3.1

Theorem 4.3.1 (PAC-Bayesian Bound Based on the Gibbs Risk). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have

$$\begin{aligned} \forall \rho \in \mathcal{M}(\mathbb{H}), \quad r_{\mathcal{D}}(\rho) &\leq \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right), \\ \text{and } \forall \rho \in \mathcal{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq 2 \left[\overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right) \right], \end{aligned} \quad (4.1)$$

where $\overline{\text{kl}}(q|\tau) \triangleq \max \left\{ p \in (0, 1) \mid \text{kl}(q\|p) \leq \tau \right\}$ (see Section 2.3.1.3).

Proof. We can apply Theorem 2.3.4 with the loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$ to obtain

$$\forall \rho \in \mathcal{M}(\mathbb{H}), \quad r_{\mathcal{D}}(\rho) \leq \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right).$$

Then, from Theorem 2.2.1, we obtain

$$\forall \rho \in \mathcal{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 2 \left[\overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right) \right].$$

■

D.2 Proof of Theorem 4.3.2

Theorem 4.3.2 (PAC-Bayesian Bound Based on the Joint Error). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$, for any

$\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad e_{\mathcal{D}}(\rho) \leq \overline{\text{kl}} \left(e_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right),$$

and $\forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 4 \left[\overline{\text{kl}} \left(e_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right].$

Proof. We consider the hypothesis set $\mathbb{H}^2 = \mathbb{H} \times \mathbb{H}$ with the distributions ρ^2 (resp. π^2) defined as $\rho^2((h, h')) = \rho(h)\rho(h')$ (resp. $\pi^2((h, h')) = \pi(h)\pi(h')$) on \mathbb{H}^2 . We can apply Theorem 2.3.4 with the loss $\ell((h, h'), (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) \neq y]$ to obtain

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad e_{\mathcal{D}}(\rho) \leq \overline{\text{kl}} \left(e_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho^2 \parallel \pi^2) + \ln \frac{2\sqrt{m}}{\delta} \right] \right. \right).$$

Then, from Theorem 2.2.2 and Lemma D.7.1, we obtain

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 4 \left[\overline{\text{kl}} \left(e_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right].$$

■

D.3 Proof of Theorem 4.3.3

Theorem 4.3.3 (PAC-Bayesian C-Bound of ROY *et al.* (2016)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, r_{\mathcal{S}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \right] \right)^2}{\underbrace{1 - 2 \max \left[0, d_{\mathcal{S}}(\rho) - \sqrt{\frac{1}{2m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]} \right]}_{\triangleq C_{\mathcal{S}}^{\text{M}}(\rho)}}. \quad (4.2)$$

D.4. Proof of Theorem 4.3.4

Proof. First, remark that from Theorem 2.2.3, we have (even for $r_{\mathcal{D}}(\rho) \geq \frac{1}{2}$)

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, r_{\mathcal{D}}(\rho) \right]\right)^2}{1 - 2 \max [0, d_{\mathcal{D}}(\rho)]}.$$

We upper-bound $r_{\mathcal{D}}(\rho)$ by applying Theorem 2.3.2 with $\frac{\delta}{2}$ and the loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$, to obtain

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad r_{\mathcal{D}}(\rho) \leq r_{\mathcal{S}}(\rho) + \sqrt{\frac{1}{2m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]}.$$

To lower-bound $d_{\mathcal{D}}(\rho)$, we have to consider the hypothesis set $\mathbb{H}^2 = \mathbb{H} \times \mathbb{H}$ with the distributions ρ^2 (*resp.* π^2) defined as $\rho^2((h, h')) = \rho(h)\rho(h')$ (*resp.* $\pi^2((h, h')) = \pi(h)\pi(h')$) on \mathbb{H}^2 . Then, we obtain the lower-bound by applying Theorem 2.3.4 with $\frac{\delta}{2}$ and the loss $\ell((h, h'), (\mathbf{x}, y)) = 2\mathbb{I}[h(\mathbf{x}) \neq y]\mathbb{I}[h'(\mathbf{x}) = y]$. We have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad d_{\mathcal{D}}(\rho) \geq d_{\mathcal{S}}(\rho) - \sqrt{\frac{1}{2m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right]},$$

where $\text{KL}(\rho^2 \parallel \pi^2) = 2 \text{KL}(\rho \parallel \pi)$ from Lemma D.7.1. ■

D.4 Proof of Theorem 4.3.4

Theorem 4.3.4 (PAC-Bayesian C-Bound of GERMAIN *et al.* (2015)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min \left[\frac{1}{2}, \overline{\text{kl}} \left(r_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right]\right)^2}{1 - 2 \max \left[0, \underline{\text{kl}} \left(d_{\mathcal{S}}(\rho) \mid \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right) \right]}. \quad (4.3)$$

$$\underbrace{\hspace{15em}}_{\triangleq C_{\mathcal{S}}^{\rho}(\rho)}$$

Proof. First, remark that from Theorem 2.2.3, we have (even for $r_{\mathcal{D}}(\rho) \geq \frac{1}{2}$)

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{\left(1 - 2 \min\left[\frac{1}{2}, r_{\mathcal{D}}(\rho)\right]\right)^2}{1 - 2 \max[0, d_{\mathcal{D}}(\rho)]}.$$

We upper-bound $r_{\mathcal{D}}(\rho)$ by applying Theorem 2.3.4 with $\frac{\delta}{2}$ and the loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[h(\mathbf{x}) \neq y]$, to obtain

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad r_{\mathcal{D}}(\rho) \leq \overline{\text{kl}}\left(r_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho \|\pi) + \ln \frac{4\sqrt{m}}{\delta} \right] \right.\right).$$

To lower-bound $d_{\mathcal{D}}(\rho)$, we have to consider the hypothesis set $\mathbb{H}^2 = \mathbb{H} \times \mathbb{H}$ with the distributions ρ^2 (*resp.* π^2) defined as $\rho^2((h, h')) = \rho(h)\rho(h')$ (*resp.* $\pi^2((h, h')) = \pi(h)\pi(h')$) on \mathbb{H}^2 . Then, we obtain the lower-bound by applying Theorem 2.3.4 with $\frac{\delta}{2}$ and the loss $\ell((h, h'), (\mathbf{x}, y)) = 2 \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) = y]$. We have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad d_{\mathcal{D}}(\rho) \geq \underline{\text{kl}}\left(d_{\mathcal{S}}(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho^2 \|\pi^2) + \ln \frac{4\sqrt{m}}{\delta} \right] \right.\right),$$

where $\text{KL}(\rho^2 \|\pi^2) = 2 \text{KL}(\rho \|\pi)$ from Lemma D.7.1. ■

D.5 Proof of Theorem 4.3.5

Theorem 4.3.5 (PAC-Bayesian C-Bound of LACASSE *et al.* (2006)). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$R_{\mathcal{D}}(\text{MV}_{\rho}) \leq \sup_{(e,d) \in \mathbb{A}_{\mathcal{S}}(\rho)} \left[1 - \frac{(1 - (2e + d))^2}{1 - 2d} \right],$$

$$\text{where } \mathbb{A}_{\mathcal{S}}(\rho) = \left\{ (e, d) \left| \begin{aligned} \text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho)) \leq \frac{1}{m} \left[2 \text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m+m}}{\delta} \right], \\ d \leq 2\sqrt{e} - 2e, \quad 2e + d < 1 \end{aligned} \right. \right\}.$$

Proof. First of all, we need to prove the following PAC-Bayesian bound:

$$\text{kl}(e_{\mathcal{S}}(\rho), d_{\mathcal{S}}(\rho) \| e_{\mathcal{D}}(\rho), d_{\mathcal{D}}(\rho)) \leq \frac{1}{m} \left[2 \text{KL}(\rho \| \pi) + \ln \frac{2\sqrt{m}+m}{\delta} \right]. \quad (\text{D.1})$$

We apply Theorem 2.3.1 with $\varphi(h, \mathcal{S}) = m \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h))$ to obtain with probability at least $1 - \delta$

$$\begin{aligned} & \mathbb{E}_{h \sim \rho} \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h)) \\ & \leq \frac{1}{m} \left[\text{KL}(\rho \| \pi) + \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h))} \right) \right], \end{aligned} \quad (\text{D.2})$$

where $\ell' : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ and $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$. Moreover, YOUNSI (2012) and YOUNSI and LACASSE (2020) proves that

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \exp \left[m \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h)) \right] \leq 2\sqrt{m} + m.$$

Hence, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h))} &= \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} e^{m \text{kl}(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h))} \\ &\leq 2\sqrt{m} + m. \end{aligned} \quad (\text{D.3})$$

Moreover, from JENSEN's inequality (Theorem A.1.1), we have

$$\begin{aligned} & \text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell}(h), \mathbb{E}_{h \sim \rho} R_{\mathcal{S}}^{\ell'}(h) \| \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h), \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell'}(h) \right) \\ & \leq \mathbb{E}_{h \sim \rho} \text{kl} \left(R_{\mathcal{S}}^{\ell}(h), R_{\mathcal{S}}^{\ell'}(h) \| R_{\mathcal{D}}^{\ell}(h), R_{\mathcal{D}}^{\ell'}(h) \right). \end{aligned} \quad (\text{D.4})$$

By combining Equations (D.2) to (D.4), we obtain Equation (D.1). Lastly the PAC-Bayesian bound in $\Delta_{\mathcal{S}}(\rho)$ is obtained by instantiating Equation (D.1) (and applying Lemma D.7.1) with

$$\begin{aligned} \ell((h, h'), (\mathbf{x}, y)) &= \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) \neq y], \\ \text{and } \ell'((h, h'), (\mathbf{x}, y)) &= 2 \mathbb{I}[h(\mathbf{x}) \neq y] \mathbb{I}[h'(\mathbf{x}) = y], \end{aligned}$$

and the distributions ρ^2 and π^2 on \mathbb{H}^2 .

Remark that the constraints involves in the set $\Delta_{\mathcal{S}}(\rho)$ do not remove the possible couple of joint error and disagreement. Indeed, since a variance is not negative,

we have

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} (\widehat{m}_\rho(\mathbf{x}, y))^2 &\geq \left(\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} [\widehat{m}_\rho(\mathbf{x}, y)] \right)^2 \\ \iff d_{\mathcal{D}'}(\rho) &\leq 2r_{\mathcal{D}'}(\rho)(1 - r_{\mathcal{D}'}(\rho)) \\ \iff d_{\mathcal{D}'}(\rho) &\leq 2 \left[\sqrt{e_{\mathcal{D}'}(\rho)} - e_{\mathcal{D}'}(\rho) \right]. \end{aligned}$$

Moreover, from Theorem 2.2.3, we know that $2e_{\mathcal{D}'}(\rho) + d_{\mathcal{D}'}(\rho) < 1$.

We can now prove that, under the constraints involved in $\mathbb{A}_S(\rho)$, we still have a valid bound on $R_{\mathcal{D}}(\text{MV}_\rho)$. To do so, we consider two cases.

Case 1: If for all $(e, d) \in \mathbb{A}_S(\rho)$ we have $2e + d < 1$.

In this case $(e_{\mathcal{D}}(\rho), d_{\mathcal{D}}(\rho)) \in \mathbb{A}_S(\rho)$ and $2e_{\mathcal{D}}(\rho) + d_{\mathcal{D}}(\rho) < 1$, thus, Theorem 2.2.3 holds with probability at least $1 - \delta$. In other words, we have $R_{\mathcal{D}}(\text{MV}_\rho) \leq 1 - \frac{[1 - (2e_{\mathcal{D}}(\rho) + d_{\mathcal{D}}(\rho))]^2}{1 - 2d_{\mathcal{D}}(\rho)} \leq \sup_{(e, d) \in \mathbb{A}'_S(\rho)} C^{\text{L}}(e, d)$ with probability at least $1 - \delta$.

Case 2: If there exists $(e, d) \in \mathbb{A}_S(\rho)$ such that $2e + d = 1$.

We have $\sup_{(e, d) \in \mathbb{A}'_S(\rho)} C^{\text{L}}(e, d) = 1$ that is a valid bound on $R_{\mathcal{D}}(\text{MV}_\rho)$. ■

D.6 Proof of Theorem 4.4.1

Theorem 4.4.1 (Reformulation of LACASSE *et al.*'s PAC-Bayesian C-Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $S \sim \mathcal{D}^m$ we have for all $\rho \in \mathbb{M}(\mathbb{H})$

$$\begin{aligned} R_{\mathcal{D}}(\text{MV}_\rho) &\leq \sup_{(e, d) \in \mathbb{A}'_S(\rho)} \left[1 - \frac{(1 - (2e + d))^2}{1 - 2d} \right], \quad (4.4) \\ \mathbb{A}'_S(\rho) &= \left\{ (e, d) \mid \text{kl}(e_S(\rho), d_S(\rho)) \mid e, d \leq \frac{1}{m} \left[2 \text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m} + m}{\delta} \right], \right. \\ &\quad \left. d \leq 2\sqrt{\min\left(e, \frac{1}{4}\right)} - 2e, \quad d < \frac{1}{2} \right\}. \end{aligned}$$

Proof. Beforehand, we explain how we fixed the constraints involved in $\mathbb{A}'_S(\rho)$. Compared to $\mathbb{A}_S(\rho)$, we add another constraint : $d < \frac{1}{2}$. Hence, we have the

D.7. Proof of Lemma D.7.1

constraints $d \leq 2\sqrt{e}-2e$, $d \leq 1-2e$, and $d < \frac{1}{2}$. We remark that when $e \leq \frac{1}{4}$, we have $2\sqrt{e}-2e \leq 1-2e$. Then, we merge $d \leq 2\sqrt{e}-2e$ and $d \leq 1-2e$ into $d \leq 2\sqrt{\min(e, \frac{1}{4})}-2e$. Indeed, we have

$$d \leq 2\sqrt{\min(e, \frac{1}{4})}-2e \iff \begin{cases} d \leq 2\sqrt{e}-2e & \text{if } e \leq \frac{1}{4}, \\ d < 1-2e & \text{if } e \geq \frac{1}{4}. \end{cases}$$

We prove now that under the constraints involved in $\mathbb{A}'_{\mathcal{S}}(\rho)$, we still have a valid bound on $R_{\mathcal{D}}(\text{MV}_{\rho})$. To do so, we consider two cases.

Case 1: If for all $(e, d) \in \mathbb{A}'_{\mathcal{S}}(\rho)$ we have $2e+d < 1$.

In this case $(e_{\mathcal{D}}(\rho), d_{\mathcal{D}}(\rho)) \in \mathbb{A}'_{\mathcal{S}}(\rho)$ and $2e_{\mathcal{D}}(\rho)+d_{\mathcal{D}}(\rho) < 1$, thus, Theorem 2.2.3 holds with probability at least $1 - \delta$. In other words, we have $R_{\mathcal{D}}(\text{MV}_{\rho}) \leq 1 - \frac{[1-(2e_{\mathcal{D}}(\rho)+d_{\mathcal{D}}(\rho))]^2}{1-2d_{\mathcal{D}}(\rho)} \leq \sup_{(e,d) \in \mathbb{A}'_{\mathcal{S}}(\rho)} C^{\text{L}}(e, d)$ with probability at least $1 - \delta$.

Case 2: If there exists $(e, d) \in \mathbb{A}'_{\mathcal{S}}(\rho)$ such that $2e+d=1$.

We have $\sup_{(e,d) \in \mathbb{A}'_{\mathcal{S}}(\rho)} C^{\text{L}}(e, d) = 1$ that is a valid bound on $R_{\mathcal{D}}(\text{MV}_{\rho})$. ■

D.7 Proof of Lemma D.7.1

Lemma D.7.1. Let the distribution ρ^2 and π^2 on $\mathbb{H}^2 = \mathbb{H} \times \mathbb{H}$ defined as

$$\rho^2((h, h')) = \rho(h)\rho(h') \quad \text{and} \quad \pi^2((h, h')) = \pi(h)\pi(h').$$

The KL divergence between ρ^2 and π^2 can be expressed *w.r.t.* ρ and π . We have

$$\text{KL}(\rho^2 \parallel \pi^2) = 2 \text{KL}(\rho \parallel \pi).$$

Proof. We develop the term $\text{KL}(\rho^2 \parallel \pi^2)$, *i.e.*, we have

$$\begin{aligned} \text{KL}(\rho^2 \parallel \pi^2) &= \mathbb{E}_{(h, h') \sim \rho^2} \ln \frac{\rho^2((h, h'))}{\pi^2((h, h'))} \\ &= \mathbb{E}_{h \sim \rho} \mathbb{E}_{h' \sim \rho} \ln \frac{\rho(h)\rho(h')}{\pi(h)\pi(h')} \\ &= \mathbb{E}_{h \sim \rho} \ln \frac{\rho(h)}{\pi(h)} + \mathbb{E}_{h' \sim \rho} \ln \frac{\rho(h')}{\pi(h')} \\ &= 2 \text{KL}(\rho \parallel \pi). \end{aligned}$$

D.8 About Danskin's Theorem

As mentioned in the context of the justification of the function `MAXIMIZE-e-d` in Algorithm 4.3, we now discuss the possible application of Danskin's Theorem (DANSKIN, 1966, Section I). The statement of the theorem is as follows.

Theorem D.8.1 (Danskin's Theorem). Let $\mathbb{A} \subset \mathbb{R}^a$ be a compact set and $\phi : \mathbb{R}^b \times \mathbb{A} \rightarrow \mathbb{R}$ s.t. for all $\mathbf{a} \in \mathbb{A}$, we have that ϕ is continuously differentiable, then $\Phi(\mathbf{x}) = \max_{\mathbf{a} \in \mathbb{A}} \phi(\mathbf{x}, \mathbf{a})$ is directionally differentiable with directional derivatives

$$\Phi'(\mathbf{x}, \mathbf{d}) = \max_{\mathbf{a} \in \mathbb{A}^*} \langle \mathbf{d}, \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}) \rangle,$$

where $\mathbb{A}^* = \{\mathbf{a}^* \mid \phi(\mathbf{x}, \mathbf{a}^*) = \max_{\mathbf{a} \in \mathbb{A}} \phi(\mathbf{x}, \mathbf{a})\}$ and $\langle \cdot, \cdot \rangle$ is the dot product.

To optimize a problem $\min_{\mathbf{x} \in \mathbb{R}^b} \Phi(\mathbf{x})$ with $\Phi(\mathbf{x}) = \max_{\mathbf{a} \in \mathbb{A}} \phi(\mathbf{x}, \mathbf{a})$, this theorem tells us that under several assumptions, if we know a maximizer $\mathbf{a} \in \mathbb{A}$, then, we have an analytical expression of the directional derivatives of $\Phi(\mathbf{x})$. Thus, from this theorem, we also know a gradient to minimize the problem $\min_{\mathbf{x} \in \mathbb{R}^b} \Phi(\mathbf{x})$; this is expressed in the following corollary.

Corollary D.8.1 (MADRY *et al.*, 2018). Assuming that the conditions of Theorem D.8.1 are fulfilled and let $\mathbf{a}^* \in \mathbb{A}^*$ be a maximizer of ϕ . If $\mathbf{d} = \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}^*)$ with $|\mathbf{d}|_2^2 > 0$ then $-\mathbf{d}$ is a descent direction for $\Phi(\mathbf{x})$, *i.e.*, $\Phi'(\mathbf{x}, \mathbf{d}) > 0$.

Proof. By definition of the directional derivative $\Phi'(\mathbf{x}, \mathbf{d})$, we have:

$$\begin{aligned} \Phi'(\mathbf{x}, \mathbf{d}) &= \max_{\mathbf{a} \in \mathbb{A}^*} \langle \mathbf{d}, \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}) \rangle \\ &= \max_{\mathbf{a} \in \mathbb{A}^*} \langle \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}^*), \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}) \rangle \geq |\nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}^*)|_2^2 > 0. \end{aligned}$$

Then, for each iteration of the min/max problem optimization, we can (i) optimize the inner maximization problem, (ii) fix the maximizer $\mathbf{a}^* \in \mathbb{A}$ and apply a gradient descent step with the derivative $\nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{a}^*)$. However, as we mentioned, the assumptions are not fulfilled in our case to apply Theorem D.8.1 since our inner objective in Equation (4.5) or its approximation is not differentiable everywhere in the compact set $[0, \frac{1}{2}]^2$. However, we never encounter problematic cases and this strategy is thus valid

D.9. Optimizing the Bound of Theorem 4.3.1

for optimizing our proposed approximation. In practice, we have found that it is indeed an efficient and sound strategy.

D.9 Optimizing the Bound of Theorem 4.3.1

To minimize the bound of Theorem 4.3.1, the objective function is defined as

$$G_U^R(\rho) = 2 \left[\overline{\text{kl}} \left(r_U(\rho) \mid \frac{1}{m} \left[\text{KL}(\rho \parallel \pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right].$$

We derive an algorithm (denoted as `Germain` in the setting description of the experiments of Section 4.5) similar to Algorithm 4.2. The algorithm is described in Algorithm D.1 below.

Algorithm D.1 Minimization of Equation (4.1) by Stochastic Gradient Descent

Given: learning sample \mathcal{S} , prior distribution π on \mathbb{H} , the objective function $G_U^R(\rho)$

Hyperparameters: number of iterations T

$\rho \leftarrow \pi$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathcal{U} \subseteq \mathcal{S}$ **do**

$\rho \leftarrow$ Update ρ with $G_U^R(\rho)$ by gradient descent

D.10 Additional Experiments

D.10.1 Details for Figures 4.2 to 4.7

We provide in Tables D.1 to D.3 the details of the results in Figures 4.2 to 4.7.

D.10.2 Experiments on the Computation Time

Figure D.1 introduces a comparison of the computation time for the different algorithms. We consider the moons binary dataset and perform, in fact, two experiments on majority votes with decision stumps.

- (i) We fix 128 stumps per feature and vary the number $m \in \{100, 500, 1000, 5000\}$ of examples in the dataset.
- (ii) We fix the number of examples $m = 5000$ and vary the number of stumps per feature in $\{32, 64, 128\}$.

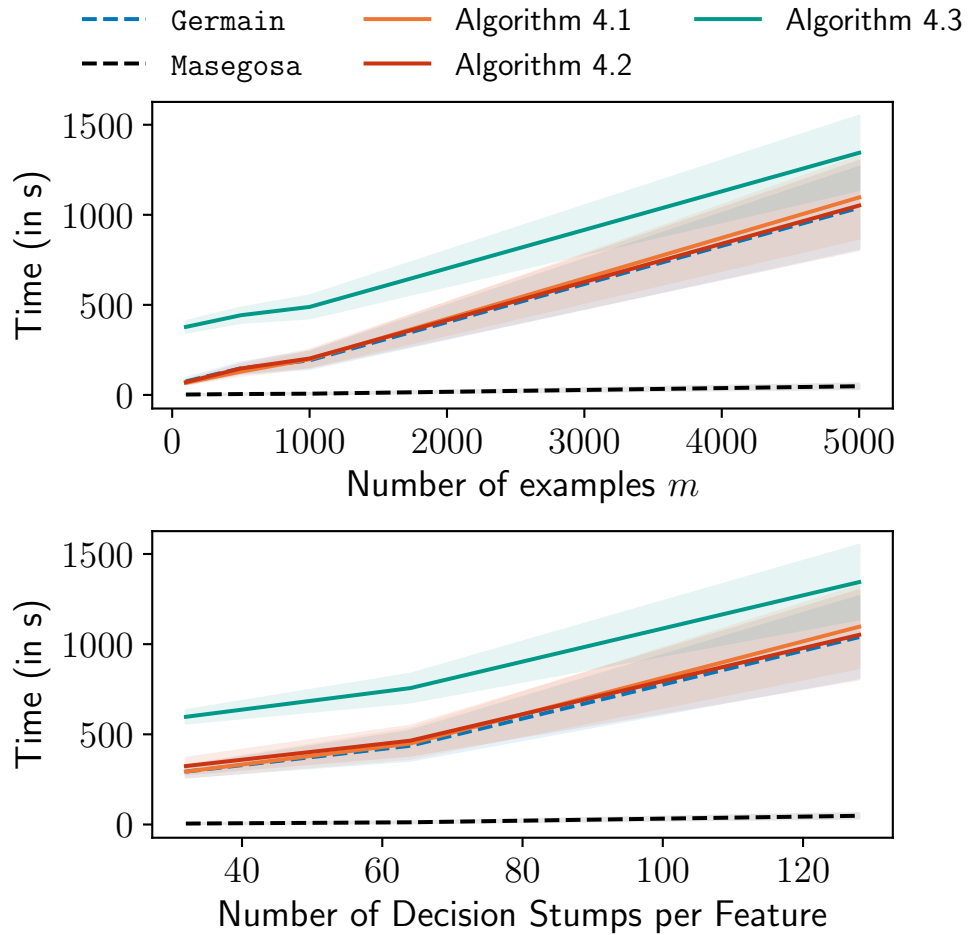


Figure D.1. In the top plot represents the evolution of the computation time (in second) with respect to the number of examples m in moons dataset. The bottom plot is the evolution of the computation time in function of the number of decision stumps per feature. For each curve, we plot the mean computation time (the plain line) and the standard deviation (the shadow) over 10 runs.

D.10.3 Details on the Empirical Joint Error and Disagreement

We report in Figures D.2 to D.7, the empirical joint error and disagreement obtained on the different datasets. These figures illustrate that the solutions found by Algorithm 4.3, Masegosa and CB-Boost are similar while MinCq and Germain provide usually very different solutions.

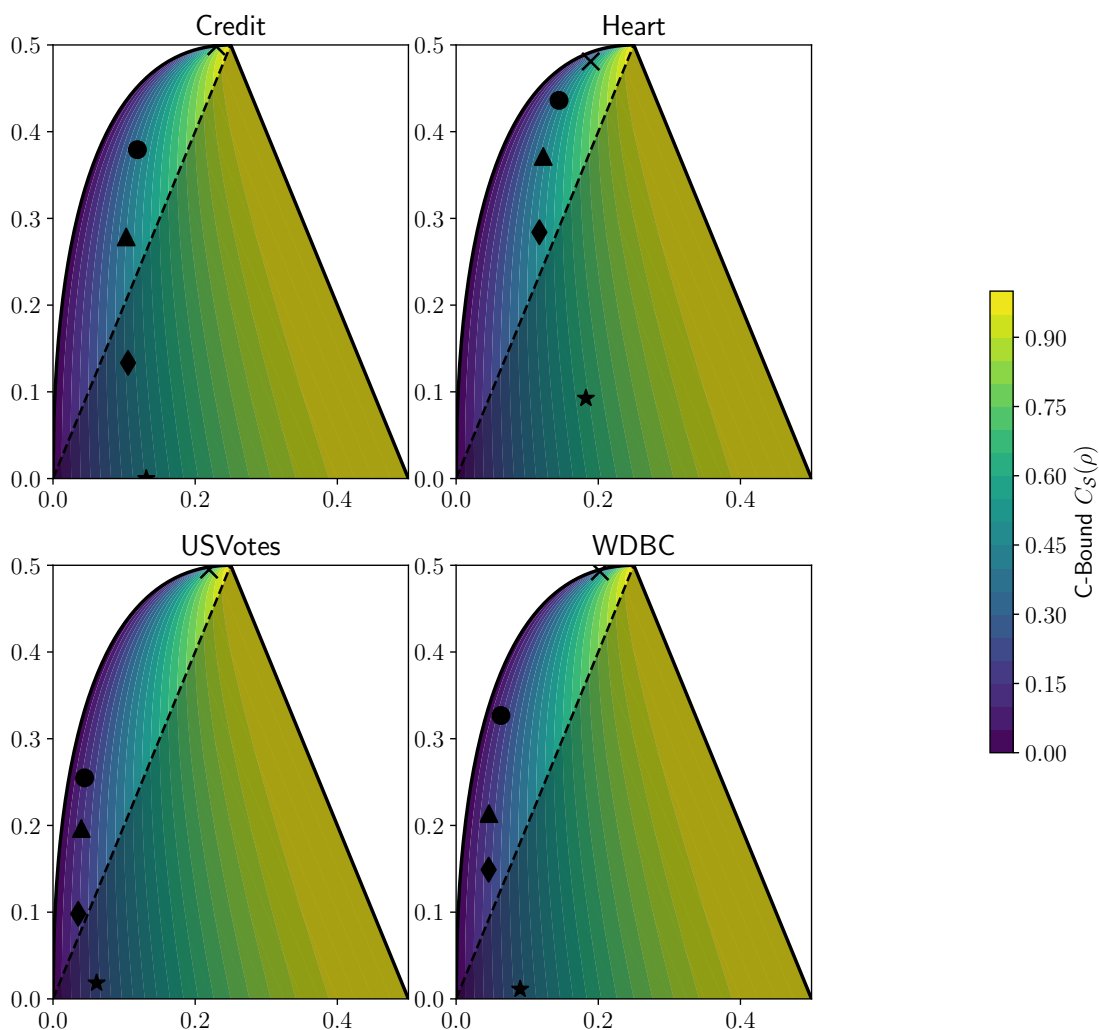


Figure D.2. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision stumps in the binary setting by Algorithm 4.3 (\diamond), Masegosa (\blacktriangle), Germain (\star), CB-Boost (\bullet), and MinCq (\times). The disagreement and the joint are averaged over the 10 runs.

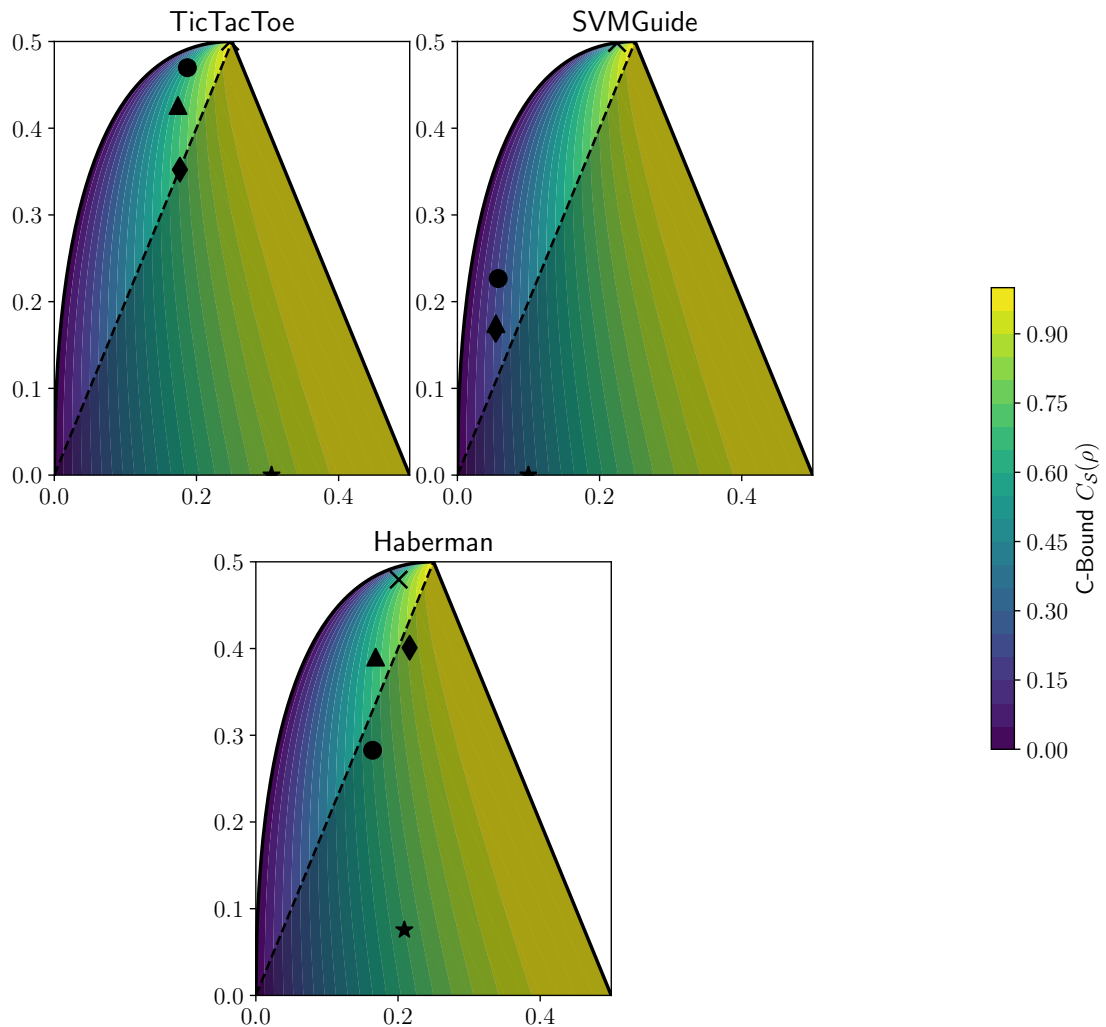


Figure D.3. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision stumps in the binary setting by Algorithm 4.3 (◆), Masegosa (▲), Germain (★), CB-Boost (●), and MinCq (×). The disagreement and the joint are averaged over the 10 runs.

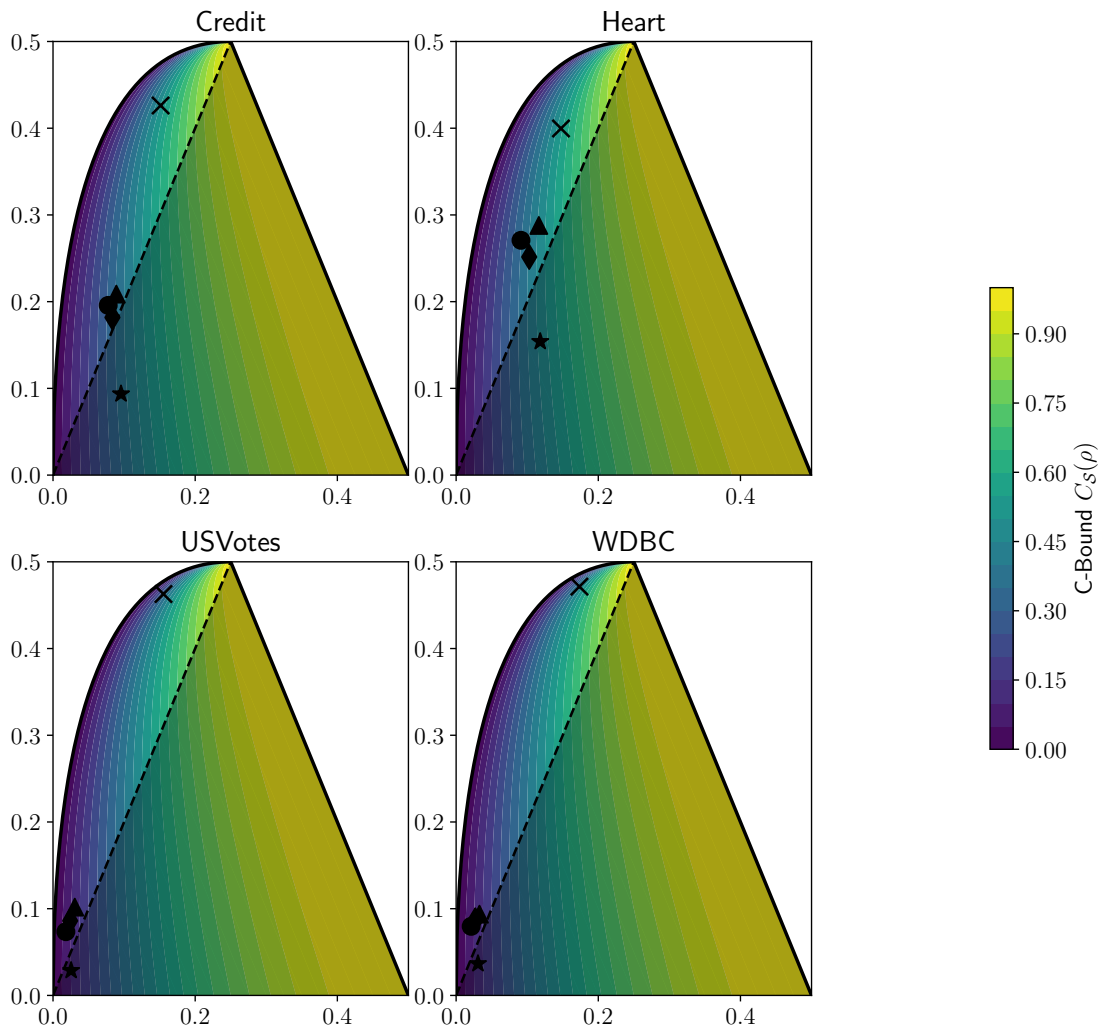


Figure D.4. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision trees in the binary setting by Algorithm 4.3 (◆), Masegosa (▲), Germain (★), CB-Boost (●), and MinCq (×). The disagreement and the joint are averaged over the 10 runs.

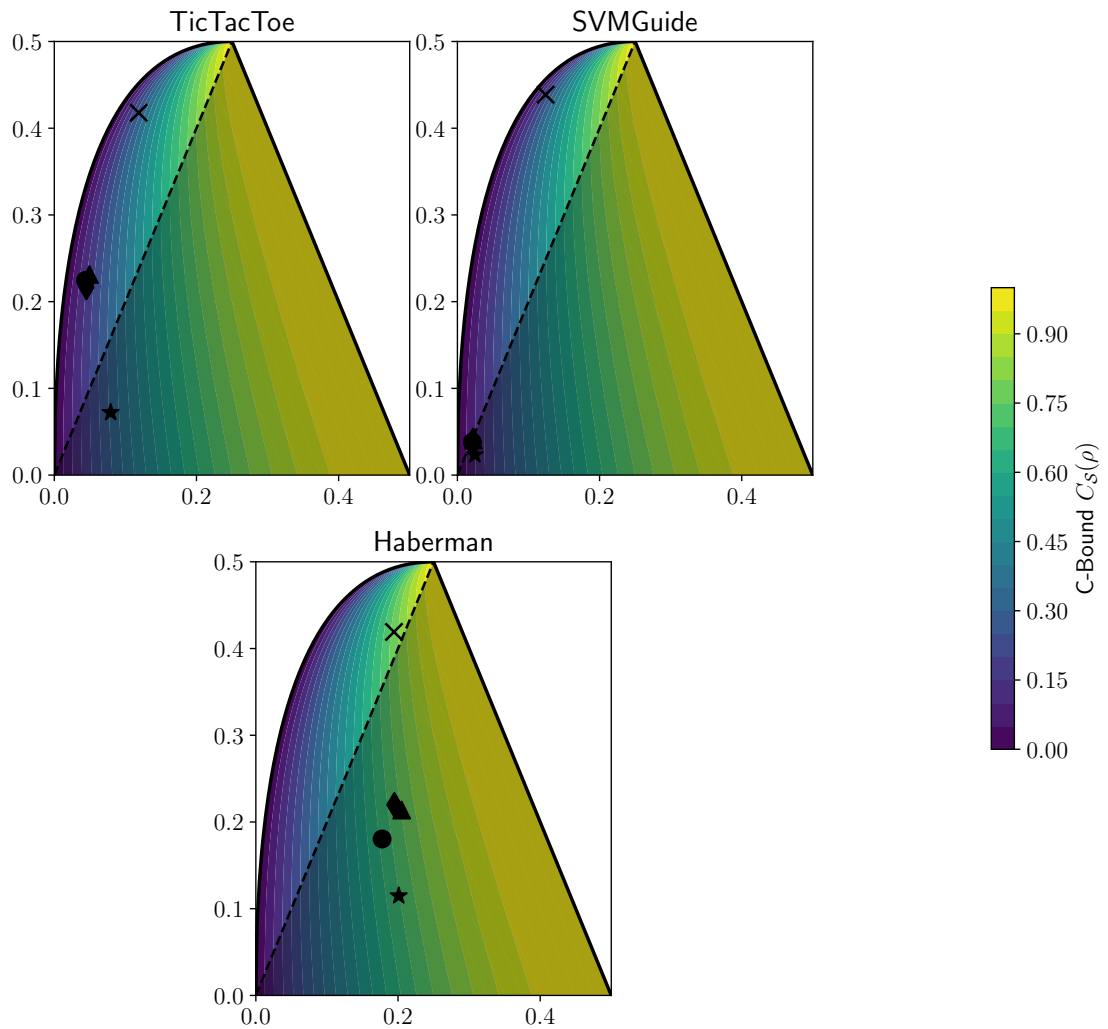


Figure D.5. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision trees in the binary setting by Algorithm 4.3 (◆), Masegosa (▲), Germain (★), CB-Boost (●), and MinCq (×). The disagreement and the joint are averaged over the 10 runs.

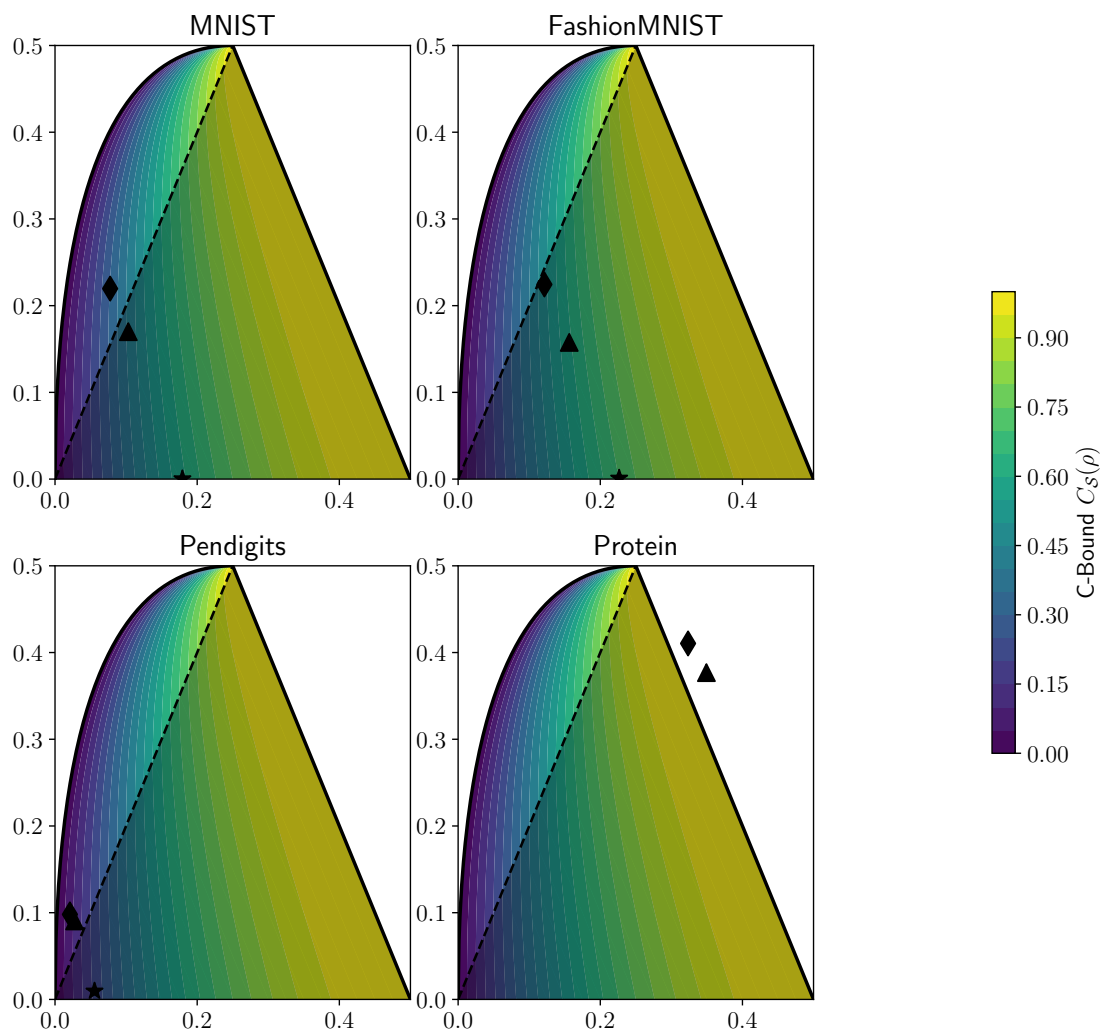


Figure D.6. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision trees in the multi-class setting by Algorithm 4.3 (◆), Masegosa (▲), Germain (★). The disagreement and the joint are averaged over the 10 runs.

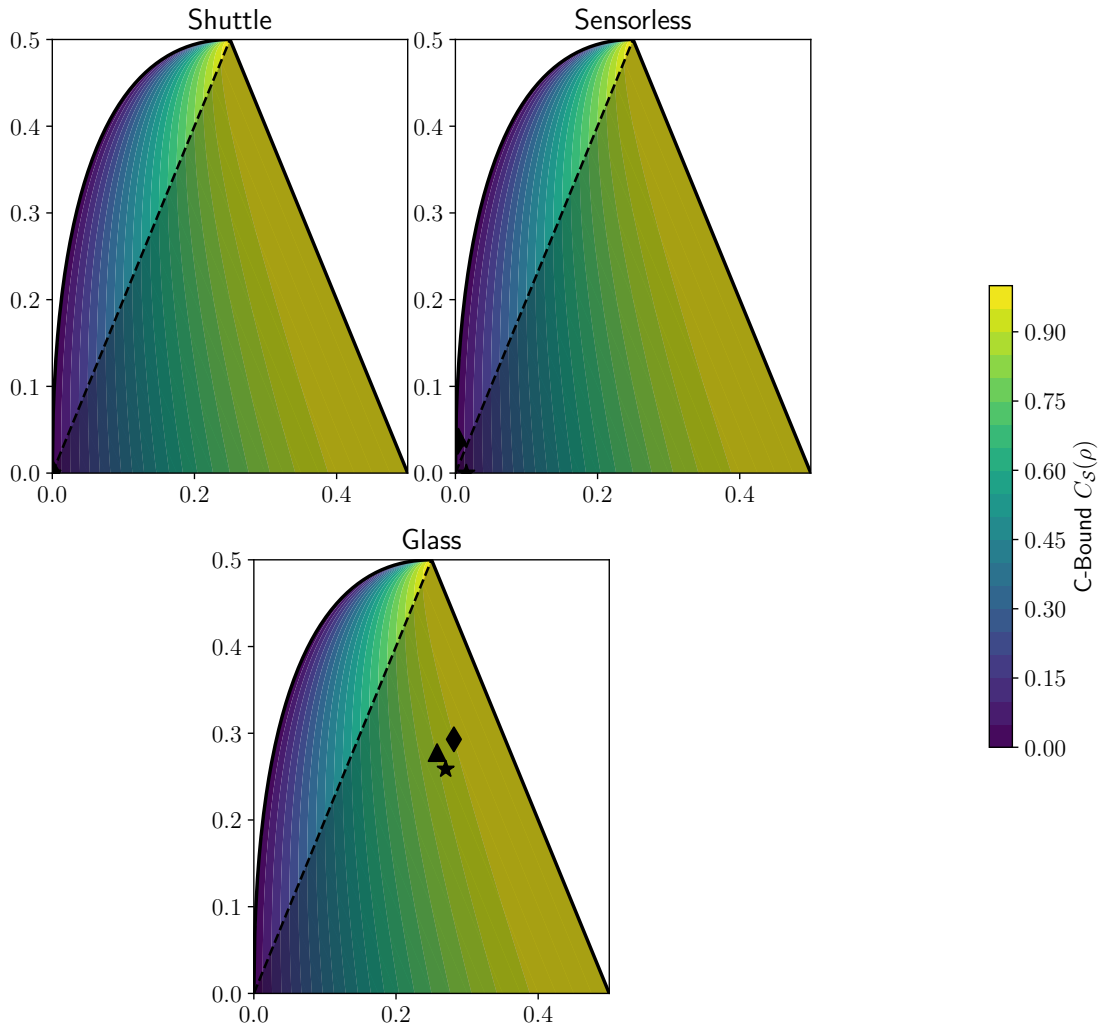


Figure D.7. Representation of all the possible values of the empirical C-Bound $C_S(\rho)$ in function of the disagreement $d_S(\rho)$ (y-axis) and joint error $e_S(\rho)$ (x-axis). We report the values obtained with the decision trees in the multi-class setting by Algorithm 4.3 (♦), Masegosa (▲), Germain (★). The disagreement and the joint are averaged over the 10 runs.

Table D.1. Comparison of the true risks $R_{\mathcal{T}}(MV_{\rho})$ and bound values obtained for each algorithm over 10 runs when the voters are decision stumps in the binary setting. More precisely, we report the mean \pm the standard deviation. “Bound” is the mean value of the bound that is optimized, excepted for MinCq and CB-Boost for which we report the bound obtained with Theorem 4.4.1 instantiated with the majority vote learned. Results in **bold** are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated to **the lowest risk** value. *Italic and underlined* results are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated respectively to the lowest bound value and the second lowest bound values.

	Algorithm 4.1		Algorithm 4.2		Algorithm 4.3		Masegosa	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
Credit	.141 \pm .014	.772 \pm .026	<u>.141 \pm .014</u>	<u>.718 \pm .036</u>	.141 \pm .014	.748 \pm .037	.141 \pm .014	.784 \pm .047
Heart	.252 \pm .034	.970 \pm .011	.252 \pm .034	.960 \pm .013	<u>.161 \pm .028</u>	<u>.937 \pm .017</u>	.163 \pm .025	1.041 \pm .034
USVotes	.043 \pm .007	.657 \pm .016	<u>.043 \pm .007</u>	<u>.494 \pm .034</u>	<u>.042 \pm .008</u>	<u>.529 \pm .029</u>	<u>.042 \pm .008</u>	<u>.520 \pm .030</u>
WDBC	.101 \pm .030	.722 \pm .016	.115 \pm .031	.675 \pm .050	.069 \pm .014	.578 \pm .024	<u>.060 \pm .010</u>	<u>.533 \pm .021</u>
TicTacToe	.296 \pm .011	.969 \pm .008	.296 \pm .011	.967 \pm .009	<u>.303 \pm .015</u>	<u>.958 \pm .009</u>	.272 \pm .024	1.021 \pm .016
SVMGuide	.085 \pm .007	.463 \pm .028	.085 \pm .007	.385 \pm .025	.081 \pm .014	.325 \pm .011	<u>.076 \pm .012</u>	<u>.313 \pm .011</u>
Haberman	.266 \pm .025	.975 \pm .012	<u>.263 \pm .025</u>	<u>.968 \pm .015</u>	.262 \pm .024	.988 \pm .020	<u>.260 \pm .019</u>	<u>1.207 \pm .050</u>

	Germain		CB-Boost		MinCq	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
Credit	<i>.141 \pm .014</i>	<i>.462 \pm .033</i>	<u>.140 \pm .015</u>	<u>.917 \pm .049</u>	.141 \pm .020	1.000 \pm .000
Heart	<i>.257 \pm .026</i>	<i>.796 \pm .036</i>	.191 \pm .016	.996 \pm .007	.185 \pm .017	1.000 \pm .000
USVotes	<i>.068 \pm .040</i>	<i>.351 \pm .108</i>	.043 \pm .007	.683 \pm .038	.048 \pm .009	1.000 \pm .000
WDBC	<i>.105 \pm .023</i>	<i>.412 \pm .032</i>	<u>.044 \pm .008</u>	<u>.687 \pm .044</u>	.044 \pm .013	1.000 \pm .000
TicTacToe	<i>.296 \pm .011</i>	<i>.812 \pm .023</i>	.202 \pm .020	1.000 \pm .000	<u>.020 \pm .003</u>	<u>1.000 \pm .000</u>
SVMGuide	<i>.102 \pm .035</i>	<i>.256 \pm .086</i>	.068 \pm .012	.334 \pm .012	<u>.048 \pm .002</u>	<u>1.000 \pm .000</u>
Haberman	<i>.265 \pm .026</i>	<i>.811 \pm .046</i>	.274 \pm .029	.988 \pm .011	.261 \pm .016	1.000 \pm .000

Table D.2. Comparison of the true risks $R_{\mathcal{T}}(MV_{\rho})$ and bound values obtained for each algorithm over 10 runs when the voters are decision trees in the binary setting. More precisely, we report the mean \pm the standard deviation. “Bound” is the mean value of the bound that is optimized, excepted for MinCq and CB-Boost for which we report the bound obtained with Theorem 4.4.1 instantiated with the majority vote learned. Results in **bold** are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated to **the lowest risk** value. *Italic and underlined* results are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated respectively to the lowest bound value and the second lowest bound values.

	Algorithm 4.1		Algorithm 4.2		Algorithm 4.3		Masegosa	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
Credit	.156 \pm .018	.861 \pm .034	.154 \pm .021	.807 \pm .048	<u>.142 \pm .015</u>	<u>.744 \pm .057</u>	.138 \pm .015	.775 \pm .081
Heart	.230 \pm .027	.984 \pm .011	.228 \pm .025	.976 \pm .018	.210 \pm .026	.956 \pm .025	.213 \pm .026	1.161 \pm .076
USVotes	.053 \pm .011	.738 \pm .031	.057 \pm .012	.571 \pm .067	.046 \pm .014	.523 \pm .059	<u>.051 \pm .023</u>	<u>.513 \pm .069</u>
WDBC	.054 \pm .010	.705 \pm .022	.061 \pm .010	.554 \pm .041	.045 \pm .007	.485 \pm .042	<u>.049 \pm .011</u>	<u>.471 \pm .050</u>
TicTacToe	.056 \pm .010	.776 \pm .025	.078 \pm .023	.685 \pm .047	.048 \pm .013	.530 \pm .052	<u>.050 \pm .010</u>	<u>.493 \pm .054</u>
SVMGuide	.033 \pm .002	.340 \pm .011	.033 \pm .001	.213 \pm .012	.032 \pm .002	.177 \pm .012	.032 \pm .002	.165 \pm .012
Haberman	.307 \pm .031	.998 \pm .004	.307 \pm .030	.997 \pm .006	<u>.295 \pm .025</u>	<u>.997 \pm .006</u>	.294 \pm .018	1.586 \pm .113

	Germain		CB-Boost		MinCq	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
Credit	.167 \pm .023	.564 \pm .059	.147 \pm .017	.772 \pm .054	.140 \pm .021	.981 \pm .025
Heart	.234 \pm .022	.840 \pm .059	.222 \pm .026	.972 \pm .019	.222 \pm .029	1.000 \pm .000
USVotes	.056 \pm .012	.334 \pm .050	.048 \pm .019	.578 \pm .050	.053 \pm .010	.989 \pm .032
WDBC	.063 \pm .010	.324 \pm .033	.044 \pm .012	.525 \pm .042	.054 \pm .013	.984 \pm .041
TicTacToe	.127 \pm .025	.457 \pm .047	.049 \pm .014	.545 \pm .052	.049 \pm .013	.901 \pm .080
SVMGuide	.038 \pm .005	.123 \pm .008	.033 \pm .002	.183 \pm .012	.033 \pm .002	.660 \pm .185
Haberman	.306 \pm .030	.970 \pm .058	.295 \pm .024	.999 \pm .002	.296 \pm .016	1.000 \pm .000

Table D.3. Comparison of the true risks $R_{\mathcal{T}}(MV_{\rho})$ and bound values obtained for each algorithm over 10 runs when the voters are decision trees in the multi-class setting. More precisely, we report the mean \pm the standard deviation. Results in **bold** are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated to **the lowest risk** value. *Italic and underlined results are the couple $(R_{\mathcal{T}}(MV_{\rho}), \text{Bound})$ associated respectively to the lowest bound value and the second lowest bound values.*

	Algorithm 4.1		Algorithm 4.2		Algorithm 4.3	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
MNIST	.034 \pm .002	.387 \pm .003	<u>.034 \pm .002</u>	<u>.371 \pm .003</u>	<i>.034 \pm .002</i>	<i>.340 \pm .004</i>
FashionMNIST	.121 \pm .004	.554 \pm .003	.121 \pm .004	.544 \pm .003	<u>.121 \pm .004</u>	<u>.523 \pm .003</u>
Pendigits	.015 \pm .007	.283 \pm .007	.015 \pm .007	.197 \pm .009	.015 \pm .007	.137 \pm .009
Protein	<i>.467 \pm .060</i>	<i>1.000 \pm .000</i>	<i>.462 \pm .057</i>	<i>1.000 \pm .000</i>	.419 \pm .038	1.000 \pm .000
Shuttle	.000 \pm .000	.061 \pm .001	.000 \pm .000	.006 \pm .001	.000 \pm .000	.005 \pm .000
Sensorless	.002 \pm .001	.097 \pm .001	.002 \pm .000	.040 \pm .001	.001 \pm .000	.022 \pm .001
Glass	<u>.328 \pm .058</u>	<u>1.000 \pm .000</u>	<u>.327 \pm .058</u>	<u>1.000 \pm .000</u>	<i>.328 \pm .060</i>	<i>1.000 \pm .000</i>

	Masegosa		Germain	
	$R_{\mathcal{T}}(MV_{\rho})$	Bound	$R_{\mathcal{T}}(MV_{\rho})$	Bound
MNIST	.095 \pm .037	.451 \pm .062	.180 \pm .005	.382 \pm .003
FashionMNIST	.195 \pm .022	.672 \pm .070	<i>.231 \pm .004</i>	<i>.479 \pm .002</i>
Pendigits	.027 \pm .017	.168 \pm .030	<u>.064 \pm .022</u>	<u>.163 \pm .010</u>
Protein	.460 \pm .025	1.503 \pm .047	<u>.529 \pm .003</u>	<u>1.098 \pm .006</u>
Shuttle	<u>.000 \pm .000</u>	<u>.005 \pm .001</u>	<i>.001 \pm .000</i>	<i>.003 \pm .001</i>
Sensorless	<u>.002 \pm .001</u>	<u>.025 \pm .001</u>	.015 \pm .001	.040 \pm .002
Glass	.324 \pm .056	1.978 \pm .140	.327 \pm .058	1.263 \pm .081

E.1 About Equation (5.1)

For the sake of completeness, we prove the bound $R_{\mathcal{D}'}(\text{MV}_\rho) \leq 2b_{\mathcal{D}'}^N(\rho)$ in the following lemma.

Lemma E.1.1. For any distribution \mathcal{D}' on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any distribution ρ on \mathbb{H} , for any $N \in \mathbb{N}_*$, we have

$$\begin{aligned} R_{\mathcal{D}'}(\text{MV}_\rho) &\leq 2b_{\mathcal{D}'}^N(\rho) \\ &= 2 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \left[\sum_{j=\lceil \frac{N}{2} \rceil}^N \binom{N}{j} \left[\frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^j \left[1 - \frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^{(N-j)} \right]. \end{aligned}$$

Proof. The proof is based on LACASSE *et al.* (2010). Note that for a given example $(\mathbf{x}, y) \sim \mathcal{D}'$ s.t. $\widehat{m}_\rho(\mathbf{x}, y) = 0$, we have

$$\underbrace{\sum_{j=\lceil \frac{N}{2} \rceil}^N \binom{N}{j} \left[\frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^j \left[1 - \frac{1}{2} (1 - \widehat{m}_\rho(\mathbf{x}, y)) \right]^{(N-j)}}_{\triangleq \heartsuit_\rho^N(\mathbf{x}, y)} \geq \frac{1}{2}.$$

Moreover, $\heartsuit_\rho^N(\mathbf{x}, y)$ is monotonically decreasing in $\widehat{m}_\rho(\mathbf{x}, y)$. From these two properties, we have for a given example $(\mathbf{x}, y) \sim \mathcal{D}'$ s.t. $\widehat{m}_\rho(\mathbf{x}, y) \leq 0$

$$\underbrace{\mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0]}_{=1} \leq 2\heartsuit_\rho^N(\mathbf{x}, y),$$

and for a given example $(\mathbf{x}, y) \sim \mathcal{D}'$ s.t. $\widehat{m}_\rho(\mathbf{x}, y) > 0$

$$\underbrace{\mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0]}_{=0} \leq 2\heartsuit_\rho^N(\mathbf{x}, y).$$

Hence, we can deduce that

$$R_{\mathcal{D}'}(\text{MV}_\rho) \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \mathbb{I}[\widehat{m}_\rho(\mathbf{x}, y) \leq 0] \leq 2 \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}'} \heartsuit_\rho^N(\mathbf{x}, y) = 2b_{\mathcal{D}'}^N(\rho). \quad \blacksquare$$

Based on Lemma E.1.1, we are now able to prove a PAC-Bayesian bound on the majority vote's true risk based on the surrogate $b_{\mathcal{D}'}^N(\rho)$. Note that LACASSE *et al.* (2010) prove a CATONI-like bound while our work is based on a SEEGER-like one. Our bound avoids using a parameter $c \geq 0$ required in CATONI's bound. We derive the bound in the following theorem.

Theorem E.1.1. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$ we have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_\rho) \leq 2 \left[\overline{\text{kl}} \left(b_{\mathcal{S}}^N(\rho) \left| \frac{1}{m} \left[N \text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right). \quad (\text{E.1})$$

Proof. First, note that we have by definition

$$b_{\mathcal{D}}^N(\rho) = \mathbb{P}_{\text{MV}_\sigma \sim \rho^N} [\widehat{m}_\sigma(\mathbf{x}, y) \leq 0].$$

We apply Theorem 2.3.4 with the loss $\ell(h, (\mathbf{x}, y)) = \mathbb{I}[\widehat{m}_\sigma(\mathbf{x}, y) \leq 0]$, the posterior distribution ρ^N and the prior distribution π^N to have

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad b_{\mathcal{D}}^N(\rho) \leq \overline{\text{kl}} \left(b_{\mathcal{S}}^N(\rho) \left| \frac{1}{m} \left[\text{KL}(\rho^N \|\pi^N) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right),$$

where $\text{KL}(\rho^N \|\pi^N) = N \text{KL}(\rho \|\pi)$. Then, from Lemma E.1.1, we obtain

$$\forall \rho \in \mathbb{M}(\mathbb{H}), \quad R_{\mathcal{D}}(\text{MV}_\rho) \leq 2 \left[\overline{\text{kl}} \left(b_{\mathcal{S}}^N(\rho) \left| \frac{1}{m} \left[N \text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right). \quad \blacksquare$$

From Theorem E.1.1, we can deduce a self-bounding algorithm denoted as Lacasse to minimize the majority vote's true risk. The objective function is defined as

$$G_{\mathcal{U}}^{\text{L}}(\rho) = 2 \left[\overline{\text{kl}} \left(b_{\mathcal{U}}^N(\rho) \left| \frac{1}{m} \left[N \text{KL}(\rho \|\pi) + \ln \frac{2\sqrt{m}}{\delta} \right] \right) \right),$$

E.2. Aggregation Property of the Dirichlet Distributions

where $b_S^N(\rho)$ is approximated through a mini-batch $\mathbb{U} \subseteq \mathbb{S}$ (and $b_{\mathbb{U}}^N(\rho)$ is computed instead). The algorithm denoted as Lacasse is described in Algorithm E.1 below.

Algorithm E.1 Minimization of Equation (E.1) by Stochastic Gradient Descent

Given: learning sample \mathbb{S} , prior distribution π on \mathbb{H} , the objective function $G_{\mathbb{U}}^L(\rho)$

Hyperparameters: number of iterations T

$\rho \leftarrow \pi$

for $t \leftarrow 1$ to T **do**

for all mini-batches $\mathbb{U} \subseteq \mathbb{S}$ **do**

$\rho \leftarrow$ Update ρ with $G_{\mathbb{U}}^L(\rho)$ by gradient descent

E.2 Aggregation Property of the Dirichlet Distributions

Lemma E.2.1. Let $\rho \sim \text{Dir}(\alpha)$ with $\alpha \in (\mathbb{R}_*^+)^K$ and $K = \text{card}(\mathbb{H})$. Then, the random variable associated to the random variable ρ where we sum the entries i and j follows a Dirichlet distribution with parameters $\alpha' = [\alpha_1, \dots, \alpha_i + \alpha_j, \dots, \alpha_K]^\top \in (\mathbb{R}_*^+)^{K-1}$, i.e., we have

$$\left(\rho(h_1), \dots, \rho(h_i) + \rho(h_j), \dots, \rho(h_K) \right) \sim \text{Dir}(\alpha').$$

Proof. Without loss of generality let $i = 1$ and $j = 2$, then, first remark that we have

$$\begin{aligned} & \int_0^{\rho(h_1) + \rho(h_2)} [x]^{\alpha_1 - 1} [\rho(h_1) + \rho(h_2) - x]^{\alpha_2 - 1} dx \\ &= \int_0^1 [(\rho(h_1) + \rho(h_2))x']^{\alpha_1 - 1} [(\rho(h_1) + \rho(h_2))(1 - x')]^{\alpha_2 - 1} (\rho(h_1) + \rho(h_2)) dx' \\ &= [\rho(h_1) + \rho(h_2)]^{\alpha_1 + \alpha_2 - 1} \int_0^1 [x']^{\alpha_1 - 1} [1 - x']^{\alpha_2 - 1} dx' \\ &\propto [\rho(h_1) + \rho(h_2)]^{\alpha_1 + \alpha_2 - 1}. \end{aligned} \tag{E.2}$$

Then, from Equation (E.2), the probability density function of $(\rho(h_1), \dots, \rho(h_i) +$

$\rho(h_j), \dots, \rho(h_K)$) can be rewritten (up to the normalization constant) as

$$\begin{aligned} & \left(\int_0^{\rho(h_1)+\rho(h_2)} [x]^{\alpha_1-1} [\rho(h_1)+\rho(h_2)-x]^{\alpha_2-1} dx \right) \left(\prod_{i=3}^K [\rho(h_i)]^{\alpha_i-1} \right) \\ & \propto [\rho(h_1)+\rho(h_2)]^{\alpha_1+\alpha_2-1} \prod_{i=3}^K [\rho(h_i)]^{\alpha_i-1}, \end{aligned}$$

which proves the desired result. ■

E.3 Proof of Lemma 5.2.1

Lemma 5.2.1 (Computation of the Stochastic Risk). For a given $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$, let

$$\mathbb{F}(\mathbf{x}, y) = \{j : h_j(\mathbf{x}) \neq y\} \quad \text{and} \quad \mathbb{T}(\mathbf{x}, y) = \{j : h_j(\mathbf{x}) = y\}$$

be respectively the set of indices of the voters that misclassify (\mathbf{x}, y) and the set of indices of the voters that correctly classify (\mathbf{x}, y) . Then, the stochastic risk $s_{\mathbb{P}}(\mathbf{x}, y)$ can be rewritten as

$$s_{\mathbb{P}}(\mathbf{x}, y) = \mathbb{E}_{\rho \sim \mathbb{P}} \mathbb{I}[\widehat{m}_{\rho}(\mathbf{x}, y) \leq 0] = I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right),$$

with $I_{0.5}(\cdot)$ the regularized incomplete beta function evaluated at 0.5. It is defined as

$$I_{0.5}(a, b) \triangleq \frac{B_{0.5}(a, b)}{B_1(a, b)}, \quad \text{where} \quad B_t(a, b) \triangleq \int_0^t x^{a-1} (1-x)^{b-1} dx$$

is the incomplete beta function.

Proof. Given an example $(\mathbf{x}, y) \in \mathbb{X} \times \mathbb{Y}$, by definition of the set $\mathbb{F}(\mathbf{x}, y)$ and

$\mathbb{T}(\mathbf{x}, y)$, we have

$$\begin{aligned} \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) \neq y] &= \sum_{j=1}^n \rho(h_j) \mathbb{I}[h_j(\mathbf{x}) \neq y] \\ &= \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \rho(h_j), \\ \text{and } \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] &= \sum_{j=1}^n \rho(h_j) \mathbb{I}[h_j(\mathbf{x}) = y] \\ &= \sum_{j \in \mathbb{T}(\mathbf{x}, y)} \rho(h_j). \end{aligned}$$

Moreover, by definition of the Dirichlet distribution (Definition 5.2.2), we have

$$\rho \sim \mathbb{P} \iff (\rho(h_1), \dots, \rho(h_n)) \sim \text{Dir}(\boldsymbol{\alpha}).$$

Then, we use the aggregation property of the Dirichlet distributions (Lemma E.2.1) to obtain

$$\begin{aligned} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \rho(h_j), \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \rho(h_j) \right) &\sim \text{Dir} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right) \\ \iff \left(\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y], \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) \neq y] \right) &\sim \text{Dir} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right). \end{aligned}$$

Thus, $(\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y], \mathbb{P}_{h \sim \rho} [h(\mathbf{x}) \neq y])$ follows a bivariate Dirichlet distribution *a.k.a.* Beta distribution. We have

$$\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] \sim \text{Beta} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right).$$

Finally, notice that the expected error is related to the cumulative probability function of the random variable $\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y]$ which is the regularized incomplete beta function $I_p : \mathbb{R}_*^+ \times \mathbb{R}_*^+ \rightarrow [0, 1]$:

$$\mathbb{P} \left[\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] \leq 0.5 \right] = I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right).$$

Then, we have

$$\begin{aligned}
 s_{\mathbb{P}}(\mathbf{x}, y) &= \mathbb{P}_{\rho \sim \mathbb{P}} \mathbb{I}[\widehat{m}_{\rho}(\mathbf{x}, y) \leq 0] \\
 &= \mathbb{P} \left[\mathbb{P}_{h \sim \rho} [h(\mathbf{x}) = y] \leq 0.5 \right] \\
 &= I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right).
 \end{aligned}$$

■

E.4 Proof of Corollary 5.2.1

Corollary 5.2.1 (Closed-form Solution of the Stochastic Risks). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any learning sample $\mathcal{S} \sim \mathcal{D}^m$, for any finite hypothesis set \mathbb{H} , for any distribution $\mathbb{P} = \text{Dir}(\boldsymbol{\alpha})$ with $\boldsymbol{\alpha} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, we have

$$\begin{aligned}
 \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{\mathbb{P}}(\mathbf{x}, y) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right), \\
 \text{and } \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{S}}(\text{MV}_{\rho}) &\leq \frac{1}{m} \sum_{i=1}^m s_{\mathbb{P}}(\mathbf{x}_i, y_i) = \frac{1}{m} \sum_{i=1}^m I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}_i, y_i)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}_i, y_i)} \alpha_j \right).
 \end{aligned}$$

Proof. From Equation (5.2) and Lemma 5.2.1, we have

$$\begin{aligned}
 \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{D}}(\text{MV}_{\rho}) &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{\mathbb{P}}(\mathbf{x}, y) \\
 &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}, y)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}, y)} \alpha_j \right).
 \end{aligned}$$

Similarly, from Equation (5.3) and Lemma 5.2.1, we have

$$\begin{aligned}
 \mathbb{E}_{\rho \sim \mathbb{P}} R_{\mathcal{S}}(\text{MV}_{\rho}) &\leq \frac{1}{m} \sum_{i=1}^m s_{\mathbb{P}}(\mathbf{x}_i, y_i) \\
 &= \frac{1}{m} \sum_{i=1}^m I_{0.5} \left(\sum_{j \in \mathbb{T}(\mathbf{x}_i, y_i)} \alpha_j, \sum_{j \in \mathbb{F}(\mathbf{x}_i, y_i)} \alpha_j \right).
 \end{aligned}$$

■

E.5 Proof of Theorem 5.3.1

Theorem 5.3.1 (PAC-Bayesian Bound for Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any finite hypothesis set \mathbb{H} , for any distribution $\Pi = \text{Dir}(\boldsymbol{\beta})$ with $\boldsymbol{\beta} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over the random choice of $\mathcal{S} \sim \mathcal{D}^m$, we have for all hyper-posterior P on \mathbb{H}

$$\mathbb{E}_{\rho \sim P} R_{\mathcal{D}}(\text{MV}_{\rho}) \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) \leq \bar{\text{kl}} \left(\frac{1}{m} \sum_{i=1}^m s_P(\mathbf{x}_i, y_i) \mid \frac{\text{KL}(P \parallel \Pi) + \ln \frac{2\sqrt{m}}{\delta}}{m} \right),$$

$$\begin{aligned} \text{with } \text{KL}(P \parallel \Pi) &= \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\beta_j)] - \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \beta_j \right) \right] - \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\alpha_j)] \\ &+ \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right] + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \left[\psi(\alpha_j) - \psi \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right], \end{aligned}$$

where $\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$ is the Gamma function and the Digamma function $\Psi(\alpha)$ is defined as the derivative of $\ln[\Gamma(\alpha)]$; these two functions are plotted in Figure 5.4.

Proof. This is a direct application of Theorem 2.3.4 (and Definition 2.3.3). Indeed, we apply Theorem 2.3.4 with the loss $\ell : \mathbb{M}(\mathbb{H}) \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ defined by $\ell(\rho, (\mathbf{x}, y)) = \mathbb{I}[\widehat{m}_{\rho}(\mathbf{x}, y) \leq 0]$.

Moreover, the closed-form solution of the KL divergence is

$$\begin{aligned}
\text{KL}(\Pi\|\text{P}) &= \mathbb{E}_{\rho\sim\text{P}} \ln \left(\frac{\text{P}(\rho)}{\Pi(\rho)} \right) \\
&= \mathbb{E}_{\rho\sim\text{Dir}(\boldsymbol{\alpha})} \ln \left(\frac{Z(\boldsymbol{\beta}) \prod_{j=1}^{\text{card}(\mathbb{H})} [\rho(h_j)]^{\alpha_j-1}}{Z(\boldsymbol{\alpha}) \prod_{j=1}^{\text{card}(\mathbb{H})} [\rho(h_j)]^{\beta_j-1}} \right) \\
&= \mathbb{E}_{\rho\sim\text{Dir}(\boldsymbol{\alpha})} \left[\ln Z(\boldsymbol{\beta}) - \ln Z(\boldsymbol{\alpha}) + \ln \left(\prod_{j=1}^{\text{card}(\mathbb{H})} [\rho(h_j)]^{\alpha_j-\beta_j} \right) \right] \\
&= \mathbb{E}_{\rho\sim\text{Dir}(\boldsymbol{\alpha})} \left[\ln Z(\boldsymbol{\beta}) - \ln Z(\boldsymbol{\alpha}) + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \ln [\rho(h_j)] \right] \\
&= \ln Z(\boldsymbol{\beta}) - \ln Z(\boldsymbol{\alpha}) + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \mathbb{E}_{\rho\sim\text{Dir}(\boldsymbol{\alpha})} \ln [\rho(h_j)] \\
&= \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\beta_j)] - \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \beta_j \right) \right] \\
&\quad - \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\alpha_j)] + \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right] \\
&\quad + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \left[\psi(\alpha_j) - \psi \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right].
\end{aligned}$$

The last equality follows by definition of Dirichlet's geometric mean

$$\mathbb{E}_{\rho\sim\text{P}} \ln [\rho(h_j)] = \psi(\alpha_j) - \psi \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right).$$

and the normalization constant

$$\ln(Z(\boldsymbol{\alpha})) = \sum_{j=1}^{\text{card}(\mathbb{H})} \ln[\Gamma(\alpha_j)] - \ln \left[\Gamma \left(\sum_{j=1}^{\text{card}(\mathbb{H})} \alpha_j \right) \right].$$

■

E.6 Proof of Theorem 5.3.2

Theorem 5.3.2 (PAC-Bayesian bound with data-dependent voters). Let Π_1 and P_1 be the hyper-prior and hyper-posterior distributions on \mathbb{H}_1 defined with \mathcal{S}_1 , and Π_2 and P_2 the prior and posterior distributions on \mathbb{H}_2 defined with \mathcal{S}_2 . For any $\lambda \in [0, 1]$ and $\delta \in (0, 1]$ with probability at least $1 - \delta$ over samples $\mathcal{S}_1 \sim \mathcal{D}^{m_1}$ and $\mathcal{S}_2 \sim \mathcal{D}^{m_2}$, we have

$$\lambda \mathbb{E}_{\rho \sim P_1} R_{\mathcal{D}}(\text{MV}_{\rho}) + (1-\lambda) \mathbb{E}_{\rho' \sim P_2} R_{\mathcal{D}}(\text{MV}_{\rho'}) \leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\lambda s_{P_1}(\mathbf{x}, y) + (1-\lambda) s_{P_2}(\mathbf{x}, y)] \leq \overline{\text{kl}} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_1} \frac{s_{P_1}(\mathbf{x}, y)}{\lambda} + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}_2} \frac{s_{P_2}(\mathbf{x}, y)}{1-\lambda} \mid \frac{\text{KL}(P_1 \parallel \Pi_1) + \ln \frac{4\sqrt{m}}{\delta}}{\frac{m}{\lambda}} + \frac{\text{KL}(P_2 \parallel \Pi_2) + \ln \frac{4\sqrt{m'}}{\delta}}{\frac{m'}{1-\lambda}} \right].$$

Proof. From the joint convexity of $\text{kl}()$, we have for any $\lambda \in [0, 1]$

$$\begin{aligned} & \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \frac{s_P(\mathbf{x}, y)}{\lambda} + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}'} \frac{s_{P'}(\mathbf{x}, y)}{1-\lambda} \mid \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{s_P(\mathbf{x}, y)}{\lambda} + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \frac{s_{P'}(\mathbf{x}, y)}{1-\lambda} \right] \\ & \leq \lambda \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} s_P(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) \right] + (1-\lambda) \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}'} s_{P'}(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{P'}(\mathbf{x}, y) \right]. \end{aligned}$$

For all hyper-prior Π on \mathbb{H} , we have with probability at least $1 - \frac{\delta}{2}$ on the random choice $\mathcal{S} \sim \mathcal{D}^m$ for all hyper-posterior P on \mathbb{H}

$$\begin{aligned} & \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} s_P(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) \right] \leq \frac{\text{KL}(P, \Pi) + \ln \frac{4\sqrt{m}}{\delta}}{m} \\ \iff & \forall \lambda \in [0, 1], \quad \lambda \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} s_P(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_P(\mathbf{x}, y) \right] \leq \lambda \frac{\text{KL}(P, \Pi) + \ln \frac{4\sqrt{m}}{\delta}}{m}. \end{aligned} \tag{E.3}$$

Similarly, for all hyper-prior Π' on \mathbb{H}' , we have with probability at least $1 - \frac{\delta}{2}$ on the random choice $\mathcal{S}' \sim \mathcal{D}^{m'}$ for all hyper-posterior P' on \mathbb{H}

$$\begin{aligned} & \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}'} s_{P'}(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{P'}(\mathbf{x}, y) \right] \leq \frac{\text{KL}(P', \Pi') + \ln \frac{4\sqrt{m'}}{\delta}}{m'} \iff \\ & \forall \lambda \in [0, 1], \quad (1-\lambda) \text{kl} \left[\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}'} s_{P'}(\mathbf{x}, y) \parallel \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} s_{P'}(\mathbf{x}, y) \right] \leq (1-\lambda) \frac{\text{KL}(P', \Pi') + \ln \frac{4\sqrt{m'}}{\delta}}{m'}. \end{aligned} \tag{E.4}$$

Combining Equation (E.3) and Equation (E.4) using the union bound, we obtain the desired result with $1 - \delta$ probability. ■

E.7 Additional Results

E.7.1 Choice of the prior

In the other experiments, we fixed the hyper-prior distribution Π (parameterized by β) to the uniform, *i.e.* $\forall j \in \{1, \dots, \text{card}(\mathbb{H})\}, \beta_j = 1$. This choice was to make the comparison with the baselines as fair as possible, as their prior was also fixed to the uniform. However, we can bias the sparsity of the posterior, or conversely its concentration, by choosing a different value for the prior distribution parameters. In some cases, tuning the prior parameters allows to obtain better performance, as reported in Figures E.1 to E.4. As in Section 5.4, the hypothesis sets \mathbb{H}_1 *resp.* \mathbb{H}_2 are composed of 50 decision trees learned with \mathcal{S}_1 *resp.* \mathcal{S}_2 with no limit on the depth. In general, these results suggest that the choice of prior distribution has a high impact on the learned model's performance and tuning its concentration parameters would be a viable option for improving the results.

E.7.2 Impact of voter strength

We report a study on the impact of voter strength on the learned models. More precisely, we provide results for additional datasets as well as the study of the expected strength of a voter as a function of the tree maximal depth. The hypothesis sets \mathbb{H}_1 *resp.* \mathbb{H}_2 are composed of 50 decision trees learned with \mathcal{S}_1 *resp.* \mathcal{S}_2 ; the prior's parameters are set to $\beta_j = 1$ for all $j \in \{1, \dots, \text{card}(\mathbb{H})\}$. The maximal depth is a values belonging to the set $\{1, 2, 4, 8, 16\}$, *i.e.*, the maximal depth varies from 1 to 16. We report in Figures E.5 to E.8 the stochastic test risks $\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{T}} s_P(\mathbf{x}, y)$ or the test risk $R_{\mathcal{T}}(\text{MV}_{\rho})$, their corresponding empirical risks and the bound values. We can see that limiting the maximal depth is an effective way for controlling the strength of the voters. Indeed, the general trend tells us that increasing the strength of the voters generally yields more powerful ensembles for all methods: the risks and the bounds are decreasing or stay constant when the tree maximal depth is increasing.

E.7. Additional Results

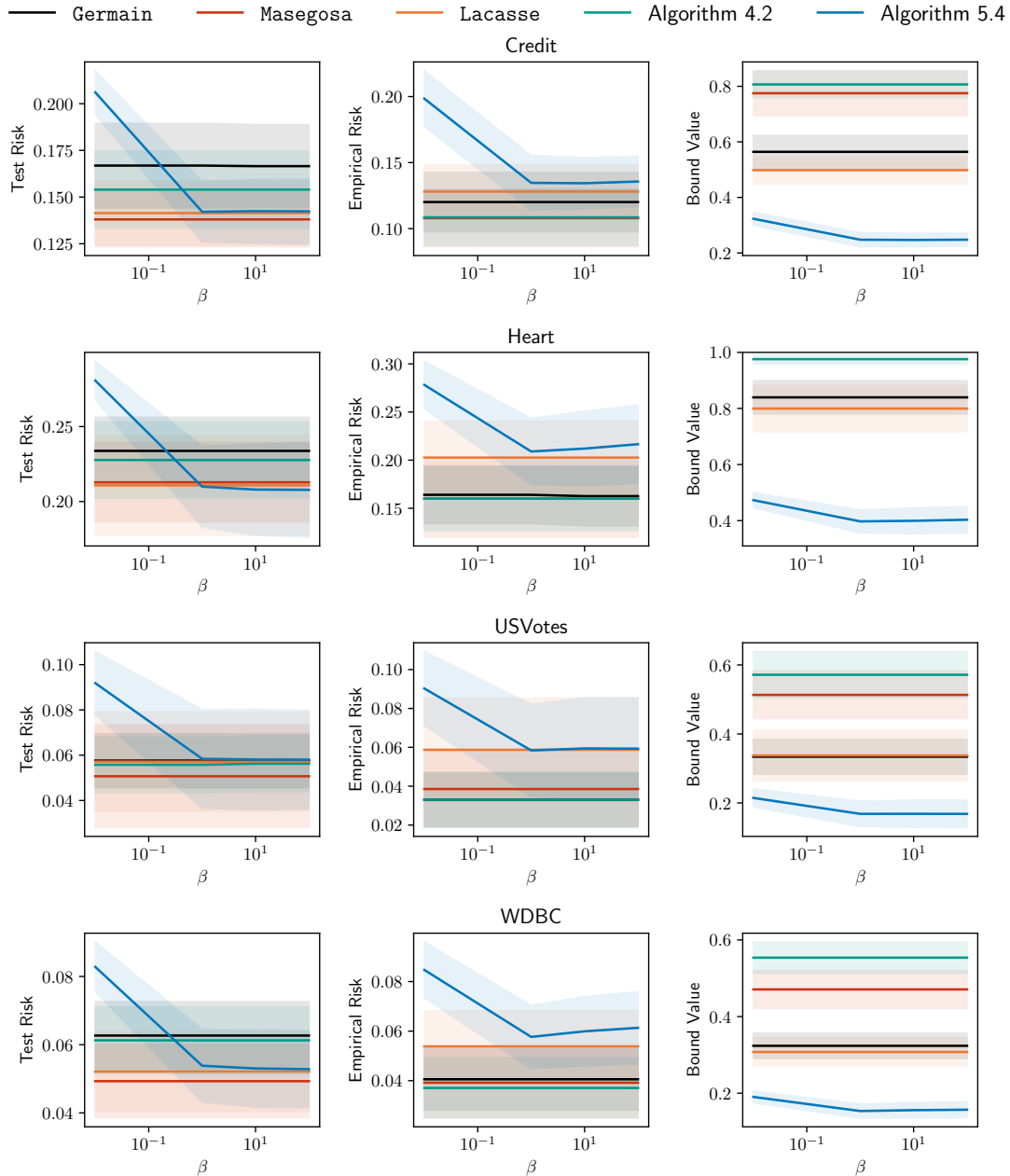


Figure E.1. Plot of the impact of the prior β on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of all the parameters β_j with $j \in \{1, \dots, \text{card}(\mathbb{H})\}$ and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

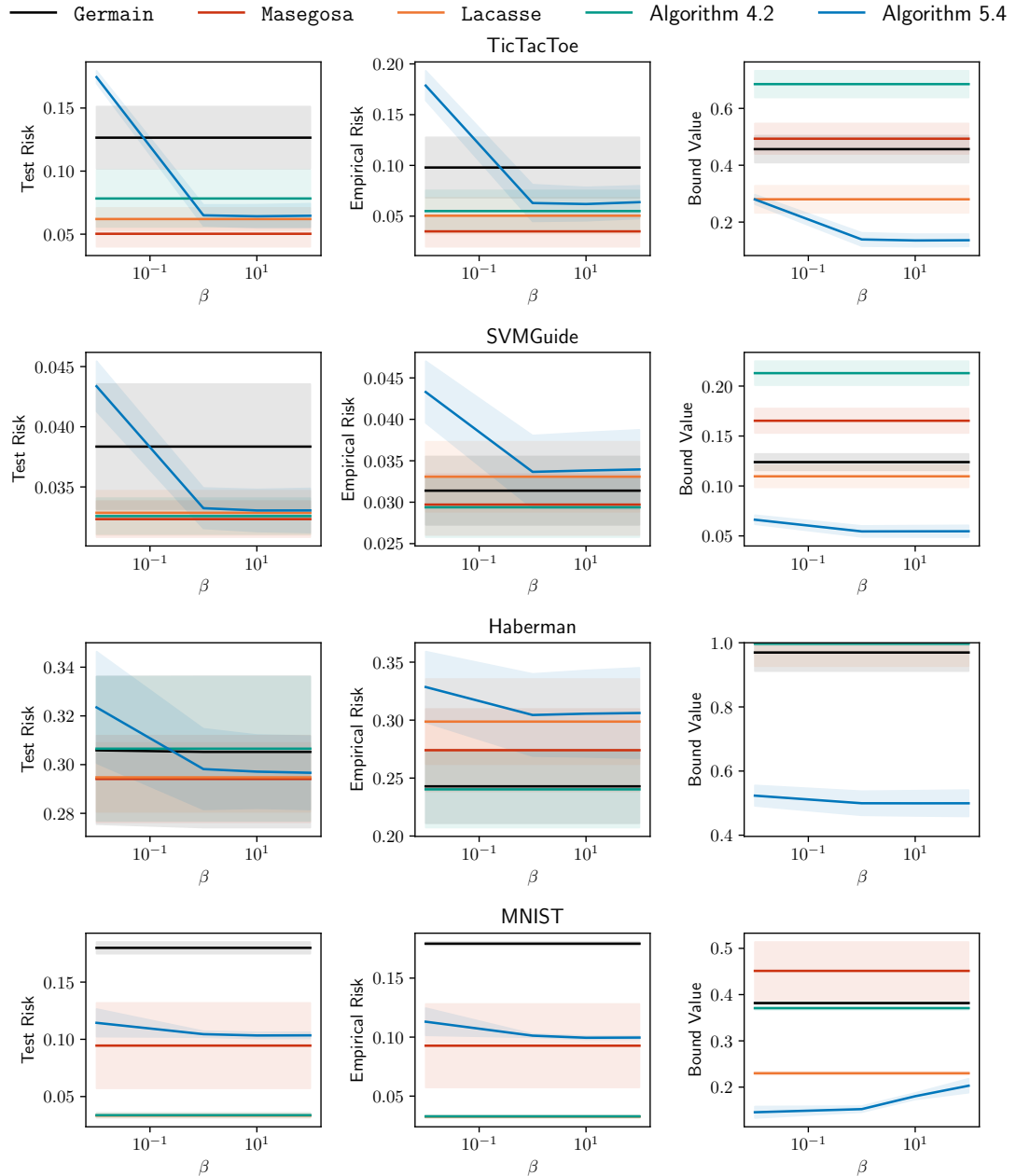


Figure E.2. Plot of the impact of the prior β on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of all the parameters β_j with $j \in \{1, \dots, \text{card}(\mathbb{H})\}$ and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

E.7. Additional Results

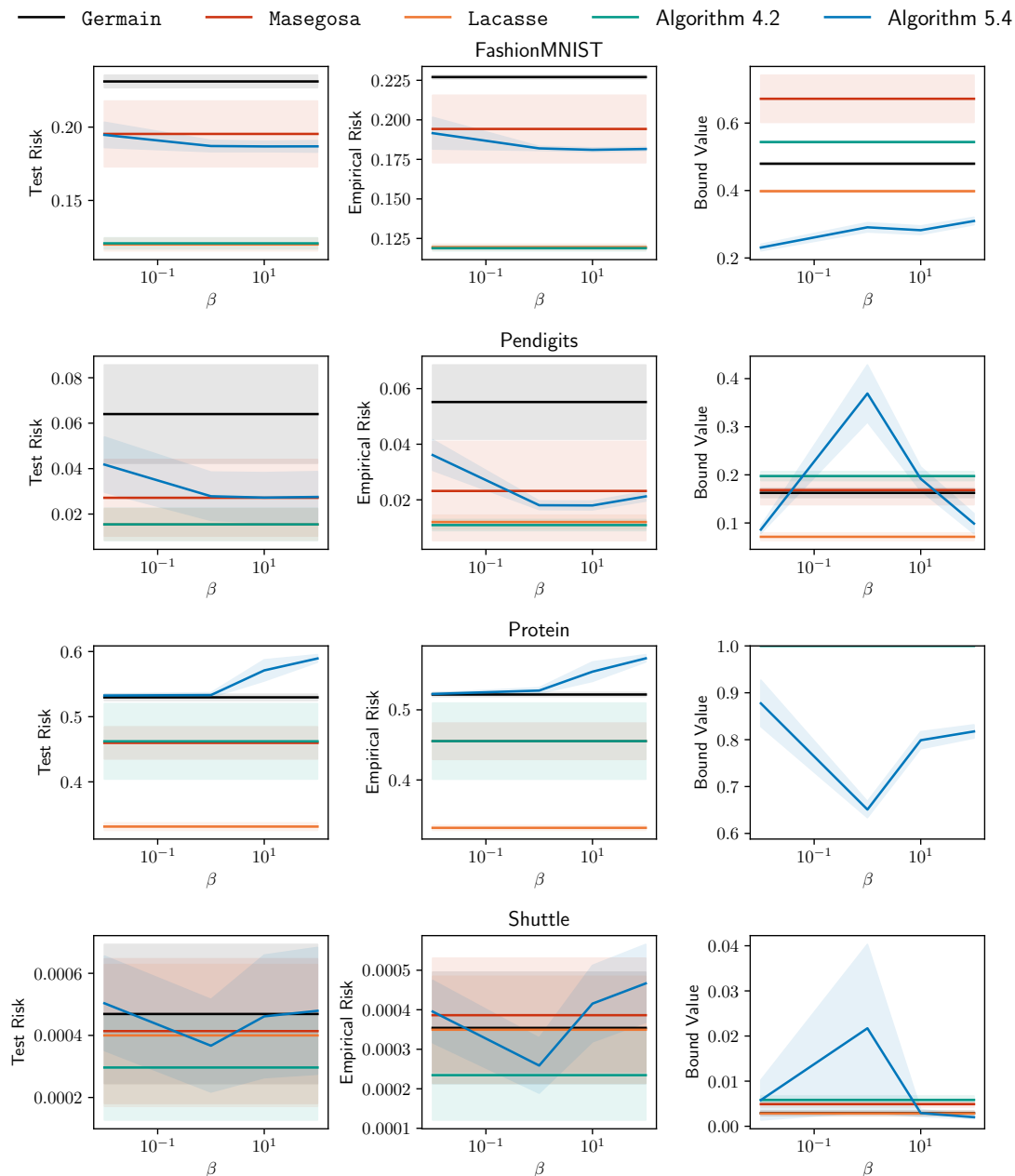


Figure E.3. Plot of the impact of the prior β on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of all the parameters β_j with $j \in \{1, \dots, \text{card}(\mathbb{H})\}$ and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

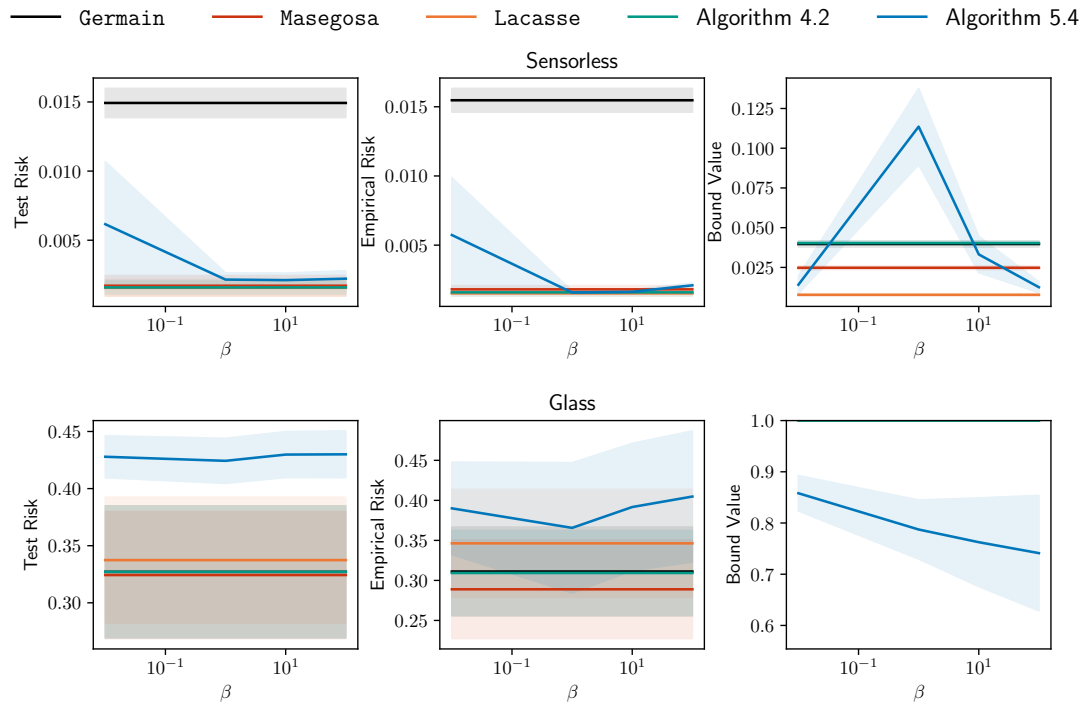


Figure E.4. Plot of the impact of the prior β on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of all the parameters β_j with $j \in \{1, \dots, \text{card}(\mathbb{H})\}$ and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

E.7. Additional Results

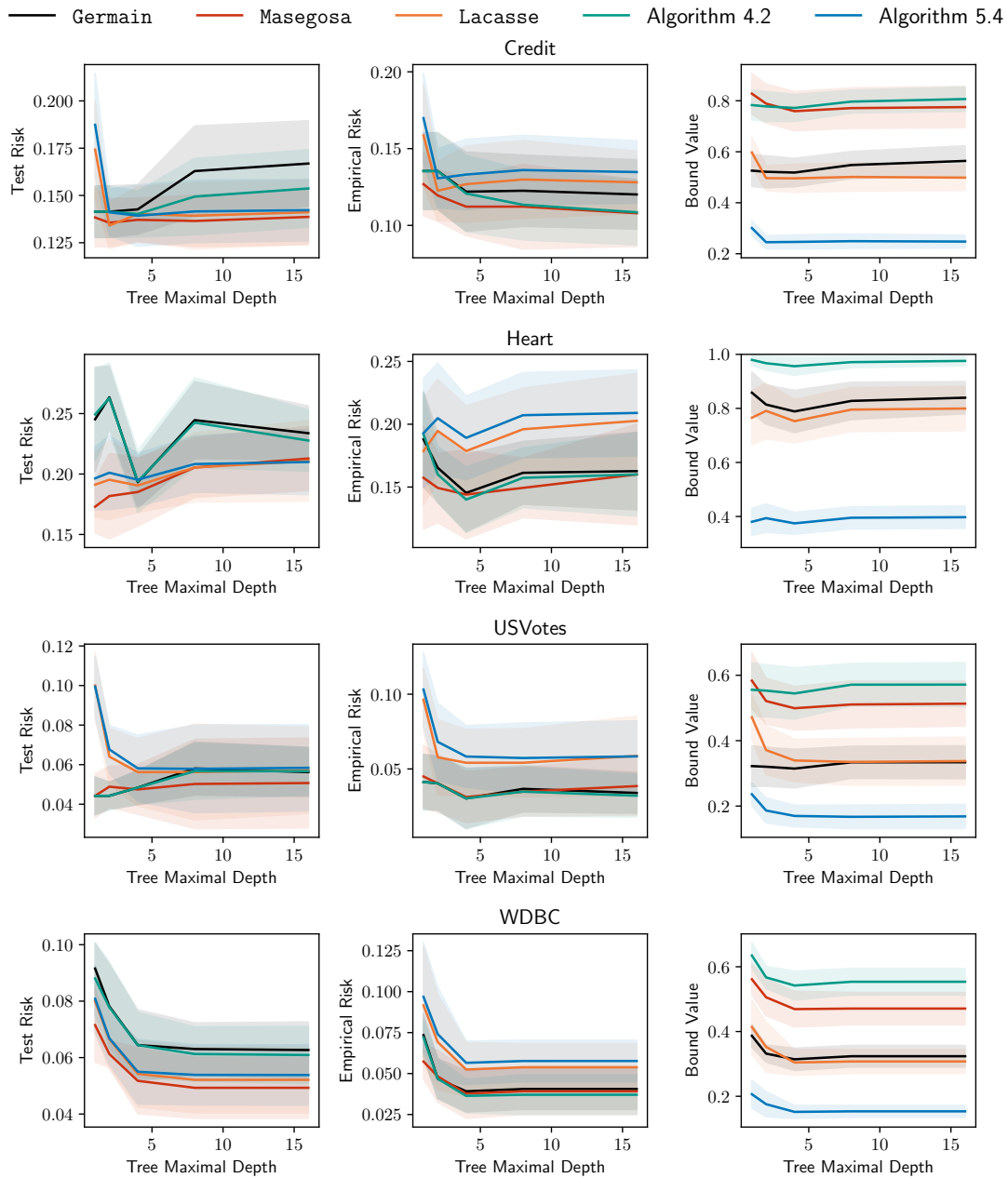


Figure E.5. Plot of the impact of tree maximal depth on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of the tree maximal depth and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

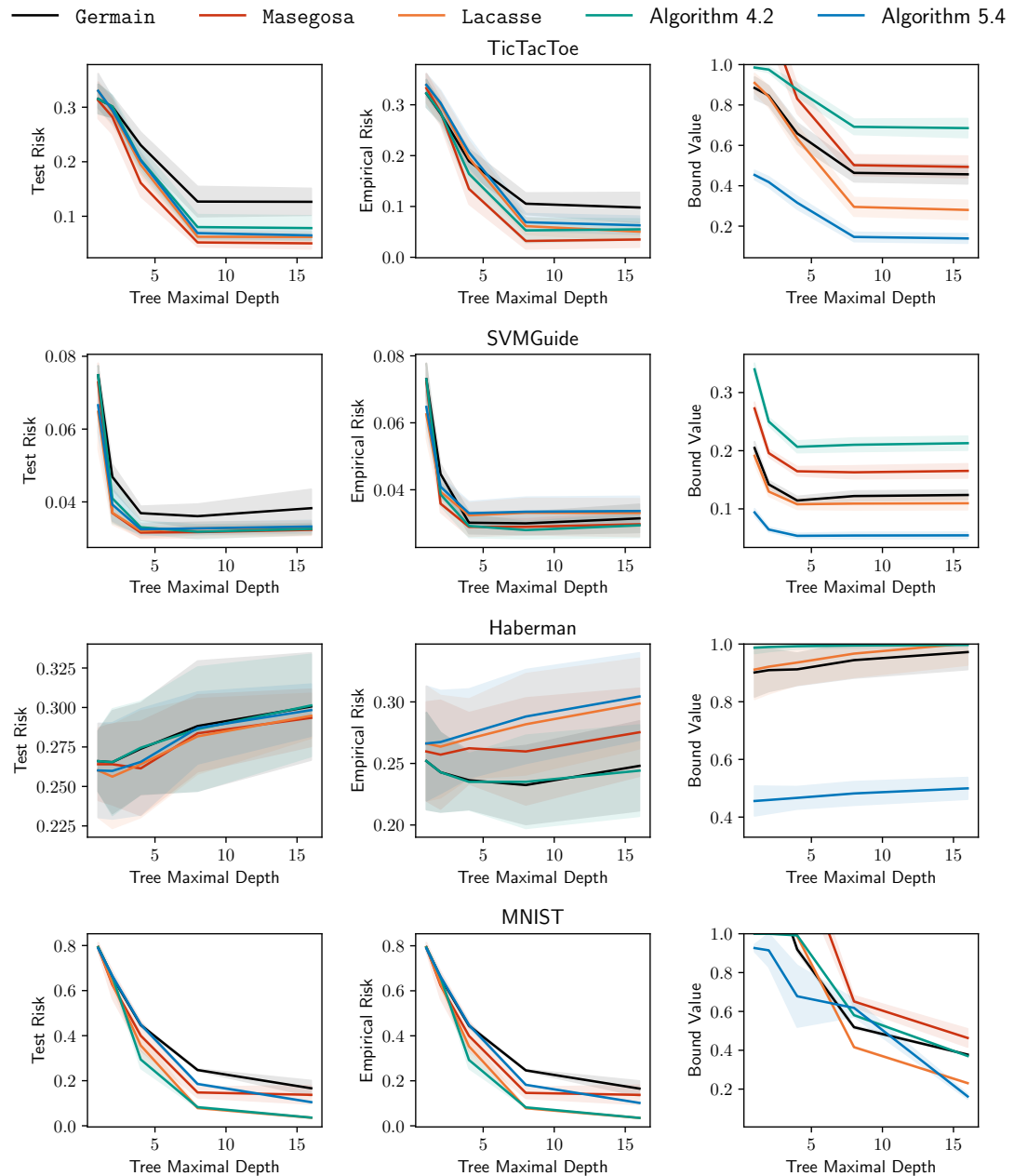


Figure E.6. Plot of the impact of tree maximal depth on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of the tree maximal depth and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

E.7. Additional Results

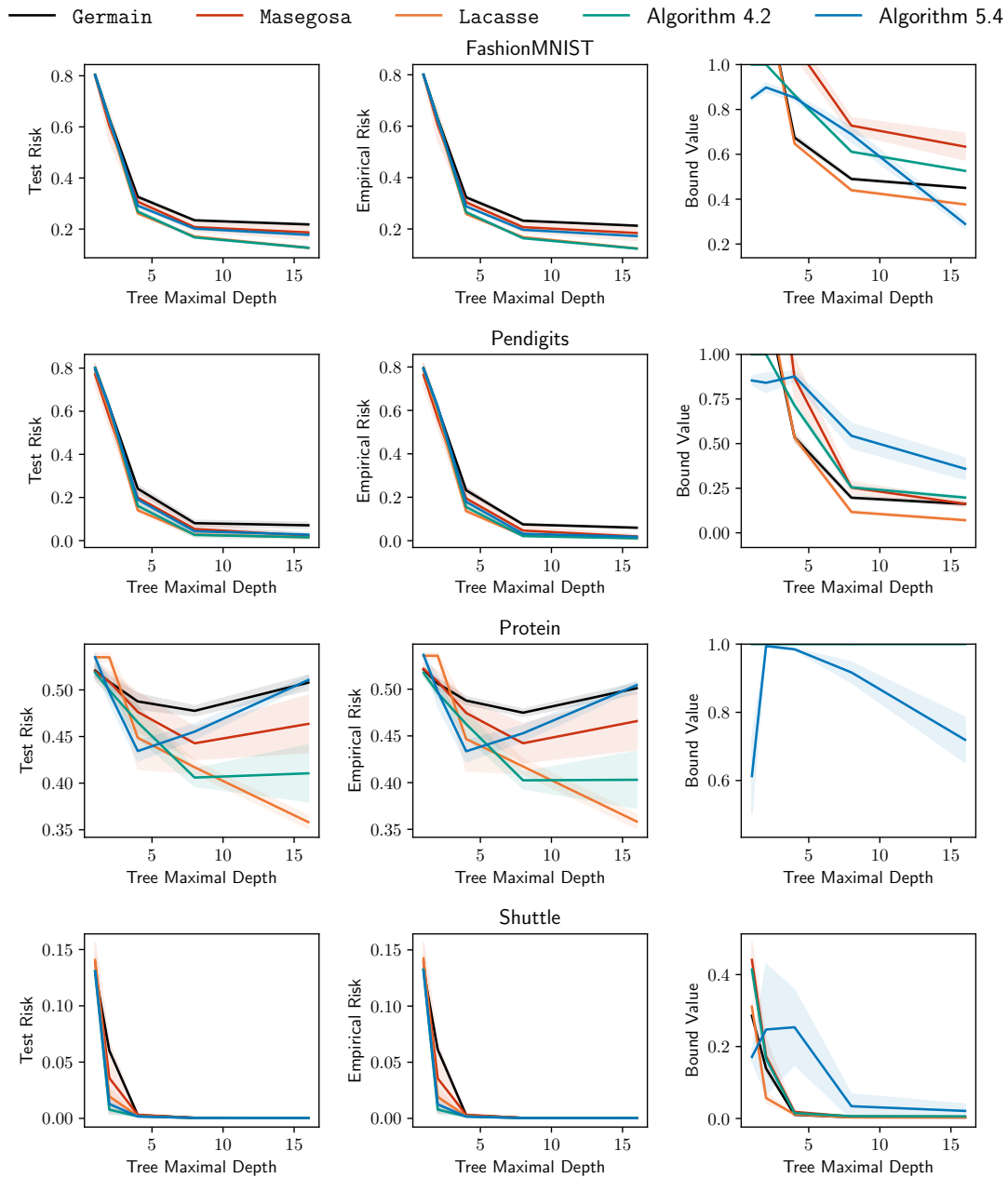


Figure E.7. Plot of the impact of tree maximal depth on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of the tree maximal depth and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

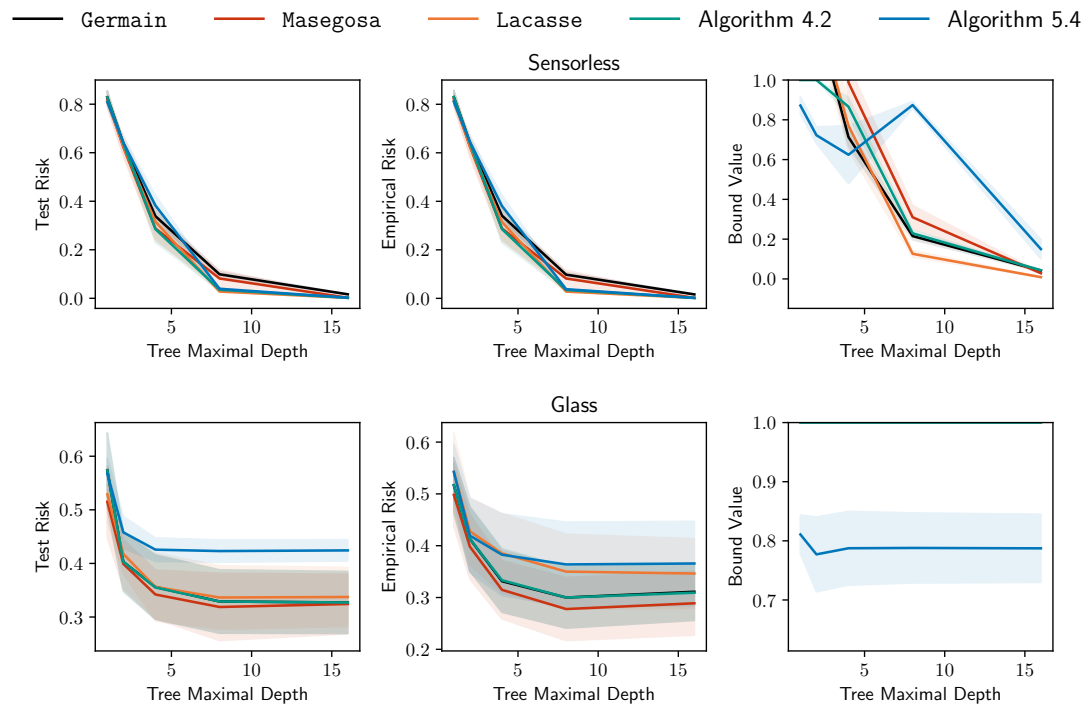


Figure E.8. Plot of the impact of tree maximal depth on the performance of the stochastic majority vote. More precisely, the x-axis represents the value of the tree maximal depth and the y-axis are the values of the test risks, the empirical risks or the bound values. The mean (plain lines) and the standard deviations (shadows) are obtained for all values on 10 runs.

APPENDIX OF CHAPTER 6

F

F.1 Proof of Theorem 6.3.1

Theorem 6.3.1 (General Disintegrated PAC-Bayes Bound). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi: \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_+$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\begin{aligned} & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \right) \\ & \leq \frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} + D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) + \ln \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \geq 1 - \delta, \end{aligned}$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Proof. For any sample $\mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m$, prior $\pi \in \mathbb{M}^*(\mathbb{H})$ and deterministic algorithm A fixed a priori, let $\rho_{\mathcal{S}} = A(\mathcal{S}, \pi)$ the distribution obtained from the algorithm A . Note that $\varphi(h, \mathcal{S})$ is a strictly-positive random variable. Hence, from MARKOV's inequality (Theorem A.2.1), we have

$$\begin{aligned} & \mathbb{P}_{h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2} \\ \iff & \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \mathbb{I} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2}. \end{aligned}$$

Taking the expectation over $\mathcal{S} \sim \mathcal{D}^m$ to both sides of the inequality gives

$$\begin{aligned} & \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \mathbb{I} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2} \\ \iff & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2}. \end{aligned}$$

Since both sides of the inequality are strictly positive, we can take the logarithm

and multiply by $\frac{\lambda}{\lambda-1} > 0$ to obtain

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq \frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right) \right] \geq 1 - \frac{\delta}{2}.$$

We develop the right-hand side of the inequality and take the expectation of the hypothesis over the prior distribution π . We have for all prior $\pi \in \mathcal{M}^*(\mathbb{H})$,

$$\begin{aligned} \frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right) &= \frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \frac{\rho_{\mathcal{S}}(h') \pi(h')}{\pi(h') \rho_{\mathcal{S}}(h')} \varphi(h', \mathcal{S}) \right) \\ &= \frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \pi} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \varphi(h', \mathcal{S}) \right), \end{aligned}$$

Remark that $\frac{1}{r} + \frac{1}{s} = 1$ with $r = \lambda$ and $s = \frac{\lambda}{\lambda-1}$. Hence, we can apply HÖLDER's inequality (Theorem A.5.1):

$$\mathbb{E}_{h' \sim \pi} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \varphi(h', \mathcal{S}) \leq \left[\mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \right]^{\lambda} \right) \right]^{\frac{1}{\lambda}} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}}.$$

Then, since both sides of the inequality are strictly positive, we take the logarithm; add $\ln(\frac{2}{\delta})$ and multiply by $\frac{\lambda}{\lambda-1} > 0$ to both sides of the inequality, to obtain

$$\begin{aligned} &\frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \pi} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \varphi(h', \mathcal{S}) \right) \\ &\leq \frac{\lambda}{\lambda-1} \ln \left(\frac{2}{\delta} \left[\mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \right]^{\lambda} \right) \right]^{\frac{1}{\lambda}} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\ &= \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{\mathcal{S}}(h')}{\pi(h')} \right]^{\lambda} \right) \right) + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left(\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right) \\ &= D_{\lambda}(\rho_{\mathcal{S}} \| \pi) + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left(\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right). \end{aligned}$$

From this inequality, we can deduce that

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\forall \pi \in \mathcal{M}^*(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq D_{\lambda}(\rho_{\mathcal{S}} \| \pi) \right. \\ \left. + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left(\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right) \right] \geq 1 - \frac{\delta}{2}. \quad (\text{F.1}) \end{aligned}$$

F.2. Proof of Corollary 6.3.1

Given a prior $\pi \in \mathbb{M}^*(\mathbb{H})$, note that $\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}}$ is a strictly positive random variable. Hence, we apply MARKOV's inequality (Theorem A.2.1) to have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \geq 1 - \frac{\delta}{2}.$$

Since the inequality does not depend on the random variable $h \sim \rho_{\mathbb{S}}$, we have

$$\begin{aligned} & \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \\ &= \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{I} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \\ &= \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathbb{S}}} \mathbb{I} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \\ &= \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]. \end{aligned}$$

Since both sides of the inequality are strictly positive, we take the logarithm to both sides of the inequality, and we add $\frac{\lambda}{\lambda-1} \ln \frac{2}{\delta}$ to have

$$\begin{aligned} & \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right] \geq 1 - \frac{\delta}{2} \iff \\ & \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left(\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^{\frac{\lambda}{\lambda-1}} \right) \right) \leq \frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} \right. \\ & \quad \left. + \ln \left(\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right) \right] \geq 1 - \frac{\delta}{2}. \end{aligned} \quad (\text{F.2})$$

Combining Equations (F.1) and (F.2) with a union bound gives us the desired result. ■

F.2 Proof of Corollary 6.3.1

Corollary 6.3.1 (Extreme Cases of Theorem 6.3.1). Under the assumptions of Theorem 6.3.1, when $\lambda \rightarrow 1^+$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left(\ln \varphi(h, \mathbb{S}) \leq \ln \frac{2}{\delta} + \ln \left[\operatorname{esssup}_{\mathbb{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathbb{H}} \varphi(h', \mathbb{S}') \right] \right) \geq 1 - \delta,$$

when $\lambda \rightarrow +\infty$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\ln \varphi(h, \mathcal{S}) \leq \ln \operatorname{esssup}_{h' \in \mathcal{H}} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} + \ln \left[\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}') \right] \right) \geq 1 - \delta,$$

where esssup is the essential supremum defined as the supremum on a set with non-zero probability measures, *i.e.*,

$$\operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathcal{H}} \varphi(h', \mathcal{S}') = \inf \left\{ \tau \in \mathbb{R}, \mathbb{P}_{\mathcal{S}' \sim \mathcal{D}^m, h' \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) > \tau \right] = 0 \right\},$$

and $\operatorname{esssup}_{h' \in \mathcal{H}} \frac{\rho_{\mathcal{S}}(h')}{\pi(h')} = \inf \left\{ \tau \in \mathbb{R}, \mathbb{P}_{h \sim \rho_{\mathcal{S}}} \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} > \tau \right] = 0 \right\}.$

Proof. Starting from Theorem 6.3.1 and rearranging, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\ln(\varphi(h, \mathcal{S})) \leq \frac{2\lambda-1}{\lambda} \ln \frac{2}{\delta} + \frac{\lambda-1}{\lambda} D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) + \ln \left(\left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \geq 1 - \delta.$$

Then, we will prove the case when $\lambda \rightarrow 1$ and $\lambda \rightarrow +\infty$ separately.

When $\lambda \rightarrow 1$. We have $\lim_{\lambda \rightarrow 1^+} \frac{2\lambda-1}{\lambda} \ln \frac{2}{\delta} = \ln \frac{2}{\delta}$ and $\lim_{\lambda \rightarrow 1^+} \frac{\lambda-1}{\lambda} D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) = 0$. Furthermore, note that

$$\|\varphi\|_{\frac{\lambda}{\lambda-1}} = \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(|\varphi(h', \mathcal{S}')|^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} = \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}}$$

is the $L^{\frac{\lambda}{\lambda-1}}$ -norm of the function $\varphi : \mathcal{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_+^*$, where $\lim_{\lambda \rightarrow 1} \|\varphi\|_{\frac{\lambda}{\lambda-1}} = \lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'}$ (since we have $\lim_{\lambda \rightarrow 1^+} \frac{\lambda}{\lambda-1} = (\lim_{\lambda \rightarrow 1} \lambda)(\lim_{\lambda \rightarrow 1} \frac{1}{\lambda-1}) = +\infty$). Then, it is well known that

$$\|\varphi\|_{\infty} = \lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'} = \operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathcal{H}} \varphi(h', \mathcal{S}').$$

Hence, we have

$$\begin{aligned}
 & \lim_{\lambda \rightarrow 1} \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow 1} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow 1} \|\varphi\|_{\frac{\lambda}{\lambda-1}} \right) = \ln \left(\lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'} \right) \\
 &= \ln (\|\varphi\|_{\infty}) = \ln \left(\operatorname{esssup}_{S' \in (\mathbb{X} \times \mathbb{Y}), h' \in \mathbb{H}} \varphi(h', S') \right).
 \end{aligned}$$

Finally, we can deduce that

$$\begin{aligned}
 & \lim_{\lambda \rightarrow 1} \left[\frac{2\lambda-1}{\lambda} \ln \frac{2}{\delta} + \frac{\lambda-1}{\lambda} D_{\lambda}(\rho_{\mathbb{S}} \|\pi) + \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \\
 &= \ln \frac{2}{\delta} + \ln \left[\operatorname{esssup}_{S' \in (\mathbb{X} \times \mathbb{Y}), h' \in \mathbb{H}} \varphi(h', S') \right].
 \end{aligned}$$

When $\lambda \rightarrow +\infty$. First, we have $\lim_{\lambda \rightarrow +\infty} \frac{2\lambda-1}{\lambda} \ln \frac{2}{\delta} = \ln \frac{2}{\delta} \left[2 - \lim_{\lambda \rightarrow +\infty} \frac{1}{\lambda} \right] = 2 \ln \frac{2}{\delta} = \ln \frac{4}{\delta^2}$ and $\lim_{\lambda \rightarrow +\infty} \|\varphi\|_{\frac{\lambda}{\lambda-1}} = \lim_{\lambda' \rightarrow 1} \|\varphi\|_{\lambda'} = \|\varphi\|_1$ (since $\lim_{\lambda \rightarrow +\infty} \frac{\lambda}{\lambda-1} = \lim_{\lambda \rightarrow +\infty} \frac{1}{1-\frac{1}{\lambda}} = 1$). Hence, we have

$$\begin{aligned}
 & \lim_{\lambda \rightarrow +\infty} \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow +\infty} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow +\infty} \|\varphi\|_{\frac{\lambda}{\lambda-1}} \right) = \ln \left(\lim_{\lambda' \rightarrow 1} \|\varphi\|_{\lambda'} \right) \\
 &= \ln (\|\varphi\|_1) = \ln \left(\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', S') \right).
 \end{aligned}$$

Moreover, by rearranging the terms in $\frac{\lambda-1}{\lambda} D_\lambda(\rho_S|\pi)$, we have

$$\begin{aligned} \frac{\lambda-1}{\lambda} D_\lambda(\rho_S|\pi) &= \frac{1}{\lambda} \ln \left(\mathbb{E}_{h \sim \pi} \left(\left[\frac{\rho_S(h)}{\pi(h)} \right]^\lambda \right) \right) = \ln \left(\left[\mathbb{E}_{h \sim \pi} \left(\left[\frac{\rho_S(h)}{\pi(h)} \right]^\lambda \right) \right]^{\frac{1}{\lambda}} \right) \\ &= \ln \left(\left[\mathbb{E}_{h \sim \pi} (\gamma(h)^\lambda) \right]^{\frac{1}{\lambda}} \right) = \ln(\|\gamma\|_\lambda), \end{aligned}$$

where $\|\gamma\|_\lambda$ is the L^λ -norm of the function γ defined as $\gamma(h) = \frac{\rho_S(h)}{\pi(h)}$. We have

$$\begin{aligned} \lim_{\lambda \rightarrow +\infty} \frac{\lambda-1}{\lambda} D_\lambda(\rho_S|\pi) &= \lim_{\lambda \rightarrow +\infty} \ln(\|\gamma\|_\lambda) = \ln \left(\lim_{\lambda \rightarrow +\infty} \|\gamma\|_\lambda \right) \\ &= \ln(\|\gamma\|_\infty) = \ln \left(\operatorname{esssup}_{h \in \mathbb{H}} \gamma(h) \right) = \ln \left(\operatorname{esssup}_{h \in \mathbb{H}} \frac{\rho_S(h)}{\pi(h)} \right). \end{aligned}$$

Finally, we can deduce that

$$\begin{aligned} &\lim_{\lambda \rightarrow +\infty} \left[\frac{2\lambda-1}{\lambda} \ln \frac{2}{\delta} + \frac{\lambda-1}{\lambda} D_\lambda(\rho_S|\pi) + \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} (\varphi(h', S')^{\frac{\lambda}{\lambda-1}}) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \\ &= \ln \operatorname{esssup}_{h' \in \mathbb{H}} \frac{\rho_S(h')}{\pi(h')} + \ln \left[\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', S') \right]. \end{aligned}$$

■

F.3 Proof of Theorem 6.3.2

For the sake of completeness, we first prove an upper bound on \sqrt{ab} (see, e.g., THIE-MANN *et al.*, 2017).

Lemma F.3.1. For any $a > 0, b > 0$, we have

$$\begin{aligned} \sqrt{\frac{a}{b}} &= \operatorname{argmin}_{\lambda > 0} \left(\frac{a}{\lambda} + \lambda b \right), \text{ and } 2\sqrt{ab} = \min_{\lambda > 0} \left(\frac{a}{\lambda} + \lambda b \right), \\ &\text{and } \forall \lambda > 0, \sqrt{ab} \leq \frac{1}{2} \left(\frac{a}{\lambda} + \lambda b \right). \end{aligned}$$

Proof. Let $f(\lambda) = \left(\frac{a}{\lambda} + \lambda b\right)$. The first derivative of f w.r.t. λ is

$$\frac{\partial f}{\partial \lambda}(\lambda) = \left(b - \frac{a}{\lambda^2}\right).$$

Moreover, from the derivative we can deduce that we have $\frac{\partial f}{\partial \lambda}(\lambda) < 0 \iff \lambda \in (0, \sqrt{\frac{a}{b}})$, and $\frac{\partial f}{\partial \lambda}(\lambda) > 0 \iff \lambda > \sqrt{\frac{a}{b}}$ and $\frac{\partial f}{\partial \lambda}(\lambda) = 0 \iff \lambda = \sqrt{\frac{a}{b}}$. It implies that the function is strictly decreasing on $\lambda \in (0, \sqrt{\frac{a}{b}})$, strictly increasing for $\lambda > \sqrt{\frac{a}{b}}$ and admit a unique minimum at $\lambda^* = \sqrt{\frac{a}{b}}$. Additionally, $f(\lambda^*) = 2\sqrt{ab}$ which proves the claim. ■

We can now prove Theorem 6.3.2 with Lemma F.3.1.

Theorem 6.3.2 (Parametrizable Disintegrated PAC-Bayes Bound). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution $\pi \in \mathcal{M}^*(\mathbb{H})$, for any measurable function $\varphi: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathbb{X} \times \mathbb{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\substack{\mathcal{S} \sim \mathcal{D}^m, \\ h \sim \rho_{\mathcal{S}}}} \left(\forall \lambda > 0, \ln(\varphi(h, \mathcal{S})) \leq \ln \left[\frac{\lambda}{2} e^{\mathcal{D}_2(\rho_{\mathcal{S}} \| \pi)} + \frac{8}{2\lambda\delta^3} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', \mathcal{S}')^2] \right] \right) \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Proof. The proof is similar to the one of Theorem 6.3.1. Since $\varphi(h, \mathcal{S})$ is a strictly positive random variable, from MARKOV's inequality (Theorem A.2.1), we have

$$\begin{aligned} & \mathbb{P}_{h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2} \\ \iff & \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \mathbb{I} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2}. \end{aligned}$$

Taking the expectation over $\mathcal{S} \sim \mathcal{D}^m$ to both sides of the inequality gives

$$\begin{aligned} & \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \mathbb{I} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2} \\ \iff & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathcal{S}}} \varphi(h', \mathcal{S}) \right] \geq 1 - \frac{\delta}{2}. \end{aligned}$$

Using Lemma F.3.1 with $a = \frac{4}{\delta^2} \varphi(h', \mathbb{S})^2$ and $b = \frac{\rho_{\mathbb{S}}(h')^2}{\pi(h')^2}$, we have for all prior $\pi \in \mathbb{M}^*(\mathbb{H})$

$$\begin{aligned} \forall \lambda > 0, \quad \frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathbb{S}}} \varphi(h', \mathbb{S}) &= \mathbb{E}_{h' \sim \pi} \sqrt{\frac{\rho_{\mathbb{S}}(h')^2}{\pi(h')^2} \frac{4}{\delta^2} \varphi(h', \mathbb{S})^2} \\ &\leq \frac{1}{2} \left[\lambda \mathbb{E}_{h' \sim \pi} \left(\frac{\rho_{\mathbb{S}}(h')}{\pi(h')} \right)^2 + \frac{4}{\lambda \delta^2} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \right]. \end{aligned}$$

Then, since both sides of the inequality are strictly positive, we take the logarithm to obtain

$$\begin{aligned} \forall \lambda > 0, \ln \left(\frac{2}{\delta} \mathbb{E}_{h' \sim \rho_{\mathbb{S}}} \varphi(h', \mathbb{S}) \right) &\leq \ln \left(\frac{1}{2} \left[\lambda \mathbb{E}_{h' \sim \pi} \left(\frac{\rho_{\mathbb{S}}(h')}{\pi(h')} \right)^2 + \frac{4}{\lambda \delta^2} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \right] \right) \\ &= \ln \left(\frac{1}{2} \left[\lambda \exp(D_2(\rho_{\mathbb{S}} \parallel \pi)) + \frac{4}{\lambda \delta^2} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \right] \right). \end{aligned}$$

Hence, we can deduce that

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\forall \pi \in \mathbb{M}^*(\mathbb{H}), \forall \lambda > 0, \ln(\varphi(h, \mathbb{S})) \right. \\ \left. \leq \ln \left(\frac{1}{2} \left[\lambda e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{4}{\lambda \delta^2} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \right] \right) \right] \geq 1 - \frac{\delta}{2}. \end{aligned} \quad (\text{F.3})$$

Given a prior $\pi \in \mathbb{M}^*(\mathbb{H})$, note that $\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^2$ is a strictly-positive random variable. Hence, we apply MARKOV's inequality (Theorem A.2.1):

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S})^2 \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}')^2 \right] \geq 1 - \frac{\delta}{2}.$$

Since the inequality does not depend on the random variable $h \sim \rho_{\mathbb{S}}$, we have

$$\begin{aligned} \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^2 \right) \right] \\ = \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S})^2 \right) \leq \frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathbb{S}')^2 \right) \right]. \end{aligned}$$

Additionally, note that multiplying by $\frac{4}{2\lambda\delta^2} > 0$, adding $\frac{\lambda}{2} \exp(D_2(\rho_{\mathbb{S}} \parallel \pi))$, and taking the logarithm to both sides of the inequality results in the same indicator

function. Indeed,

$$\begin{aligned}
 & \mathbb{I} \left[\mathbb{E}_{h' \sim \pi} (\varphi(h', \mathbb{S})^2) \leq \frac{2}{\delta} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} (\varphi(h', S')^2) \right] \\
 &= \mathbb{I} \left[\forall \lambda > 0, \frac{4}{2\lambda\delta^2} \mathbb{E}_{h' \sim \pi} (\varphi(h', \mathbb{S})^2) \leq \frac{8}{2\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} (\varphi(h', S')^2) \right] \\
 &= \mathbb{I} \left[\forall \lambda > 0, \ln \left(\frac{\lambda}{2} \exp(D_2(\rho_{\mathbb{S}} \parallel \pi)) + \frac{4}{2\lambda\delta^2} \mathbb{E}_{h' \sim \pi} (\varphi(h', \mathbb{S})^2) \right) \right. \\
 &\quad \left. \leq \ln \left(\frac{\lambda}{2} \exp(D_2(\rho_{\mathbb{S}} \parallel \pi)) + \frac{8}{2\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} (\varphi(h', S')^2) \right) \right].
 \end{aligned}$$

Hence, we can deduce that

$$\begin{aligned}
 & \mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\forall \lambda > 0, \ln \left(\frac{1}{2} \left[\lambda \exp(D_2(\rho_{\mathbb{S}} \parallel \pi)) + \frac{4}{\lambda\delta^2} \mathbb{E}_{h' \sim \pi} (\varphi(h', \mathbb{S})^2) \right] \right) \right. \\
 &\quad \left. \leq \ln \left(\frac{1}{2} \left[\lambda \exp(D_2(\rho_{\mathbb{S}} \parallel \pi)) + \frac{8}{\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} (\varphi(h', S')^2) \right] \right) \right] \geq 1 - \frac{\delta}{2}. \quad (\text{F.4})
 \end{aligned}$$

Combining Equations (F.3) and (F.4) with a union bound gives us the desired result. \blacksquare

F.4 Proof of Proposition 6.3.1

Proposition 6.3.1 (Optimal Bound of Theorem 6.3.2). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution π on \mathbb{H} , for any $\delta \in (0, 1]$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_+^*$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$, let

$$\lambda^* = \operatorname{argmin}_{\lambda > 0} \ln \left[\frac{\lambda}{2} e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [8\varphi(h', S')^2]}{2\lambda\delta^3} \right],$$

$$\begin{aligned}
 \text{then, we have} \quad & 2 \ln \left[\frac{\lambda^*}{2} e^{D_2(\rho_{\mathbb{S}} \parallel \pi)} + \frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\frac{8\varphi(h', S')^2}{2\lambda^*\delta^3} \right) \right] \\
 &= \underbrace{D_2(\rho_{\mathbb{S}} \parallel \pi) + \ln \left[\frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\frac{8\varphi(h', S')^2}{\delta^3} \right) \right]}_{\text{Theorem 6.3.1 with } \lambda = 2},
 \end{aligned}$$

$$\text{where } \lambda^* = \sqrt{\frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [8\varphi(h', S')^2]}{\delta^3 \exp(D_2(\rho_S \|\pi))}}.$$

Put into words: the optimal λ^* gives the same bound for Theorem 6.3.1 and Theorem 6.3.2.

Proof. We consider the right-hand side of the inequality of Theorem 6.3.2 (which is strictly positive): we have

$$\ln \left[\frac{\lambda}{2} e^{D_2(\rho_S \|\pi)} + \frac{8}{2\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', S')^2] \right]. \quad (\text{F.5})$$

Since \ln is a strictly increasing function, we have

$$\begin{aligned} & \min_{\lambda > 0} \left\{ \ln \left[\frac{\lambda}{2} e^{D_2(\rho_S \|\pi)} + \frac{8}{2\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', S')^2] \right] \right\} \\ &= \ln \left[\min_{\lambda > 0} \left\{ \frac{\lambda}{2} e^{D_2(\rho_S \|\pi)} + \frac{8}{2\lambda\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', S')^2] \right\} \right]. \end{aligned}$$

Then, we apply Lemma F.3.1 by taking $a = \frac{8}{2\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', S')^2]$ and $b = \frac{1}{2} e^{D_2(\rho_S \|\pi)}$ to obtain $\lambda^* = \sqrt{\frac{a}{b}} = \sqrt{\frac{\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [8\varphi(h', S')^2]}{\delta^3 \exp(D_2(\rho_S \|\pi))}}$. Finally, by substituting λ^* into Equation (F.5), we obtain

$$\begin{aligned} & \ln \left[\frac{\lambda^*}{2} e^{D_2(\rho_S \|\pi)} + \frac{8}{2\lambda^*\delta^3} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} [\varphi(h', S')^2] \right] \\ &= \frac{1}{2} \left(D_2(\rho_S \|\pi) + \ln \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\frac{8\varphi(h', S')^2}{\delta^3} \right) \right] \right), \end{aligned}$$

which is the desired result. ■

F.5 Proof of Corollary 6.4.1

We introduce Theorem F.5.1 which takes into account a set of priors \mathbb{P} while Theorem 6.3.1 handles a unique prior π .

Theorem F.5.1. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any priors set $\mathbb{P} = \{\pi_t\}_{t=1}^T$ of T prior $\pi \in \mathbb{M}^*(\mathbb{H})$, for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm A :

F.5. Proof of Corollary 6.4.1

$(\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\forall \pi_t \in \mathbb{P}, \frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \frac{2T}{\delta} + \ln \left(\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right) \right] \geq 1 - \delta,$$

where $\rho_{\mathcal{S}} \triangleq A(\mathcal{S}, \pi)$ is output by the deterministic algorithm A .

Proof. The proof is mainly the same as Theorem 6.3.1. Indeed, we first derive the same equation as Equation (F.1), we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\forall \pi \in \mathbb{M}^*(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq D_{\lambda}(\rho_{\mathcal{S}} \parallel \pi) + \frac{\lambda}{\lambda-1} \ln \frac{2}{\delta} + \ln \left(\mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right) \right) \right] \geq 1 - \frac{\delta}{2}.$$

Then, we apply MARKOV's inequality (as in Theorem 6.3.1) T times with the T priors π_t belonging to \mathbb{P} , however, we set the confidence to $\frac{\delta}{2T}$ instead of $\frac{\delta}{2}$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\ln \left(\mathbb{E}_{h' \sim \pi_t} \left[\varphi(h', \mathcal{S})^{\frac{\lambda}{\lambda-1}} \right] \right) \leq \ln \frac{2T}{\delta} + \ln \left(\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi_t} \left[\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right] \right) \right] \geq 1 - \frac{\delta}{2T}.$$

Finally, combining the $T + 1$ bounds with a union bound gives us the desired result. \blacksquare

We now prove Corollary 6.4.1 from Theorem F.5.1.

Corollary 6.4.1 (Instantiation of Theorem 6.3.1 for Neural Networks). For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$, for any algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, for any loss $\ell : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left(\forall \pi_t \in \mathbb{P}, \text{kl}(\mathbf{R}_{\mathcal{S}}^{\ell}(h) \parallel \mathbf{R}_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{\sigma^2} + \ln \frac{16T\sqrt{m}}{\delta^3} \right] \right) \geq 1 - \delta,$$

where $\text{kl}(a||b) = a \ln \frac{a}{b} + (1-a) \ln \frac{1-a}{1-b}$, $\rho_{\mathcal{S}} = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$, and the hypothesis $h \sim \rho_{\mathcal{S}}$ is parameterized by $\mathbf{w} + \boldsymbol{\epsilon}$.

Proof. We instantiate Theorem F.5.1 with $\varphi(h, \mathcal{S}) = \exp\left[\frac{\lambda-1}{\lambda} m \text{kl}(\mathbf{R}_{\mathcal{S}}^{\ell}(h) || \mathbf{R}_{\mathcal{D}}^{\ell}(h))\right]$ and $\lambda = 2$. We have with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$, for all prior $\pi_t \in \mathbb{P}$

$$\text{kl}(\mathbf{R}_{\mathcal{S}}^{\ell}(h) || \mathbf{R}_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[D_2(\rho_{\mathcal{S}} || \pi_t) + \ln \left(\frac{8T}{\delta^3} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi_t} e^{m \text{kl}(\mathbf{R}_{\mathcal{S}'}^{\ell}(h') || \mathbf{R}_{\mathcal{D}}^{\ell}(h'))} \right) \right].$$

From MAURER (2004) we upper-bound $\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi_t} e^{m \text{kl}(\mathbf{R}_{\mathcal{S}'}^{\ell}(h') || \mathbf{R}_{\mathcal{D}}^{\ell}(h'))}$ by $2\sqrt{m}$ for each prior π_t (Lemma B.16.1). Hence, we have, for all prior $\pi_t \in \mathbb{P}$

$$\text{kl}(\mathbf{R}_{\mathcal{S}}^{\ell}(h) || \mathbf{R}_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[D_2(\rho_{\mathcal{S}} || \pi_t) + \ln \left(\frac{16T\sqrt{m}}{\delta^3} \right) \right].$$

Additionally, the RÉNYI divergence $D_2(\rho_{\mathcal{S}} || \pi_t)$ between two multivariate Gaussians $\rho_{\mathcal{S}} = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ and $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$ is well known: its closed-form solution is $D_2(\rho_{\mathcal{S}} || \pi_t) = \frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{\sigma^2}$ (see, for example, (GIL *et al.*, 2013)). ■

F.6 Proof of Corollary 6.4.2

We first prove the following lemma in order to prove Corollary 6.4.2.

Lemma F.6.1. If $\rho_{\mathcal{S}} = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ and $\pi = \mathcal{N}(\mathbf{v}, \sigma^2 \mathbf{I}_D)$, we have

$$\ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} = \frac{1}{2\sigma^2} \left[\|\mathbf{w} + \boldsymbol{\epsilon} - \mathbf{v}\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2 \right],$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$ is a Gaussian noise such that $\mathbf{w} + \boldsymbol{\epsilon}$ are the weights of $h \sim \rho_{\mathcal{S}}$ with $\rho_{\mathcal{S}} = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$.

Proof. The probability density functions of $\rho_{\mathcal{S}}$ and π for $h \sim \rho_{\mathcal{S}}$ (with the weights

$\mathbf{w}+\boldsymbol{\epsilon}$) can be rewritten as

$$\rho_{\mathbb{S}}(h) = \left[\frac{1}{\sigma\sqrt{2\pi}} \right]^D \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{w}\|_2^2\right) = \left[\frac{1}{\sigma\sqrt{2\pi}} \right]^D \exp\left(-\frac{1}{2\sigma^2}\|\boldsymbol{\epsilon}\|_2^2\right)$$

and $\pi(h) = \left[\frac{1}{\sigma\sqrt{2\pi}} \right]^D \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}\|_2^2\right)$.

We can derive a closed-form expression of $\ln\left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)}\right]$. Indeed, we have

$$\begin{aligned} \ln\left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)}\right] &= \ln[\rho_{\mathbb{S}}(h)] - \ln[\pi(h)] \\ &= \ln\left(\left[\frac{1}{\sigma\sqrt{2\pi}}\right]^D \exp\left(-\frac{1}{2\sigma^2}\|\boldsymbol{\epsilon}\|_2^2\right)\right) \\ &\quad - \ln\left(\left[\frac{1}{\sigma\sqrt{2\pi}}\right]^D \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}\|_2^2\right)\right) \\ &= -\frac{1}{2\sigma^2}\|\boldsymbol{\epsilon}\|_2^2 + \frac{1}{2\sigma^2}\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}\|_2^2 = \frac{1}{2\sigma^2}\left[\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2\right]. \end{aligned}$$

■

We can now prove Corollary 6.4.2.

Corollary 6.4.2 (Instantiation of Known Bounds for Neural Networks). For any distribution \mathcal{D} on $\mathbb{X}\times\mathbb{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2\mathbf{I}_D)$, for any algorithm $A : (\mathbb{X}\times\mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, for any loss $\ell : \mathbb{H}\times(\mathbb{X}\times\mathbb{Y}) \rightarrow \{0, 1\}$, for any $\delta \in (0, 1]$, with probability at least $1-\delta$ over the learning sample $\mathbb{S} \sim \mathcal{D}^m$ and the hypothesis $h \sim \rho_{\mathbb{S}}$ parameterized by $\mathbf{w}+\boldsymbol{\epsilon}$, we have $\forall \pi_t \in \mathbb{P}$

$$\text{kl}(\mathbb{R}_{\mathbb{S}}^\ell(h) \parallel \mathbb{R}_{\mathcal{D}}^\ell(h)) \leq \frac{1}{m} \left[\frac{\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} + \ln \frac{2T\sqrt{m}}{\delta} \right], \quad (6.1)$$

$$\forall b \in \mathbb{b}, \text{kl}_+(\mathbb{R}_{\mathbb{S}}^\ell(h) \parallel \mathbb{R}_{\mathcal{D}}^\ell(h)) \leq \frac{1}{m} \left[\frac{b+1}{b} \left[\frac{\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} \right]_+ + \ln \frac{(b+1)T \text{card}(\mathbb{b})}{\delta} \right], \quad (6.2)$$

$$\forall c \in \mathbb{c}, \mathbb{R}_{\mathcal{D}}^\ell(h) \leq \frac{1 - \exp\left(-c\mathbb{R}_{\mathbb{S}}^\ell(h) - \frac{1}{m} \left[\frac{\|\mathbf{w}+\boldsymbol{\epsilon}-\mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} + \ln \frac{T \text{card}(\mathbb{c})}{\delta} \right]\right)}{1 - e^{-c}}, \quad (6.3)$$

with $[x]_+ = \max(x, 0)$, and $\text{kl}_+(\mathbf{R}_S^\ell(h) \|\mathbf{R}_D^\ell(h)) = \text{kl}(\mathbf{R}_S^\ell(h) \|\mathbf{R}_D^\ell(h))$ if $\mathbf{R}_S^\ell(h) < \mathbf{R}_D^\ell(h)$ and 0 otherwise. Moreover, $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$ is a Gaussian noise such that $\mathbf{w} + \epsilon$ are the weights of $h \sim \rho_S$ with $\rho_S = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$, and $\mathfrak{c}, \mathfrak{b}$ are two sets of hyperparameters fixed a priori.

Proof. We will prove the three bounds separately.

Equation (6.1). We instantiate Theorem 1(i) of RIVASPLATA *et al.* (2020) (proved in Theorem 2.4.1) with $\varphi(h, \mathbb{S}) = \exp\left[m \text{kl}(\mathbf{R}_S^\ell(h) \|\mathbf{R}_D^\ell(h))\right]$, however, we apply the theorem T times for each prior $\pi_t \in \mathbb{P}$ (with a confidence $\frac{\delta}{T}$ instead of δ). Hence, for each prior $\pi_t \in \mathbb{P}$, we have with probability at least $1 - \frac{\delta}{T}$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$ and $h \sim \rho_S$

$$\text{kl}(\mathbf{R}_S^\ell(h) \|\mathbf{R}_D^\ell(h)) \leq \frac{1}{m} \left[\ln \left[\frac{\rho_S(h)}{\pi_t(h)} \right] + \ln \left(\frac{T}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbf{R}_{\mathbb{S}'}^\ell(h') \|\mathbf{R}_D^\ell(h'))} \right) \right].$$

From MAURER (2004), we upper-bound $\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl}(\mathbf{R}_{\mathbb{S}'}^\ell(h') \|\mathbf{R}_D^\ell(h'))}$ by $2\sqrt{m}$ (Lemma B.16.1) and using Lemma F.6.1 we rewrite the disintegrated KL divergence. Finally, a union bound argument gives us the claim.

Equation (6.2). We apply $T \text{card}(\mathfrak{b})$ times Proposition 3.1 of BLANCHARD and FLEURET (2007) (proved in Theorem 2.4.3) with a confidence $\frac{\delta}{T \text{card}(\mathfrak{b})}$ instead of δ . For each prior $\pi_t \in \mathbb{P}$ and hyperparameters $b \in \mathfrak{b}$, we have with probability at least $1 - \frac{\delta}{T \text{card}(\mathfrak{b})}$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$ and $h \sim \rho_S$

$$\text{kl}(\mathbf{R}_S^\ell(h) \|\mathbf{R}_D^\ell(h)) \leq \frac{1}{m} \left[\frac{b+1}{b} \left[\ln \left[\frac{\rho_S(h)}{\pi_t(h)} \right] \right]_+ + \ln \left(\frac{T \text{card}(\mathfrak{b})(b+1)}{\delta} \right) \right].$$

From Lemma F.6.1 and a union bound argument, we obtain the claim.

Equation (6.3). We apply $T \text{card}(\mathfrak{c})$ times Theorem 1.2.7 of CATONI (2007) (proved in Theorem 2.4.2) with a confidence $\frac{\delta}{T \text{card}(\mathfrak{c})}$ instead of δ . For each prior $\pi_t \in \mathbb{P}$ and hyperparameter $c \in \mathfrak{c}$, we have with probability at least $1 - \frac{\delta}{T \text{card}(\mathfrak{c})}$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$ and $h \sim \rho_S$

$$\mathbf{R}_D^\ell(h) \leq \frac{1}{1-e^{-c}} \left[1 - \exp \left(-c \mathbf{R}_S^\ell(h) - \frac{1}{m} \left[\ln \left[\frac{\rho_S(h)}{\pi_t(h)} \right] + \ln \frac{T \text{card}(\mathfrak{c})}{\delta} \right] \right) \right].$$

From Lemma F.6.1 and a union bound argument, we obtain the claim. ■

F.7 Proof of Corollary 6.4.3

Corollary 6.4.3 (PAC-Bayesian Bound for Stochastic Neural Networks). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any set $\mathbb{P} = \{\pi_1, \dots, \pi_T\}$ of T priors on \mathbb{H} where $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$, for any loss $\ell: \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow \{0, 1\}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over $\mathbb{S} \sim \mathcal{D}^m$ and $\{h_1, \dots, h_K\} \sim \rho^K$, we have simultaneously $\forall \pi_t \in \mathbb{P}$,

$$\text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{2\sigma^2} + \ln \frac{4T\sqrt{m}}{\delta} \right], \quad (6.4)$$

$$\text{and} \quad \text{kl} \left(\frac{1}{K} \sum_{i=1}^K R_{\mathbb{S}}(h_i) \parallel \mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \right) \leq \frac{1}{n} \ln \frac{4}{\delta}, \quad (6.5)$$

where $\rho = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ and the hypothesis h sampled from ρ is parameterized by $\mathbf{w} + \boldsymbol{\epsilon}$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)$.

Proof. We instantiate Theorem 2.3.4 and apply JENSEN's inequality (Theorem A.1.1) on the left-hand side of the inequation for each prior $\pi_t = \mathcal{N}(\mathbf{v}_t, \sigma^2 \mathbf{I}_D)$ with the posterior $\rho = \mathcal{N}(\mathbf{w}, \sigma^2 \mathbf{I}_D)$ with a confidence $\frac{\delta}{2T}$ instead of δ . Indeed, for each prior π_t , with probability at least $1 - \frac{\delta}{2T}$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$, we have for all posterior ρ on \mathbb{H} ,

$$\text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) \leq \frac{1}{m} \left[\text{KL}(\rho \parallel \pi_t) + \ln \frac{4T\sqrt{m}}{\delta} \right].$$

Note that the closed-form solution of the KL divergence between the Gaussian distributions ρ and π_t is well known, we have $\text{KL}(\rho \parallel \pi_t) = \frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{2\sigma^2}$. Then, by applying a union bound argument over the T bounds obtained with the T priors π_t , we have with probability at least $1 - \frac{\delta}{2}$ over the random choice of $\mathbb{S} \sim \mathcal{D}^m$, for all prior $\pi_t \in \mathbb{P}$, for all posterior ρ

$$\text{kl} \left(\mathbb{E}_{h \sim \rho} R_{\mathbb{S}}^{\ell}(h) \parallel \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \right) \leq \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{2\sigma^2} + \ln \frac{4T\sqrt{m}}{\delta} \right]. \quad (\text{Equation (6.4)})$$

Additionally, we obtained Equation (6.5) by a direct application the Theorem 2.2 of DZIUGAITE and ROY (2017) (with confidence $\frac{\delta}{2}$ instead of δ). Finally, from a union bound of the two bounds in Equations (6.4) and (6.5) gives the result. ■

F.8 Evaluation and Minimization of the Bounds of Corollaries 6.4.1 to 6.4.3

We optimize and evaluate the bounds of the corollaries (except Equation (6.3)) thanks to the inverting functions of $\text{kl}()$ defined in Definition 2.3.3. Indeed, for the different corollaries, the PAC-Bayesian generalization bounds become

$$R_{\mathcal{D}}^{\ell}(h) \leq \underbrace{\overline{\text{kl}}\left(R_{\mathcal{S}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{\sigma^2} + \ln \frac{16T\sqrt{m}}{\delta^3} \right] \right.\right)}_{\text{Corollary 6.4.1}},$$

$$R_{\mathcal{D}}^{\ell}(h) \leq \underbrace{\overline{\text{kl}}\left(R_{\mathcal{S}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{\|\mathbf{w} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} + \ln \frac{2T\sqrt{m}}{\delta} \right] \right.\right)}_{\text{Equation (6.1)}},$$

$$R_{\mathcal{D}}^{\ell}(h) \leq \underbrace{\overline{\text{kl}}\left(R_{\mathcal{S}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{b+1}{b} \left[\frac{\|\mathbf{w} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} \right] + \ln \frac{(b+1)T \text{card}(\mathcal{b})}{\delta} \right] \right.\right)}_{\text{Equation (6.2)}},$$

$$\text{and } \mathbb{E}_{h \sim \rho} R_{\mathcal{D}}^{\ell}(h) \leq \underbrace{\overline{\text{kl}}\left(\spadesuit \left| \frac{1}{m} \left[\frac{\|\mathbf{w} - \mathbf{v}_t\|_2^2}{2\sigma^2} + \ln \frac{4T\sqrt{m}}{\delta} \right] \right.\right)}_{\text{Corollary 6.4.3}},$$

$$\text{where } \spadesuit = \overline{\text{kl}}\left(\frac{1}{K} \sum_{i=1}^K R_{\mathcal{S}}(h_i) \left| \frac{1}{K} \ln \frac{4}{\delta} \right.\right).$$

Based on these bounds, we can deduce some objective functions that is approximated on a mini-batch $\mathcal{U} \subseteq \mathcal{S}$. Indeed, at each iteration in phase **2**), after sampling the noise $\boldsymbol{\epsilon}$, the algorithm updates the weights $\boldsymbol{\omega}$ (*i.e.*, the hypothesis h) by optimizing

$$\underbrace{\overline{\text{kl}}\left(R_{\mathcal{U}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{\|\boldsymbol{\omega} - \mathbf{v}_t\|_2^2}{\sigma^2} + \ln \frac{16T\sqrt{m}}{\delta^3} \right] \right.\right)}_{\text{Objective function for Corollary 6.4.1}},$$

$$\underbrace{\overline{\text{kl}}\left(R_{\mathcal{U}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{\|\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} + \ln \frac{2T\sqrt{m}}{\delta} \right] \right.\right)}_{\text{Objective function for Equation (6.1)}},$$

$$\underbrace{\overline{\text{kl}}\left(R_{\mathcal{U}}^{\ell}(h) \left| \frac{1}{m} \left[\frac{b+1}{b} \left[\frac{\|\boldsymbol{\omega} + \boldsymbol{\epsilon} - \mathbf{v}_t\|_2^2 - \|\boldsymbol{\epsilon}\|_2^2}{2\sigma^2} \right] + \ln \frac{(b+1)T \text{card}(\mathcal{b})}{\delta} \right] \right.\right)}_{\text{Objective function for Equation (6.2)}},$$

where the loss $\ell()$ is the bounded cross-entropy loss of DZIUGAITE and ROY (2018), *i.e.*, $\ell(h, (\mathbf{x}, y)) = -\frac{1}{2} \ln[e^{-Z} + (1 - 2e^{-Z})h[y]]$.

Concerning the optimization of the hyperparameters $c \in \mathfrak{c}$ and $b \in \mathfrak{b}$ for Equations (6.2) and (6.3), we (i) initialize $b \in \mathfrak{b}$ or $c \in \mathfrak{c}$ with the one that performs best on the first mini-batch and (ii) optimize by gradient descent the hyperparameter. To evaluate Equations (6.2) and (6.3), we take $b \in \mathfrak{b}$ and $c \in \mathfrak{c}$ that leads to the tightest bound.

F.9 Disintegrated Information-theoretic Bounds

We discuss in this section another interpretation of the disintegration procedure through Theorems F.9.1 and F.9.2 below. Actually, the RÉNYI divergence between π and ρ is sensitive to the choice of the learning sample \mathbb{S} : when the posterior ρ learned from \mathbb{S} differs greatly from the prior π the divergence is high. To avoid such a behavior, we consider mutual information which is a measure of dependence between the random variables $\mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m$ and $h \in \mathbb{H}$. More formally, the mutual information is defined as

$$\text{MI}(h; \mathbb{S}) = \min_{\pi \in \mathfrak{M}^*(\mathbb{H})} \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathbb{S}} \| \pi).$$

From this quantity, we can derive the generalization bound introduced in the following theorem.

Theorem F.9.1. For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [1, +\infty[$, for any $\delta \in (0, 1]$, for any deterministic algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathfrak{M}^*(\mathbb{H}) \rightarrow \mathfrak{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\ln \varphi(h, \mathbb{S}) \leq \frac{1}{\delta} \left[\text{MI}(h; \mathbb{S}) + \ln \left(\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \varphi(h, \mathbb{S}) \right) \right] \right] \geq 1 - \delta,$$

where π^* is defined such that $\pi^*(h) = \mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \rho_{\mathbb{S}}(h)$.

Proof. Deferred to Appendix F.11. ■

As for the disintegrated bounds introduced in Section 6.3, the bound on $\ln \varphi(h, \mathbb{S})$ depends on mainly two terms: a term (i.e., $\text{MI}(h; \mathbb{S})$) that measures the dependence of $h \in \mathbb{H}$ on the learning sample \mathbb{S} and $\ln \left(\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \varphi(h, \mathbb{S}) \right)$ that must be upper-bounded to obtain a computable bound. However, the bound has a polynomial dependence of δ , i.e., we have $\frac{1}{\delta}$ instead of $\ln \frac{1}{\delta}$. To improve such dependence, we consider Sibson's mutual information (VERDÚ, 2015). It involves an expectation over the learning samples of a given size m and is defined for a given $\lambda > 1$ by

$$\text{MI}_{\lambda}(h; \mathbb{S}) \triangleq \min_{\pi \in \mathfrak{M}^*(\mathbb{H})} \frac{1}{\lambda - 1} \ln \left[\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right]^{\lambda} \right].$$

The higher $\text{MI}_\lambda(h; \mathcal{S})$, the higher the correlation is, meaning that the sampling of h is highly dependent on the choice of \mathcal{S} . This measure has two interesting properties: it generalizes the mutual information (VERDÚ, 2015), and it can be related to the RÉNYI divergence. Indeed, let $\rho(h, \mathcal{S}) = \rho_{\mathcal{S}}(h) \mathcal{D}^m(\mathcal{S})$, resp. $\pi(h, \mathcal{S}) = \pi(h) \mathcal{D}^m(\mathcal{S})$, be the probability of sampling both $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$, resp. $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \pi$. Then we can write:

$$\begin{aligned} \text{MI}_\lambda(h; \mathcal{S}) &= \min_{\pi \in \mathcal{M}^*(\mathbb{H})} \frac{1}{\lambda-1} \ln \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \left[\frac{\rho_{\mathcal{S}}(h) \mathcal{D}^m(\mathcal{S})}{\pi(h) \mathcal{D}^m(\mathcal{S})} \right]^\lambda \right] \\ &= \min_{\pi \in \mathcal{M}^*(\mathbb{H})} D_\lambda(\rho \| \pi). \end{aligned} \quad (\text{F.6})$$

From VERDÚ (2015) the optimal prior π^* minimizing Equation (F.6) is a *distribution-dependent* prior:

$$\pi^*(h) = \frac{\left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)^\lambda \right]^{\frac{1}{\lambda}}}{\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h')^\lambda \right]^{\frac{1}{\lambda}}}.$$

This leads to an *Information-Theoretic generalization bound*.

Theorem F.9.2 (Disintegrated Information-Theoretic Bound). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\substack{\mathcal{S} \sim \mathcal{D}^m, \\ h \sim \rho_{\mathcal{S}}}} \left(\frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq \text{MI}_\lambda(h'; \mathcal{S}') + \ln \left[\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left[\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right] \right] \right) \geq 1 - \delta.$$

Proof. Deferred to Appendix F.11. ■

We can remark that Theorem F.9.2 is tighter than Theorem F.9.1. For example, when we instantiate Theorem F.9.1 with $\varphi(h, \mathcal{S}) = \exp \left[m \text{kl}(\mathbb{R}_{\mathcal{S}}^\ell(h) \| \mathbb{R}_{\mathcal{D}}^\ell(h)) \right]$, the bound will be multiplied by $\frac{1}{\delta^{\frac{1}{m}}}$, while the bound of Theorem F.9.2 is only multiplied by $\frac{1}{m}$ (but we add the term $\frac{1}{m} \ln \frac{1}{\delta}$ to the bound which is small even for small m).

For the sake of comparison, we introduce the following corollary of Theorem F.9.2.

Corollary F.9.1. Under the assumptions of Theorem F.9.2, when $\lambda \rightarrow 1^+$, with probability at least $1 - \delta$ we have

$$\ln \varphi(h, \mathcal{S}) \leq \ln \frac{1}{\delta} + \ln \left[\operatorname{esssup}_{\mathcal{S}' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathbb{H}} \varphi(h', \mathcal{S}') \right].$$

When $\lambda \rightarrow +\infty$, with probability at least $1 - \delta$ we have

$$\ln \varphi(h, \mathcal{S}) \leq \ln \left(\operatorname{esssup}_{\mathcal{S} \in \mathcal{S}, h \in \mathbb{H}} \frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)} \right) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathcal{S}') \right].$$

As for Theorem 6.3.1, this corollary illustrates a trade-off introduced by λ between the Sibson's mutual information $\operatorname{MI}_{\lambda}(h'; \mathcal{S}')$ and the term $\ln \left(\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right)$.

Furthermore, ESPOSITO *et al.* (2020, Cor.4) introduced a bound involving Sibson's mutual information. Their bound holds with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$:

$$2(\mathcal{R}_{\mathcal{S}}^{\ell}(h) - \mathcal{R}_{\mathcal{D}}^{\ell}(h))^2 \leq \frac{1}{m} \left[\operatorname{MI}_{\lambda}(h'; \mathcal{S}') + \ln \frac{2}{\delta^{\frac{\lambda}{\lambda-1}}} \right]. \quad (\text{F.7})$$

Hence, we compare Equation (F.7) with the equations of the following corollary.

Corollary F.9.2. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm $A: (\mathcal{X} \times \mathcal{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, with probability at least $1 - \delta$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$, we have

$$\operatorname{kl}(\mathcal{R}_{\mathcal{S}}^{\ell}(h) \| \mathcal{R}_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[\operatorname{MI}_{\lambda}(h'; \mathcal{S}') + \ln \frac{2\sqrt{m}}{\delta^{\frac{\lambda}{\lambda-1}}} \right] \quad (\text{F.8})$$

$$\text{and } 2(\mathcal{R}_{\mathcal{S}}^{\ell}(h) - \mathcal{R}_{\mathcal{D}}^{\ell}(h))^2 \leq \frac{1}{m} \left[\operatorname{MI}_{\lambda}(h'; \mathcal{S}') + \ln \frac{2\sqrt{m}}{\delta^{\frac{\lambda}{\lambda-1}}} \right]. \quad (\text{F.9})$$

Proof. First of all, we instantiate Theorem F.9.2 with the function $\varphi(h, \mathcal{S}) = \exp \left[\frac{\lambda-1}{\lambda} m \operatorname{kl}(\mathcal{R}_{\mathcal{S}}^{\ell}(h) \| \mathcal{R}_{\mathcal{D}}^{\ell}(h)) \right]$, we have (by rearranging the terms)

$$\operatorname{kl}(\mathcal{R}_{\mathcal{S}}^{\ell}(h) \| \mathcal{R}_{\mathcal{D}}^{\ell}(h)) \leq \frac{1}{m} \left[\operatorname{MI}_{\lambda}(h'; \mathcal{S}') + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \operatorname{kl}(\mathcal{R}_{\mathcal{S}'}^{\ell}(h') \| \mathcal{R}_{\mathcal{D}}^{\ell}(h'))} \right) \right].$$

Then, from MAURER (2004), we upper-bound $\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \operatorname{kl}(\mathcal{R}_{\mathcal{S}'}^{\ell}(h') \| \mathcal{R}_{\mathcal{D}}^{\ell}(h'))}$ by $2\sqrt{m}$ (Lemma B.16.1) to obtain Equation (F.8). Finally, to obtain Equ-

tion (F.9), we apply PINSKER's inequality (Theorem B.5.1), *i.e.*, we have the inequality $2(\mathbb{R}_S^\ell(h) - \mathbb{R}_D^\ell(h))^2 \leq \text{kl}(\mathbb{R}_S^\ell(h) \| \mathbb{R}_D^\ell(h))$ on Equation (F.8). ■

Equation (F.9) is slightly looser than Equation (F.7) since it involves an extra term of $\frac{1}{m} \ln \sqrt{m}$. However, Equation (F.8) is tighter than Equation (F.7) when we have $\text{kl}(\mathbb{R}_S^\ell(h) \| \mathbb{R}_D^\ell(h)) - 2(\mathbb{R}_S^\ell(h) - \mathbb{R}_D^\ell(h))^2 \geq \frac{1}{m} \ln \sqrt{m}$ (which becomes more frequent as m grows). Moreover, from a theoretical view, Theorem F.9.2 brings a different philosophy than the disintegrated PAC-Bayes bounds. Indeed, in Theorems 6.3.1 and 6.3.2, given \mathbb{S} , the RÉNYI divergence $D_\lambda(\rho_S \| \pi)$ suggests that the learned posterior ρ_S should be close enough to the prior π to get a low bound. While in Theorem F.9.2, the Sibson's mutual information $\text{MI}_\lambda(h'; S')$ suggests that the random variable h has to be *not too much correlated* to \mathbb{S} . However, the bound of Theorem F.9.2 is not computable in practice due notably to the sample expectation over the unknown distribution \mathcal{D} in $\text{MI}_\lambda()$. An exciting line of future works could be to study how we can make use of Theorem F.9.2 in practice.

F.10 Proof of Theorem F.9.1

In order to prove Theorem F.9.1, we need to prove Lemma F.10.1.

Lemma F.10.1. For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [1, +\infty[$, for any $\delta \in (0, 1]$, for any deterministic algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_S} \left[\forall \pi \in \mathbb{M}^*(\mathbb{H}), \ln \varphi(h, \mathbb{S}) \leq \frac{1}{\delta} \left[\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \text{KL}(\rho_S \| \pi) + \ln \left(\mathbb{E}_{\mathbb{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \varphi(h, \mathbb{S}) \right) \right] \right] \geq 1 - \delta.$$

Proof. By developing $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \varphi(h, \mathcal{S})$, we have for all prior $\pi \in \mathbb{M}^*(\mathbb{H})$

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \varphi(h, \mathcal{S}) &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\rho_{\mathcal{S}}(h) \pi(h)}{\pi(h) \rho_{\mathcal{S}}(h)} \varphi(h, \mathcal{S}) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\pi(h)}{\rho_{\mathcal{S}}(h)} \varphi(h, \mathcal{S}) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \| \pi) + \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\pi(h)}{\rho_{\mathcal{S}}(h)} \varphi(h, \mathcal{S}) \right]. \end{aligned}$$

From JENSEN's inequality (Theorem A.1.1), we have for all prior $\pi \in \mathbb{M}^*(\mathbb{H})$

$$\begin{aligned} &\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \| \pi) + \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\pi(h)}{\rho_{\mathcal{S}}(h)} \varphi(h, \mathcal{S}) \right] \\ &\leq \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \| \pi) + \ln \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \frac{\pi(h)}{\rho_{\mathcal{S}}(h)} \varphi(h, \mathcal{S}) \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \| \pi) + \ln \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \varphi(h, \mathcal{S}) \right]. \end{aligned} \quad (\text{F.10})$$

Since we assume in this case that $\varphi(h, \mathcal{S}) \geq 1$ for all $h \in \mathbb{H}$ and $\mathcal{S} \in (\mathbb{X} \times \mathbb{Y})^m$, we have $\ln \varphi(h, \mathcal{S}) \geq 0$; we can apply MARKOV's inequality (Theorem A.2.1) to obtain

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\ln \varphi(h, \mathcal{S}) \leq \frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \ln \varphi(h', \mathcal{S}') \right] \geq 1 - \delta. \quad (\text{F.11})$$

Then, from Equations (F.10) and (F.11), we can deduce the stated result. \blacksquare

We are now ready to prove Theorem F.9.1.

Theorem F.9.1. For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow [1, +\infty[$, for any $\delta \in (0, 1]$, for any deterministic algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathbb{M}^*(\mathbb{H}) \rightarrow \mathbb{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\ln \varphi(h, \mathcal{S}) \leq \frac{1}{\delta} \left[\text{MI}(h; \mathcal{S}) + \ln \left(\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \varphi(h, \mathcal{S}) \right) \right] \right] \geq 1 - \delta,$$

where π^* is defined such that $\pi^*(h) = \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)$.

Proof. Note that the mutual information is $\text{MI}(h; \mathcal{S}) = \min_{\pi \in \mathbb{M}^*(\mathbb{H})} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \| \pi)$. Hence, to prove Theorem F.9.1, we have to instantiate Lemma F.10.1 with the

optimal prior, *i.e.*, the prior π which minimizes $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \parallel \pi)$. The optimal prior is well-known (see, *e.g.*, CATONI, 2007; LEVER *et al.*, 2013): for the sake of completeness, we derive it. First, we have

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \parallel \pi) &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\rho_{\mathcal{S}}(h) [\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)]}{\pi(h) [\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)]} \right] \\ &= \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\rho_{\mathcal{S}}(h)}{\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)} \right] + \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)}{\pi(h)} \right]. \end{aligned}$$

Hence,

$$\begin{aligned} \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \parallel \pi) &= \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\rho_{\mathcal{S}}(h)}{\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)} \right] \right. \\ &\quad \left. + \mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)}{\pi(h)} \right] \right] \\ &= \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} \left[\mathbb{E}_{h \sim \rho_{\mathcal{S}}} \ln \left[\frac{\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)}{\pi(h)} \right] \right] = \pi^*, \end{aligned}$$

where $\pi^*(h) = \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \rho_{\mathcal{S}'}(h)$. Note that π^* is defined from the data distribution \mathcal{D} , hence, π^* is a valid prior when instantiating Lemma F.10.1 with π^* . Then, we have with probability at least $1-\delta$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$

$$\begin{aligned} \ln \varphi(h, \mathcal{S}) &\leq \frac{1}{\delta} \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \text{KL}(\rho_{\mathcal{S}} \parallel \pi^*) + \ln \left(\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \varphi(h, \mathcal{S}) \right) \right] \\ &= \frac{1}{\delta} \left[\text{MI}(h; \mathcal{S}) + \ln \left(\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \varphi(h, \mathcal{S}) \right) \right]. \end{aligned}$$

■

F.11 Proof of Theorem F.9.2

We first introduce Lemma F.11.1 in order to prove Theorem F.9.2.

Lemma F.11.1. For any distribution \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any hypothesis set \mathbb{H} , for any prior distribution π on \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any deterministic algorithm $A : (\mathcal{X} \times \mathcal{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$,

we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\forall \pi \in \mathcal{M}^*(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq D_{\lambda}(\rho \| \pi) + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right) \right] \geq 1 - \delta.$$

where $\rho(h, \mathcal{S}) = \rho_{\mathcal{S}}(h) \mathcal{D}^m(\mathcal{S})$; $\pi(h, \mathcal{S}) = \pi(h) \mathcal{D}^m(\mathcal{S})$.

Proof. Note that $\varphi(h, \mathcal{S})$ is a non-negative random variable. From MARKOV's inequality (Theorem A.2.1), we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\varphi(h, \mathcal{S}) \leq \frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \varphi(h', \mathcal{S}') \right] \geq 1 - \delta.$$

Then, since both sides of the inequality are strictly positive, we take the logarithm to both sides of the equality and multiply by $\frac{\lambda}{\lambda-1} > 0$ to obtain

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) \leq \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \varphi(h', \mathcal{S}') \right) \right] \geq 1 - \delta.$$

We develop the right-hand side of the inequality in the indicator function and make the expectation of the hypothesis over the distribution π appear. We have for all priors $\pi \in \mathcal{M}^*(\mathbb{H})$,

$$\begin{aligned} & \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \varphi(h', \mathcal{S}') \right) \\ &= \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \frac{\rho_{\mathcal{S}'}(h')}{\pi(h')} \frac{\pi(h')}{\rho_{\mathcal{S}'}(h')} \varphi(h', \mathcal{S}') \right) \\ &= \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \frac{\rho_{\mathcal{S}'}(h')}{\pi(h')} \varphi(h', \mathcal{S}') \right). \end{aligned}$$

Then, since $\frac{1}{r} + \frac{1}{s} = 1$ where $r = \lambda$ and $s = \frac{\lambda}{\lambda-1}$. Hence, HÖLDER's inequality (Theorem A.5.1) gives

$$\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{\mathcal{S}'}} \varphi(h', \mathcal{S}') \leq \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{\mathcal{S}'}(h')}{\pi(h')} \right]^{\lambda} \right) \right]^{\frac{1}{\lambda}} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}}.$$

Since both sides of the inequality are positive, we take the logarithm. Moreover, we add $\ln(\frac{1}{\delta})$, and we multiply by $\frac{\lambda}{\lambda-1} > 0$ to both sides of the inequality. We have

$$\begin{aligned} & \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \rho_{S'}} \varphi(h', S') \right) \\ & \leq \frac{\lambda}{\lambda-1} \ln \left(\frac{1}{\delta} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{S'}(h')}{\pi(h')} \right]^\lambda \right) \right]^{\frac{1}{\lambda}} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda-1}{\lambda}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\ & = \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{S'}(h')}{\pi(h')} \right]^\lambda \right) \right) + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda-1}{\lambda}} \right) \right). \end{aligned}$$

Hence, we can deduce that

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m, h \sim \rho_S} \left[\forall \pi \in \mathcal{M}^*(\mathbb{H}), \frac{\lambda}{\lambda-1} \ln(\varphi(h, S)) \leq \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{S'}(h')}{\pi(h')} \right]^\lambda \right) \right) \right. \\ \left. + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', S')^{\frac{\lambda-1}{\lambda}} \right) \right) \right] \geq 1 - \delta, \end{aligned}$$

where by definition we have $D_\lambda(\rho \parallel \pi) = \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\left[\frac{\rho_{S'}(h')}{\pi(h')} \right]^\lambda \right) \right)$. ■

From Lemma F.11.1, we prove Theorem F.9.2.

Theorem F.9.2 (Disintegrated Information-Theoretic Bound). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any measurable function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm $A : (\mathbb{X} \times \mathbb{Y})^m \times \mathcal{M}^*(\mathbb{H}) \rightarrow \mathcal{M}(\mathbb{H})$, we have

$$\mathbb{P}_{\substack{S \sim \mathcal{D}^m, \\ h \sim \rho_S}} \left(\frac{\lambda}{\lambda-1} \ln(\varphi(h, S)) \leq \text{MI}_\lambda(h'; S') + \ln \left[\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left[\varphi(h', S')^{\frac{\lambda-1}{\lambda}} \right] \right] \right) \geq 1 - \delta.$$

Proof. Sibson's mutual information is $\text{MI}_\lambda(h; S) = \min_{\pi \in \mathcal{M}^*(\mathbb{H})} D_\lambda(\rho \parallel \pi)$. Hence, in order to prove Theorem F.9.2, we have to instantiate Lemma F.11.1 with the optimal prior, i.e., the prior π which minimizes $D_\lambda(\rho \parallel \pi)$. Actually, this optimal prior has a closed-form solution (VERDÚ, 2015). For the sake of completeness,

we derive it. First, we have

$$\begin{aligned}
 & D_\lambda(\rho\|\pi) \\
 &= \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi} \left(\left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right]^\lambda \right) \right) \\
 &= \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} (\rho_{\mathcal{S}}(h)^\lambda) \right] (\pi(h)^{-\lambda}) \right) \\
 &= \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} (\rho_{\mathcal{S}}(h)^\lambda) \right] (\pi(h)^{-\lambda}) \left[\frac{\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}}}{\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}}} \right]^\lambda \right) \\
 &= \frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}} \right) \\
 &\quad + \frac{1}{\lambda-1} \ln \left(\mathbb{E}_{h \sim \pi} \frac{1}{\pi(h)^\lambda} \left[\frac{\left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} (\rho_{\mathcal{S}}(h)^\lambda) \right]^{\frac{1}{\lambda}}}{\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}}} \right]^\lambda \right) \\
 &= \frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}} \right) + D_\lambda(\pi^* \|\pi),
 \end{aligned}$$

where $\pi^*(h) = \left[\frac{\left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} (\rho_{\mathcal{S}}(h)^\lambda) \right]^{\frac{1}{\lambda}}}{\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}}} \right]$.

From these equalities and using the fact that $D_\lambda(\pi^* \|\pi)$ is minimal (*i.e.*, equal to zero) when $\pi^* = \pi$, we can deduce that

$$\begin{aligned}
 & \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} D_\lambda(\rho\|\pi) \\
 &= \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} \left[\frac{\lambda}{\lambda-1} \ln \left(\mathbb{E}_{h' \sim \pi} \frac{1}{\pi(h')} \left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} (\rho_{\mathcal{S}'}(h')^\lambda) \right]^{\frac{1}{\lambda}} \right) + D_\lambda(\pi^* \|\pi) \right] \\
 &= \operatorname{argmin}_{\pi \in \mathcal{M}^*(\mathbb{H})} D_\lambda(\pi^* \|\pi) = \pi^*.
 \end{aligned}$$

Note that π^* is defined from the data distribution \mathcal{D} , hence, π^* is a valid prior when instantiating Lemma F.11.1 with π^* . Then, we have with probability at least $1-\delta$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$

$$\begin{aligned}
 \frac{\lambda}{\lambda-1} \ln(\varphi(h, \mathcal{S})) &\leq D_\lambda(\rho\|\pi^*) + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right) \\
 &= \operatorname{MI}_\lambda(h'; \mathcal{S}') + \ln \left(\frac{1}{\delta^{\frac{\lambda}{\lambda-1}}} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \left(\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}} \right) \right).
 \end{aligned}$$

where $\pi^*(h, \mathbb{S}) = \pi^*(h)\mathcal{D}^m(\mathbb{S})$. ■

F.12 Proof of Corollary F.9.1

Corollary F.9.1. Under the assumptions of Theorem F.9.2, when $\lambda \rightarrow 1^+$, with probability at least $1 - \delta$ we have

$$\ln \varphi(h, \mathbb{S}) \leq \ln \frac{1}{\delta} + \ln \left[\operatorname{esssup}_{\mathbb{S}' \in (\mathbb{X} \times \mathbb{Y}), h' \in \mathbb{H}} \varphi(h', \mathbb{S}') \right].$$

When $\lambda \rightarrow +\infty$, with probability at least $1 - \delta$ we have

$$\ln \varphi(h, \mathbb{S}) \leq \ln \left(\operatorname{esssup}_{\mathbb{S} \in \mathbb{S}, h \in \mathbb{H}} \frac{\rho_{\mathbb{S}}(h)}{\pi^*(h)} \right) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \varphi(h', \mathbb{S}') \right].$$

Proof. The proof is similar to Corollary 6.3.1. Starting from Theorem F.9.2 and rearranging, we have

$$\begin{aligned} \mathbb{P}_{\substack{\mathbb{S} \sim \mathcal{D}^m \\ h \sim \rho_{\mathbb{S}}}} \left[\ln(\varphi(h, \mathbb{S})) \leq \frac{\lambda-1}{\lambda} \operatorname{MI}_{\lambda}(h'; \mathbb{S}') \right. \\ \left. + \ln \frac{1}{\delta} + \ln \left(\left[\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \geq 1 - \delta, \end{aligned}$$

Then, we will prove separately the case when $\lambda \rightarrow 1$ and $\lambda \rightarrow +\infty$.

When $\lambda \rightarrow 1$. We have $\lim_{\lambda \rightarrow 1^+} \frac{\lambda-1}{\lambda} \operatorname{MI}_{\lambda}(h'; \mathbb{S}') = 0$. Furthermore, note that

$$\|\varphi\|_{\frac{\lambda}{\lambda-1}} = \left[\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(|\varphi(h', \mathbb{S}')|^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} = \left[\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', \mathbb{S}')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}}$$

is the $L^{\frac{\lambda}{\lambda-1}}$ -norm of the function $\varphi : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}_+^*$, where $\lim_{\lambda \rightarrow 1} \|\varphi\|_{\frac{\lambda}{\lambda-1}} = \lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'}$ (since we have $\lim_{\lambda \rightarrow 1^+} \frac{\lambda}{\lambda-1} = (\lim_{\lambda \rightarrow 1} \lambda)(\lim_{\lambda \rightarrow 1} \frac{1}{\lambda-1}) = +\infty$). Then, it is well known that

$$\|\varphi\|_{\infty} = \lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'} = \operatorname{esssup}_{\mathbb{S}' \in (\mathbb{X} \times \mathbb{Y}), h' \in \mathbb{H}} \varphi(h', \mathbb{S}').$$

Hence, we have

$$\begin{aligned}
 & \lim_{\lambda \rightarrow 1} \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow 1} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow 1} \|\varphi\|_{\frac{\lambda}{\lambda-1}} \right) = \ln \left(\lim_{\lambda' \rightarrow +\infty} \|\varphi\|_{\lambda'} \right) \\
 &= \ln (\|\varphi\|_{\infty}) = \ln \left(\operatorname{esssup}_{S' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathcal{H}} \varphi(h', S') \right).
 \end{aligned}$$

Finally, we can deduce that

$$\begin{aligned}
 & \lim_{\lambda \rightarrow 1} \left[\frac{\lambda-1}{\lambda} \operatorname{MI}_{\lambda}(h'; S') + \ln \frac{1}{\delta} + \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \\
 &= \ln \frac{1}{\delta} + \ln \left[\operatorname{esssup}_{S' \in (\mathcal{X} \times \mathcal{Y}), h' \in \mathcal{H}} \varphi(h', S') \right].
 \end{aligned}$$

When $\lambda \rightarrow +\infty$. First, we have $\lim_{\lambda \rightarrow +\infty} \|\varphi\|_{\frac{\lambda}{\lambda-1}} = \lim_{\lambda' \rightarrow 1} \|\varphi\|_{\lambda'} = \|\varphi\|_1$. Hence, we have

$$\begin{aligned}
 & \lim_{\lambda \rightarrow +\infty} \ln \left(\left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow +\infty} \left[\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \left(\varphi(h', S')^{\frac{\lambda}{\lambda-1}} \right) \right]^{\frac{\lambda-1}{\lambda}} \right) \\
 &= \ln \left(\lim_{\lambda \rightarrow +\infty} \|\varphi\|_{\frac{\lambda}{\lambda-1}} \right) = \ln \left(\lim_{\lambda' \rightarrow 1} \|\varphi\|_{\lambda'} \right) \\
 &= \ln (\|\varphi\|_1) = \ln \left(\mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \varphi(h', S') \right).
 \end{aligned}$$

Moreover, by rearranging the terms in $\frac{\lambda-1}{\lambda} \text{MI}_\lambda(h'; \mathcal{S}')$, we have

$$\begin{aligned} \frac{\lambda-1}{\lambda} \text{MI}_\lambda(h'; \mathcal{S}') &= \frac{1}{\lambda} \ln \left(\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \left(\left[\frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)} \right]^\lambda \right) \right) \\ &= \ln \left(\left[\mathbb{E}_{\mathcal{S} \sim \mathcal{D}^m} \mathbb{E}_{h \sim \pi^*} \left(\left[\frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)} \right]^\lambda \right) \right]^{\frac{1}{\lambda}} \right) \\ &= \ln \left(\left[\mathbb{E}_{h \sim \pi^*} (\gamma(h)^\lambda) \right]^{\frac{1}{\lambda}} \right) = \ln(\|\gamma\|_\lambda), \end{aligned}$$

where $\|\gamma\|_\lambda$ is the L^λ -norm of the function γ defined as $\gamma(h) = \frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)}$. We have

$$\begin{aligned} \lim_{\lambda \rightarrow +\infty} \frac{\lambda-1}{\lambda} \text{MI}_\lambda(h'; \mathcal{S}') &= \lim_{\lambda \rightarrow +\infty} \ln(\|\gamma\|_\lambda) = \ln \left(\lim_{\lambda \rightarrow +\infty} \|\gamma\|_\lambda \right) \\ &= \ln(\|\gamma\|_\infty) = \ln \left(\text{esssup}_{\mathcal{S} \in \mathcal{S}, h \in \mathcal{H}} \gamma(h) \right) = \ln \left(\text{esssup}_{\mathcal{S} \in \mathcal{S}, h \in \mathcal{H}} \frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)} \right). \end{aligned}$$

Finally, we can deduce that

$$\begin{aligned} &\lim_{\lambda \rightarrow 1} \left[\frac{\lambda-1}{\lambda} \text{MI}_\lambda(h'; \mathcal{S}') + \ln \frac{1}{\delta} + \ln \left(\left[\mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} (\varphi(h', \mathcal{S}')^{\frac{\lambda}{\lambda-1}}) \right]^{\frac{\lambda-1}{\lambda}} \right) \right] \\ &= \ln \left(\text{esssup}_{\mathcal{S} \in \mathcal{S}, h \in \mathcal{H}} \frac{\rho_{\mathcal{S}}(h)}{\pi^*(h)} \right) + \ln \left[\frac{1}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi^*} \varphi(h', \mathcal{S}') \right]. \end{aligned}$$

■

F.13 Proof of Corollary 6.6.1

Corollary 6.6.1 (Instantiation of Theorem 2.4.1 to Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any finite set of voters \mathbb{H} , for any hyper-prior distribution $\Pi = \text{Dir}(\beta)$ on \mathbb{H} with $\beta \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $\delta \in (0, 1]$, for any algorithm A that outputs a hyper-posterior given a learning sample and a hyper-prior, with probability at least $1-\delta$ over the learning sample $\mathcal{S} \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim P_{\mathcal{S}} = \text{Dir}(\alpha)$ with $\alpha \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$ we have

$$\text{kl}(\mathbb{R}_{\mathcal{S}}^\ell(\text{MV}_\rho) \| \mathbb{R}_{\mathcal{D}}^\ell(\text{MV}_\rho)) \leq \frac{1}{m} \left[\ln \frac{Z(\beta)}{Z(\alpha)} + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \ln(\rho(h_j)) + \ln \frac{2\sqrt{m}}{\delta} \right],$$

where $P_S \triangleq A(S, \Pi)$ is output by the deterministic algorithm A .

Proof. We apply Theorem 2.4.1 with $\phi(\rho, S) = m \text{kl}(R_S^\ell(MV_\rho) \| R_D^\ell(MV_\rho))$ to obtain with probability at least $1-\delta$ over the learning sample $S \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim P$, we have

$$\text{kl}(R_S^\ell(MV_\rho) \| R_D^\ell(MV_\rho)) \leq \frac{1}{m} \left[\ln \frac{P(\rho)}{\Pi(\rho)} + \ln \left[\frac{1}{\delta} \mathbb{E}_{\rho' \sim \Pi} e^{m \text{kl}(R_S^\ell(MV_{\rho'}) \| R_D^\ell(MV_{\rho'}))} \right] \right]. \quad (\text{F.12})$$

Moreover, the closed form solution of the disintegrated KL divergence $\ln \frac{P(\rho)}{\Pi(\rho)}$ is

$$\begin{aligned} \ln \frac{P(\rho)}{\Pi(\rho)} &= \ln(P(\rho)) - \ln(\Pi(\rho)) \\ &= \ln \left(\frac{1}{Z(\boldsymbol{\alpha})} \prod_{j=1}^{\text{card}(\mathbb{H})} [\rho(h_j)]^{\alpha_j - 1} \right) - \ln \left(\frac{1}{Z(\boldsymbol{\beta})} \prod_{j=1}^{\text{card}(\mathbb{H})} [\rho(h_j)]^{\beta_j - 1} \right) \\ &= \ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - 1) \ln(\rho(h_j)) - \sum_{j=1}^{\text{card}(\mathbb{H})} (\beta_j - 1) \ln(\rho(h_j)) \\ &= \ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \sum_{j=1}^{\text{card}(\mathbb{H})} (\alpha_j - \beta_j) \ln(\rho(h_j)). \end{aligned} \quad (\text{F.13})$$

Additionally, from Lemma B.16.1, we have

$$\mathbb{E}_{\rho' \sim \Pi} e^{m \text{kl}(R_S^\ell(MV_{\rho'}) \| R_D^\ell(MV_{\rho'}))} \leq 2\sqrt{m}. \quad (\text{F.14})$$

Lastly, by merging Equations (F.12) to (F.14) we obtain the claim. \blacksquare

F.14 Proof of Corollary 6.6.2

Corollary 6.6.2 (Instantiation of Theorem 6.3.1 to Stochastic Majority Votes). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any finite set of voters \mathbb{H} , for any hyper-prior distribution $\Pi = \text{Dir}(\boldsymbol{\beta})$ on \mathbb{H} with $\boldsymbol{\beta} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$, for any loss $\ell : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y}) \rightarrow [0, 1]$, for any $\lambda > 1$, for any $\delta \in (0, 1]$, for any algorithm A that outputs a hyper-posterior given a learning sample and a hyper-prior, with probability at least $1-\delta$ over the learning sample $S \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim P_S = \text{Dir}(\boldsymbol{\alpha})$

with $\boldsymbol{\alpha} \in (\mathbb{R}_*^+)^{\text{card}(\mathbb{H})}$ we have

$$\begin{aligned} \text{kl}(\mathbb{R}_S^\ell(\text{MV}_\rho) \parallel \mathbb{R}_D^\ell(\text{MV}_\rho)) &\leq \frac{1}{m} \left[\frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} + \ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} \right. \\ &\quad \left. + \frac{1}{\lambda-1} \ln \frac{Z(\lambda\boldsymbol{\alpha}+(1-\lambda)\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \ln(2\sqrt{m}) \right], \end{aligned}$$

where $P_S \triangleq A(S, \Pi)$ is output by the deterministic algorithm A .

Proof. We apply Theorem 6.3.1 with $\phi(\rho, S) = \exp \left[\frac{\lambda-1}{\lambda} m \text{kl}(\mathbb{R}_S^\ell(\text{MV}_\rho) \parallel \mathbb{R}_D^\ell(\text{MV}_\rho)) \right]$ to obtain with probability at least $1-\delta$ over the learning sample $S \sim \mathcal{D}^m$ and the posterior distribution $\rho \sim P$, we have

$$\begin{aligned} \text{kl}(\mathbb{R}_S^\ell(\text{MV}_\rho) \parallel \mathbb{R}_D^\ell(\text{MV}_\rho)) &\leq \frac{1}{m} \left[\frac{2\lambda-1}{\lambda-1} \ln \frac{2}{\delta} + D_\lambda(P \parallel \Pi) \right. \\ &\quad \left. + \ln \left(\mathbb{E}_{\rho' \sim \Pi} e^{m \text{kl}(\mathbb{R}_S^\ell(\text{MV}_{\rho'}) \parallel \mathbb{R}_D^\ell(\text{MV}_{\rho'}))} \right) \right]. \quad (\text{F.15}) \end{aligned}$$

Moreover, the closed form solution of the RÉNYI divergence $D_\lambda(P \parallel \Pi)$ (GIL *et al.*,

2013) is given by

$$\begin{aligned}
 D_\lambda(P\|\Pi) &= \frac{1}{\lambda-1} \ln \left(\int_{\mathcal{M}(\mathbb{H})} P(\rho)^\lambda \Pi(\rho)^{1-\lambda} d\xi(\rho) \right) \\
 &= \frac{1}{\lambda-1} \ln \left(\int_{\mathcal{M}(\mathbb{H})} \frac{1}{Z(\boldsymbol{\alpha})^\lambda} \prod_{j=1}^{\text{card}(\mathbb{H})} (\rho(h_j))^{\lambda(\alpha_j-1)} \right. \\
 &\quad \left. \frac{1}{Z(\boldsymbol{\beta})^{1-\lambda}} \prod_{j=1}^{\text{card}(\mathbb{H})} (\rho(h_j))^{(1-\lambda)(\beta_j-1)} d\xi(\rho) \right) \\
 &= \frac{1}{\lambda-1} \ln \left(\frac{Z(\boldsymbol{\beta})^{\lambda-1}}{Z(\boldsymbol{\alpha})^\lambda} \right) \\
 &\quad + \frac{1}{\lambda-1} \ln \left(\int_{\mathcal{M}(\mathbb{H})} \prod_{i=1}^{\text{card}(\mathbb{H})} (\rho(h_i))^{\lambda\alpha_i+(1-\lambda)\beta_i-1} d\xi(\rho) \right) \\
 &= \frac{1}{\lambda-1} \ln \left(\frac{Z(\boldsymbol{\beta})^{\lambda-1}}{Z(\boldsymbol{\alpha})^\lambda} \right) + \ln Z(\lambda\boldsymbol{\alpha}+(1-\lambda)\boldsymbol{\beta}) \\
 &= \frac{1}{\lambda-1} \ln \left(\frac{Z(\boldsymbol{\beta})^{\lambda-1}}{Z(\boldsymbol{\alpha})^{\lambda-1+1}} \right) + \ln Z(\lambda\boldsymbol{\alpha}+(1-\lambda)\boldsymbol{\beta}) \\
 &= \ln \frac{Z(\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})} + \frac{1}{\lambda-1} \ln \frac{Z(\lambda\boldsymbol{\alpha}+(1-\lambda)\boldsymbol{\beta})}{Z(\boldsymbol{\alpha})}, \tag{F.16}
 \end{aligned}$$

where ξ is the reference measure on $\mathcal{M}(\mathbb{H})$. Additionally, from Lemma B.16.1, we have

$$\mathbb{E}_{\rho' \sim \Pi} e^{m \text{kl}(\mathbb{R}_S^\ell(\text{MV}_{\rho'}) \|\mathbb{R}_D^\ell(\text{MV}_{\rho'}))} \leq 2\sqrt{m}. \tag{F.17}$$

Lastly, by merging Equations (F.15) to (F.17) we obtain the claim. ■

F.15 Details of the Results

Table F.2 to Table F.10 report empirical results for split ratios going from 0.0 to 0.9. Table F.11 to Table F.13 report the performances of the prior before applying Step 2).

For the split 0.0, since Step 1) is skipped, the prior distribution π is only initialized as introduced in Section 6.5.3.2. Note that in this case, $T = 1$ since we have only one prior. To do the same number of epochs compared to the other splits, we perform 11 epochs (instead of 1) for MNIST and Fashion-MNIST and 110 epochs (instead of 10) for CIFAR-10 during Step 2). The other parameters are not changed.

Table F.1. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the R ENYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this figure, that the split ratio is 0.0.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.901 \pm .002	.908 \pm .002	.901 \pm .002	.005	.897 \pm .013	.904 \pm .012	.897 \pm .012	.009	.898 \pm .017	.905 \pm .016	.898 \pm .016	.027	.902 \pm .015	.908 \pm .014	.901 \pm .015	.671
	blanchard	.901 \pm .002	.926 \pm .002	.901 \pm .002	122.846 \pm 15.952	.897 \pm .013	.912 \pm .012	.897 \pm .013	39.350 \pm 8.999	.898 \pm .017	.907 \pm .016	.898 \pm .017	13.023 \pm 4.818	.901 \pm .015	.907 \pm .014	.901 \pm .014	3.041 \pm 2.459
	catoni	.901 \pm .002	.926 \pm .003	.901 \pm .002	121.860 \pm 15.930	.897 \pm .013	.909 \pm .012	.897 \pm .013	38.552 \pm 8.872	.898 \pm .017	.905 \pm .016	.898 \pm .017	12.474 \pm 4.774	.901 \pm .014	.906 \pm .013	.901 \pm .014	3.088 \pm 2.379
	rivasplata	.901 \pm .002	.920 \pm .002	.901 \pm .002	123.301 \pm 15.941	.896 \pm .014	.908 \pm .012	.896 \pm .013	39.195 \pm 8.959	.897 \pm .017	.905 \pm .016	.897 \pm .017	12.827 \pm 4.858	.902 \pm .015	.907 \pm .014	.901 \pm .015	3.232 \pm 2.454
	stochastic	—	.944	—	.002	—	.941	—	.004	—	.941	—	.014	—	.944	—	.336
Fashion	ours	.970 \pm .028	.972 \pm .025	.970 \pm .027	.016	.944 \pm .038	.949 \pm .035	.944 \pm .037	.046	.910 \pm .027	.917 \pm .026	.910 \pm .027	.140	.901 \pm .026	.909 \pm .025	.901 \pm .026	1.255
	blanchard	.970 \pm .029	.978 \pm .019	.970 \pm .028	122.508 \pm 16.085	.942 \pm .038	.952 \pm .032	.943 \pm .038	39.957 \pm 8.610	.910 \pm .031	.919 \pm .029	.910 \pm .031	12.649 \pm 4.846	.899 \pm .028	.905 \pm .027	.899 \pm .028	3.206 \pm 2.566
	catoni	.970 \pm .028	.983 \pm .017	.970 \pm .027	122.364 \pm 15.860	.945 \pm .038	.954 \pm .036	.945 \pm .037	38.555 \pm 8.873	.912 \pm .032	.919 \pm .031	.912 \pm .032	12.167 \pm 4.762	.899 \pm .027	.905 \pm .026	.899 \pm .027	3.122 \pm 2.392
	rivasplata	.970 \pm .028	.977 \pm .021	.971 \pm .027	123.328 \pm 15.929	.943 \pm .038	.950 \pm .033	.943 \pm .038	39.300 \pm 8.991	.908 \pm .031	.916 \pm .029	.908 \pm .031	12.627 \pm 4.890	.899 \pm .028	.905 \pm .027	.899 \pm .028	3.591 \pm 2.610
	stochastic	—	.990	—	.008	—	.975	—	.023	—	.950	—	.070	—	.944	—	.627
CIFAR-10	ours	.899 \pm .000	.907 \pm .000	.899 \pm .000	3.113	.896 \pm .002	.914 \pm .002	.894 \pm .002	107.797	.826 \pm .011	.885 \pm .009	.825 \pm .010	76.475	.786 \pm .019	.851 \pm .015	.788 \pm .018	714.351
	blanchard	.899 \pm .000	.940 \pm .001	.898 \pm .000	314.983 \pm 26.377	.888 \pm .004	.927 \pm .002	.885 \pm .003	28.250 \pm 25.255	.823 \pm .010	.885 \pm .008	.822 \pm .010	422.401 \pm 29.323	.798 \pm .019	.856 \pm .015	.799 \pm .018	292.706 \pm 25.318
	catoni	.899 \pm .000	.941 \pm .000	.898 \pm .000	285.415 \pm 25.085	.894 \pm .002	.930 \pm .004	.892 \pm .002	169.713 \pm 19.543	.857 \pm .010	.915 \pm .009	.856 \pm .010	273.554 \pm 23.212	.815 \pm .019	.864 \pm .017	.816 \pm .018	209.069 \pm 21.230
	rivasplata	.899 \pm .001	.930 \pm .001	.898 \pm .000	362.070 \pm 28.420	.864 \pm .004	.933 \pm .002	.862 \pm .004	1568.007 \pm 55.492	.748 \pm .010	.837 \pm .007	.750 \pm .009	1219.178 \pm 49.610	.769 \pm .018	.828 \pm .015	.771 \pm .017	526.068 \pm 33.837
	stochastic	—	.942	—	1.557	—	.945	—	53.898	—	.914	—	38.237	—	.884	—	357.175
		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.901 \pm .002	.909 \pm .002	.901 \pm .002	3.767	.896 \pm .014	.904 \pm .013	.896 \pm .014	.835	.898 \pm .016	.905 \pm .015	.898 \pm .016	1.062	.901 \pm .015	.909 \pm .014	.901 \pm .015	6.022
	blanchard	.900 \pm .003	.990 \pm .000	.900 \pm .003	12004.196 \pm 152.632	.894 \pm .017	.986 \pm .006	.894 \pm .016	3837.785 \pm 93.560	.888 \pm .021	.957 \pm .013	.888 \pm .020	1221.198 \pm 49.920	.898 \pm .015	.939 \pm .012	.897 \pm .015	391.343 \pm 28.182
	catoni	.900 \pm .003	.997 \pm .002	.900 \pm .003	5694.194 \pm 102.906	.889 \pm .020	.967 \pm .012	.889 \pm .019	3331.617 \pm 78.945	.879 \pm .025	.941 \pm .016	.880 \pm .025	1481.726 \pm 53.973	.888 \pm .023	.937 \pm .015	.888 \pm .023	567.893 \pm 33.441
	rivasplata	.900 \pm .004	.990 \pm .000	.900 \pm .003	1199.818 \pm 152.557	.892 \pm .017	.970 \pm .009	.892 \pm .016	3846.699 \pm 84.643	.886 \pm .020	.940 \pm .015	.886 \pm .020	1224.463 \pm 49.970	.897 \pm .018	.928 \pm .015	.897 \pm .018	393.757 \pm 29.158
	stochastic	—	.944	—	1.884	—	.940	—	.417	—	.941	—	.531	—	.944	—	3.011
Fashion	ours	.977 \pm .024	.979 \pm .021	.977 \pm .023	3.926	.947 \pm .038	.951 \pm .035	.947 \pm .038	1.623	.907 \pm .030	.914 \pm .029	.907 \pm .030	2.947	.900 \pm .026	.910 \pm .025	.900 \pm .026	15.978
	blanchard	.984 \pm .015	.990 \pm .000	.984 \pm .015	12019.121 \pm 166.251	.912 \pm .029	.988 \pm .004	.911 \pm .029	3846.861 \pm 84.568	.883 \pm .029	.953 \pm .019	.883 \pm .029	1232.645 \pm 5.285	.403 \pm .041	.648 \pm .038	.399 \pm .041	3853.231 \pm 87.867
	catoni	.983 \pm .018	1.000 \pm .000	.983 \pm .017	5654.642 \pm 114.040	.903 \pm .021	.985 \pm .012	.902 \pm .021	4354.538 \pm 94.427	.751 \pm .033	.867 \pm .023	.750 \pm .033	2702.652 \pm 76.863	.504 \pm .041	.673 \pm .037	.502 \pm .041	3172.609 \pm 78.698
	rivasplata	.983 \pm .016	.990 \pm .000	.983 \pm .016	11976.720 \pm 165.964	.905 \pm .023	.975 \pm .007	.905 \pm .023	3855.872 \pm 84.676	.855 \pm .035	.916 \pm .027	.855 \pm .035	125.110 \pm 51.837	.365 \pm .032	.559 \pm .032	.359 \pm .033	4823.725 \pm 103.813
	stochastic	—	.990	—	1.963	—	.977	—	.812	—	.948	—	1.473	—	.944	—	7.989
CIFAR-10	ours	.899 \pm .000	.915 \pm .000	.899 \pm .000	63.416	.890 \pm .003	.932 \pm .003	.886 \pm .003	68.353	.786 \pm .011	.888 \pm .008	.787 \pm .010	2072.610	.769 \pm .017	.859 \pm .013	.770 \pm .017	1406.824
	blanchard	.869 \pm .002	.990 \pm .000	.866 \pm .001	27237.938 \pm 251.770	.813 \pm .004	.990 \pm .000	.812 \pm .003	12052.733 \pm 159.732	.697 \pm .011	.920 \pm .005	.700 \pm .009	5137.799 \pm 103.680	.674 \pm .020	.861 \pm .014	.675 \pm .020	2814.450 \pm 76.004
	catoni	.928 \pm .001	1.000 \pm .000	.925 \pm .001	2145276.795 \pm 2095.160	.821 \pm .002	1.000 \pm .000	.821 \pm .002	375019.277 \pm 896.780	.689 \pm .011	.870 \pm .007	.692 \pm .010	5292.535 \pm 106.380	.629 \pm .019	.805 \pm .015	.628 \pm .019	4159.131 \pm 96.763
	rivasplata	.867 \pm .002	.990 \pm .000	.864 \pm .001	35956.152 \pm 268.304	.812 \pm .004	.976 \pm .001	.811 \pm .003	12135.134 \pm 157.621	.698 \pm .010	.874 \pm .006	.701 \pm .009	5191.665 \pm 102.712	.677 \pm .020	.819 \pm .015	.678 \pm .019	2839.514 \pm 81.432
	stochastic	—	.947	—	31.708	—	.954	—	34.176	—	.908	—	1036.305	—	.886	—	703.412

Table F.2. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.1.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
$\eta = 10^{-6}$		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.035 \pm .000	.044 \pm .000	.039 \pm .000	.622	.024 \pm .000	.034 \pm .000	.029 \pm .000	2.122	.029 \pm .002	.040 \pm .002	.034 \pm .002	12.754	.034 \pm .004	.044 \pm .004	.038 \pm .004	7.303
	blanchard	.034 \pm .000	.058 \pm .002	.038 \pm .000	99.876 \pm 14.858	.024 \pm .000	.038 \pm .001	.030 \pm .000	21.775 \pm 6.848	.034 \pm .002	.043 \pm .002	.038 \pm .002	3.949 \pm 2.877	.039 \pm .005	.047 \pm .005	.043 \pm .005	.590 \pm 1.085
	catoni	.035 \pm .000	.064 \pm .001	.039 \pm .000	119.663 \pm 15.854	.024 \pm .000	.038 \pm .001	.030 \pm .000	26.277 \pm 7.490	.033 \pm .002	.041 \pm .002	.037 \pm .002	4.067 \pm 2.882	.038 \pm .005	.045 \pm .005	.042 \pm .004	.759 \pm 1.217
	rivasplata	.034 \pm .000	.052 \pm .001	.038 \pm .000	104.880 \pm 15.268	.024 \pm .000	.036 \pm .001	.029 \pm .000	23.007 \pm 7.187	.033 \pm .002	.042 \pm .002	.037 \pm .002	4.116 \pm 2.845	.038 \pm .005	.046 \pm .004	.042 \pm .004	.775 \pm 1.231
	stochastic	—	.080	—	.311	—	.067	—	1.061	—	.074	—	6.377	—	.079	—	3.651
Fashion	ours	.166 \pm .001	.169 \pm .000	.159 \pm .000	.580	.157 \pm .001	.160 \pm .001	.150 \pm .001	2.128	.160 \pm .002	.161 \pm .003	.151 \pm .002	3.503	.176 \pm .006	.179 \pm .006	.168 \pm .005	1.268
	blanchard	.165 \pm .001	.192 \pm .002	.159 \pm .000	96.822 \pm 14.116	.157 \pm .001	.166 \pm .002	.150 \pm .001	21.592 \pm 6.681	.163 \pm .003	.162 \pm .003	.153 \pm .003	3.846 \pm 2.660	.178 \pm .005	.178 \pm .005	.170 \pm .005	.463 \pm .954
	catoni	.165 \pm .001	.190 \pm .003	.159 \pm .000	119.927 \pm 15.938	.157 \pm .001	.163 \pm .002	.150 \pm .001	26.363 \pm 7.355	.162 \pm .003	.161 \pm .003	.152 \pm .003	4.152 \pm 2.945	.177 \pm .006	.178 \pm .006	.169 \pm .006	.548 \pm 1.032
	rivasplata	.165 \pm .001	.183 \pm .002	.158 \pm .000	101.954 \pm 14.463	.157 \pm .001	.163 \pm .002	.150 \pm .001	23.098 \pm 6.977	.162 \pm .003	.161 \pm .003	.153 \pm .003	3.852 \pm 2.798	.177 \pm .006	.177 \pm .006	.169 \pm .006	.516 \pm .985
	stochastic	—	.227	—	.290	—	.216	—	1.064	—	.218	—	1.751	—	.237	—	.634
CIFAR-10	ours	.479 \pm .000	.487 \pm .000	.472 \pm .000	.052	.479 \pm .000	.493 \pm .000	.477 \pm .000	.065	.458 \pm .001	.479 \pm .000	.463 \pm .000	.299	.480 \pm .002	.495 \pm .001	.480 \pm .001	.793
	blanchard	.479 \pm .000	.550 \pm .003	.472 \pm .000	27.644 \pm 22.868	.479 \pm .000	.522 \pm .003	.477 \pm .000	85.476 \pm 12.781	.458 \pm .001	.489 \pm .003	.463 \pm .000	24.608 \pm 7.136	.481 \pm .002	.495 \pm .002	.480 \pm .001	5.093 \pm 3.299
	catoni	.479 \pm .000	.546 \pm .005	.472 \pm .000	269.855 \pm 22.883	.479 \pm .000	.511 \pm .003	.477 \pm .000	85.113 \pm 12.806	.458 \pm .001	.483 \pm .002	.463 \pm .000	25.453 \pm 7.155	.480 \pm .002	.495 \pm .001	.480 \pm .001	5.468 \pm 3.315
	rivasplata	.479 \pm .000	.528 \pm .002	.472 \pm .000	27.588 \pm 22.859	.479 \pm .000	.511 \pm .002	.477 \pm .000	85.745 \pm 13.357	.458 \pm .001	.484 \pm .002	.463 \pm .001	25.051 \pm 7.005	.481 \pm .002	.494 \pm .001	.480 \pm .001	5.155 \pm 3.260
	stochastic	—	.558	—	.026	—	.564	—	.032	—	.550	—	.150	—	.566	—	.397
		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
$\eta = 10^{-4}$		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.035 \pm .000	.048 \pm .000	.039 \pm .000	35.348	.024 \pm .000	.037 \pm .001	.029 \pm .000	3.753	.022 \pm .001	.042 \pm .001	.027 \pm .001	153.773	.025 \pm .002	.041 \pm .002	.029 \pm .002	97.840
	blanchard	.032 \pm .000	.442 \pm .003	.036 \pm .000	1181.482 \pm 14.449	.022 \pm .000	.206 \pm .003	.027 \pm .000	3851.110 \pm 84.274	.019 \pm .001	.102 \pm .002	.023 \pm .001	1306.371 \pm 51.396	.024 \pm .002	.065 \pm .003	.027 \pm .002	411.772 \pm 29.458
	catoni	.035 \pm .000	.362 \pm .003	.039 \pm .000	11925.734 \pm 145.511	.024 \pm .000	.152 \pm .002	.029 \pm .000	3841.248 \pm 84.033	.027 \pm .002	.084 \pm .002	.032 \pm .001	1235.287 \pm 49.454	.027 \pm .002	.059 \pm .003	.030 \pm .002	403.300 \pm 28.587
	rivasplata	.030 \pm .000	.289 \pm .002	.034 \pm .000	12022.576 \pm 151.157	.021 \pm .000	.134 \pm .002	.026 \pm .000	3912.803 \pm 85.146	.018 \pm .000	.072 \pm .001	.022 \pm .000	1348.169 \pm 53.400	.023 \pm .002	.051 \pm .002	.026 \pm .001	424.971 \pm 29.301
	stochastic	—	.084	—	17.674	—	.069	—	15.376	—	.072	—	76.887	—	.072	—	48.920
Fashion	ours	.166 \pm .001	.172 \pm .000	.159 \pm .000	13.084	.157 \pm .001	.163 \pm .001	.150 \pm .001	16.513	.159 \pm .002	.164 \pm .002	.149 \pm .002	2.344	.176 \pm .005	.181 \pm .005	.168 \pm .005	11.331
	blanchard	.160 \pm .001	.588 \pm .003	.153 \pm .000	1089.829 \pm 137.125	.150 \pm .001	.379 \pm .003	.141 \pm .001	3744.491 \pm 83.656	.155 \pm .002	.271 \pm .003	.145 \pm .002	1221.062 \pm 49.548	.173 \pm .005	.233 \pm .006	.165 \pm .004	369.721 \pm 27.211
	catoni	.165 \pm .001	.500 \pm .003	.159 \pm .000	11954.591 \pm 141.463	.156 \pm .001	.311 \pm .002	.148 \pm .001	3826.848 \pm 86.111	.158 \pm .002	.248 \pm .003	.148 \pm .002	1226.282 \pm 5.332	.174 \pm .005	.252 \pm .006	.166 \pm .004	393.542 \pm 27.890
	rivasplata	.158 \pm .001	.459 \pm .002	.151 \pm .000	11541.128 \pm 14.706	.149 \pm .001	.302 \pm .002	.140 \pm .001	3878.145 \pm 85.782	.154 \pm .002	.230 \pm .002	.144 \pm .001	1244.035 \pm 49.268	.172 \pm .005	.212 \pm .005	.164 \pm .004	378.990 \pm 27.559
	stochastic	—	.229	—	6.542	—	.219	—	8.257	—	.219	—	1.172	—	.239	—	5.666
CIFAR-10	ours	.479 \pm .000	.489 \pm .000	.472 \pm .000	4.882	.479 \pm .000	.496 \pm .000	.477 \pm .000	9.273	.458 \pm .001	.480 \pm .000	.463 \pm .000	4.988	.480 \pm .002	.497 \pm .001	.479 \pm .001	8.681
	blanchard	.479 \pm .000	.957 \pm .001	.471 \pm .000	22201.935 \pm 218.369	.479 \pm .000	.854 \pm .002	.477 \pm .000	8777.551 \pm 125.716	.457 \pm .001	.699 \pm .003	.461 \pm .000	2758.075 \pm 77.155	.474 \pm .001	.613 \pm .003	.472 \pm .001	903.948 \pm 4.742
	catoni	.479 \pm .000	.995 \pm .000	.471 \pm .000	26347.736 \pm 225.908	.479 \pm .000	.771 \pm .002	.477 \pm .000	8566.272 \pm 124.834	.455 \pm .001	.650 \pm .002	.459 \pm .000	3117.566 \pm 75.178	.468 \pm .001	.621 \pm .001	.466 \pm .001	1481.520 \pm 52.533
	rivasplata	.479 \pm .000	.915 \pm .001	.471 \pm .000	29489.241 \pm 241.010	.479 \pm .000	.765 \pm .002	.477 \pm .000	867.264 \pm 126.038	.456 \pm .001	.633 \pm .002	.460 \pm .000	2776.052 \pm 72.901	.472 \pm .001	.572 \pm .002	.470 \pm .001	937.091 \pm 42.116
	stochastic	—	.559	—	2.441	—	.566	—	4.637	—	.551	—	2.494	—	.567	—	4.340

Table F.3. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the R NYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.2.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
$\eta = 10^{-6}$		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.016 \pm .000	.023 \pm .000	.019 \pm .000	.336	.015 \pm .000	.023 \pm .000	.019 \pm .000	.748	.014 \pm .001	.020 \pm .001	.016 \pm .000	2.096	.019 \pm .002	.024 \pm .002	.020 \pm .002	2.244
	blanchard	.016 \pm .000	.034 \pm .001	.019 \pm .000	97.590 \pm 14.260	.015 \pm .000	.026 \pm .001	.019 \pm .000	21.153 \pm 6.514	.015 \pm .001	.020 \pm .001	.016 \pm .001	3.362 \pm 2.569	.020 \pm .002	.024 \pm .002	.021 \pm .002	.371 \pm .875
	catoni	.016 \pm .000	.034 \pm .001	.019 \pm .000	116.744 \pm 15.447	.015 \pm .000	.027 \pm .002	.019 \pm .000	24.135 \pm 7.075	.015 \pm .001	.020 \pm .001	.016 \pm .001	3.352 \pm 2.667	.020 \pm .002	.024 \pm .002	.021 \pm .002	.410 \pm .890
	rivasplata	.016 \pm .000	.030 \pm .001	.019 \pm .000	101.334 \pm 14.728	.015 \pm .000	.024 \pm .001	.019 \pm .000	21.663 \pm 6.603	.015 \pm .001	.020 \pm .001	.016 \pm .001	3.409 \pm 2.666	.020 \pm .002	.024 \pm .002	.021 \pm .002	.446 \pm .927
	stochastic	—	.052	—	.168	—	.051	—	.374	—	.047	—	1.048	—	.053	—	1.122
	Fashion	ours	.165 \pm .002	.169 \pm .001	.157 \pm .001	4.811	.148 \pm .003	.155 \pm .002	.143 \pm .002	1.856	.145 \pm .005	.153 \pm .006	.139 \pm .005	15.453	.160 \pm .005	.166 \pm .005	.155 \pm .005
blanchard		.163 \pm .002	.190 \pm .003	.155 \pm .001	96.264 \pm 14.472	.152 \pm .003	.163 \pm .003	.147 \pm .003	21.099 \pm 6.507	.155 \pm .007	.160 \pm .007	.151 \pm .007	3.929 \pm 2.841	.163 \pm .006	.165 \pm .006	.158 \pm .006	.340 \pm .885
catoni		.163 \pm .002	.190 \pm .004	.156 \pm .001	121.542 \pm 16.499	.150 \pm .002	.158 \pm .003	.144 \pm .002	27.241 \pm 7.318	.151 \pm .006	.155 \pm .006	.146 \pm .006	5.120 \pm 3.150	.162 \pm .005	.165 \pm .005	.157 \pm .005	.444 \pm .968
rivasplata		.161 \pm .001	.180 \pm .002	.153 \pm .001	106.403 \pm 14.044	.150 \pm .002	.158 \pm .003	.145 \pm .003	23.134 \pm 7.064	.153 \pm .006	.157 \pm .006	.148 \pm .007	4.439 \pm 2.924	.162 \pm .006	.165 \pm .005	.157 \pm .005	.417 \pm .928
stochastic		—	.226	—	2.405	—	.210	—	5.428	—	.207	—	7.727	—	.223	—	.816
CIFAR-10		ours	.390 \pm .000	.407 \pm .000	.391 \pm .000	.040	.404 \pm .000	.414 \pm .000	.398 \pm .000	.070	.396 \pm .001	.411 \pm .000	.395 \pm .000	.155	.416 \pm .002	.432 \pm .001	.415 \pm .001
	blanchard	.390 \pm .000	.473 \pm .004	.391 \pm .000	271.616 \pm 23.555	.404 \pm .000	.445 \pm .003	.398 \pm .000	84.868 \pm 13.050	.396 \pm .001	.422 \pm .003	.395 \pm .000	23.962 \pm 7.208	.416 \pm .002	.432 \pm .002	.416 \pm .001	4.496 \pm 3.018
	catoni	.390 \pm .000	.473 \pm .006	.391 \pm .000	27.502 \pm 23.371	.404 \pm .000	.434 \pm .003	.398 \pm .000	84.848 \pm 12.992	.396 \pm .001	.415 \pm .002	.395 \pm .000	24.505 \pm 6.942	.416 \pm .002	.431 \pm .001	.415 \pm .001	4.859 \pm 3.176
	rivasplata	.390 \pm .000	.450 \pm .002	.391 \pm .000	271.700 \pm 23.586	.403 \pm .000	.433 \pm .002	.398 \pm .000	85.027 \pm 13.047	.396 \pm .001	.416 \pm .002	.395 \pm .000	23.955 \pm 7.093	.416 \pm .002	.431 \pm .002	.416 \pm .001	4.610 \pm 3.084
	stochastic	—	.477	—	.020	—	.485	—	.035	—	.482	—	.077	—	.503	—	.485
			$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$		
$\eta = 10^{-4}$		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.016 \pm .000	.025 \pm .000	.019 \pm .000	14.490	.015 \pm .000	.024 \pm .000	.019 \pm .000	8.583	.014 \pm .000	.021 \pm .001	.016 \pm .000	13.055	.016 \pm .001	.023 \pm .001	.017 \pm .001	25.556
	blanchard	.016 \pm .000	.430 \pm .004	.018 \pm .000	11405.062 \pm 153.554	.014 \pm .000	.200 \pm .003	.018 \pm .000	3799.912 \pm 89.585	.013 \pm .000	.086 \pm .002	.014 \pm .000	1187.859 \pm 48.700	.015 \pm .001	.049 \pm .002	.016 \pm .001	38.983 \pm 27.857
	catoni	.016 \pm .000	.355 \pm .002	.019 \pm .000	11954.106 \pm 15.709	.015 \pm .000	.149 \pm .003	.019 \pm .000	3828.342 \pm 83.937	.014 \pm .001	.064 \pm .002	.016 \pm .001	1218.708 \pm 48.514	.017 \pm .001	.041 \pm .002	.018 \pm .001	389.726 \pm 29.076
	rivasplata	.015 \pm .000	.272 \pm .002	.018 \pm .000	1173.953 \pm 149.364	.013 \pm .000	.122 \pm .002	.017 \pm .000	3691.345 \pm 82.512	.012 \pm .000	.056 \pm .001	.013 \pm .000	1206.615 \pm 5.381	.015 \pm .001	.037 \pm .001	.015 \pm .001	391.881 \pm 28.344
	stochastic	—	.053	—	7.245	—	.052	—	4.292	—	.048	—	6.528	—	.051	—	12.778
	Fashion	ours	.165 \pm .002	.172 \pm .001	.157 \pm .001	23.705	.141 \pm .002	.156 \pm .002	.137 \pm .002	52.736	.131 \pm .003	.147 \pm .003	.126 \pm .003	7.515	.156 \pm .004	.165 \pm .004	.151 \pm .003
blanchard		.136 \pm .001	.598 \pm .003	.130 \pm .001	11334.327 \pm 145.083	.125 \pm .001	.379 \pm .003	.121 \pm .001	3998.068 \pm 88.992	.124 \pm .001	.247 \pm .003	.117 \pm .001	126.184 \pm 48.814	.152 \pm .003	.216 \pm .004	.147 \pm .003	364.531 \pm 28.029
catoni		.162 \pm .001	.525 \pm .004	.154 \pm .001	11965.668 \pm 152.681	.141 \pm .002	.309 \pm .003	.137 \pm .002	384.802 \pm 84.123	.132 \pm .003	.224 \pm .004	.127 \pm .002	1239.918 \pm 49.594	.155 \pm .004	.232 \pm .005	.150 \pm .004	394.607 \pm 28.146
rivasplata		.131 \pm .001	.455 \pm .002	.127 \pm .001	1193.209 \pm 155.390	.123 \pm .001	.290 \pm .002	.119 \pm .001	4005.169 \pm 89.793	.123 \pm .001	.204 \pm .002	.116 \pm .001	1294.726 \pm 49.874	.152 \pm .004	.195 \pm .004	.146 \pm .003	378.905 \pm 27.422
stochastic		—	.228	—	11.853	—	.209	—	26.368	—	.198	—	35.258	—	.221	—	8.477
CIFAR-10		ours	.390 \pm .000	.411 \pm .000	.391 \pm .000	13.286	.404 \pm .000	.415 \pm .000	.398 \pm .000	3.305	.396 \pm .001	.412 \pm .000	.395 \pm .000	3.136	.415 \pm .001	.433 \pm .001	.415 \pm .001
	blanchard	.389 \pm .000	.990 \pm .000	.391 \pm .000	75424.764 \pm 397.521	.403 \pm .000	.820 \pm .002	.397 \pm .000	8815.324 \pm 126.764	.395 \pm .001	.651 \pm .003	.394 \pm .000	2738.066 \pm 75.053	.408 \pm .001	.557 \pm .003	.405 \pm .001	918.500 \pm 42.347
	catoni	.390 \pm .000	.990 \pm .000	.391 \pm .000	26434.787 \pm 228.500	.403 \pm .000	.726 \pm .003	.397 \pm .000	8651.380 \pm 126.473	.394 \pm .001	.620 \pm .002	.393 \pm .000	4178.302 \pm 9.315	.401 \pm .001	.556 \pm .001	.396 \pm .001	1462.235 \pm 55.526
	rivasplata	.389 \pm .000	.902 \pm .001	.391 \pm .000	31497.669 \pm 249.683	.403 \pm .000	.715 \pm .002	.397 \pm .000	8707.893 \pm 133.239	.394 \pm .001	.578 \pm .003	.393 \pm .000	2741.257 \pm 74.942	.405 \pm .001	.512 \pm .002	.402 \pm .001	967.818 \pm 43.629
	stochastic	—	.480	—	6.643	—	.486	—	1.653	—	.483	—	1.568	—	.503	—	3.032

Table F.4. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.3.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
$\eta = 10^{-6}$		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.012 \pm .000	.017 \pm .000	.013 \pm .000	.181	.009 \pm .000	.015 \pm .000	.011 \pm .000	.155	.012 \pm .000	.020 \pm .000	.016 \pm .000	1.655	.013 \pm .001	.019 \pm .001	.015 \pm .001	.615
	blanchard	.012 \pm .000	.027 \pm .001	.013 \pm .000	93.915 \pm 14.109	.012 \pm .000	.021 \pm .001	.014 \pm .000	19.292 \pm 6.037	.012 \pm .000	.020 \pm .001	.016 \pm .000	3.023 \pm 2.430	.014 \pm .001	.018 \pm .001	.015 \pm .001	.368 \pm .831
	catoni	.012 \pm .000	.025 \pm .001	.013 \pm .000	113.574 \pm 15.436	.012 \pm .000	.023 \pm .002	.014 \pm .000	22.347 \pm 6.877	.012 \pm .000	.020 \pm .001	.016 \pm .000	2.918 \pm 2.341	.013 \pm .001	.018 \pm .001	.015 \pm .001	.336 \pm .807
	rivasplata	.012 \pm .000	.023 \pm .001	.013 \pm .000	96.392 \pm 14.300	.012 \pm .000	.020 \pm .001	.014 \pm .000	19.905 \pm 6.254	.012 \pm .000	.020 \pm .001	.016 \pm .000	2.931 \pm 2.446	.013 \pm .001	.018 \pm .001	.015 \pm .001	.355 \pm .813
	stochastic	—	.042	—	.091	—	.039	—	.077	—	.047	—	.827	—	.045	—	.308
	Fashion	ours	.126 \pm .000	.134 \pm .000	.124 \pm .000	.328	.126 \pm .001	.130 \pm .001	.119 \pm .001	1.692	.122 \pm .002	.126 \pm .002	.115 \pm .002	4.617	.139 \pm .005	.145 \pm .005	.133 \pm .005
blanchard		.126 \pm .000	.157 \pm .003	.124 \pm .000	88.034 \pm 13.485	.126 \pm .001	.136 \pm .002	.120 \pm .001	18.852 \pm 6.115	.124 \pm .002	.127 \pm .002	.118 \pm .002	3.014 \pm 2.395	.142 \pm .006	.144 \pm .006	.137 \pm .006	.370 \pm .819
catoni		.126 \pm .000	.159 \pm .004	.124 \pm .000	114.259 \pm 15.300	.126 \pm .001	.133 \pm .002	.120 \pm .001	22.607 \pm 6.871	.124 \pm .002	.126 \pm .002	.118 \pm .002	3.100 \pm 2.513	.141 \pm .006	.144 \pm .006	.136 \pm .006	.390 \pm .898
rivasplata		.126 \pm .000	.148 \pm .002	.124 \pm .000	93.107 \pm 13.630	.126 \pm .001	.133 \pm .002	.120 \pm .001	19.724 \pm 6.320	.124 \pm .002	.126 \pm .002	.118 \pm .002	2.980 \pm 2.451	.142 \pm .006	.144 \pm .006	.136 \pm .006	.371 \pm .869
stochastic		—	.187	—	.164	—	.182	—	.846	—	.178	—	2.309	—	.199	—	1.212
CIFAR-10		ours	.369 \pm .000	.375 \pm .000	.358 \pm .000	.028	.351 \pm .000	.368 \pm .000	.352 \pm .000	.041	.359 \pm .001	.377 \pm .000	.360 \pm .000	.183	.419 \pm .001	.433 \pm .001	.416 \pm .001
	blanchard	.369 \pm .000	.446 \pm .004	.358 \pm .000	269.789 \pm 22.724	.351 \pm .000	.401 \pm .004	.352 \pm .000	84.113 \pm 12.530	.359 \pm .001	.388 \pm .003	.360 \pm .000	22.878 \pm 6.728	.419 \pm .001	.432 \pm .003	.416 \pm .001	4.089 \pm 2.818
	catoni	.369 \pm .000	.450 \pm .007	.358 \pm .000	269.843 \pm 24.225	.351 \pm .000	.390 \pm .004	.352 \pm .000	84.500 \pm 12.608	.359 \pm .001	.381 \pm .002	.360 \pm .000	23.567 \pm 7.181	.419 \pm .001	.432 \pm .001	.416 \pm .001	4.285 \pm 2.942
	rivasplata	.369 \pm .000	.421 \pm .003	.358 \pm .000	27.224 \pm 24.187	.351 \pm .000	.388 \pm .002	.352 \pm .000	84.250 \pm 13.274	.359 \pm .001	.382 \pm .002	.360 \pm .000	23.053 \pm 6.724	.419 \pm .001	.431 \pm .002	.416 \pm .001	4.141 \pm 2.985
	stochastic	—	.445	—	.014	—	.438	—	.020	—	.447	—	.092	—	.504	—	.380
	$\eta = 10^{-4}$		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$		
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.012 \pm .000	.019 \pm .000	.013 \pm .000	24.837	.012 \pm .000	.020 \pm .000	.014 \pm .000	12.358	.012 \pm .000	.021 \pm .000	.015 \pm .000	13.908	.013 \pm .001	.019 \pm .001	.014 \pm .001	16.179
	blanchard	.012 \pm .000	.467 \pm .004	.013 \pm .000	11819.223 \pm 154.992	.011 \pm .000	.211 \pm .003	.014 \pm .000	3808.981 \pm 86.014	.010 \pm .000	.094 \pm .003	.014 \pm .000	121.397 \pm 51.944	.012 \pm .001	.046 \pm .002	.013 \pm .001	372.832 \pm 26.602
	catoni	.012 \pm .000	.339 \pm .002	.013 \pm .000	1196.394 \pm 15.704	.012 \pm .000	.159 \pm .003	.014 \pm .000	3838.459 \pm 88.155	.012 \pm .000	.070 \pm .002	.016 \pm .000	1218.505 \pm 51.783	.013 \pm .001	.037 \pm .001	.014 \pm .001	386.824 \pm 28.233
	rivasplata	.012 \pm .000	.289 \pm .003	.013 \pm .000	1191.037 \pm 152.759	.011 \pm .000	.128 \pm .002	.014 \pm .000	3768.785 \pm 9.947	.010 \pm .000	.061 \pm .001	.013 \pm .000	1231.638 \pm 49.362	.011 \pm .001	.033 \pm .001	.012 \pm .000	382.225 \pm 28.481
	stochastic	—	.044	—	12.418	—	.046	—	6.179	—	.047	—	6.954	—	.045	—	8.089
	Fashion	ours	.126 \pm .000	.137 \pm .000	.124 \pm .000	12.401	.125 \pm .001	.132 \pm .001	.119 \pm .001	14.631	.120 \pm .002	.128 \pm .002	.113 \pm .001	26.499	.133 \pm .003	.143 \pm .003	.127 \pm .003
blanchard		.123 \pm .000	.602 \pm .003	.121 \pm .000	10558.872 \pm 139.107	.119 \pm .001	.383 \pm .004	.112 \pm .001	3893.091 \pm 86.176	.113 \pm .001	.239 \pm .003	.106 \pm .001	1204.211 \pm 5.815	.132 \pm .003	.195 \pm .004	.125 \pm .003	362.146 \pm 27.801
catoni		.126 \pm .000	.531 \pm .004	.124 \pm .000	11966.223 \pm 148.195	.125 \pm .001	.299 \pm .003	.118 \pm .001	3829.806 \pm 85.864	.119 \pm .002	.209 \pm .002	.113 \pm .001	1225.310 \pm 48.090	.134 \pm .004	.202 \pm .005	.127 \pm .003	395.243 \pm 29.182
rivasplata		.123 \pm .000	.458 \pm .003	.120 \pm .000	11209.156 \pm 143.319	.118 \pm .001	.287 \pm .002	.111 \pm .001	3815.804 \pm 85.091	.112 \pm .001	.196 \pm .002	.105 \pm .001	126.956 \pm 49.255	.130 \pm .003	.173 \pm .004	.124 \pm .003	376.904 \pm 27.549
stochastic		—	.189	—	6.200	—	.184	—	7.316	—	.195	—	13.250	—	.189	—	11.851
CIFAR-10		ours	.369 \pm .000	.379 \pm .000	.358 \pm .000	11.657	.351 \pm .000	.369 \pm .000	.352 \pm .000	2.267	.359 \pm .001	.378 \pm .000	.360 \pm .000	2.616	.418 \pm .001	.434 \pm .001	.415 \pm .001
	blanchard	.369 \pm .000	.990 \pm .000	.358 \pm .000	40152.974 \pm 291.721	.351 \pm .000	.809 \pm .003	.351 \pm .000	8753.816 \pm 136.801	.358 \pm .001	.635 \pm .004	.359 \pm .000	2728.436 \pm 73.835	.412 \pm .001	.568 \pm .004	.407 \pm .001	91.026 \pm 44.096
	catoni	.369 \pm .000	.986 \pm .000	.358 \pm .000	24477.984 \pm 223.367	.351 \pm .000	.708 \pm .003	.351 \pm .000	8463.452 \pm 135.001	.357 \pm .001	.578 \pm .002	.357 \pm .000	3401.221 \pm 84.878	.405 \pm .001	.561 \pm .002	.399 \pm .001	1354.100 \pm 51.315
	rivasplata	.369 \pm .000	.868 \pm .001	.358 \pm .000	24424.968 \pm 223.601	.351 \pm .000	.694 \pm .002	.351 \pm .000	8665.339 \pm 136.361	.358 \pm .001	.555 \pm .003	.358 \pm .000	274.651 \pm 74.784	.409 \pm .001	.521 \pm .003	.403 \pm .001	955.211 \pm 44.609
	stochastic	—	.448	—	5.829	—	.439	—	1.134	—	.448	—	1.308	—	.504	—	2.838

Table F.5. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.4.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$				
		$\eta = 10^{-6}$	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.010 \pm .000	.017 \pm .000	.013 \pm .000	.194	.012 \pm .000	.018 \pm .000	.014 \pm .000	.138	.009 \pm .000	.015 \pm .000	.011 \pm .000	.235	.014 \pm .001	.020 \pm .001	.015 \pm .001	1.111	
	blanchard	.010 \pm .000	.028 \pm .001	.013 \pm .000	88.323 \pm 13.740	.012 \pm .000	.021 \pm .001	.014 \pm .000	16.792 \pm 5.702	.009 \pm .000	.014 \pm .001	.011 \pm .000	2.449 \pm 2.313	.014 \pm .001	.019 \pm .001	.016 \pm .001	.244 \pm .765	
	catoni	.010 \pm .000	.026 \pm .001	.013 \pm .000	109.202 \pm 15.634	.012 \pm .000	.023 \pm .002	.014 \pm .000	19.918 \pm 6.526	.009 \pm .000	.015 \pm .001	.011 \pm .000	2.486 \pm 2.362	.014 \pm .001	.019 \pm .001	.016 \pm .001	.298 \pm .762	
	rivasplata	.010 \pm .000	.024 \pm .001	.013 \pm .000	91.872 \pm 14.470	.012 \pm .000	.019 \pm .001	.014 \pm .000	17.002 \pm 5.882	.009 \pm .000	.014 \pm .000	.011 \pm .000	2.529 \pm 2.251	.014 \pm .001	.019 \pm .001	.016 \pm .001	.308 \pm .778	
	stochastic	—	.043	—	.097	—	.044	—	.069	—	.039	—	.117	—	.047	—	.555	
Fashion	ours	.118 \pm .001	.123 \pm .000	.112 \pm .000	.269	.113 \pm .001	.118 \pm .001	.107 \pm .001	.743	.117 \pm .002	.121 \pm .002	.110 \pm .002	2.600	.131 \pm .004	.138 \pm .004	.126 \pm .004	1.229	
	blanchard	.118 \pm .001	.145 \pm .003	.112 \pm .000	82.403 \pm 13.230	.113 \pm .001	.123 \pm .002	.107 \pm .001	16.836 \pm 5.583	.119 \pm .002	.121 \pm .003	.112 \pm .003	2.641 \pm 2.369	.133 \pm .004	.136 \pm .004	.128 \pm .004	.297 \pm .731	
	catoni	.118 \pm .001	.151 \pm .004	.112 \pm .000	109.988 \pm 15.347	.113 \pm .001	.120 \pm .002	.107 \pm .001	19.889 \pm 6.689	.118 \pm .002	.120 \pm .003	.112 \pm .003	2.615 \pm 2.234	.132 \pm .004	.136 \pm .004	.128 \pm .004	.300 \pm .811	
	rivasplata	.118 \pm .001	.137 \pm .002	.112 \pm .000	87.804 \pm 13.640	.113 \pm .001	.120 \pm .002	.107 \pm .001	17.491 \pm 6.144	.118 \pm .002	.121 \pm .003	.112 \pm .003	2.549 \pm 2.175	.133 \pm .005	.137 \pm .004	.128 \pm .004	.322 \pm .794	
	stochastic	—	.174	—	.135	—	.168	—	.372	—	.172	—	1.300	—	.191	—	.615	
CIFAR-10	ours	.334 \pm .000	.346 \pm .000	.328 \pm .000	.025	.322 \pm .000	.331 \pm .000	.313 \pm .000	.050	.323 \pm .001	.334 \pm .000	.316 \pm .000	.160	.333 \pm .001	.341 \pm .001	.323 \pm .001	.461	
	blanchard	.334 \pm .000	.421 \pm .004	.328 \pm .000	269.875 \pm 23.982	.322 \pm .000	.364 \pm .004	.313 \pm .000	83.082 \pm 13.029	.323 \pm .001	.345 \pm .004	.316 \pm .000	21.614 \pm 6.670	.333 \pm .001	.340 \pm .002	.323 \pm .001	3.630 \pm 2.750	
	catoni	.334 \pm .000	.433 \pm .008	.328 \pm .000	27.270 \pm 24.201	.322 \pm .000	.355 \pm .005	.313 \pm .000	84.148 \pm 13.578	.323 \pm .001	.338 \pm .002	.316 \pm .000	22.547 \pm 6.801	.333 \pm .001	.338 \pm .001	.323 \pm .001	3.831 \pm 2.801	
	rivasplata	.334 \pm .000	.394 \pm .003	.328 \pm .000	27.133 \pm 24.109	.322 \pm .000	.351 \pm .003	.313 \pm .000	83.438 \pm 13.033	.323 \pm .001	.339 \pm .002	.316 \pm .000	21.688 \pm 6.718	.333 \pm .001	.339 \pm .002	.323 \pm .001	3.667 \pm 2.757	
	stochastic	—	.414	—	.013	—	.399	—	.025	—	.403	—	.080	—	.409	—	.230	
		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$				
		$\eta = 10^{-4}$	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.010 \pm .000	.019 \pm .000	.013 \pm .000	23.992	.012 \pm .000	.019 \pm .000	.014 \pm .000	7.767	.009 \pm .000	.015 \pm .000	.011 \pm .000	3.165	.012 \pm .001	.019 \pm .001	.013 \pm .001	18.413	
	blanchard	.010 \pm .000	.500 \pm .004	.013 \pm .000	1123.328 \pm 151.115	.012 \pm .000	.236 \pm .004	.014 \pm .000	381.840 \pm 94.218	.009 \pm .000	.096 \pm .003	.011 \pm .000	1184.214 \pm 47.208	.011 \pm .001	.048 \pm .002	.012 \pm .001	363.194 \pm 26.547	
	catoni	.010 \pm .000	.369 \pm .003	.013 \pm .000	1191.598 \pm 154.180	.012 \pm .000	.180 \pm .003	.014 \pm .000	3826.581 \pm 85.362	.009 \pm .000	.070 \pm .002	.011 \pm .000	1217.723 \pm 49.984	.012 \pm .001	.039 \pm .002	.014 \pm .001	384.476 \pm 29.126	
	rivasplata	.010 \pm .000	.316 \pm .003	.013 \pm .000	11557.703 \pm 151.498	.012 \pm .000	.142 \pm .002	.014 \pm .000	3751.391 \pm 84.542	.009 \pm .000	.061 \pm .002	.011 \pm .000	1172.156 \pm 46.933	.010 \pm .001	.035 \pm .001	.012 \pm .001	373.003 \pm 27.844	
	stochastic	—	.045	—	11.996	—	.045	—	3.884	—	.040	—	1.583	—	.045	—	9.207	
Fashion	ours	.118 \pm .000	.127 \pm .000	.112 \pm .000	17.987	.113 \pm .001	.119 \pm .001	.107 \pm .001	6.361	.114 \pm .002	.123 \pm .002	.107 \pm .002	22.582	.125 \pm .003	.137 \pm .003	.122 \pm .003	16.872	
	blanchard	.115 \pm .001	.659 \pm .004	.110 \pm .000	11835.780 \pm 161.816	.110 \pm .001	.395 \pm .004	.104 \pm .000	3828.562 \pm 94.279	.108 \pm .001	.244 \pm .004	.102 \pm .001	1185.882 \pm 5.575	.123 \pm .003	.192 \pm .004	.119 \pm .002	346.265 \pm 27.827	
	catoni	.118 \pm .001	.566 \pm .004	.112 \pm .000	11921.114 \pm 153.739	.113 \pm .001	.304 \pm .003	.107 \pm .000	3822.647 \pm 85.225	.114 \pm .002	.208 \pm .003	.107 \pm .002	1217.879 \pm 52.353	.125 \pm .003	.196 \pm .004	.121 \pm .002	388.473 \pm 29.475	
	rivasplata	.114 \pm .000	.476 \pm .003	.109 \pm .000	11206.239 \pm 149.549	.110 \pm .001	.292 \pm .003	.103 \pm .000	3745.930 \pm 84.367	.106 \pm .001	.197 \pm .003	.101 \pm .001	1229.005 \pm 51.052	.122 \pm .003	.170 \pm .004	.118 \pm .003	361.652 \pm 28.452	
	stochastic	—	.177	—	8.994	—	.169	—	1.129	—	.172	—	1.191	—	.189	—	8.436	
CIFAR-10	ours	.334 \pm .000	.350 \pm .000	.328 \pm .000	12.067	.322 \pm .000	.332 \pm .000	.313 \pm .000	4.172	.323 \pm .001	.336 \pm .000	.316 \pm .000	3.382	.332 \pm .001	.343 \pm .001	.322 \pm .001	6.855	
	blanchard	.334 \pm .000	.977 \pm .001	.328 \pm .000	28565.558 \pm 245.568	.322 \pm .000	.803 \pm .003	.313 \pm .000	8479.553 \pm 126.804	.321 \pm .001	.614 \pm .004	.315 \pm .000	2727.786 \pm 7.572	.327 \pm .001	.487 \pm .004	.317 \pm .001	887.578 \pm 42.449	
	catoni	.334 \pm .000	.983 \pm .000	.328 \pm .000	24136.528 \pm 211.963	.322 \pm .000	.694 \pm .004	.313 \pm .000	7928.671 \pm 122.159	.320 \pm .001	.515 \pm .002	.314 \pm .000	237.703 \pm 65.952	.323 \pm .001	.468 \pm .002	.312 \pm .001	1157.073 \pm 47.283	
	rivasplata	.334 \pm .000	.922 \pm .001	.328 \pm .000	33282.032 \pm 246.654	.322 \pm .000	.680 \pm .003	.312 \pm .000	8493.458 \pm 128.894	.320 \pm .001	.527 \pm .003	.314 \pm .000	2739.108 \pm 7.556	.325 \pm .001	.436 \pm .003	.314 \pm .001	91.066 \pm 43.389	
	stochastic	—	.417	—	6.033	—	.400	—	2.086	—	.403	—	1.691	—	.410	—	3.427	

Table F.6. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the R NYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.5.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$\eta = 10^{-6}$															
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.008 \pm .000	.015 \pm .000	.010 \pm .000	.084	.006 \pm .000	.012 \pm .000	.009 \pm .000	.053	.008 \pm .000	.014 \pm .000	.010 \pm .000	.179	.014 \pm .001	.019 \pm .001	.014 \pm .001	.576
	blanchard	.008 \pm .000	.025 \pm .001	.010 \pm .000	81.167 \pm 12.801	.006 \pm .000	.014 \pm .001	.009 \pm .000	15.518 \pm 5.438	.009 \pm .000	.014 \pm .001	.010 \pm .000	2.140 \pm 2.072	.015 \pm .001	.018 \pm .001	.015 \pm .001	.284 \pm .649
	catoni	.008 \pm .000	.022 \pm .001	.010 \pm .000	104.063 \pm 14.662	.006 \pm .000	.015 \pm .000	.009 \pm .000	17.676 \pm 5.963	.008 \pm .000	.014 \pm .001	.010 \pm .000	2.152 \pm 2.085	.015 \pm .001	.018 \pm .001	.015 \pm .001	.252 \pm .680
	rivasplata	.008 \pm .000	.021 \pm .001	.010 \pm .000	84.581 \pm 13.035	.006 \pm .000	.013 \pm .001	.009 \pm .000	15.545 \pm 5.594	.008 \pm .000	.014 \pm .000	.010 \pm .000	2.185 \pm 1.992	.015 \pm .001	.018 \pm .001	.015 \pm .001	.276 \pm .693
	stochastic	—	.039	—	.042	—	.035	—	.026	—	.038	—	.090	—	.045	—	.288
	Fashion	ours	.106 \pm .000	.113 \pm .000	.101 \pm .000	.133	.104 \pm .001	.110 \pm .000	.099 \pm .000	.327	.108 \pm .002	.112 \pm .001	.101 \pm .001	.903	.120 \pm .004	.127 \pm .003	.115 \pm .003
blanchard		.106 \pm .000	.136 \pm .003	.101 \pm .000	77.573 \pm 12.564	.104 \pm .001	.115 \pm .003	.099 \pm .000	15.278 \pm 5.599	.109 \pm .002	.111 \pm .002	.102 \pm .001	2.153 \pm 2.081	.122 \pm .004	.126 \pm .004	.117 \pm .004	.248 \pm .715
catoni		.106 \pm .000	.145 \pm .005	.101 \pm .000	104.356 \pm 14.712	.104 \pm .001	.112 \pm .002	.099 \pm .000	17.566 \pm 5.996	.109 \pm .002	.110 \pm .001	.102 \pm .001	2.217 \pm 2.084	.122 \pm .004	.125 \pm .004	.117 \pm .004	.262 \pm .699
rivasplata		.106 \pm .000	.127 \pm .002	.101 \pm .000	82.150 \pm 12.955	.104 \pm .001	.112 \pm .001	.099 \pm .000	15.509 \pm 5.629	.109 \pm .002	.111 \pm .001	.102 \pm .001	2.178 \pm 2.060	.122 \pm .004	.126 \pm .004	.117 \pm .004	.264 \pm .704
stochastic		—	.162	—	.066	—	.159	—	.164	—	.162	—	.451	—	.179	—	.434
CIFAR-10		ours	.312 \pm .000	.323 \pm .000	.304 \pm .000	.027	.281 \pm .000	.304 \pm .000	.285 \pm .000	.035	.298 \pm .001	.310 \pm .000	.291 \pm .000	.101	.315 \pm .001	.329 \pm .001	.309 \pm .001
	blanchard	.312 \pm .000	.405 \pm .004	.304 \pm .000	268.149 \pm 22.835	.281 \pm .000	.339 \pm .004	.285 \pm .000	8.690 \pm 12.628	.298 \pm .001	.320 \pm .004	.291 \pm .000	19.648 \pm 6.249	.315 \pm .001	.327 \pm .003	.310 \pm .001	3.213 \pm 2.590
	catoni	.312 \pm .000	.428 \pm .009	.304 \pm .000	269.415 \pm 22.884	.281 \pm .000	.333 \pm .005	.285 \pm .000	83.414 \pm 13.018	.298 \pm .001	.314 \pm .003	.291 \pm .000	2.711 \pm 6.481	.315 \pm .001	.326 \pm .001	.310 \pm .001	3.273 \pm 2.597
	rivasplata	.312 \pm .000	.375 \pm .003	.304 \pm .000	268.589 \pm 22.845	.281 \pm .000	.325 \pm .003	.285 \pm .000	81.532 \pm 12.712	.298 \pm .001	.315 \pm .002	.291 \pm .000	19.813 \pm 6.288	.315 \pm .001	.327 \pm .002	.310 \pm .001	3.233 \pm 2.599
	stochastic	—	.391	—	.013	—	.370	—	.017	—	.377	—	.050	—	.397	—	.184
	MNIST	ours	.008 \pm .000	.017 \pm .000	.010 \pm .000	29.993	.006 \pm .000	.013 \pm .000	.009 \pm .000	3.162	.008 \pm .000	.015 \pm .000	.010 \pm .000	1.418	.013 \pm .001	.019 \pm .001	.013 \pm .001
blanchard		.008 \pm .000	.574 \pm .005	.010 \pm .000	11894.556 \pm 155.958	.006 \pm .000	.256 \pm .004	.009 \pm .000	3826.515 \pm 86.973	.008 \pm .000	.108 \pm .003	.010 \pm .000	1184.777 \pm 48.158	.010 \pm .001	.052 \pm .002	.011 \pm .000	36.865 \pm 28.054
catoni		.008 \pm .000	.396 \pm .003	.010 \pm .000	11986.455 \pm 15.722	.006 \pm .000	.192 \pm .002	.009 \pm .000	3824.971 \pm 85.072	.008 \pm .000	.079 \pm .002	.010 \pm .000	1213.611 \pm 48.751	.013 \pm .001	.042 \pm .002	.014 \pm .001	384.275 \pm 28.556
rivasplata		.008 \pm .000	.362 \pm .003	.010 \pm .000	11905.971 \pm 15.609	.006 \pm .000	.148 \pm .003	.009 \pm .000	377.259 \pm 84.127	.008 \pm .000	.067 \pm .002	.010 \pm .000	118.841 \pm 5.043	.010 \pm .001	.036 \pm .001	.011 \pm .000	369.675 \pm 27.947
stochastic		—	.041	—	14.996	—	.035	—	1.581	—	.039	—	.709	—	.045	—	6.116
Fashion		ours	.106 \pm .000	.114 \pm .000	.101 \pm .000	6.310	.103 \pm .001	.113 \pm .000	.099 \pm .000	9.312	.106 \pm .002	.115 \pm .001	.100 \pm .001	14.924	.115 \pm .003	.126 \pm .003	.110 \pm .002
	blanchard	.105 \pm .000	.674 \pm .004	.101 \pm .000	10795.464 \pm 143.426	.102 \pm .000	.412 \pm .004	.098 \pm .000	3685.940 \pm 82.481	.103 \pm .001	.253 \pm .004	.097 \pm .001	1178.401 \pm 48.359	.113 \pm .002	.186 \pm .004	.108 \pm .002	338.697 \pm 27.104
	catoni	.106 \pm .000	.623 \pm .005	.101 \pm .000	11971.564 \pm 15.589	.104 \pm .001	.321 \pm .004	.099 \pm .000	3825.370 \pm 87.728	.107 \pm .002	.208 \pm .003	.100 \pm .001	1214.976 \pm 48.846	.116 \pm .003	.184 \pm .004	.111 \pm .003	388.197 \pm 27.580
	rivasplata	.105 \pm .000	.503 \pm .003	.100 \pm .000	11139.304 \pm 15.540	.102 \pm .000	.307 \pm .003	.097 \pm .000	381.075 \pm 87.924	.102 \pm .001	.201 \pm .003	.096 \pm .001	1201.832 \pm 48.877	.112 \pm .002	.161 \pm .003	.107 \pm .002	349.146 \pm 27.482
	stochastic	—	.163	—	3.155	—	.161	—	4.656	—	.163	—	7.462	—	.176	—	9.182
	CIFAR-10	ours	.312 \pm .000	.328 \pm .000	.304 \pm .000	12.006	.281 \pm .000	.304 \pm .000	.285 \pm .000	1.802	.297 \pm .001	.311 \pm .000	.291 \pm .000	2.056	.314 \pm .001	.330 \pm .001	.309 \pm .001
blanchard		.312 \pm .000	.990 \pm .000	.304 \pm .000	48007.471 \pm 31.730	.280 \pm .000	.825 \pm .003	.284 \pm .000	8824.774 \pm 134.331	.296 \pm .001	.617 \pm .004	.290 \pm .000	2723.775 \pm 66.832	.309 \pm .001	.490 \pm .004	.303 \pm .001	888.277 \pm 41.530
catoni		.312 \pm .000	.980 \pm .000	.304 \pm .000	21278.808 \pm 207.839	.280 \pm .000	.681 \pm .004	.284 \pm .000	6951.932 \pm 118.540	.296 \pm .001	.496 \pm .003	.290 \pm .000	2145.470 \pm 6.045	.305 \pm .001	.457 \pm .002	.299 \pm .001	103.494 \pm 47.021
rivasplata		.312 \pm .000	.964 \pm .001	.304 \pm .000	42834.626 \pm 284.116	.280 \pm .000	.690 \pm .003	.284 \pm .000	8675.531 \pm 136.658	.296 \pm .001	.521 \pm .003	.290 \pm .000	2718.415 \pm 66.664	.307 \pm .001	.434 \pm .003	.301 \pm .001	921.068 \pm 42.158
stochastic		—	.394	—	6.003	—	.371	—	.901	—	.378	—	1.028	—	.397	—	2.391

Table F.7. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the R NYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.6.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$				
		$\eta = 10^{-6}$	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.008 \pm .000	.014 \pm .000	.010 \pm .000	.040	.007 \pm .000	.014 \pm .000	.009 \pm .000	.068	.008 \pm .000	.013 \pm .000	.009 \pm .000	.092	.008 \pm .000	.014 \pm .001	.009 \pm .000	.128	
	blanchard	.008 \pm .000	.026 \pm .002	.010 \pm .000	75.043 \pm 11.586	.007 \pm .000	.016 \pm .001	.009 \pm .000	13.220 \pm 4.956	.008 \pm .000	.012 \pm .001	.009 \pm .000	1.774 \pm 1.772	.008 \pm .000	.012 \pm .001	.009 \pm .000	.190 \pm .594	
	catoni	.008 \pm .000	.022 \pm .001	.010 \pm .000	96.561 \pm 13.980	.007 \pm .000	.016 \pm .000	.009 \pm .000	15.107 \pm 5.370	.008 \pm .000	.013 \pm .001	.009 \pm .000	1.835 \pm 1.837	.008 \pm .000	.013 \pm .000	.009 \pm .000	.219 \pm .619	
	rivasplata	.008 \pm .000	.021 \pm .001	.010 \pm .000	76.898 \pm 12.301	.007 \pm .000	.014 \pm .001	.009 \pm .000	13.370 \pm 4.931	.008 \pm .000	.013 \pm .000	.009 \pm .000	1.695 \pm 1.741	.008 \pm .000	.013 \pm .001	.009 \pm .000	.183 \pm .580	
	stochastic	—	.038	—	.020	—	.037	—	.034	—	.037	—	.046	—	.037	—	.064	
	Fashion	ours	.109 \pm .000	.115 \pm .000	.102 \pm .000	.128	.114 \pm .001	.117 \pm .001	.104 \pm .001	.436	.101 \pm .001	.108 \pm .001	.096 \pm .001	.452	.110 \pm .003	.116 \pm .003	.103 \pm .003	.438
blanchard		.109 \pm .000	.139 \pm .003	.102 \pm .000	7.878 \pm 11.599	.114 \pm .001	.121 \pm .003	.104 \pm .001	13.041 \pm 5.012	.102 \pm .001	.106 \pm .002	.096 \pm .001	1.840 \pm 1.864	.111 \pm .003	.113 \pm .003	.104 \pm .002	.184 \pm .600	
catoni		.109 \pm .000	.152 \pm .006	.102 \pm .000	96.732 \pm 13.464	.114 \pm .001	.119 \pm .002	.104 \pm .001	15.103 \pm 5.363	.102 \pm .001	.105 \pm .001	.096 \pm .001	1.825 \pm 1.886	.111 \pm .003	.112 \pm .003	.104 \pm .003	.224 \pm .610	
rivasplata		.109 \pm .000	.129 \pm .002	.102 \pm .000	75.029 \pm 11.918	.114 \pm .001	.118 \pm .002	.104 \pm .001	13.495 \pm 5.112	.102 \pm .001	.106 \pm .001	.096 \pm .001	1.798 \pm 1.859	.111 \pm .003	.114 \pm .003	.104 \pm .002	.219 \pm .610	
stochastic		—	.164	—	.064	—	.167	—	.218	—	.157	—	.226	—	.165	—	.219	
CIFAR-10		ours	.277 \pm .000	.297 \pm .000	.276 \pm .000	.021	.288 \pm .000	.307 \pm .000	.286 \pm .000	.027	.273 \pm .001	.284 \pm .000	.263 \pm .000	.079	.281 \pm .001	.302 \pm .001	.281 \pm .001	.227
	blanchard	.277 \pm .000	.386 \pm .005	.276 \pm .000	262.952 \pm 24.385	.288 \pm .000	.346 \pm .005	.286 \pm .000	76.609 \pm 12.923	.273 \pm .001	.293 \pm .004	.263 \pm .000	17.724 \pm 6.241	.281 \pm .001	.299 \pm .002	.281 \pm .001	2.580 \pm 2.299	
	catoni	.277 \pm .000	.398 \pm .001	.276 \pm .000	268.083 \pm 24.567	.288 \pm .000	.343 \pm .007	.286 \pm .000	82.887 \pm 13.493	.273 \pm .001	.287 \pm .003	.263 \pm .000	18.978 \pm 6.437	.281 \pm .001	.297 \pm .001	.281 \pm .001	2.661 \pm 2.317	
	rivasplata	.277 \pm .000	.354 \pm .004	.276 \pm .000	263.581 \pm 24.435	.288 \pm .000	.330 \pm .003	.286 \pm .000	77.488 \pm 12.464	.273 \pm .001	.288 \pm .002	.263 \pm .000	17.704 \pm 5.927	.281 \pm .001	.299 \pm .002	.281 \pm .001	2.619 \pm 2.297	
	stochastic	—	.363	—	.010	—	.374	—	.014	—	.349	—	.040	—	.368	—	.113	
			$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$\eta = 10^{-4}$	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.008 \pm .000	.016 \pm .000	.010 \pm .000	9.520	.007 \pm .000	.014 \pm .000	.009 \pm .000	3.594	.008 \pm .000	.014 \pm .000	.009 \pm .000	1.877	.008 \pm .000	.014 \pm .001	.009 \pm .000	6.589	
	blanchard	.008 \pm .000	.657 \pm .005	.010 \pm .000	1209.158 \pm 157.539	.007 \pm .000	.304 \pm .005	.009 \pm .000	3795.285 \pm 88.141	.008 \pm .000	.124 \pm .004	.009 \pm .000	1183.704 \pm 5.113	.007 \pm .000	.052 \pm .002	.009 \pm .000	347.860 \pm 25.275	
	catoni	.008 \pm .000	.452 \pm .004	.010 \pm .000	12032.708 \pm 157.184	.007 \pm .000	.225 \pm .003	.009 \pm .000	3834.246 \pm 89.809	.008 \pm .000	.093 \pm .003	.009 \pm .000	1225.575 \pm 51.027	.007 \pm .000	.039 \pm .002	.008 \pm .000	39.374 \pm 26.987	
	rivasplata	.008 \pm .000	.423 \pm .004	.010 \pm .000	11943.688 \pm 156.365	.007 \pm .000	.179 \pm .003	.009 \pm .000	3787.407 \pm 87.968	.008 \pm .000	.075 \pm .002	.009 \pm .000	1173.457 \pm 49.846	.007 \pm .000	.035 \pm .002	.008 \pm .000	348.717 \pm 26.495	
	stochastic	—	.039	—	4.760	—	.038	—	1.797	—	.037	—	.938	—	.038	—	3.294	
	Fashion	ours	.109 \pm .000	.119 \pm .000	.102 \pm .000	16.776	.114 \pm .001	.119 \pm .001	.104 \pm .001	7.869	.101 \pm .001	.111 \pm .001	.095 \pm .001	14.224	.109 \pm .002	.116 \pm .002	.101 \pm .002	9.187
blanchard		.108 \pm .000	.743 \pm .004	.101 \pm .000	11048.501 \pm 146.969	.112 \pm .001	.468 \pm .005	.101 \pm .001	3798.865 \pm 87.270	.099 \pm .001	.268 \pm .005	.093 \pm .001	1144.740 \pm 49.199	.106 \pm .002	.183 \pm .004	.099 \pm .002	328.466 \pm 24.435	
catoni		.109 \pm .000	.712 \pm .005	.102 \pm .000	1191.096 \pm 15.212	.114 \pm .001	.367 \pm .005	.104 \pm .001	3831.104 \pm 88.371	.101 \pm .001	.216 \pm .003	.095 \pm .001	1221.392 \pm 5.970	.108 \pm .002	.175 \pm .003	.101 \pm .002	386.528 \pm 26.498	
rivasplata		.108 \pm .000	.557 \pm .003	.101 \pm .000	11148.085 \pm 145.818	.111 \pm .001	.340 \pm .003	.100 \pm .001	3757.976 \pm 83.965	.098 \pm .001	.209 \pm .003	.092 \pm .001	1176.081 \pm 49.829	.106 \pm .002	.156 \pm .003	.098 \pm .002	34.716 \pm 24.874	
stochastic		—	.168	—	8.388	—	.168	—	3.935	—	.159	—	7.112	—	.165	—	4.594	
CIFAR-10		ours	.277 \pm .000	.301 \pm .000	.276 \pm .000	8.466	.288 \pm .000	.308 \pm .000	.286 \pm .000	2.415	.273 \pm .001	.285 \pm .000	.263 \pm .000	2.256	.280 \pm .001	.303 \pm .001	.280 \pm .001	2.747
	blanchard	.277 \pm .000	.990 \pm .000	.276 \pm .000	58878.209 \pm 356.845	.288 \pm .000	.868 \pm .003	.286 \pm .000	8858.838 \pm 134.545	.272 \pm .001	.625 \pm .005	.262 \pm .000	2709.659 \pm 76.197	.278 \pm .001	.480 \pm .005	.277 \pm .001	86.940 \pm 43.864	
	catoni	.277 \pm .000	.974 \pm .000	.276 \pm .000	17581.286 \pm 185.476	.288 \pm .000	.662 \pm .005	.286 \pm .000	5118.582 \pm 105.636	.272 \pm .001	.456 \pm .003	.262 \pm .000	1548.107 \pm 58.565	.277 \pm .001	.426 \pm .002	.274 \pm .001	783.103 \pm 41.593	
	rivasplata	.277 \pm .000	.990 \pm .000	.276 \pm .000	82459.214 \pm 398.763	.288 \pm .000	.733 \pm .003	.286 \pm .000	8674.850 \pm 13.468	.272 \pm .001	.518 \pm .004	.262 \pm .000	2709.173 \pm 77.205	.277 \pm .001	.418 \pm .004	.275 \pm .001	874.307 \pm 44.089	
	stochastic	—	.366	—	4.233	—	.374	—	1.207	—	.350	—	1.128	—	.369	—	1.374	

Table F.8. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.7.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.011 \pm .000	.019 \pm .000	.013 \pm .000	.047	.010 \pm .000	.018 \pm .000	.012 \pm .000	.125	.010 \pm .000	.017 \pm .000	.012 \pm .000	.116	.010 \pm .001	.018 \pm .001	.012 \pm .001	.132
	blanchard	.011 \pm .000	.032 \pm .002	.013 \pm .000	65.017 \pm 11.099	.010 \pm .000	.019 \pm .001	.012 \pm .000	1.819 \pm 4.995	.010 \pm .000	.016 \pm .001	.012 \pm .000	1.551 \pm 1.635	.010 \pm .001	.016 \pm .001	.012 \pm .001	1.15 \pm .560
	catoni	.011 \pm .000	.028 \pm .001	.013 \pm .000	84.529 \pm 13.023	.010 \pm .000	.021 \pm .000	.012 \pm .000	11.910 \pm 5.053	.010 \pm .000	.017 \pm .001	.012 \pm .000	1.228 \pm 1.637	.010 \pm .001	.017 \pm .001	.012 \pm .001	.173 \pm .512
	rivasplata	.011 \pm .000	.026 \pm .001	.013 \pm .000	68.055 \pm 11.606	.010 \pm .000	.018 \pm .001	.012 \pm .000	1.637 \pm 4.962	.010 \pm .000	.016 \pm .000	.012 \pm .000	1.408 \pm 1.639	.010 \pm .001	.016 \pm .001	.012 \pm .001	.160 \pm .529
	stochastic	—	.044	—	.023	—	.043	—	.062	—	.042	—	.058	—	.043	—	.066
	Fashion	ours	.099 \pm .000	.112 \pm .000	.098 \pm .000	.067	.107 \pm .001	.115 \pm .001	.100 \pm .001	.542	.098 \pm .002	.107 \pm .001	.093 \pm .001	.353	.108 \pm .003	.117 \pm .002	.102 \pm .002
blanchard		.099 \pm .000	.138 \pm .004	.098 \pm .000	61.733 \pm 1.862	.107 \pm .001	.119 \pm .003	.101 \pm .001	1.651 \pm 4.230	.099 \pm .001	.104 \pm .002	.094 \pm .001	1.342 \pm 1.664	.108 \pm .003	.113 \pm .003	.103 \pm .002	.143 \pm .534
catoni		.099 \pm .000	.155 \pm .007	.098 \pm .000	83.929 \pm 12.212	.107 \pm .001	.116 \pm .003	.101 \pm .001	11.543 \pm 4.870	.099 \pm .002	.103 \pm .002	.094 \pm .001	1.437 \pm 1.594	.108 \pm .003	.112 \pm .003	.103 \pm .002	.153 \pm .545
rivasplata		.099 \pm .000	.128 \pm .002	.098 \pm .000	65.737 \pm 11.733	.107 \pm .001	.116 \pm .002	.101 \pm .001	1.958 \pm 4.794	.099 \pm .002	.105 \pm .002	.094 \pm .001	1.491 \pm 1.618	.108 \pm .003	.114 \pm .003	.103 \pm .002	.155 \pm .546
stochastic		—	.161	—	.034	—	.164	—	.271	—	.155	—	.177	—	.166	—	.156
CIFAR-10		ours	.277 \pm .000	.296 \pm .000	.272 \pm .000	.016	.266 \pm .000	.281 \pm .000	.257 \pm .000	.022	.253 \pm .001	.272 \pm .000	.248 \pm .000	.069	.236 \pm .001	.258 \pm .001	.235 \pm .001
	blanchard	.277 \pm .000	.399 \pm .006	.272 \pm .000	257.371 \pm 23.327	.266 \pm .000	.322 \pm .005	.257 \pm .000	7.190 \pm 11.685	.253 \pm .001	.281 \pm .005	.248 \pm .000	15.214 \pm 5.838	.236 \pm .001	.255 \pm .002	.235 \pm .001	2.223 \pm 2.016
	catoni	.277 \pm .000	.399 \pm .002	.272 \pm .000	269.048 \pm 23.489	.266 \pm .000	.328 \pm .009	.257 \pm .000	81.217 \pm 13.476	.253 \pm .001	.275 \pm .004	.248 \pm .000	16.576 \pm 6.087	.236 \pm .001	.253 \pm .002	.235 \pm .001	2.248 \pm 2.082
	rivasplata	.277 \pm .000	.362 \pm .004	.272 \pm .000	258.993 \pm 23.725	.266 \pm .000	.305 \pm .004	.257 \pm .000	72.737 \pm 12.750	.253 \pm .001	.275 \pm .003	.248 \pm .000	15.342 \pm 5.948	.236 \pm .001	.255 \pm .002	.235 \pm .001	2.220 \pm 2.180
	stochastic	—	.362	—	.008	—	.345	—	.011	—	.336	—	.034	—	.322	—	.059
			$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$		
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.011 \pm .000	.025 \pm .000	.013 \pm .000	45.094	.010 \pm .000	.019 \pm .000	.012 \pm .000	7.479	.010 \pm .000	.018 \pm .000	.012 \pm .000	5.269	.010 \pm .000	.018 \pm .001	.011 \pm .001	6.510
	blanchard	.011 \pm .000	.737 \pm .004	.013 \pm .000	11285.050 \pm 147.363	.010 \pm .000	.381 \pm .006	.012 \pm .000	3785.071 \pm 85.889	.010 \pm .000	.160 \pm .005	.011 \pm .000	1181.043 \pm 46.219	.009 \pm .000	.067 \pm .003	.011 \pm .001	34.267 \pm 26.244
	catoni	.011 \pm .000	.547 \pm .004	.013 \pm .000	11965.668 \pm 153.481	.010 \pm .000	.283 \pm .004	.012 \pm .000	3811.642 \pm 88.111	.010 \pm .000	.120 \pm .004	.011 \pm .000	1212.373 \pm 48.835	.009 \pm .000	.050 \pm .002	.010 \pm .000	383.387 \pm 27.059
	rivasplata	.011 \pm .000	.509 \pm .004	.013 \pm .000	11555.623 \pm 15.287	.010 \pm .000	.226 \pm .004	.012 \pm .000	3695.054 \pm 9.289	.009 \pm .000	.096 \pm .003	.011 \pm .000	1171.892 \pm 47.812	.009 \pm .000	.044 \pm .002	.010 \pm .000	343.025 \pm 25.804
	stochastic	—	.050	—	22.547	—	.044	—	3.740	—	.043	—	2.634	—	.043	—	3.255
	Fashion	ours	.099 \pm .000	.116 \pm .000	.098 \pm .000	11.922	.107 \pm .001	.117 \pm .001	.101 \pm .001	6.556	.097 \pm .001	.109 \pm .001	.092 \pm .001	9.235	.105 \pm .002	.118 \pm .002	.100 \pm .002
blanchard		.098 \pm .000	.795 \pm .004	.098 \pm .000	10179.790 \pm 138.889	.101 \pm .001	.524 \pm .006	.096 \pm .001	3752.748 \pm 9.952	.095 \pm .001	.291 \pm .005	.090 \pm .001	1091.018 \pm 47.577	.104 \pm .002	.195 \pm .005	.098 \pm .002	309.857 \pm 24.422
catoni		.099 \pm .000	.808 \pm .002	.098 \pm .000	11999.071 \pm 158.418	.107 \pm .001	.425 \pm .006	.100 \pm .001	3817.800 \pm 91.674	.098 \pm .001	.235 \pm .004	.093 \pm .001	1216.042 \pm 5.641	.106 \pm .002	.182 \pm .004	.101 \pm .002	376.493 \pm 27.018
rivasplata		.098 \pm .000	.619 \pm .004	.097 \pm .000	10768.160 \pm 146.634	.099 \pm .001	.369 \pm .004	.094 \pm .001	3565.270 \pm 88.164	.094 \pm .001	.224 \pm .004	.089 \pm .001	1137.876 \pm 48.421	.103 \pm .002	.164 \pm .003	.097 \pm .002	318.512 \pm 24.741
stochastic		—	.164	—	5.961	—	.166	—	3.278	—	.166	—	4.618	—	.166	—	5.181
CIFAR-10		ours	.277 \pm .000	.303 \pm .000	.272 \pm .000	12.803	.266 \pm .000	.282 \pm .000	.257 \pm .000	2.312	.253 \pm .001	.272 \pm .000	.248 \pm .000	1.641	.236 \pm .001	.259 \pm .001	.235 \pm .001
	blanchard	.277 \pm .000	.990 \pm .000	.272 \pm .000	2577.092 \pm 236.075	.266 \pm .000	.901 \pm .003	.257 \pm .000	8788.732 \pm 134.680	.253 \pm .001	.662 \pm .005	.247 \pm .000	2683.054 \pm 73.139	.235 \pm .001	.464 \pm .006	.233 \pm .001	85.586 \pm 41.917
	catoni	.277 \pm .000	1.000 \pm .000	.272 \pm .000	177807.417 \pm 546.892	.266 \pm .000	.601 \pm .005	.257 \pm .000	331.757 \pm 83.561	.253 \pm .001	.416 \pm .003	.247 \pm .000	85.973 \pm 4.961	.234 \pm .001	.369 \pm .003	.233 \pm .001	485.863 \pm 31.335
	rivasplata	.277 \pm .000	.990 \pm .000	.272 \pm .000	48522.489 \pm 309.735	.266 \pm .000	.762 \pm .003	.257 \pm .000	850.968 \pm 131.507	.252 \pm .001	.542 \pm .004	.247 \pm .000	2696.074 \pm 73.062	.234 \pm .001	.393 \pm .004	.232 \pm .001	858.936 \pm 41.972
	stochastic	—	.366	—	6.401	—	.346	—	1.156	—	.336	—	.821	—	.322	—	.965

Table F.9. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.8.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$			
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.011 \pm .000	.020 \pm .000	.013 \pm .000	.064	.008 \pm .000	.017 \pm .000	.010 \pm .000	.050	.011 \pm .000	.018 \pm .000	.011 \pm .000	.112	.010 \pm .001	.016 \pm .001	.009 \pm .001	.073
	blanchard	.011 \pm .000	.034 \pm .003	.013 \pm .000	49.248 \pm 1.541	.008 \pm .000	.018 \pm .001	.010 \pm .000	8.031 \pm 3.654	.011 \pm .000	.016 \pm .001	.011 \pm .000	810 \pm 1.248	.010 \pm .001	.014 \pm .001	.010 \pm .001	102 \pm 448
	catoni	.011 \pm .000	.030 \pm .002	.013 \pm .000	66.244 \pm 11.961	.008 \pm .000	.018 \pm .001	.010 \pm .000	8.685 \pm 3.987	.011 \pm .000	.019 \pm .001	.011 \pm .000	1.011 \pm 1.283	.010 \pm .001	.016 \pm .001	.010 \pm .001	131 \pm 422
	rivasplata	.011 \pm .000	.028 \pm .002	.013 \pm .000	5.344 \pm 1.600	.008 \pm .000	.017 \pm .001	.010 \pm .000	7.757 \pm 4.187	.011 \pm .000	.017 \pm .001	.011 \pm .000	861 \pm 1.361	.010 \pm .001	.014 \pm .001	.010 \pm .001	090 \pm 460
	stochastic	—	.046	—	.032	—	.041	—	.025	—	.043	—	.056	—	.040	—	.037
	Fashion	ours	.103 \pm .000	.117 \pm .000	.099 \pm .000	.068	.098 \pm .001	.114 \pm .001	.096 \pm .001	.178	.104 \pm .001	.117 \pm .002	.099 \pm .002	.587	.107 \pm .004	.119 \pm .004	.101 \pm .003
blanchard		.103 \pm .000	.144 \pm .004	.099 \pm .000	5.069 \pm 9.537	.098 \pm .001	.116 \pm .003	.096 \pm .001	8.105 \pm 3.874	.104 \pm .001	.113 \pm .002	.100 \pm .002	990 \pm 1.435	.109 \pm .004	.114 \pm .004	.102 \pm .004	102 \pm 444
catoni		.103 \pm .000	.168 \pm .009	.099 \pm .000	66.761 \pm 1.939	.098 \pm .001	.115 \pm .004	.096 \pm .001	8.698 \pm 3.974	.104 \pm .001	.112 \pm .002	.100 \pm .002	934 \pm 1.413	.109 \pm .004	.113 \pm .004	.102 \pm .004	100 \pm 457
rivasplata		.103 \pm .000	.132 \pm .003	.099 \pm .000	52.096 \pm 1.745	.098 \pm .001	.113 \pm .002	.096 \pm .001	7.820 \pm 4.154	.104 \pm .001	.114 \pm .002	.100 \pm .002	939 \pm 1.417	.108 \pm .004	.115 \pm .004	.102 \pm .004	100 \pm 464
stochastic		—	.165	—	.034	—	.162	—	.089	—	.166	—	.294	—	.168	—	.164
CIFAR-10		ours	.249 \pm .000	.265 \pm .000	.237 \pm .000	.014	.247 \pm .000	.271 \pm .000	.243 \pm .000	.018	.259 \pm .001	.281 \pm .001	.252 \pm .001	.055	.249 \pm .001	.274 \pm .001	.245 \pm .001
	blanchard	.249 \pm .000	.384 \pm .007	.237 \pm .000	24.108 \pm 22.114	.247 \pm .000	.316 \pm .006	.243 \pm .000	59.096 \pm 1.459	.259 \pm .001	.289 \pm .006	.252 \pm .001	11.804 \pm 5.001	.249 \pm .001	.269 \pm .003	.245 \pm .001	1.578 \pm 1.705
	catoni	.249 \pm .000	.368 \pm .003	.237 \pm .000	27.284 \pm 23.618	.247 \pm .000	.337 \pm .012	.243 \pm .000	73.833 \pm 11.692	.259 \pm .001	.285 \pm .005	.252 \pm .001	12.808 \pm 5.089	.249 \pm .001	.266 \pm .002	.245 \pm .001	1.635 \pm 1.773
	rivasplata	.249 \pm .000	.341 \pm .005	.237 \pm .000	244.258 \pm 22.339	.247 \pm .000	.298 \pm .004	.243 \pm .000	61.907 \pm 11.135	.259 \pm .001	.283 \pm .003	.252 \pm .001	11.818 \pm 4.923	.249 \pm .001	.269 \pm .002	.245 \pm .001	1.629 \pm 1.763
	stochastic	—	.328	—	.007	—	.334	—	.009	—	.344	—	.028	—	.337	—	.036
			$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$		
		$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.011 \pm .000	.030 \pm .000	.013 \pm .000	53.875	.008 \pm .000	.018 \pm .000	.010 \pm .000	4.369	.011 \pm .000	.019 \pm .000	.011 \pm .000	5.063	.009 \pm .001	.016 \pm .001	.009 \pm .001	4.854
	blanchard	.011 \pm .000	.828 \pm .004	.013 \pm .000	10014.066 \pm 14.769	.008 \pm .000	.491 \pm .007	.010 \pm .000	3707.758 \pm 86.461	.011 \pm .000	.211 \pm .007	.011 \pm .000	1151.660 \pm 47.238	.009 \pm .001	.076 \pm .005	.009 \pm .001	303.402 \pm 25.072
	catoni	.011 \pm .000	.684 \pm .004	.013 \pm .000	12238.359 \pm 158.595	.008 \pm .000	.343 \pm .005	.010 \pm .000	3834.114 \pm 88.516	.011 \pm .000	.168 \pm .006	.011 \pm .000	121.777 \pm 48.460	.009 \pm .001	.060 \pm .003	.008 \pm .001	356.740 \pm 25.649
	rivasplata	.011 \pm .000	.662 \pm .005	.013 \pm .000	1207.265 \pm 161.842	.008 \pm .000	.305 \pm .005	.010 \pm .000	3785.930 \pm 87.976	.011 \pm .000	.125 \pm .004	.010 \pm .000	1141.437 \pm 46.910	.009 \pm .001	.048 \pm .002	.008 \pm .000	305.573 \pm 23.629
	stochastic	—	.055	—	26.937	—	.042	—	2.185	—	.044	—	2.532	—	.040	—	2.427
	Fashion	ours	.102 \pm .000	.121 \pm .000	.099 \pm .000	1.120	.098 \pm .001	.115 \pm .001	.096 \pm .001	3.956	.102 \pm .001	.118 \pm .002	.098 \pm .001	7.830	.103 \pm .003	.118 \pm .003	.097 \pm .002
blanchard		.101 \pm .000	.990 \pm .000	.098 \pm .000	27936.970 \pm 235.840	.096 \pm .001	.585 \pm .007	.094 \pm .001	321.105 \pm 81.006	.098 \pm .001	.348 \pm .007	.094 \pm .001	1045.641 \pm 44.087	.101 \pm .002	.208 \pm .006	.095 \pm .002	273.641 \pm 24.046
catoni		.103 \pm .000	.865 \pm .002	.099 \pm .000	12143.837 \pm 161.857	.098 \pm .001	.536 \pm .008	.096 \pm .001	3802.871 \pm 87.750	.103 \pm .001	.286 \pm .006	.098 \pm .001	1202.907 \pm 47.928	.105 \pm .003	.191 \pm .005	.098 \pm .003	354.246 \pm 25.507
rivasplata		.102 \pm .000	.746 \pm .004	.098 \pm .000	11305.448 \pm 149.693	.096 \pm .001	.438 \pm .005	.093 \pm .001	3458.977 \pm 83.715	.097 \pm .001	.264 \pm .004	.093 \pm .001	1101.567 \pm 44.816	.099 \pm .002	.172 \pm .004	.094 \pm .002	285.588 \pm 24.451
stochastic		—	.168	—	5.060	—	.163	—	1.978	—	.166	—	3.915	—	.166	—	4.399
CIFAR-10		ours	.249 \pm .000	.274 \pm .000	.237 \pm .000	14.083	.247 \pm .000	.273 \pm .000	.243 \pm .000	1.770	.259 \pm .001	.282 \pm .001	.252 \pm .001	1.098	.248 \pm .001	.275 \pm .001	.245 \pm .001
	blanchard	.249 \pm .000	.990 \pm .000	.237 \pm .000	26575.507 \pm 218.278	.247 \pm .000	.925 \pm .002	.243 \pm .000	7135.143 \pm 117.030	.259 \pm .001	.739 \pm .006	.251 \pm .001	2581.211 \pm 74.799	.247 \pm .001	.526 \pm .007	.243 \pm .001	831.790 \pm 4.592
	catoni	.249 \pm .000	1.000 \pm .000	.237 \pm .000	154168.585 \pm 539.590	.247 \pm .000	.677 \pm .008	.243 \pm .000	3148.174 \pm 83.069	.259 \pm .001	.549 \pm .006	.252 \pm .001	1735.530 \pm 57.888	.248 \pm .001	.425 \pm .005	.244 \pm .001	675.780 \pm 38.306
	rivasplata	.249 \pm .000	.990 \pm .000	.237 \pm .000	35062.089 \pm 246.257	.247 \pm .000	.824 \pm .003	.243 \pm .000	8092.236 \pm 125.162	.259 \pm .001	.610 \pm .005	.251 \pm .001	2652.857 \pm 75.369	.247 \pm .001	.441 \pm .005	.242 \pm .001	84.056 \pm 4.952
	stochastic	—	.334	—	7.041	—	.335	—	.885	—	.345	—	.549	—	.337	—	.731

Table F.10. Comparison of ours, rivasplata, blanchard and catoni based on the disintegrated bounds, and stochastic based on the randomized bounds learned with two learning rates $\eta \in \{10^{-4}, 10^{-6}\}$ and different variances $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$. We report the test risk ($R_T(h)$), the bound value (Bnd), the empirical risk ($R_S(h)$), and the divergence (Div) associated with each bound (the RÉNYI divergence for ours, the KL divergence for stochastic, and the disintegrated KL divergence for rivasplata, blanchard and catoni). More precisely, we report the mean \pm the standard deviation for 400 neural networks sampled from ρ_S for ours, rivasplata, blanchard, and catoni. We consider, in this table, that the split ratio is 0.9.

		$\sigma^2 = 10^{-6}$				$\sigma^2 = 10^{-5}$				$\sigma^2 = 10^{-4}$				$\sigma^2 = 10^{-3}$				
		$\eta = 10^{-6}$	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div	$R_T(h)$	Bnd	$R_S(h)$	Div
MNIST	ours	.008 \pm .000	.018 \pm .000	.008 \pm .000	.029	.011 \pm .000	.020 \pm .000	.010 \pm .000	.052	.009 \pm .000	.018 \pm .001	.009 \pm .000	.059	.008 \pm .000	.019 \pm .001	.009 \pm .001	.023	
	blanchard	.008 \pm .000	.033 \pm .004	.009 \pm .000	35.446 \pm 8.610	.011 \pm .000	.020 \pm .002	.010 \pm .000	4.933 \pm 2.958	.009 \pm .000	.015 \pm .001	.009 \pm .000	.490 \pm .960	.008 \pm .001	.016 \pm .001	.009 \pm .001	.059 \pm .299	
	catoni	.008 \pm .000	.026 \pm .002	.009 \pm .000	41.267 \pm 9.234	.011 \pm .000	.019 \pm .001	.010 \pm .000	4.564 \pm 3.263	.009 \pm .000	.016 \pm .001	.009 \pm .000	.581 \pm .989	.008 \pm .001	.017 \pm .001	.009 \pm .001	.078 \pm .320	
	rivasplata	.008 \pm .000	.025 \pm .002	.009 \pm .000	35.856 \pm 8.648	.011 \pm .000	.019 \pm .001	.010 \pm .000	4.620 \pm 2.983	.009 \pm .000	.015 \pm .001	.009 \pm .000	.448 \pm 1.045	.008 \pm .000	.016 \pm .001	.009 \pm .001	.041 \pm .330	
	stochastic	—	.041	—	.014	—	.045	—	.026	—	.042	—	.030	—	.043	—	.012	
	ours	.094 \pm .000	.113 \pm .000	.089 \pm .000	.029	.091 \pm .001	.119 \pm .001	.095 \pm .001	.107	.092 \pm .002	.113 \pm .001	.089 \pm .001	.097	.103 \pm .003	.124 \pm .003	.099 \pm .003	.045	
blanchard	.094 \pm .000	.140 \pm .006	.089 \pm .000	32.563 \pm 8.007	.091 \pm .001	.119 \pm .004	.095 \pm .001	4.567 \pm 2.912	.092 \pm .002	.106 \pm .002	.089 \pm .001	.468 \pm 1.101	.104 \pm .003	.116 \pm .003	.099 \pm .003	.063 \pm .300		
catoni	.094 \pm .000	.146 \pm .002	.089 \pm .000	4.355 \pm 9.121	.091 \pm .001	.120 \pm .005	.095 \pm .001	4.895 \pm 3.064	.092 \pm .002	.106 \pm .002	.089 \pm .001	.473 \pm 1.052	.103 \pm .003	.117 \pm .003	.099 \pm .003	.079 \pm .319		
rivasplata	.094 \pm .000	.127 \pm .004	.089 \pm .000	33.175 \pm 8.710	.091 \pm .001	.117 \pm .002	.095 \pm .001	4.774 \pm 3.003	.092 \pm .002	.107 \pm .002	.089 \pm .001	.479 \pm .924	.103 \pm .003	.118 \pm .003	.099 \pm .002	.045 \pm .330		
stochastic	—	.159	—	.015	—	.166	—	.053	—	.159	—	.048	—	.172	—	.023		
CIFAR-10	ours	.231 \pm .000	.268 \pm .000	.228 \pm .000	.011	.235 \pm .000	.267 \pm .000	.227 \pm .000	.009	.218 \pm .001	.253 \pm .001	.214 \pm .001	.024	.231 \pm .001	.264 \pm .002	.224 \pm .002	.036	
	blanchard	.231 \pm .000	.418 \pm .010	.228 \pm .000	193.922 \pm 19.216	.235 \pm .000	.312 \pm .009	.227 \pm .000	39.705 \pm 8.929	.218 \pm .000	.256 \pm .007	.214 \pm .001	6.919 \pm 3.722	.231 \pm .001	.255 \pm .003	.224 \pm .002	.878 \pm 1.248	
	catoni	.231 \pm .000	.388 \pm .005	.228 \pm .000	255.538 \pm 22.306	.235 \pm .000	.337 \pm .003	.227 \pm .000	53.736 \pm 1.302	.218 \pm .000	.257 \pm .007	.214 \pm .001	7.060 \pm 3.626	.231 \pm .001	.255 \pm .003	.224 \pm .002	.857 \pm 1.264	
	rivasplata	.231 \pm .000	.364 \pm .007	.228 \pm .000	202.026 \pm 19.688	.235 \pm .000	.293 \pm .006	.227 \pm .000	42.458 \pm 9.250	.218 \pm .001	.251 \pm .004	.214 \pm .001	6.780 \pm 3.575	.231 \pm .001	.256 \pm .002	.224 \pm .002	.854 \pm 1.275	
	stochastic	—	.328	—	.005	—	.327	—	.005	—	.312	—	.012	—	.324	—	.018	
	ours	.094 \pm .000	.115 \pm .000	.089 \pm .000	2.501	.091 \pm .001	.121 \pm .001	.095 \pm .001	2.925	.092 \pm .002	.114 \pm .001	.088 \pm .001	3.069	.102 \pm .002	.125 \pm .003	.098 \pm .002	3.159	
blanchard	.094 \pm .000	.990 \pm .000	.089 \pm .000	19455.864 \pm 19.460	.089 \pm .001	.792 \pm .007	.093 \pm .001	3402.546 \pm 86.590	.090 \pm .001	.461 \pm .010	.087 \pm .001	1002.861 \pm 44.393	.102 \pm .002	.244 \pm .009	.098 \pm .002	206.177 \pm 2.051		
catoni	.094 \pm .000	1.000 \pm .000	.089 \pm .000	60888.029 \pm 346.501	.091 \pm .001	.813 \pm .012	.095 \pm .001	3756.375 \pm 9.419	.092 \pm .002	.390 \pm .010	.089 \pm .001	1161.884 \pm 52.073	.103 \pm .003	.215 \pm .007	.099 \pm .002	277.284 \pm 25.479		
rivasplata	.094 \pm .000	.990 \pm .000	.089 \pm .000	27137.315 \pm 227.934	.088 \pm .001	.597 \pm .007	.093 \pm .001	3371.321 \pm 86.352	.090 \pm .001	.331 \pm .007	.086 \pm .001	1003.481 \pm 48.362	.101 \pm .002	.195 \pm .006	.097 \pm .002	207.442 \pm 21.896		
stochastic	—	.160	—	1.250	—	.167	—	1.463	—	.160	—	1.535	—	.172	—	1.579		
CIFAR-10	ours	.231 \pm .000	.279 \pm .000	.228 \pm .000	12.925	.235 \pm .000	.268 \pm .000	.227 \pm .000	1.371	.218 \pm .001	.254 \pm .001	.214 \pm .001	.715	.231 \pm .001	.264 \pm .002	.224 \pm .002	1.019	
	blanchard	.231 \pm .000	.990 \pm .000	.228 \pm .000	26032.808 \pm 222.475	.235 \pm .000	.986 \pm .001	.227 \pm .000	6875.633 \pm 112.137	.217 \pm .000	.831 \pm .006	.214 \pm .001	2292.053 \pm 68.347	.230 \pm .001	.606 \pm .010	.222 \pm .001	76.644 \pm 39.246	
	catoni	.231 \pm .000	1.000 \pm .000	.228 \pm .000	17684.651 \pm 576.711	.235 \pm .000	.980 \pm .000	.227 \pm .000	8265.727 \pm 123.941	.218 \pm .000	.834 \pm .011	.214 \pm .001	2664.069 \pm 73.915	.231 \pm .001	.517 \pm .009	.224 \pm .002	85.593 \pm 41.022	
	rivasplata	.231 \pm .000	.988 \pm .001	.228 \pm .000	14284.846 \pm 169.166	.235 \pm .000	.919 \pm .002	.227 \pm .000	7121.350 \pm 114.645	.217 \pm .000	.699 \pm .006	.213 \pm .001	2502.412 \pm 68.728	.229 \pm .001	.494 \pm .007	.221 \pm .001	776.237 \pm 39.540	
	stochastic	—	.335	—	6.462	—	.328	—	.685	—	.313	—	.358	—	.324	—	.510	

Table F.11. Comparison of the bound values before performing Step 2) of our Training Method for ours, rivasplata, blanchard and catoni. More precisely, for each split and each variance $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$, we report the mean \pm the standard deviation (for 400 neural networks sampled from π) of the test risk ($R_T(h)$), the empirical risk ($R_S(h)$), and the value of the bounds of Corollaries 6.4.1 and 6.4.2. We consider in this table that the dataset is MNIST.

	Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)		Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)
$9 \cdot 10^{-6} = \sigma^2$.0	.901 \pm .002	.901 \pm .002	.908 \pm .002	.906 \pm .002	.905 \pm .002	.906 \pm .002	$1 \cdot 10^{-3} = \sigma^2$.0	.898 \pm .017	.898 \pm .017	.905 \pm .016	.903 \pm .016	.902 \pm .016	.903 \pm .016
	.1	.035 \pm .000	.039 \pm .000	.045 \pm .000	.043 \pm .000	.043 \pm .000	.042 \pm .000		.1	.035 \pm .003	.039 \pm .002	.045 \pm .002	.044 \pm .002	.043 \pm .002	.043 \pm .002
	.2	.016 \pm .000	.019 \pm .000	.023 \pm .000	.022 \pm .000	.022 \pm .000	.022 \pm .000		.2	.015 \pm .001	.016 \pm .001	.020 \pm .001	.019 \pm .001	.019 \pm .001	.019 \pm .001
	.3	.012 \pm .000	.013 \pm .000	.017 \pm .000	.016 \pm .000	.015 \pm .000	.015 \pm .000		.3	.012 \pm .000	.016 \pm .000	.020 \pm .001	.019 \pm .001	.019 \pm .001	.019 \pm .001
	.4	.010 \pm .000	.013 \pm .000	.017 \pm .000	.016 \pm .000	.016 \pm .000	.016 \pm .000		.4	.009 \pm .000	.011 \pm .000	.015 \pm .000	.014 \pm .000	.014 \pm .000	.014 \pm .000
	.5	.008 \pm .000	.010 \pm .000	.015 \pm .000	.013 \pm .000	.013 \pm .000	.014 \pm .000		.5	.008 \pm .000	.010 \pm .000	.015 \pm .000	.013 \pm .000	.013 \pm .000	.014 \pm .000
	.6	.008 \pm .000	.010 \pm .000	.014 \pm .000	.013 \pm .000	.013 \pm .000	.014 \pm .000		.6	.008 \pm .000	.009 \pm .000	.013 \pm .000	.012 \pm .000	.012 \pm .000	.013 \pm .000
	.7	.011 \pm .000	.013 \pm .000	.019 \pm .000	.017 \pm .000	.017 \pm .000	.018 \pm .000		.7	.010 \pm .000	.012 \pm .000	.017 \pm .000	.016 \pm .000	.015 \pm .000	.016 \pm .000
	.8	.011 \pm .000	.013 \pm .000	.020 \pm .000	.018 \pm .000	.018 \pm .000	.020 \pm .000		.8	.011 \pm .000	.011 \pm .000	.018 \pm .000	.016 \pm .000	.016 \pm .000	.018 \pm .000
.9	.008 \pm .000	.009 \pm .000	.018 \pm .000	.015 \pm .000	.014 \pm .000	.015 \pm .000	.9	.009 \pm .000	.009 \pm .000	.018 \pm .001	.015 \pm .001	.015 \pm .001	.016 \pm .001		
$9 \cdot 10^{-5} = \sigma^2$.0	.897 \pm .013	.897 \pm .012	.904 \pm .012	.902 \pm .012	.902 \pm .012	.903 \pm .012	$9 \cdot 10^{-3} = \sigma^2$.0	.903 \pm .014	.902 \pm .014	.909 \pm .013	.907 \pm .013	.907 \pm .013	.907 \pm .013
	.1	.024 \pm .000	.030 \pm .001	.035 \pm .001	.034 \pm .001	.033 \pm .001	.033 \pm .001		.1	.041 \pm .005	.045 \pm .005	.050 \pm .005	.049 \pm .005	.048 \pm .005	.048 \pm .005
	.2	.015 \pm .000	.019 \pm .000	.023 \pm .000	.022 \pm .000	.021 \pm .000	.021 \pm .000		.2	.020 \pm .002	.022 \pm .002	.026 \pm .002	.025 \pm .002	.025 \pm .002	.024 \pm .002
	.3	.009 \pm .000	.011 \pm .000	.015 \pm .000	.014 \pm .000	.013 \pm .000	.013 \pm .000		.3	.014 \pm .001	.015 \pm .001	.019 \pm .001	.018 \pm .001	.018 \pm .001	.018 \pm .001
	.4	.012 \pm .000	.014 \pm .000	.018 \pm .000	.017 \pm .000	.017 \pm .000	.017 \pm .000		.4	.015 \pm .001	.016 \pm .001	.021 \pm .001	.020 \pm .001	.019 \pm .001	.019 \pm .001
	.5	.006 \pm .000	.009 \pm .000	.012 \pm .000	.011 \pm .000	.011 \pm .000	.012 \pm .000		.5	.015 \pm .001	.015 \pm .001	.020 \pm .001	.019 \pm .001	.018 \pm .001	.018 \pm .001
	.6	.007 \pm .000	.009 \pm .000	.014 \pm .000	.013 \pm .000	.012 \pm .000	.013 \pm .000		.6	.008 \pm .000	.010 \pm .000	.014 \pm .001	.013 \pm .001	.012 \pm .000	.013 \pm .000
	.7	.010 \pm .000	.012 \pm .000	.018 \pm .000	.016 \pm .000	.016 \pm .000	.017 \pm .000		.7	.010 \pm .001	.012 \pm .001	.018 \pm .001	.016 \pm .001	.016 \pm .001	.017 \pm .001
	.8	.008 \pm .000	.010 \pm .000	.017 \pm .000	.015 \pm .000	.014 \pm .000	.017 \pm .000		.8	.010 \pm .001	.010 \pm .001	.016 \pm .001	.014 \pm .001	.014 \pm .001	.016 \pm .001
.9	.011 \pm .000	.010 \pm .000	.020 \pm .000	.017 \pm .000	.017 \pm .000	.018 \pm .000	.9	.008 \pm .000	.009 \pm .001	.019 \pm .001	.016 \pm .001	.015 \pm .001	.017 \pm .001		

Table F.12. Comparison of the bound values before performing Step 2) of our Training Method for ours, rivasplata, blanchard and catoni. More precisely, for each split and each variance $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$, we report the mean \pm the standard deviation (for 400 neural networks sampled from π) of the test risk ($R_T(h)$), the empirical risk ($R_S(h)$), and the value of the bounds of Corollaries 6.4.1 and 6.4.2. We consider in this table that the dataset is Fashion-MNIST.

	Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)		Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)
$10^{-1} = \sigma^2$.0	.970 \pm .028	.970 \pm .027	.972 \pm .025	.971 \pm .025	.971 \pm .026	.972 \pm .026	$10^{-1} = \sigma^2$.0	.912 \pm .027	.912 \pm .027	.918 \pm .026	.916 \pm .027	.916 \pm .027	.916 \pm .026
	.1	.166 \pm .001	.159 \pm .000	.169 \pm .000	.167 \pm .000	.166 \pm .000	.167 \pm .000		.1	.164 \pm .003	.154 \pm .003	.164 \pm .003	.162 \pm .003	.161 \pm .003	.162 \pm .004
	.2	.168 \pm .002	.160 \pm .001	.170 \pm .001	.168 \pm .001	.167 \pm .001	.168 \pm .001		.2	.164 \pm .009	.160 \pm .009	.170 \pm .010	.168 \pm .010	.167 \pm .010	.168 \pm .010
	.3	.126 \pm .000	.124 \pm .000	.134 \pm .000	.132 \pm .000	.131 \pm .000	.131 \pm .000		.3	.125 \pm .002	.119 \pm .002	.129 \pm .002	.126 \pm .002	.126 \pm .002	.126 \pm .002
	.4	.118 \pm .001	.112 \pm .000	.123 \pm .000	.120 \pm .000	.119 \pm .000	.119 \pm .000		.4	.119 \pm .003	.113 \pm .003	.124 \pm .003	.121 \pm .003	.120 \pm .003	.120 \pm .003
	.5	.106 \pm .000	.101 \pm .000	.113 \pm .000	.110 \pm .000	.109 \pm .000	.109 \pm .000		.5	.109 \pm .002	.102 \pm .001	.113 \pm .001	.110 \pm .001	.109 \pm .001	.109 \pm .001
	.6	.109 \pm .000	.102 \pm .000	.115 \pm .000	.112 \pm .000	.110 \pm .000	.110 \pm .000		.6	.102 \pm .001	.096 \pm .001	.109 \pm .001	.105 \pm .001	.105 \pm .001	.104 \pm .001
	.7	.099 \pm .000	.098 \pm .000	.112 \pm .000	.109 \pm .000	.108 \pm .000	.107 \pm .000		.7	.099 \pm .002	.094 \pm .001	.108 \pm .001	.104 \pm .001	.103 \pm .001	.102 \pm .001
	.8	.103 \pm .000	.099 \pm .000	.117 \pm .000	.112 \pm .000	.111 \pm .000	.110 \pm .000		.8	.104 \pm .001	.100 \pm .002	.118 \pm .002	.113 \pm .002	.112 \pm .002	.111 \pm .002
.9	.094 \pm .000	.089 \pm .000	.113 \pm .000	.107 \pm .000	.105 \pm .000	.106 \pm .000	.9	.092 \pm .002	.089 \pm .001	.113 \pm .001	.107 \pm .001	.105 \pm .001	.106 \pm .001		
$10^{-2} = \sigma^2$.0	.945 \pm .038	.945 \pm .037	.949 \pm .035	.948 \pm .035	.948 \pm .036	.948 \pm .036	$10^{-2} = \sigma^2$.0	.899 \pm .026	.899 \pm .027	.906 \pm .026	.904 \pm .026	.904 \pm .026	.905 \pm .025
	.1	.158 \pm .001	.151 \pm .001	.161 \pm .001	.159 \pm .001	.158 \pm .001	.159 \pm .001		.1	.178 \pm .006	.170 \pm .006	.181 \pm .006	.178 \pm .006	.177 \pm .006	.179 \pm .006
	.2	.157 \pm .003	.151 \pm .003	.162 \pm .003	.159 \pm .003	.158 \pm .003	.159 \pm .003		.2	.164 \pm .006	.159 \pm .006	.169 \pm .006	.167 \pm .006	.166 \pm .006	.167 \pm .006
	.3	.126 \pm .001	.121 \pm .001	.131 \pm .001	.128 \pm .001	.127 \pm .001	.128 \pm .001		.3	.143 \pm .007	.138 \pm .007	.148 \pm .007	.146 \pm .007	.145 \pm .007	.145 \pm .007
	.4	.114 \pm .001	.107 \pm .001	.118 \pm .001	.115 \pm .001	.114 \pm .001	.114 \pm .001		.4	.133 \pm .005	.129 \pm .005	.140 \pm .005	.137 \pm .005	.137 \pm .005	.137 \pm .005
	.5	.104 \pm .001	.099 \pm .000	.110 \pm .000	.108 \pm .000	.107 \pm .000	.106 \pm .000		.5	.122 \pm .004	.117 \pm .004	.129 \pm .004	.126 \pm .004	.125 \pm .004	.125 \pm .004
	.6	.115 \pm .001	.104 \pm .001	.117 \pm .001	.114 \pm .001	.113 \pm .001	.112 \pm .001		.6	.111 \pm .003	.104 \pm .003	.117 \pm .003	.114 \pm .003	.113 \pm .003	.112 \pm .003
	.7	.107 \pm .001	.101 \pm .001	.115 \pm .001	.111 \pm .001	.110 \pm .001	.109 \pm .001		.7	.109 \pm .003	.103 \pm .003	.118 \pm .003	.114 \pm .003	.113 \pm .003	.112 \pm .003
	.8	.098 \pm .001	.096 \pm .001	.114 \pm .001	.109 \pm .001	.108 \pm .001	.107 \pm .001		.8	.108 \pm .004	.102 \pm .004	.120 \pm .004	.115 \pm .004	.114 \pm .004	.113 \pm .004
.9	.091 \pm .001	.095 \pm .001	.119 \pm .001	.113 \pm .001	.111 \pm .001	.112 \pm .001	.9	.103 \pm .003	.099 \pm .002	.124 \pm .003	.118 \pm .003	.116 \pm .003	.116 \pm .003		

Table F.13. Comparison of the bound values before performing Step 2) of our Training Method for ours, rivasplata, blanchard and catoni. More precisely, for each split and each variance $\sigma^2 \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$, we report the mean \pm the standard deviation (for 400 neural networks sampled from π) of the test risk ($R_T(h)$), the empirical risk ($R_S(h)$), and the value of the bounds of Corollaries 6.4.1 and 6.4.2. We consider in this table that the dataset is CIFAR-10.

	Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)		Split	$R_T(h)$	$R_S(h)$	Cor. 6.4.1	Eq. (6.1)	Eq. (6.2)	Eq. (6.3)
$10^{-1} = \tau^2$.0	.899 \pm .000	.899 \pm .000	.906 \pm .000	.904 \pm .000	.903 \pm .000	.904 \pm .000	$10^{-1} = \tau^2$.0	.900 \pm .004	.900 \pm .003	.907 \pm .003	.905 \pm .003	.905 \pm .003	.905 \pm .003
	.1	.476 \pm .000	.470 \pm .000	.486 \pm .000	.482 \pm .000	.481 \pm .000	.485 \pm .000		.1	.458 \pm .001	.464 \pm .001	.479 \pm .001	.476 \pm .001	.475 \pm .001	.478 \pm .001
	.2	.390 \pm .000	.389 \pm .000	.406 \pm .000	.402 \pm .000	.401 \pm .000	.404 \pm .000		.2	.395 \pm .001	.396 \pm .000	.412 \pm .000	.409 \pm .000	.408 \pm .000	.411 \pm .000
	.3	.370 \pm .000	.358 \pm .000	.374 \pm .000	.371 \pm .000	.370 \pm .000	.372 \pm .000		.3	.361 \pm .001	.361 \pm .000	.378 \pm .000	.375 \pm .000	.373 \pm .000	.376 \pm .000
	.4	.334 \pm .000	.328 \pm .000	.346 \pm .000	.342 \pm .000	.341 \pm .000	.342 \pm .000		.4	.323 \pm .001	.316 \pm .000	.334 \pm .000	.330 \pm .000	.329 \pm .000	.331 \pm .000
	.5	.307 \pm .000	.302 \pm .000	.321 \pm .000	.317 \pm .000	.316 \pm .000	.317 \pm .000		.5	.296 \pm .001	.291 \pm .000	.310 \pm .000	.306 \pm .000	.304 \pm .000	.305 \pm .000
	.6	.274 \pm .000	.276 \pm .000	.297 \pm .000	.293 \pm .000	.291 \pm .000	.291 \pm .000		.6	.271 \pm .001	.263 \pm .000	.284 \pm .000	.279 \pm .000	.278 \pm .000	.278 \pm .000
	.7	.275 \pm .000	.272 \pm .000	.296 \pm .000	.290 \pm .000	.289 \pm .000	.288 \pm .000		.7	.253 \pm .001	.246 \pm .000	.270 \pm .000	.265 \pm .000	.263 \pm .000	.262 \pm .000
	.8	.249 \pm .000	.237 \pm .000	.265 \pm .000	.259 \pm .000	.257 \pm .000	.256 \pm .000		.8	.259 \pm .001	.252 \pm .001	.281 \pm .001	.275 \pm .001	.273 \pm .001	.272 \pm .001
.9	.227 \pm .000	.230 \pm .000	.269 \pm .000	.260 \pm .000	.258 \pm .000	.258 \pm .000	.9	.217 \pm .000	.216 \pm .001	.255 \pm .001	.246 \pm .001	.243 \pm .001	.244 \pm .001		
$10^{-1} = \tau^2$.0	.899 \pm .001	.899 \pm .000	.906 \pm .000	.904 \pm .000	.904 \pm .000	.904 \pm .000	$10^{-1} = \tau^2$.0	.905 \pm .012	.904 \pm .012	.911 \pm .011	.909 \pm .011	.909 \pm .011	.909 \pm .011
	.1	.476 \pm .000	.478 \pm .000	.494 \pm .000	.490 \pm .000	.489 \pm .000	.493 \pm .000		.1	.479 \pm .002	.480 \pm .001	.496 \pm .001	.493 \pm .001	.491 \pm .001	.495 \pm .001
	.2	.403 \pm .000	.398 \pm .000	.414 \pm .000	.410 \pm .000	.409 \pm .000	.412 \pm .000		.2	.415 \pm .002	.415 \pm .001	.432 \pm .001	.428 \pm .001	.427 \pm .001	.430 \pm .001
	.3	.349 \pm .000	.350 \pm .000	.367 \pm .000	.363 \pm .000	.362 \pm .000	.364 \pm .000		.3	.417 \pm .001	.416 \pm .001	.434 \pm .001	.430 \pm .001	.429 \pm .001	.431 \pm .001
	.4	.322 \pm .000	.313 \pm .000	.330 \pm .000	.327 \pm .000	.326 \pm .000	.327 \pm .000		.4	.333 \pm .001	.323 \pm .001	.341 \pm .001	.337 \pm .001	.336 \pm .001	.338 \pm .001
	.5	.281 \pm .000	.283 \pm .000	.302 \pm .000	.298 \pm .000	.297 \pm .000	.297 \pm .000		.5	.316 \pm .001	.311 \pm .001	.331 \pm .001	.327 \pm .001	.325 \pm .001	.326 \pm .001
	.6	.290 \pm .000	.286 \pm .000	.307 \pm .000	.303 \pm .000	.301 \pm .000	.301 \pm .000		.6	.280 \pm .001	.281 \pm .001	.302 \pm .001	.298 \pm .001	.296 \pm .001	.296 \pm .001
	.7	.266 \pm .000	.257 \pm .000	.281 \pm .000	.276 \pm .000	.274 \pm .000	.274 \pm .000		.7	.239 \pm .001	.234 \pm .001	.257 \pm .001	.252 \pm .001	.250 \pm .001	.250 \pm .001
	.8	.247 \pm .000	.243 \pm .000	.271 \pm .000	.265 \pm .000	.263 \pm .000	.262 \pm .000		.8	.249 \pm .001	.245 \pm .001	.274 \pm .001	.268 \pm .001	.266 \pm .001	.264 \pm .001
.9	.236 \pm .000	.227 \pm .000	.266 \pm .000	.257 \pm .000	.255 \pm .000	.255 \pm .000	.9	.233 \pm .001	.232 \pm .002	.272 \pm .002	.263 \pm .002	.260 \pm .002	.260 \pm .002		

APPENDIX OF CHAPTER 7

G.1 Proof of Theorem 7.3.1

Theorem 7.3.1 (Generalization Bound with Complexity Measures). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap. For any \mathcal{D} on $\mathcal{X} \times \mathcal{Y}$, for any \mathbb{H} , for any distribution $\pi \in \mathcal{M}^*(\mathbb{H})$ on \mathbb{H} , for any $\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h' \sim \pi, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \left[\alpha \mathbf{R}_{\mathcal{S}}(h') + \mu(h', \mathcal{S}) \right] - \left[\alpha \mathbf{R}_{\mathcal{S}}(h) + \mu(h, \mathcal{S}) \right] \right. \\ \left. + \ln \frac{\pi(h')}{\pi(h)} + \ln \left(\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp[\phi(\mathbf{R}_{\mathcal{D}}(h'), \mathbf{R}_{\mathcal{S}'}(h'))] \right) \right] \geq 1 - \delta. \end{aligned}$$

Proof. First of all, we denote as $Z = \int_{\mathbb{H}} \exp[-\alpha \mathbf{R}_{\mathcal{S}}(h') - \mu(h', \mathcal{S})] d\xi(h')$, the normalization constant of the Gibbs distribution $\rho_{\mathcal{S}}$ and ξ the reference measure on \mathbb{H} . Moreover, we have

$$\rho_{\mathcal{S}}(h) = \frac{1}{Z} \exp[-\alpha \mathbf{R}_{\mathcal{S}}(h) - \mu(h, \mathcal{S})] \propto \exp[-\alpha \mathbf{R}_{\mathcal{S}}(h) - \mu(h, \mathcal{S})].$$

We apply Theorem 2.4.1 with $\frac{\delta}{2}$ instead of δ and with the function $\varphi(h, \mathcal{S}) = \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h))$ to obtain with probability at least $1 - \frac{\delta}{2}$ over $\mathcal{S} \sim \mathcal{D}^m$ and $h \sim \rho_{\mathcal{S}}$

$$\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \ln \left[\frac{\rho_{\mathcal{S}}(h)}{\pi(h)} \right] + \ln \left[\frac{2}{\delta} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\phi(\mathbf{R}_{\mathcal{D}}(h'), \mathbf{R}_{\mathcal{S}'}(h'))} \right].$$

We develop the term $\ln \left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right]$ in Theorem 2.4.1. We have

$$\begin{aligned}
\ln \left[\frac{\rho_{\mathbb{S}}(h)}{\pi(h)} \right] &= \ln \left(\frac{\exp[-\alpha R_{\mathbb{S}}(h) - \mu(h, \mathbb{S})]}{Z} \frac{1}{\pi(h)} \right) \\
&= \ln(\exp[-\alpha R_{\mathbb{S}}(h) - \mu(h, \mathbb{S})]) \\
&\quad - \ln \left(\pi(h) \int_{\mathbb{H}} \exp[-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})] d\xi(h') \right) \\
&= -\alpha R_{\mathbb{S}}(h) - \mu(h, \mathbb{S}) \\
&\quad - \ln \left(\pi(h) \int_{\mathbb{H}} \frac{\pi(h')}{\pi(h')} \exp[-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})] d\xi(h') \right) \\
&= -\alpha R_{\mathbb{S}}(h) - \mu(h, \mathbb{S}) - \ln \left(\mathbb{E}_{h' \sim \pi} \frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right).
\end{aligned}$$

Hence, we obtain the following

$$\begin{aligned}
\mathbb{P}_{\mathbb{S} \sim \mathcal{D}^m, h \sim \rho_{\mathbb{S}}} \left[\phi(R_{\mathcal{D}}(h), R_{\mathbb{S}}(h)) \leq \ln \left[\frac{2}{\delta} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{\phi(R_{\mathcal{D}}(h), R_{\mathbb{S}'}(h'))} \right] \right. \\
\left. - \alpha R_{\mathbb{S}}(h) - \mu(h, \mathbb{S}) - \ln \left(\mathbb{E}_{h' \sim \pi} \frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right) \right] \geq 1 - \frac{\delta}{2}. \quad (\text{G.1})
\end{aligned}$$

We can now upper-bound the term $-\ln \left(\mathbb{E}_{h' \sim \pi} \frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right)$. To do so, since $\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} > 0$ for all $h \in \mathbb{H}$ and $\mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m$, we apply MARKOV's inequality (Theorem A.2.1) to obtain for all $h \in \mathbb{H}$ and $\mathbb{S} \in (\mathbb{X} \times \mathbb{Y})^m$ with probability at least $1 - \frac{\delta}{2}$ over $h' \sim \pi$

$$\begin{aligned}
\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} &\leq \frac{2}{\delta} \mathbb{E}_{h' \sim \pi} \left(\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right) \\
\iff -\ln \left(\mathbb{E}_{h' \sim \pi} \left[\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right] \right) &\leq \ln \frac{2}{\delta} - \ln \left(\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right).
\end{aligned}$$

Moreover, by simplifying the right-hand side of the inequality, we have

$$-\ln \left(\frac{\pi(h)}{\pi(h')} e^{-\alpha R_{\mathbb{S}}(h') - \mu(h', \mathbb{S})} \right) = \ln \frac{\pi(h')}{\pi(h)} + \alpha R_{\mathbb{S}}(h') + \mu(h', \mathbb{S}).$$

G.2. Proof of Corollary 7.3.1

Hence, we obtain the following inequality:

$$\mathbb{P}_{h' \sim \pi} \left[-\ln \left(\frac{\pi(h)}{\pi(h')} e^{-\alpha R_S(h') - \mu(h', \mathbb{S})} \right) \leq \ln \frac{2}{\delta} + \ln \frac{\pi(h')}{\pi(h)} + \alpha R_S(h') + \mu(h', \mathbb{S}) \right] \geq 1 - \frac{\delta}{2}. \quad (\text{G.2})$$

By using an union bound on Equations (G.1) and (G.2) an rearranging the terms, we obtain the claimed result. ■

G.2 Proof of Corollary 7.3.1

Corollary 7.3.1 (Practical Generalization Bound with Complexity Measures). For any distribution \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any bounded hypothesis set \mathbb{H} , given the uniform prior distribution π on \mathbb{H} , for any $\mu : \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m \rightarrow \mathbb{R}$, for any $\delta \in (0, 1]$, with probability at least $1 - \delta$ over $\mathbb{S} \sim \mathcal{D}^m$, $h' \sim \pi$, $h \sim \rho_{\mathbb{S}}$ we have

$$\text{kl} [R_S(h) \| R_{\mathcal{D}}(h)] \leq \frac{1}{m} \left[\left[\alpha R_S(h') + \mu(h', \mathbb{S}) \right] - \left[\alpha R_S(h) + \mu(h, \mathbb{S}) \right] + \frac{8\sqrt{m}}{\delta^2} \right]_+, \quad (\text{7.5})$$

$$\left| R_{\mathcal{D}}(h) - R_S(h) \right| \leq \sqrt{\frac{1}{2m} \left[\left[\alpha R_S(h') + \mu(h', \mathbb{S}) \right] - \left[\alpha R_S(h) + \mu(h, \mathbb{S}) \right] + \frac{8\sqrt{m}}{\delta^2} \right]_+}, \quad (\text{7.6})$$

where $[a]_+ = \max(0, a)$, and $\rho_{\mathbb{S}}$ is the Gibbs distribution defined by Equation (7.2).

Proof. Since π is the uniform distribution we have: $\mathbb{E}_{h' \sim \pi} \ln \frac{\pi(h)}{\pi(h')} = 0$. We instantiate Th. 7.3.1 with $\phi(R_{\mathcal{D}}(h), R_S(h)) = m \text{kl} [R_S(h) \| R_{\mathcal{D}}(h)]$. It remains to upper-bound $\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(m \text{kl} [R_{\mathbb{S}'}(h') \| R_{\mathcal{D}}(h')])$. We have

$$\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl} [R_{\mathbb{S}'}(h') \| R_{\mathcal{D}}(h')]} = \mathbb{E}_{h' \sim \pi} \mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} e^{m \text{kl} [R_{\mathbb{S}'}(h') \| R_{\mathcal{D}}(h')]} \quad (\text{G.3})$$

$$\mathbb{E}_{\mathbb{S}' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} e^{m \text{kl} [R_{\mathbb{S}'}(h') \| R_{\mathcal{D}}(h')]} \leq 2\sqrt{m}, \quad (\text{G.4})$$

where Equation (G.3) is due to FUBINI's theorem (*i.e.*, we can exchange the two expectations), and Equation (G.4) is due to MAURER (2004) (see Lemma B.16.1). By rearranging the terms, with probability at least $1 - \delta$ over $\mathbb{S} \sim \mathcal{D}^m$, $h \sim \rho_{\mathbb{S}}$, and

$h' \sim \pi$ we have

$$\text{kl} [R_S(h) \| R_D(h)] \leq \frac{1}{m} \left[[\alpha R_S(h') + \mu(h', S)] - [\alpha R_S(h) + \mu(h, S)] + \ln \frac{8\sqrt{m}}{\delta^2} \right].$$

Hence, by definition of $[a]_+$, we can deduce Equation (7.5). From PINSKER's inequality (Theorem B.5.1), we have

$$2(R_D(h) - R_S(h))^2 \leq \text{kl} [R_S(h) \| R_D(h)].$$

Hence, thanks to this inequality and by rearranging the terms, we obtain Equation (7.6). \blacksquare

G.3 Proof of Proposition 7.5.1

Proposition 7.5.1 (Set-theoretic view of Definition 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_\mu : \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m \times (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 7.3.1. Under these conditions, with $\mathbb{Z}_d = \left\{ (h, S) \in \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m : \phi(R_D(h), R_S(h)) \leq \Phi_\mu(h, S, \delta) \right\}$, and $\mathbb{P}_{S \sim \mathcal{D}^m, h \sim \rho_S} [(h, S) \in \mathbb{Z}_d] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (7.1)} &\iff \forall (h, S) \in \mathbb{Z}_d, \phi(R_D(h), R_S(h)) \leq \Phi_\mu(h, S, \delta) \\ &\iff \sup_{(h, S) \in \mathbb{Z}_d} \left\{ \phi(R_D(h), R_S(h)) - \Phi_\mu(h, S, \delta) \right\} \leq 0. \end{aligned}$$

Proof. First of all, by definition of the supremum, we have

$$\begin{aligned} \forall \delta \in (0, 1], \quad \forall (h, S) \in \mathbb{Z}_d, \quad \phi(R_D(h), R_S(h)) &\leq \Phi_\mu(h, S, \delta) \\ \forall (h, S) \in \mathbb{Z}_d, \quad \phi(R_D(h), R_S(h)) - \Phi_\mu(h, S, \delta) &\leq 0 \\ \iff \sup_{(h, S) \in \mathbb{Z}_d} \left\{ \phi(R_D(h), R_S(h)) - \Phi_\mu(h, S, \delta) \right\} &\leq 0. \end{aligned}$$

It remains to prove that

$$\underbrace{\text{Equation (7.1)}}_{(A)} \iff \underbrace{\forall (h, S) \in \mathbb{Z}_d, \phi(R_D(h), R_S(h)) \leq \Phi_\mu(h, S, \delta)}_{(B)}$$

to complete the proof.

Step 1 ((A) \Rightarrow (B)). Assuming that (A) holds, by definition of \mathbb{Z}_d , we have

$$\begin{aligned} & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mu}(h, \mathcal{S}, \delta) \right] \geq 1 - \delta \\ \iff & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[(h, \mathcal{S}) \in \mathbb{Z}_d \right] \geq 1 - \delta, \end{aligned}$$

since $\mathbb{I}[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mu}(h, \mathcal{S}, \delta)] = \mathbb{I}[(h, \mathcal{S}) \in \mathbb{Z}_d]$.

Additionally, by definition of \mathbb{Z}_d , we know that

$$\forall (h, \mathcal{S}) \in \mathbb{Z}_d, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mu}(h, \mathcal{S}, \delta),$$

where $\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [(h, \mathcal{S}) \in \mathbb{Z}_d] \geq 1 - \delta$.

Step 2 ((A) \Leftarrow (B)). Note that from the definition of \mathbb{Z}_d we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [(h, \mathcal{S}) \in \mathbb{Z}_d] \geq 1 - \delta.$$

Additionally, we can deduce that

$$\begin{aligned} & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [(h, \mathcal{S}) \in \mathbb{Z}_d] \geq 1 - \delta \\ \iff & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} [\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mu}(h, \mathcal{S}, \delta)] \geq 1 - \delta. \end{aligned}$$

■

G.4 Proof of Proposition 7.5.2

Proposition 7.5.2 (Set-theoretic View of Uniform Convergence Bounds). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.2. Under these conditions, with $\mathbb{Z}_u = \left\{ \mathcal{S} \in (\mathcal{X} \times \mathcal{Y})^m : \forall h \in \mathbb{H}, \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_u(\delta) \right\}$, and $\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m} [\mathcal{S} \in \mathbb{Z}_u] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (1.1)} & \iff \forall \mathcal{S} \in \mathbb{Z}_u, \forall h \in \mathbb{H}, \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_u(\delta) \\ & \iff \sup_{\mathcal{S} \in \mathbb{Z}_u} \sup_{h \in \mathbb{H}} \left\{ \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \right\} \leq \Phi_u(\delta). \end{aligned}$$

Proof. First of all, by definition of the supremum, we have

$$\begin{aligned} \forall \delta \in (0, 1], \quad \forall S \in \mathbb{Z}_u, \quad \forall h \in \mathbb{H}, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta) \\ \iff \sup_{S \in \mathbb{Z}_u} \sup_{h \in \mathbb{H}} \left\{ \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \right\} \leq \Phi_u(\delta). \end{aligned}$$

It remains to prove that

$$\underbrace{\text{Equation (1.1)}}_{(A)} \iff \underbrace{\forall S \in \mathbb{Z}_u, \quad \forall h \in \mathbb{H}, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta)}_{(B)}$$

to complete the proof.

Step 1 ((A) \Rightarrow (B)). Assuming that (A) holds, by definition of \mathbb{Z}_u , we have

$$\begin{aligned} \mathbb{P}_{S \sim \mathcal{D}^m} \left[\forall h \in \mathbb{H}, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta \\ \iff \mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_u \right] \geq 1 - \delta, \end{aligned}$$

since $I[\forall h \in \mathbb{H}, \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta)] = I[S \in \mathbb{Z}_u]$.

Additionally, by definition of \mathbb{Z}_u , we know that

$$\forall S \in \mathbb{Z}_u, \quad \forall h \in \mathbb{H}, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta),$$

where $\mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_u \right] \geq 1 - \delta$.

Step 2 ((A) \Leftarrow (B)). Note that from the definition of \mathbb{Z}_u we have

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_u \right] \geq 1 - \delta.$$

Additionally, we can deduce that

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_u \right] \geq 1 - \delta \iff \mathbb{P}_{S \sim \mathcal{D}^m} \left[\forall h \in \mathbb{H}, \quad \phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta.$$

■

G.5 Proof of Corollary 7.5.1

Corollary 7.5.1 (Uniform Convergence Bound from Theorem 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap and assume that there exists a function $\Phi_u : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.2 such that $\Phi_u(\delta) \geq$

G.6. Proof of Proposition 7.5.3

$\ln \left[\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\phi(\mathbf{R}_{\mathcal{D}}(h'), \mathbf{R}_{S'}(h'))) \right]$. For any \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hypothesis set \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{S \sim \mathcal{D}^m, h \sim \rho_S} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \Phi_u(\delta) \right] \geq 1 - \delta.$$

Proof. Let the parametric function $\mu(\cdot)$ defined as

$$\forall (h, S) \in \mathbb{H} \times (\mathbb{X} \times \mathbb{Y})^m, \quad \mu(h, S) = -\alpha \mathbf{R}_S(h) - \ln \pi(h) + \Phi_u(\delta).$$

Given the definition of ρ_S (with the parametric function $\mu(\cdot)$ defined above), we can deduce from Theorem 7.3.1 that

$$\begin{aligned} & \mathbb{P}_{S \sim \mathcal{D}^m, h' \sim \pi, h \sim \rho_S} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \ln \left(\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp[\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{S'}(g))] \right) \right] \\ &= \mathbb{P}_{S \sim \mathcal{D}^m, h \sim \rho_S} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_S(h)) \leq \ln \left(\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp[\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{S'}(g))] \right) \right] \geq 1 - \delta. \end{aligned}$$

Note that the equality holds since $h' \sim \pi$ does not appear in the bound. If the assumption $\Phi_u(\delta) \geq \ln \left[\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp(\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{S'}(g))) \right]$ is satisfied, then, we can deduce Corollary 7.5.1. \blacksquare

G.6 Proof of Proposition 7.5.3

Proposition 7.5.3 (Set-theoretic View of Algorithmic Dependent Bounds). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be a generalization gap and assume that there exists a function $\Phi_a : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.5. Under these conditions, with $\mathbb{Z}_a = \{S \in (\mathbb{X} \times \mathbb{Y})^m : \phi(\mathbf{R}_{\mathcal{D}}(h_S), \mathbf{R}_S(h_S)) \leq \Phi_a(\delta)\}$ and $\mathbb{P}_{S \sim \mathcal{D}^m} [S \in \mathbb{Z}_a] \geq 1 - \delta$, we have

$$\begin{aligned} \text{Equation (1.2)} & \iff \forall S \in \mathbb{Z}_a, \phi(\mathbf{R}_{\mathcal{D}}(h_S), \mathbf{R}_S(h_S)) \leq \Phi_a(\delta) \\ & \iff \sup_{S \in \mathbb{Z}_a} \phi(\mathbf{R}_{\mathcal{D}}(h_S), \mathbf{R}_S(h_S)) \leq \Phi_a(\delta). \end{aligned}$$

Proof. First of all, by definition of the supremum, we have

$$\forall \delta \in (0, 1], \quad \forall S \in \mathbb{Z}_a, \phi(h_S, S) \leq \Phi_a(\delta) \iff \sup_{S \in \mathbb{Z}_a} \phi(h_S, S) \leq \Phi_a(\delta).$$

It remains to prove that

$$\underbrace{\text{Equation (1.2)}}_{(A)} \iff \underbrace{\forall S \in \mathbb{Z}_a, \phi(R_{\mathcal{D}}(h_S), R_S(h_S)) \leq \Phi_a(\delta)}_{(B)}$$

to complete the proof.

Step 1 ((A) \Rightarrow (B)). Assuming that (A) holds, by definition of \mathbb{Z}_a , we have

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[\phi(h_S, S) \leq \Phi_a(\delta) \right] \geq 1 - \delta \iff \mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_a \right] \geq 1 - \delta.$$

since $\mathbb{I}[\phi(h_S, S) \leq \Phi_a(\delta)] = \mathbb{I}[S \in \mathbb{Z}_a]$.

Additionally, by definition of \mathbb{Z}_a , we know that

$$\forall S \in \mathbb{Z}_a, \quad \phi(h_S, S) \leq \Phi_a(\delta), \quad \text{where} \quad \mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_a \right] \geq 1 - \delta.$$

Step 2 ((A) \Leftarrow (B)). Note that from the definition of \mathbb{Z}_a we have

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_a \right] \geq 1 - \delta.$$

Additionally, we can deduce that

$$\mathbb{P}_{S \sim \mathcal{D}^m} \left[S \in \mathbb{Z}_a \right] \geq 1 - \delta \iff \mathbb{P}_{S \sim \mathcal{D}^m} \left[\phi(h_S, S) \leq \Phi_a(\delta) \right] \geq 1 - \delta.$$

■

G.7 Proof of Corollary 7.5.2

Corollary 7.5.2 (Algorithmic-dependent Bound from Theorem 7.3.1). Let $\phi : [0, 1]^2 \rightarrow \mathbb{R}$ be the generalization gap and assume that there exists a function $\Phi_a : (0, 1] \rightarrow \mathbb{R}$ fulfilling Definition 1.3.5 such that $\Phi_a(\delta) \geq \ln \left[\frac{4}{\delta^2} \mathbb{E}_{S' \sim \mathcal{D}^m} \mathbb{E}_{h' \sim \pi} \exp(\phi(R_{\mathcal{D}}(h'), R_{S'}(h'))) \right]$. For any \mathcal{D} on $\mathbb{X} \times \mathbb{Y}$, for any hy-

pothesis set \mathbb{H} , for any $\delta \in (0, 1]$, we have

$$\mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \Phi_{\mathbf{a}}(\delta) \right] \geq 1 - \delta.$$

Proof. Let the parametric function $\mu(\cdot)$ defined as

$$\forall (h, \mathcal{S}) \in \mathbb{H} \times (\mathcal{X} \times \mathcal{Y})^m, \quad \mu(h, \mathcal{S}) = -\alpha \mathbf{R}_{\mathcal{S}}(h) - \ln \pi(h) + \Phi_{\mathbf{a}}(\delta).$$

Given the definition of $\rho_{\mathcal{S}}$ (with the parametric function $\mu(\cdot)$ defined above), we can deduce from Theorem 7.3.1 that

$$\begin{aligned} & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h' \sim \pi, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \ln \left(\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp [\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{\mathcal{S}'}(g))] \right) \right] \\ = & \mathbb{P}_{\mathcal{S} \sim \mathcal{D}^m, h \sim \rho_{\mathcal{S}}} \left[\phi(\mathbf{R}_{\mathcal{D}}(h), \mathbf{R}_{\mathcal{S}}(h)) \leq \ln \left(\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp [\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{\mathcal{S}'}(g))] \right) \right] \geq 1 - \delta. \end{aligned}$$

Note that the equality holds since $h' \sim \pi$ does not appear in the bound. If the assumption $\Phi_{\mathbf{a}}(\delta) \geq \ln \left[\frac{4}{\delta^2} \mathbb{E}_{\mathcal{S}' \sim \mathcal{D}^m} \mathbb{E}_{g \sim \pi} \exp (\phi(\mathbf{R}_{\mathcal{D}}(g), \mathbf{R}_{\mathcal{S}'}(g))) \right]$ is satisfied, then, we can deduce Corollary 7.5.2. \blacksquare

G.8 Details on the Experiments

In this section, we introduce additional figures concerning the tightness, the influence of α , and the influence of the number of parameters. Additionally, we provide more experiments with data-dependent complexity measures that we present in Appendix G.8.1.

G.8.1 Data-dependent Complexity Measures $\Phi_\mu(h, \mathbb{S}, \delta)$

As we have pointed out in the paper, the parametric function $\mu()$ depends on the learning sample \mathbb{S} . We illustrate this dependence with other parametric functions defined as

$$\begin{aligned} \text{DIST_FRO-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{DIST_FRO}(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \\ \text{DIST_L}_2\text{-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{DIST_L}_2(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \\ \text{PARAM_NORM-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{PARAM_NORM}(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \\ \text{PATH_NORM-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{PATH_NORM}(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \\ \text{SUM_FRO-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{SUM_FRO}(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \\ \text{ZERO-AUG}(h_{\mathbf{w}}, \mathbb{S}) &= \text{ZERO-AUG}(h_{\mathbf{w}}) + \text{AUG}(h_{\mathbf{w}}, \mathbb{S}), \end{aligned}$$

where

$$\text{AUG}(h, \mathbb{S}) = -\frac{1}{2}R_{\mathbb{S}}(h) + \frac{1}{2}R_{\widehat{\mathbb{S}}}(h),$$

and $\widehat{\mathbb{S}}$ is a data-augmented learning sample. More precisely, we apply to each example $(\mathbf{x}, y) \in \mathbb{S}$ (a) a random rotation (with a maximum angle set to 20°) and (b) a random translation (with a maximum of 3 translated pixels per dimension).

G.8.2 Tightness of the Bounds

Figures G.1 and G.2 report the tightness of the bounds for the data-dependent complexity measures introduced in Appendix G.8.1.

G.8.3 Influence of the Parameter α

Figures G.3 to G.6 shows the influence of the parameter α for all parametric functions.

G.8.4 Influence of the Depth/Width

Figures G.7 to G.10 shows the influence of the depth and the width for all parametric functions.

G.8. Details on the Experiments

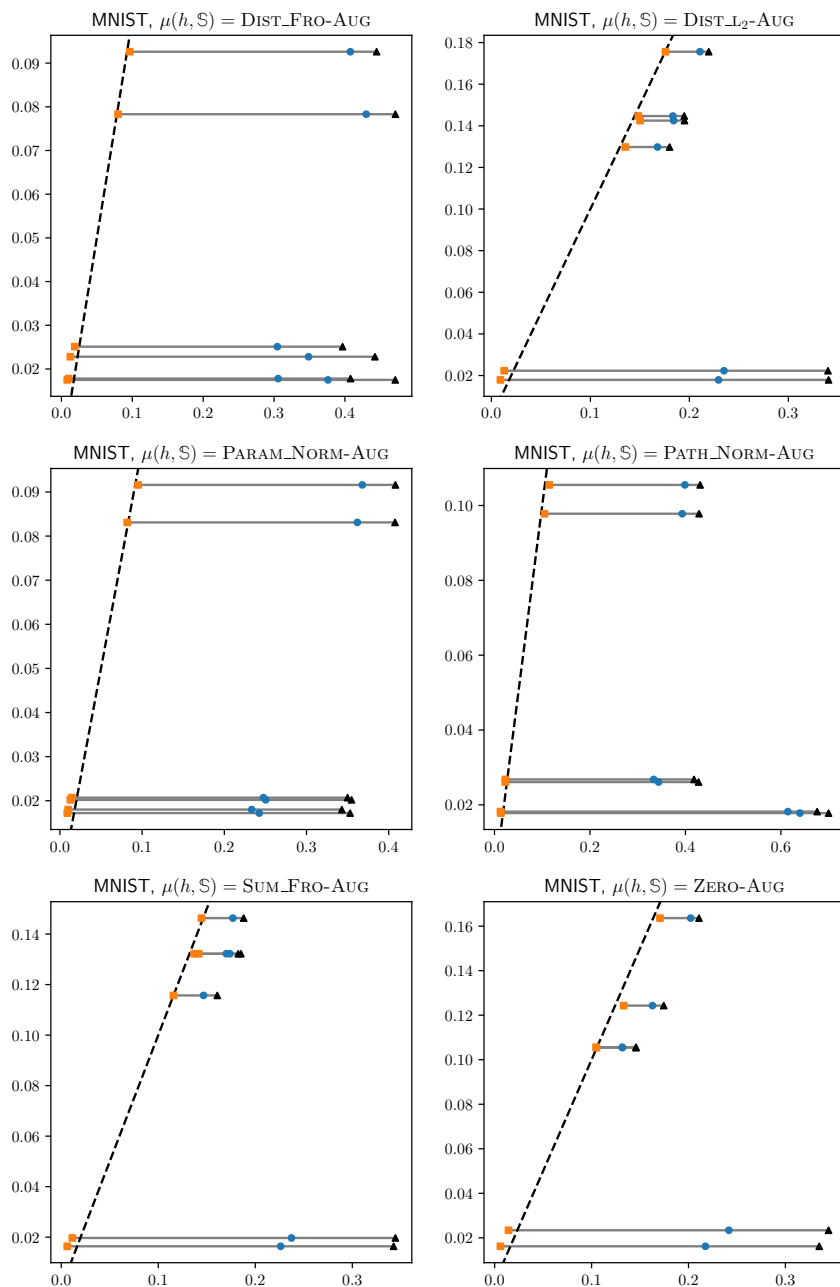


Figure G.1. Scatter plot given a parametric function $\mu(h, \mathcal{S})$, where each segment represents a neural network $h_{\mathbf{w}}$ learned with a given α , width H and depth L . For each segment, there is a corresponding orange square and a blue circle. The orange squares corresponds to the empirical risk $R_{\mathcal{S}}(h)$ (x-axis) and test risk $R_{\mathcal{T}}(h)$ (y-axis). The blue circle resp. the black triangle represents Equation (7.7) resp. Equation (7.8) in the x-axis and the test risk $R_{\mathcal{T}}(h)$ in the y-axis. The dashed line is the identity function.

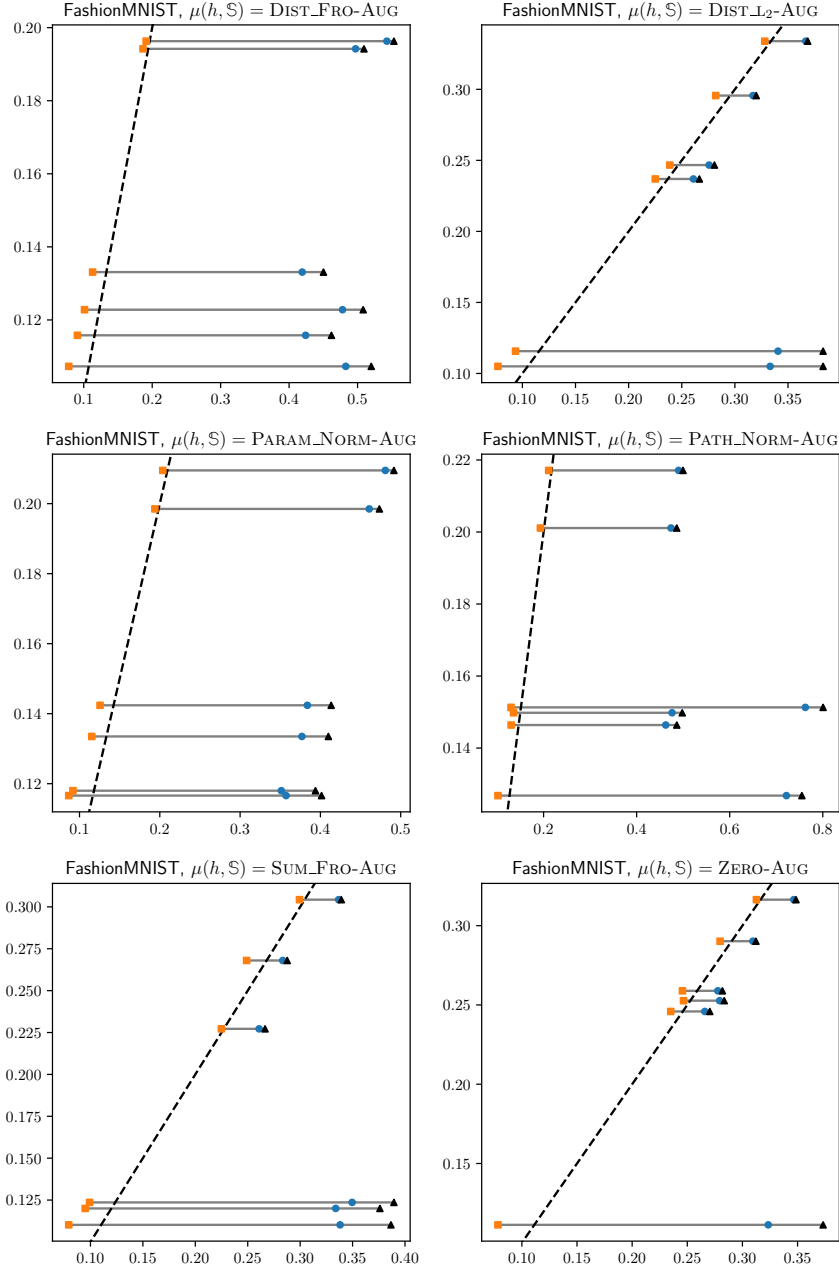


Figure G.2. Scatter plot given a parametric function $\mu(h, \mathcal{S})$, where each segment represents a neural network $h_{\mathcal{W}}$ learned with a given α , width H and depth L . For each segment, there is a corresponding orange square and a blue circle. The orange squares corresponds to the empirical risk $R_{\mathcal{S}}(h)$ (x-axis) and test risk $R_{\mathcal{T}}(h)$ (y-axis). The blue circle resp. the black triangle represents Equation (7.7) resp. Equation (7.8) in the x-axis and the test risk $R_{\mathcal{T}}(h)$ in the y-axis. The dashed line is the identity function.

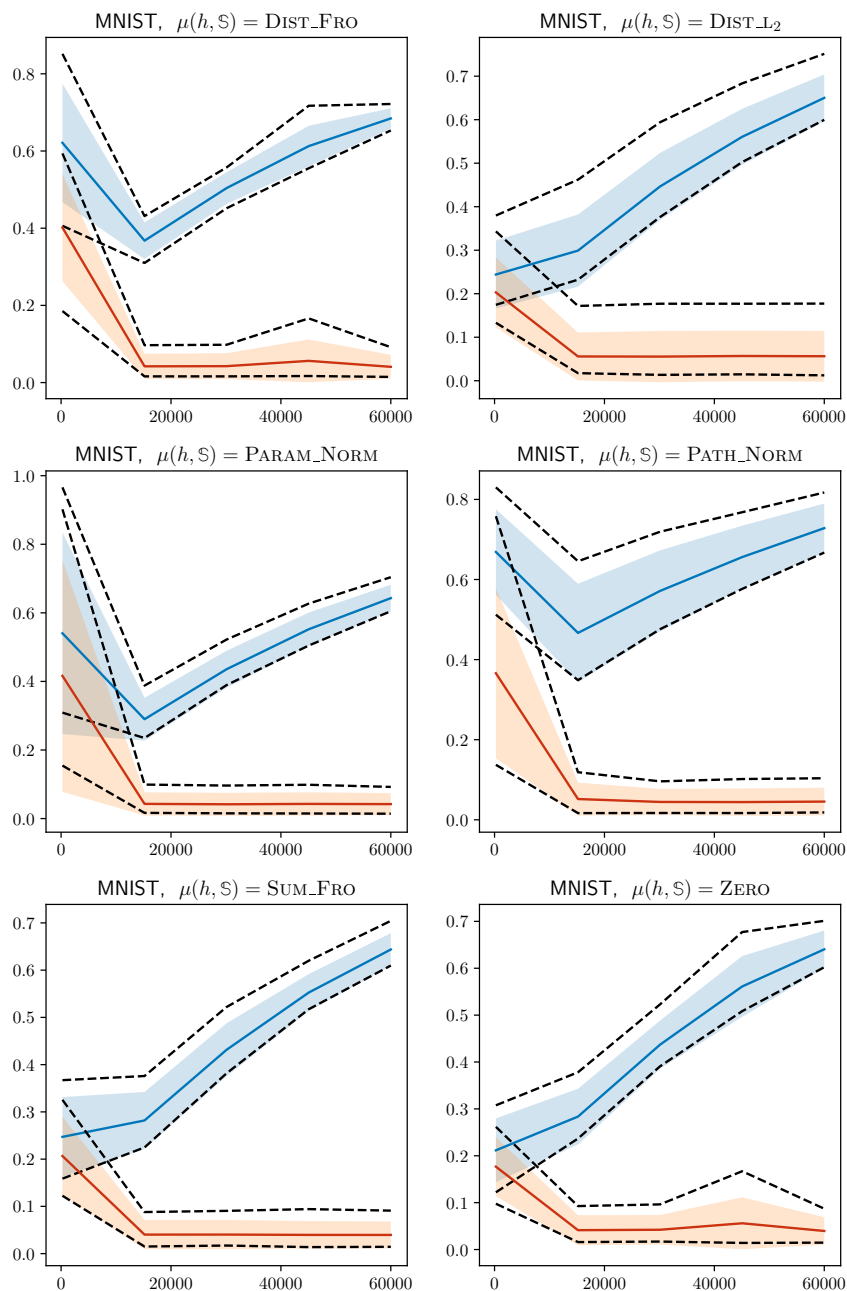


Figure G.3. Influence of the parameter α in the x-axis. The bound values are represented in blue and the test risk in red. The two (solid) lines are the mean values computed on the depths and widths; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

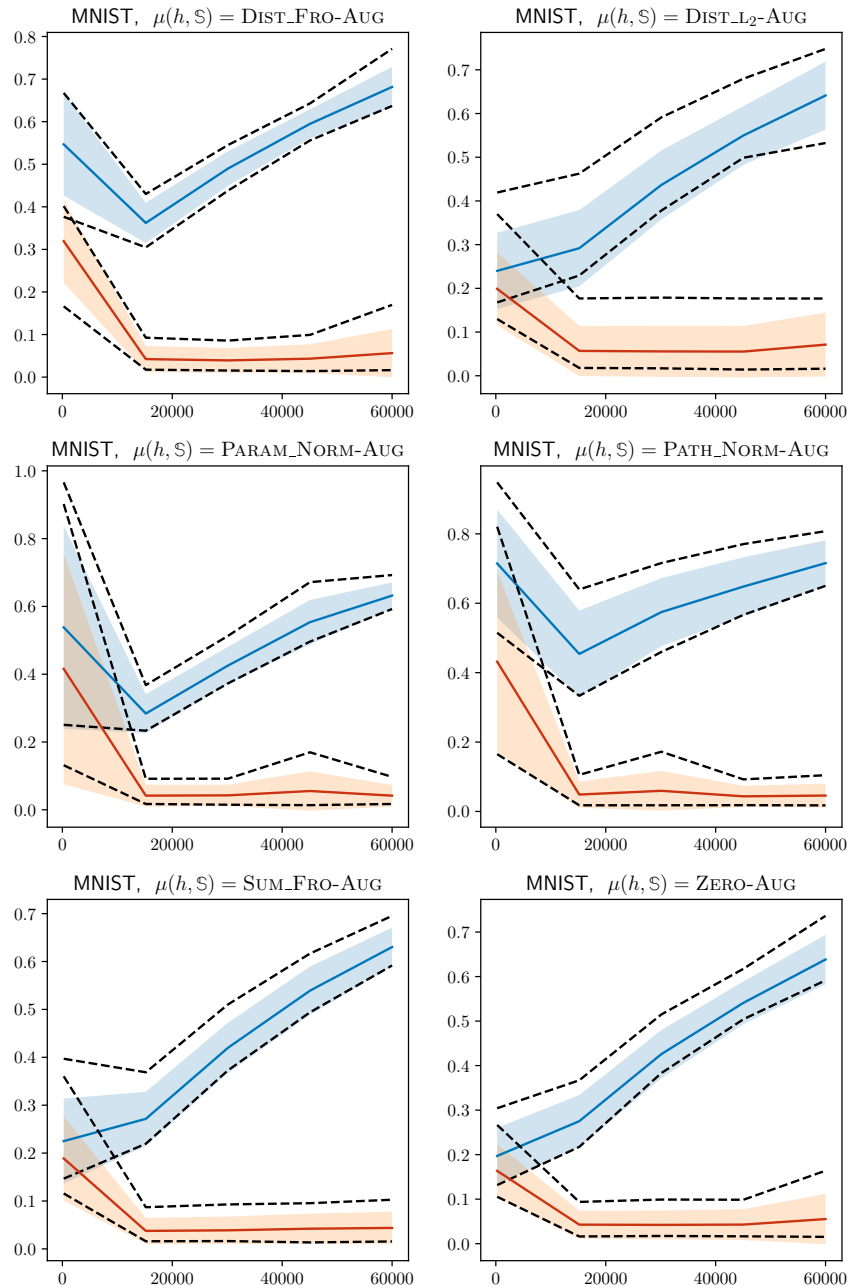


Figure G.4. Influence of the parameter α in the x-axis. The bound values are represented in blue and the test risk in red. The two (solid) lines are the mean values computed on the depths and widths; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

G.8. Details on the Experiments

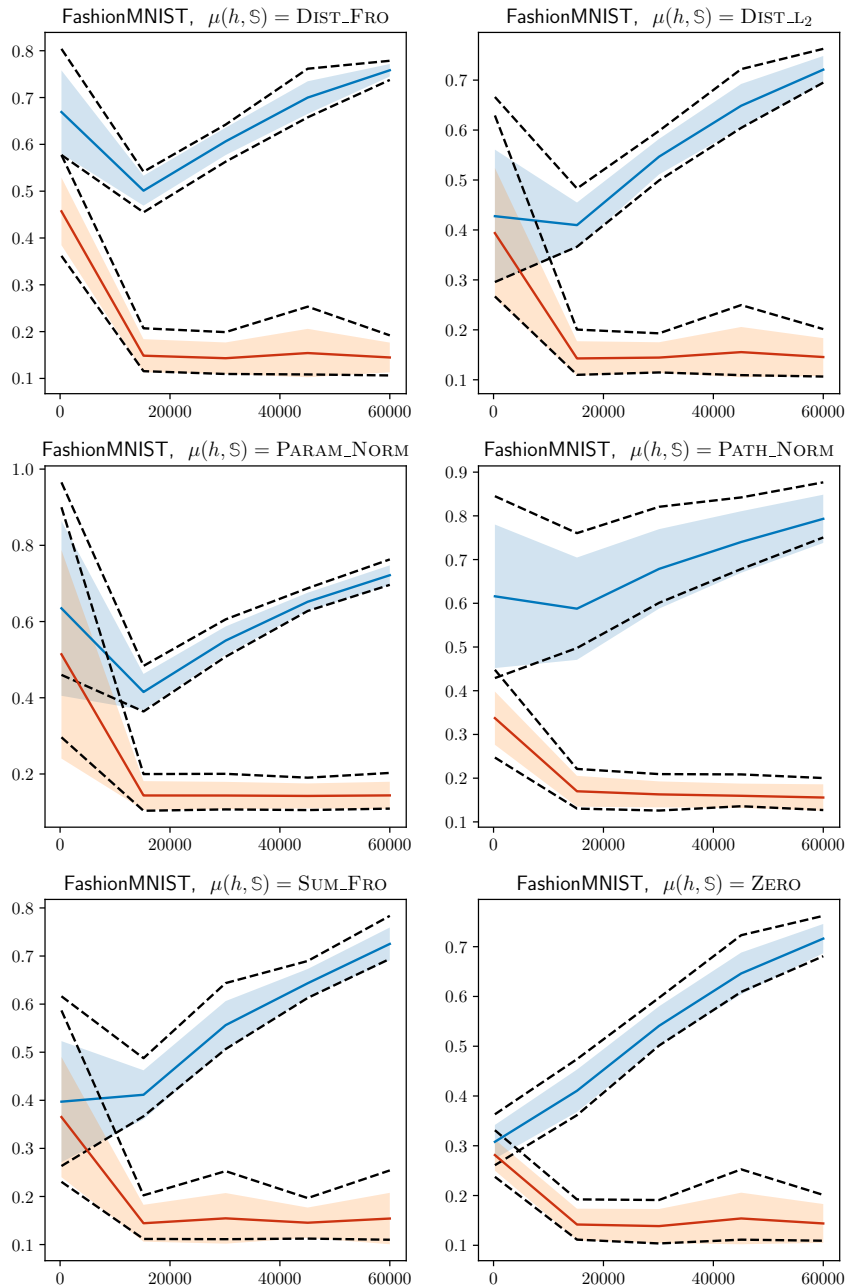


Figure G.5. Influence of the parameter α in the x-axis. The bound values are represented in blue and the test risk in red. The two (solid) lines are the mean values computed on the depths and widths; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

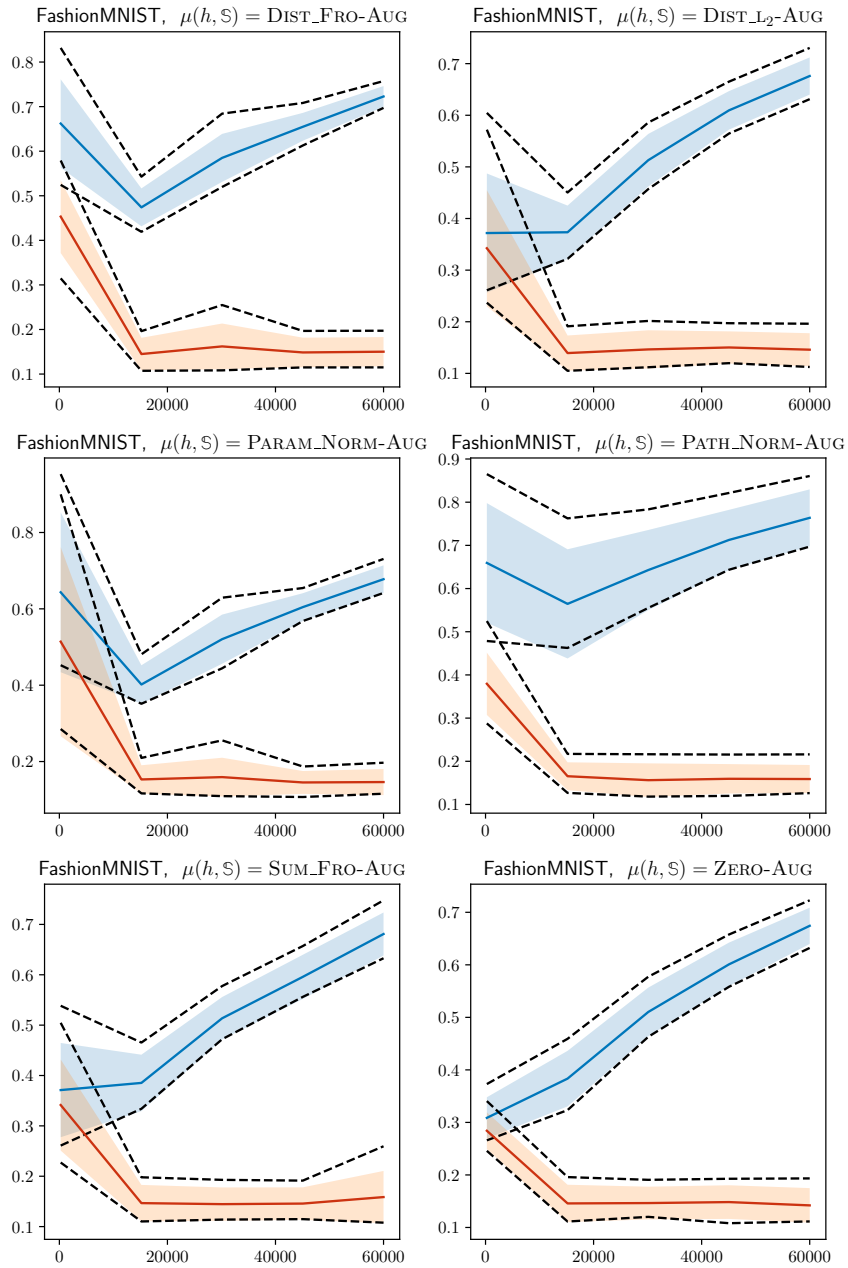


Figure G.6. Influence of the parameter α in the x -axis. The bound values are represented in blue and the test risk in red. The two (solid) lines are the mean values computed on the depths and widths; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

G.8. Details on the Experiments

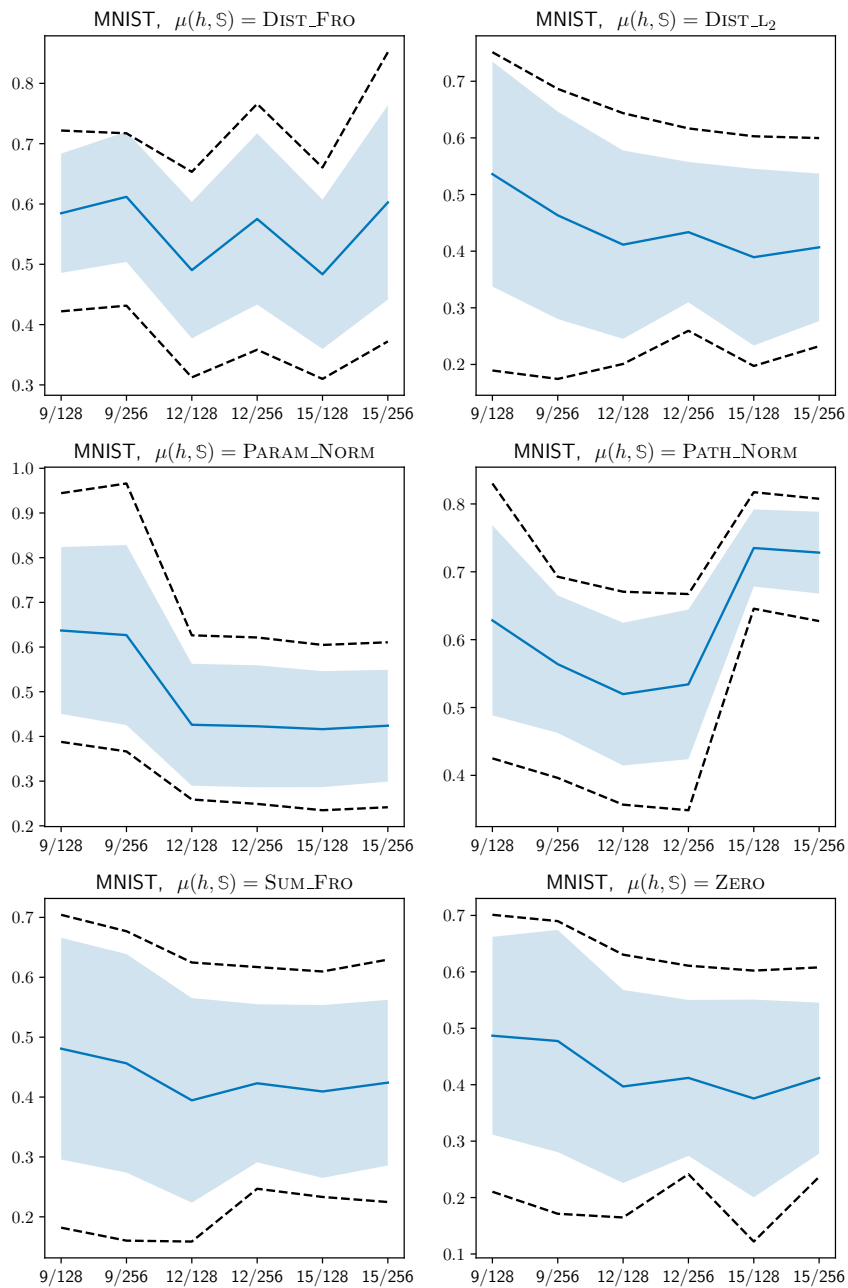


Figure G.7. Influence of the depth and the width in the x-axis as “depth/width”. The (solid) lines are the mean values computed on the different values of α ; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

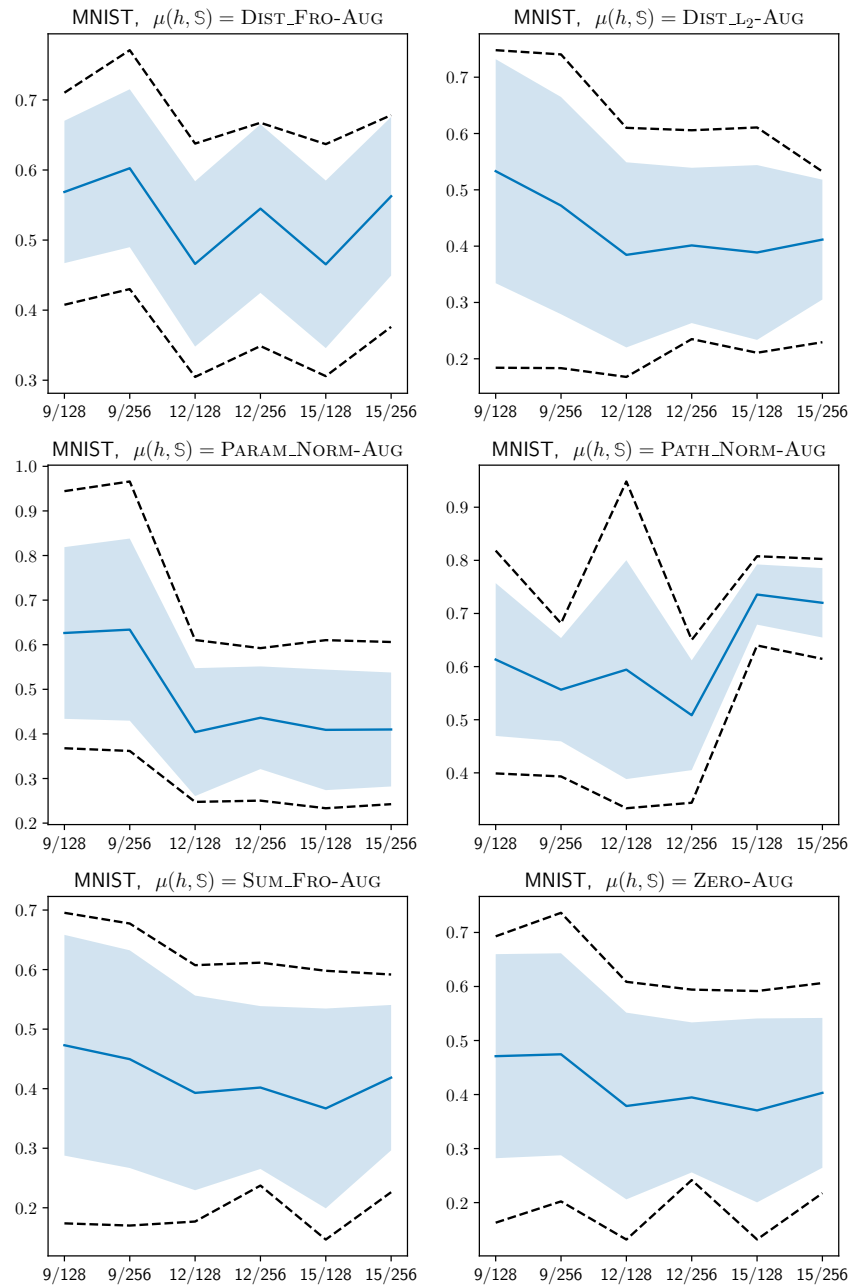


Figure G.8. Influence of the depth and the width in the x-axis as “depth/width”. The (solid) lines are the mean values computed on the different values of α ; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

G.8. Details on the Experiments

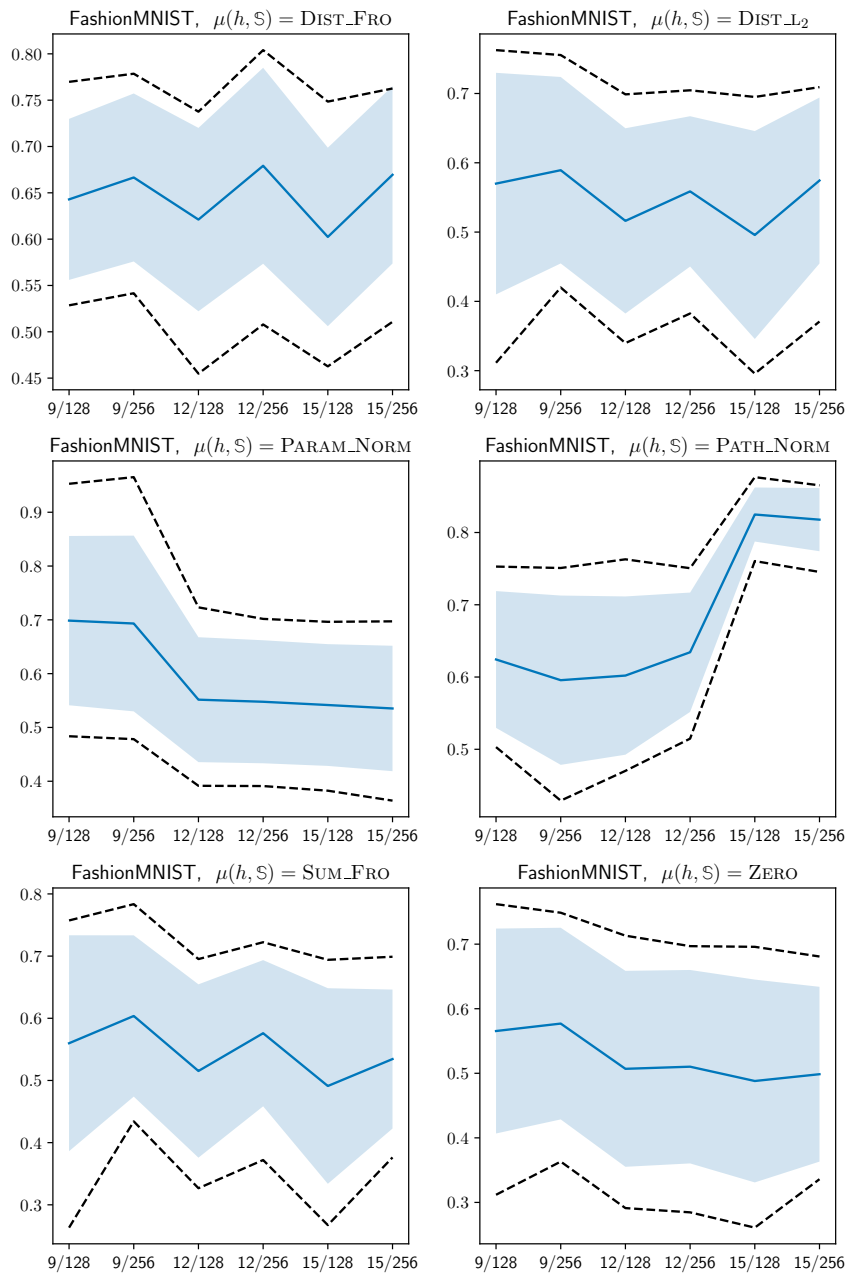


Figure G.9. Influence of the depth and the width in the x-axis as “depth/width”. The (solid) lines are the mean values computed on the different values of α ; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

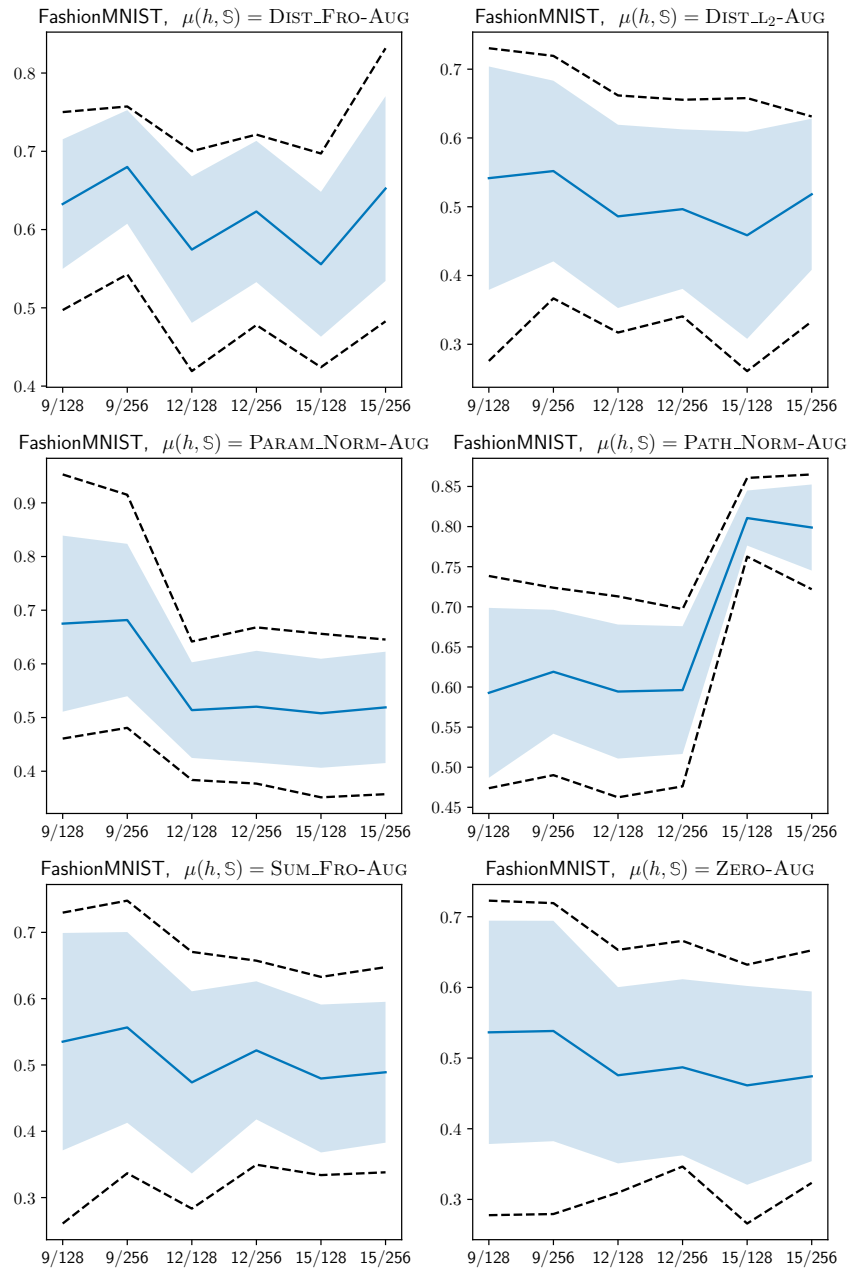


Figure G.10. Influence of the depth and the width in the x-axis as “depth/width”. The (solid) lines are the mean values computed on the different values of α ; the shadows are the standard deviation. The dashed lines are the minimum and the maximum values.

REFERENCES

MARTÍN ABADI *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. (2015)

—— Cited on page 35.

JEAN-BAPTISTE ALAYRAC, JONATHAN UESATO, PO-SEN HUANG, ALHUSSEIN FAWZI, ROBERT STANFORTH, and PUSHMEET KOHLI. Are Labels Required for Improving Adversarial Robustness? *NeurIPS*. (2019)

—— Cited on page 202.

PIERRE ALQUIER. User-friendly introduction to PAC-Bayes bounds. *CoRR*. abs/2110.11216. (2021)

—— Cited on page 64.

PIERRE ALQUIER, JAMES RIDGWAY, and NICOLAS CHOPIN. On the properties of variational approximations of Gibbs posteriors. *Journal of Machine Learning Research*. (2016)

—— Cited on page 181.

AMIRAN AMBROLADZE, EMILIO PARRADO-HERNÁNDEZ, and JOHN SHAWE-TAYLOR. Tighter PAC-Bayes Bounds. *NIPS*. (2006)

—— Cited on page 158.

PETER BARTLETT and SHAHAR MENDELSON. Rademacher and Gaussian Complexities: Risk Bounds and Structural Results. *Journal of Machine Learning Research*. (2002)

—— Cited on pages 19, 46, 47, 193.

BAPTISTE BAUVIN, CÉCILE CAPPONI, JEAN-FRANCIS ROY, and FRANÇOIS LAVIOLETTE. Fast greedy C -bound minimization with guarantees. *Machine Learning*. (2020)

—— Cited on pages 102, 106, 121, 122.

LUC BÉGIN, PASCAL GERMAIN, FRANÇOIS LAVIOLETTE, and JEAN-FRANCIS ROY. PAC-Bayesian Theory for Transductive Learning. *AISTATS*. (2014)

—— Cited on page 65.

LUC BÉGIN, PASCAL GERMAIN, FRANÇOIS LAVIOLETTE, and JEAN-FRANCIS ROY. PAC-Bayesian Bounds based on the Rényi Divergence. *AISTATS*. (2016)

—— Cited on pages 73, 74, 151, 154, 225.

AURÉLIEN BELLET, AMAURY HABRARD, EMILIE MORVANT, and MARC SEBBAN. Learning A Priori Constrained Weighted Majority Votes. *Machine Learning*. (2014)

—— Cited on pages 19, 57, 102.

SHAI BEN-DAVID, JOHN BLITZER, KOBY CRAMMER, ALEX KULESZA, FERNANDO PEREIRA, and JENNIFER WORTMAN VAUGHAN. A theory of learning from different domains. *Machine Learning*. (2010)

—— Cited on page 202.

JULIAN BESAG. Comments on “Representations of knowledge in complex systems” by U. Grenander and MI Miller. *Journal of the Royal Statistical Society, Series B*. (1994)

—— Cited on page 185.

BATTISTA BIGGIO, IGINO CORONA, DAVIDE MAIORCA, BLAINE NELSON, NEDIM SRNDIC, PAVEL LASKOV, GIORGIO GIACINTO, and FABIO ROLI. Evasion Attacks against Machine Learning at Test Time. *ECML PKDD*. (2013)

—— Cited on page 82.

FELIX BIGGS and BENJAMIN GUEDJ. Differentiable PAC-Bayes Objectives with Partially Aggregated Neural Networks. *Entropy*. (2021)

—— Cited on page 150.

FELIX BIGGS and BENJAMIN GUEDJ. On Margins and Derandomisation in PAC-Bayes. *AISTATS*. (2022)

—— Cited on pages 150, 152.

FELIX BIGGS, VALENTINA ZANTEDESCHI, and BENJAMIN GUEDJ. On Margins and Generalisation for Voting Classifiers. *NeurIPS*. (2020)

—— Cited on pages 171, 172.

CHRISTOPHER BISHOP. Pattern recognition and machine learning (5th Edition). *Information science and statistics*. Springer. (2007)

—— Cited on pages 17, 65.

References

GILLES BLANCHARD and FRANÇOIS FLEURET. Occam's Hammer. *COLT*. (2007)
—— Cited on pages 20, 56, 75, 77, 141, 150, 155, 156, 159, 162, 173, 178, 201, 232, 236, 316.

ROBERT BOIK and JAMES ROBINSON-COX. Derivatives of the incomplete beta function. *Journal of Statistical Software*. (1999)
—— Cited on page 138.

BERNHARD BOSER, ISABELLE GUYON, and VLADIMIR VAPNIK. A training algorithm for optimal margin classifiers. *COLT*. (1992)
—— Cited on page 35.

STÉPHANE BOUCHERON, GÁBOR LUGOSI, and PASCAL MASSART. Concentration Inequalities - A Nonasymptotic Theory of Independence. *Oxford University Press*. (2013)
—— Cited on page 43.

OLIVIER BOUSQUET and ANDRÉ ELISSEEFF. Stability and Generalization. *Journal of Machine Learning Research*. (2002)
—— Cited on pages 19, 48, 49, 195.

OLIVIER BOUSQUET, YEGOR KLOCHKOV, and NIKITA ZHIVOTOVSKIY. Sharper Bounds for Uniformly Stable Algorithms. *COLT*. (2020)
—— Cited on page 50.

STEPHEN BOYD and LIEVEN VANDENBERGHE. Convex Optimization. *Cambridge University Press*. (2004)
—— Cited on pages 109, 114, 228, 229, 242.

JAMES BRADBURY *et al.* JAX: composable transformations of Python+NumPy programs. (2018)
—— Cited on page 35.

LEO BREIMAN. Bagging Predictors. *Machine Learning*. (1996)
—— Cited on pages 19, 20, 56, 57, 203.

LEO BREIMAN. Random Forests. *Machine Learning*. (2001)
—— Cited on pages 19, 20, 57, 60, 62, 102, 105, 203.

LEO BREIMAN, JEROME FRIEDMAN, RICHARD OLSHEN, and CHARLES STONE. Classification and Regression Trees. *Wadsworth*. (1984)

—— Cited on pages 32, 34.

CLÉMENT CANONNE. A short note on an inequality between KL and TV. *CoRR*. abs/2202.07198. (2022)

—— Cited on page 219.

NICHOLAS CARLINI and DAVID WAGNER. Towards Evaluating the Robustness of Neural Networks. *IEEE Symposium on Security and Privacy*. (2017)

—— Cited on pages 82, 85.

YAIR CARMON, ADITI RAGHUNATHAN, LUDWIG SCHMIDT, JOHN DUCHI, and PERCY LIANG. Unlabeled Data Improves Adversarial Robustness. *NeurIPS*. (2019)

—— Cited on page 202.

OLIVIER CATONI. Statistical learning theory and stochastic optimization: Ecole d'Eté de Probabilités de Saint-Flour, XXXI-2001. *Springer Science & Business Media*. (2004)

—— Cited on page 181.

OLIVIER CATONI. Pac-Bayesian Supervised Classification: The Thermodynamics of Statistical Learning. *CoRR*. abs/0712.0248. (2007)

—— Cited on pages 20, 52, 56, 66–68, 70, 74–76, 78, 134, 141, 150, 155, 156, 159, 162, 173, 178, 181, 201, 222, 227, 228, 231, 286, 316, 324.

HERMAN CHERNOFF. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*. (1952)

—— Cited on pages 234, 235.

TZUU-SHUH CHIANG, CHII-RUEY HWANG, and SHUENN JYI SHEU. Diffusion for global optimization in \mathbb{R}^n . *SIAM Journal on Control and Optimization*. (1987)

—— Cited on page 181.

SIDDHARTHA CHIB and EDWARD GREENBERG. Understanding the Metropolis-Hastings Algorithm. *The american statistician*. (1995)

—— Cited on page 185.

JEREMY COHEN, ELAN ROSENFELD, and ZICO KOLTER. Certified Adversarial Robustness via Randomized Smoothing. *ICML*. (2019)

—— Cited on page 85.

References

RONAN COLLOBERT, KORAY KAVUKCUOGLU, and CLÉMENT FARABET. Torch7: A matlab-like environment for machine learning. *NIPS 2011 BigLearn Workshop*. (2011)

—— Cited on page 35.

NICOLAS DE CONDORCET. Essai sur l'Application de l'Analyse à la Probabilité des Décisions Rendues à la Pluralité des Voix. *Imprimerie Royale*. (1785)

—— Cited on page 19.

CORINNA CORTES and VLADIMIR VAPNIK. Support-vector networks. *Machine Learning*. (1995)

—— Cited on pages 35, 57.

THOMAS COVER and PETER HART. Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*. (1967)

—— Cited on page 57.

THOMAS COVER and JOY THOMAS. Elements of Information Theory (2nd edition). *Wiley*. (2006)

—— Cited on pages 218, 219.

JOHN DANSKIN. The Theory of Max-Min, with Applications. *SIAM Journal on Applied Mathematics*. (1966)

—— Cited on pages 114, 272.

ZHUN DENG, LINJUN ZHANG, AMIRATA GHORBANI, and JAMES ZOU. Improving Adversarial Robustness via Unlabeled Out-of-Domain Data. *AISTATS*. (2021)

—— Cited on page 202.

PHILIP DERBEKO, RAN EL-YANIV, and RON MEIR. Explicit Learning Curves for Transduction and Application to Clustering and Compression Algorithms. *Journal of Artificial Intelligence Research*. (2004)

—— Cited on page 65.

STEVEN DIAMOND and STEPHEN BOYD. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*. (2016)

—— Cited on page 113.

THOMAS DIETTERICH. Ensemble methods in machine learning. *International workshop on multiple classifier systems*. (2000)

—— Cited on pages 19, 20, 102.

MONROE DONSKER and SRINIVASA VARADHAN. Asymptotic evaluation of certain Markov process expectations for large time - III. *Communications on pure and applied Mathematics*. (1976)

—— Cited on page 73.

DHEERU DUA and CASEY GRAFF. UCI Machine Learning Repository. (2017). URL: <http://archive.ics.uci.edu/ml>

—— Cited on pages 121, 139.

GINTARE KAROLINA DZIUGAITE, ALEXANDRE DROUIN, BRADY NEAL, NITARSHAN RAJKUMAR, ETHAN CABALLERO, LINBO WANG, IOANNIS MITLIAGKAS, and DANIEL M. ROY. In search of robust measures of generalization. *NeurIPS*. (2020)

—— Cited on pages 176, 186.

GINTARE KAROLINA DZIUGAITE, KYLE HSU, WASEEM GHARBIH, GABRIEL ARPINO, and DANIEL ROY. On the role of data in PAC-Bayes. *AISTATS*. (2021)

—— Cited on pages 66, 93, 135, 184.

GINTARE KAROLINA DZIUGAITE and DANIEL ROY. Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data. *UAI*. (2017)

—— Cited on pages 23, 158, 161, 317.

GINTARE KAROLINA DZIUGAITE and DANIEL ROY. Data-dependent PAC-Bayes priors via differential privacy. *NeurIPS*. (2018)

—— Cited on pages 93, 161, 162, 169, 187, 318.

TIM VAN ERVEN and PETER HARREMOËS. Rényi Divergence and Kullback-Leibler Divergence. *IEEE Transactions on Information Theory*. (2014)

—— Cited on pages 74, 156.

AMEDEO ROBERTO ESPOSITO, MICHAEL GASTPAR, and IBRAHIM ISSA. Robust Generalization via α -Mutual Information. *International Zurich Seminar on Information and Communication*. (2020)

—— Cited on page 321.

FARZAN FARNIA, JESSE ZHANG, and DAVID TSE. Generalizable Adversarial Training via Spectral Normalization. *ICLR*. (2019)

—— Cited on page 85.

References

VITALY FELDMAN and JAN VONDRÁK. Generalization Bounds for Uniformly Stable Algorithms. *NeurIPS*. (2018)

—— Cited on page 50.

VITALY FELDMAN and JAN VONDRÁK. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. *COLT*. (2019)

—— Cited on page 50.

MIKHAIL FIGURNOV, SHAKIR MOHAMED, and ANDRIY MNIH. Implicit Reparameterization Gradients. *NeurIPS*. (2018)

—— Cited on page 138.

YOAV FREUND. Self Bounding Learning Algorithms. *COLT*. (1998)

—— Cited on pages 19, 23, 60, 64, 83, 92, 102, 123, 126, 136, 153.

YOAV FREUND, RAJ IYER, ROBERT SCHAPIRE, and YORAM SINGER. An Efficient Boosting Algorithm for Combining Preferences. *ICML*. (1998)

—— Cited on page 57.

YOAV FREUND, RAJ IYER, ROBERT SCHAPIRE, and YORAM SINGER. An Efficient Boosting Algorithm for Combining Preferences. *Journal of Machine Learning Research*. (2003)

—— Cited on page 57.

YOAV FREUND and ROBERT SCHAPIRE. Experiments with a New Boosting Algorithm. *ICML*. (1996)

—— Cited on pages 19, 20, 56, 57, 203.

JEROME FRIEDMAN. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*. (2001)

—— Cited on pages 57, 203.

ROY FROSTIG, MATTHEW JAMES JOHNSON, and CHRIS LEARY. Compiling machine learning programs via high-level tracing. *Systems for Machine Learning*. (2018)

—— Cited on page 35.

TOMER GALANTI, LIOR WOLF, and TAMIR HAZAN. A theoretical framework for deep transfer learning. *Information and Inference: A Journal of the IMA*. (2016)

—— Cited on page 202.

PASCAL GERMAIN, FRANCIS BACH, ALEXANDRE LACOSTE, and SIMON LACOSTE-JULIEN. PAC-Bayesian Theory Meets Bayesian Inference. *NIPS*. (2016)

—— Cited on page 65.

PASCAL GERMAIN, AMAURY HABRARD, FRANÇOIS LAVIOLETTE, and EMILIE MORVANT. PAC-Bayes and domain adaptation. *Neurocomputing*. (2020)

—— Cited on pages 65, 158, 202.

PASCAL GERMAIN, ALEXANDRE LACASSE, FRANÇOIS LAVIOLETTE, and MARIO MARCHAND. PAC-Bayesian Learning of Linear Classifiers. *ICML*. (2009)

—— Cited on pages 22, 66, 67, 70, 73, 75, 76, 152, 158, 162, 164, 220, 238.

PASCAL GERMAIN, ALEXANDRE LACASSE, FRANÇOIS LAVIOLETTE, MARIO MARCHAND, and JEAN-FRANCIS ROY. Risk Bounds for the Majority Vote: From a PAC-Bayesian Analysis to a Learning Algorithm. *Journal of Machine Learning Research*. (2015)

—— Cited on pages 60, 61, 64, 104–107, 113, 122, 138, 213, 242, 251, 267.

MANUEL GIL, FADY ALAJAJI, and TAMÁS LINDER. Rényi divergence measures for commonly used univariate continuous distributions. *Information Sciences*. (2013)

—— Cited on pages 314, 332.

XAVIER GLOROT and YOSHUA BENGIO. Understanding the difficulty of training deep feedforward neural networks. *AISTATS*. (2010)

—— Cited on pages 96, 163.

IAN GOODFELLOW, YOSHUA BENGIO, and AARON COURVILLE. Deep Learning. *Adaptive computation and machine learning*. MIT Press. (2016)

—— Cited on pages 35, 162.

IAN GOODFELLOW, JONATHON SHLENS, and CHRISTIAN SZEGEDY. Explaining and Harnessing Adversarial Examples. *ICLR*. (2015)

—— Cited on pages 23, 82, 84.

THORE GRAEPEL, RALF HERBRICH, and JOHN SHAWE-TAYLOR. PAC-Bayesian Compression Bounds on the Prediction Error of Learning Algorithms for Classification. *Machine Learning*. (2005)

—— Cited on pages 19, 57, 203.

References

MICHAEL GRANT, STEPHEN BOYD, and YINYU YE. Disciplined Convex Programming. *Global optimization*. (2006)

—— Cited on page 113.

BENJAMIN GUEDJ. A Primer on PAC-Bayesian Learning. *CoRR*. abs/1901.05353. (2019)

—— Cited on page 64.

KAIMING HE, XIANGYU ZHANG, SHAOQING REN, and JIAN SUN. Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification. *IEEE ICCV*. (2015)

—— Cited on pages 163, 186.

KAIMING HE, XIANGYU ZHANG, SHAOQING REN, and JIAN SUN. Deep Residual Learning for Image Recognition. *CVPR*. (2016)

—— Cited on page 163.

DAN HENDRYCKS and THOMAS DIETTERICH. Benchmarking Neural Network Robustness to Common Corruptions and Perturbations. *ICLR*. (2019)

—— Cited on pages 85, 86.

WAYNE IBA and PAT LANGLEY. Induction of One-Level Decision Trees. *Ninth International Workshop on Machine Learning*. (1992)

—— Cited on pages 32–34.

MARTIN JANKOWIAK and FRITZ OBERMEYER. Pathwise Derivatives Beyond the Reparameterization Trick. *ICML*. (2018)

—— Cited on page 138.

YIDING JIANG, BEHNAM NEYSHABUR, HOSSEIN MOBAHI, DILIP KRISHNAN, and SAMY BENGIO. Fantastic Generalization Measures and Where to Find Them. *ICLR*. (2019)

—— Cited on pages 175, 176, 186.

YIDING JIANG *et al.* Methods and Analysis of The First Competition in Predicting Generalization of Deep Learning. *NeurIPS 2020 Competition and Demonstration Track*. (2021)

—— Cited on page 176.

PETER KAIROUZ *et al.* Advances and Open Problems in Federated Learning. *Foundations and Trends in Machine Learning*. (2021)

—— Cited on page 203.

SHAM KAKADE, KARTHIK SRIDHARAN, and AMBUJ TEWARI. On the Complexity of Linear Prediction: Risk Bounds, Margin Bounds, and Regularization. *NIPS*. (2008)

—— Cited on pages 47, 180.

KENJI KAWAGUCHI, LESLIE PACK KAEHLING, and YOSHUA BENGIO. Generalization in Deep Learning. *CoRR*. abs/1710.05468. (2017)

—— Cited on pages 19, 57.

MAURICE KENDALL. A new measure of rank correlation. *Biometrika*. (1938)

—— Cited on page 176.

HOEL KERVADEC, JOSE DOLZ, JING YUAN, CHRISTIAN DESROSIERS, ERIC GRANGER, and ISMAIL BEN AYED. Constrained Deep Networks: Lagrangian Optimization via Log-Barrier Extensions. *CoRR*. abs/1904.04205. (2019)

—— Cited on page 109.

JUSTIN KHIM and PO-LING LOH. Adversarial Risk Bounds for Binary Classification via Function Transformation. *CoRR*. abs/1810.09519. (2018)

—— Cited on page 85.

DIEDERIK KINGMA and JIMMY BA. Adam: A Method for Stochastic Optimization. *ICLR*. (2015)

—— Cited on pages 96, 108, 163.

PETER KONTSCHIEDER, MADALINA FITERAU, ANTONIO CRIMINISI, and SAMUEL ROTA BULÒ. Deep Neural Decision Forests. *IJCAI*. (2016)

—— Cited on pages 95, 253, 254.

ALEX KRIZHEVSKY. Learning Multiple Layers of Features from Tiny Images. MA thesis. University of Toronto, (2009)

—— Cited on page 163.

ALEX KRIZHEVSKY, ILYA SUTSKEVER, and GEOFFREY HINTON. ImageNet Classification with Deep Convolutional Neural Networks. *NIPS*. (2012)

—— Cited on page 35.

References

LUDMILA KUNCHEVA. Combining pattern classifiers: methods and algorithms. *John Wiley & Sons*. (2014)

—— Cited on pages 20, 102.

ALEXEY KURAKIN, IAN GOODFELLOW, and SAMY BENGIO. Adversarial Machine Learning at Scale. *ICLR*. (2017)

—— Cited on pages 82, 84, 95.

ALEXANDRE LACASSE. Bornes PAC-Bayes et algorithmes d'apprentissage. PhD thesis. (2010)

—— Cited on pages 70, 127.

ALEXANDRE LACASSE, FRANÇOIS LAVIOLETTE, MARIO MARCHAND, PASCAL GERMAIN, and NICOLAS USUNIER. PAC-Bayes Bounds for the Risk of the Majority Vote and the Variance of the Gibbs Classifier. *NIPS*. (2006)

—— Cited on pages 60–63, 100, 102, 105, 107, 111, 122, 268, 270.

ALEXANDRE LACASSE, FRANÇOIS LAVIOLETTE, MARIO MARCHAND, and FRANCIS TURGEON-BOUTIN. Learning with Randomized Majority Votes. *ECML PKDD*. (2010)

—— Cited on pages 126, 127, 138, 285, 286.

JOHN LANGFORD. Tutorial on Practical Prediction Theory for Classification. *Journal of Machine Learning Research*. (2005)

—— Cited on pages 75, 234, 235.

JOHN LANGFORD and JOHN SHAWE-TAYLOR. PAC-Bayes & Margins. *NIPS*. (2002)

—— Cited on pages 22, 39, 60, 75, 150.

FRANÇOIS LAVIOLETTE, EMILIE MORVANT, LIVA RALAIVOLA, and JEAN-FRANCIS ROY. Risk upper bounds for general ensemble methods with an application to multi-class classification. *Neurocomputing*. (2017)

—— Cited on pages 58, 63, 65, 100, 104–106, 129, 140, 141, 172.

YANN LECUN, CORINNA CORTES, and CHRISTOPHER BURGES. THE MNIST DATASET of handwritten digits. (1998). URL: <http://yann.lecun.com/exdb/mnist/>

—— Cited on pages 31, 96, 121, 139, 163, 186.

GAËL LETARTE, PASCAL GERMAIN, BENJAMIN GUEDJ, and FRANÇOIS LAVIOLETTE. Dichotomize and Generalize: PAC-Bayesian Binary Activated Deep Neural Networks. *NeurIPS*. (2019)

—— Cited on pages 22, 150, 158, 164.

GUY LEVER, FRANÇOIS LAVIOLETTE, and JOHN SHAWE-TAYLOR. Tighter PAC-Bayes bounds through distribution-dependent priors. *Theoretical Computer Science*. (2013)

—— Cited on pages 93, 324.

MIN LIN, QIANG CHEN, and SHUICHENG YAN. Network in network. *CoRR*. abs/1312.4400. (2013)

—— Cited on page 186.

BEN LONDON. A PAC-Bayesian Analysis of Randomized Learning with Application to Stochastic Gradient Descent. *NIPS*. (2017)

—— Cited on page 65.

STEPHAN SLOTH LORENZEN, CHRISTIAN IGEL, and YEVGENY SELDIN. On PAC-Bayesian bounds for random forests. *Machine Learning*. (2019)

—— Cited on pages 102, 201.

ALEKSANDER MADRY, ALEKSANDAR MAKELOV, LUDWIG SCHMIDT, DIMITRIS TSIPRAS, and ADRIAN VLADU. Towards Deep Learning Models Resistant to Adversarial Attacks. *ICLR*. (2018)

—— Cited on pages 82, 85, 95, 272.

ANDRÉS MASEGOSA, STEPHAN SLOTH LORENZEN, CHRISTIAN IGEL, and YEVGENY SELDIN. Second Order PAC-Bayesian Bounds for the Weighted Majority Vote. *NeurIPS*. (2020)

—— Cited on pages 22, 39, 60, 61, 64, 92, 100, 102, 105, 121, 122, 128, 138, 201, 214.

ANDREAS MAURER. A Note on the PAC Bayesian Theorem. *CoRR*. cs.LG/0411099. (2004)

—— Cited on pages 52, 70, 238, 242, 314, 316, 321, 349.

DAVID MCALLESTER. Some PAC-Bayesian Theorems. *COLT*. (1998)

—— Cited on pages 19, 56, 90, 134.

References

- DAVID MCALLESTER. Some PAC-Bayesian Theorems. *Machine Learning*. (1999)
—— Cited on pages 52, 64.
- DAVID MCALLESTER. PAC-Bayesian Stochastic Model Selection. *Machine Learning*. (2003)
—— Cited on pages 60, 66, 67, 74, 102, 106, 107, 110, 221, 227.
- WARREN MCCULLOCH and WALTER PITTS. A Logical Calculus of Ideas Immanent in Nervous Activity. *Bulletin of Mathematical Biophysics*. (1943)
—— Cited on page 35.
- COLIN MCDIARMID. On the method of bounded differences. *Surveys in Combinatorics*. (1989)
—— Cited on pages 47, 49.
- BRENDAN MCMAHAN, EIDER MOORE, DANIEL RAMAGE, SETH HAMPSON, and BLAISE AGÜERA Y ARCAS. Communication-Efficient Learning of Deep Networks from Decentralized Data. *AISTATS*. (2017)
—— Cited on page 203.
- DANIEL MCNAMARA and MARIA-FLORINA BALCAN. Risk Bounds for Transferring Representations With and Without Fine-Tuning. *ICML*. (2017)
—— Cited on page 202.
- ZAKARIA MHAMMEDI, PETER GRÜNWARD, and BENJAMIN GUEDJ. PAC-Bayes Un-Expected Bernstein Inequality. *NeurIPS*. (2019)
—— Cited on pages 135, 136.
- MARVIN MINSKY and SEYMOUR PAPERT. Perceptrons: an introduction to computational geometry. *The MIT Press*. (1972)
—— Cited on pages 35, 36.
- MEHRYAR MOHRI, AFSHIN ROSTAMIZADEH, and AMEET TALWALKAR. Foundations of Machine Learning. *Adaptive computation and machine learning*. MIT Press. (2012)
—— Cited on pages 17, 44–46.
- OMAR MONTASSER, STEVE HANNEKE, and NATHAN SREBRO. VC Classes are Adversarially Robustly Learnable, but Only Improperly. *COLT*. (2019)
—— Cited on page 85.

OMAR MONTASSER, STEVE HANNEKE, and NATHAN SREBRO. Reducing Adversarially Robust Learning to Non-Robust PAC Learning. *NeurIPS*. (2020)

—— Cited on page 85.

VAISHNAVH NAGARAJAN and ZICO KOLTER. Deterministic PAC-Bayesian generalization bounds for deep networks via generalizing noise-resilience. *ICLR*. (2019a)

—— Cited on page 152.

VAISHNAVH NAGARAJAN and ZICO KOLTER. Uniform convergence may be unable to explain generalization in deep learning. *NeurIPS*. (2019b)

—— Cited on pages 85, 150, 194.

YURII NESTEROV. Smooth minimization of non-smooth functions. *Mathematical Programming*. (2005)

—— Cited on page 131.

BEHNAM NEYSHABUR, SRINADH BHOJANAPALLI, and NATHAN SREBRO. A PAC-Bayesian Approach to Spectrally-Normalized Margin Bounds for Neural Networks. *ICLR*. (2018)

—— Cited on pages 150, 152.

YUKI OHNISHI and JEAN HONORIO. Novel Change of Measure Inequalities with Applications to PAC-Bayesian Bounds and Monte Carlo Estimation. *AISTATS*. (2021)

—— Cited on pages 73, 248, 251.

FRANCESCO ORABONA and TATIANA TOMMASI. Training Deep Networks without Learning Rates Through Coin Betting. *NIPS*. (2017)

—— Cited on pages 108, 121, 139, 203.

NICOLAS PAPERNOT, PATRICK MCDANIEL, SOMESH JHA, MATT FREDRIKSON, Z. BERKAY CELIK, and ANANTHRAM SWAMI. The Limitations of Deep Learning in Adversarial Settings. *IEEE EuroS&P*. (2016)

—— Cited on page 82.

EMILIO PARRADO-HERNÁNDEZ, AMIRAN AMBROLADZE, JOHN SHAWE-TAYLOR, and SHILIANG SUN. PAC-Bayes Bounds with Data Dependent Priors. *Journal of Machine Learning Research*. (2012)

—— Cited on pages 66, 93, 135, 184.

References

ADAM PASZKE *et al.* PyTorch: An Imperative Style, High-Performance Deep Learning Library. *NeurIPS*. (2019)

—— Cited on pages 35, 111, 138.

FABIAN PEDREGOSA *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*. (2011)

—— Cited on page 35.

ANASTASIA PENTINA and CHRISTOPH H. LAMPERT. A PAC-Bayesian bound for Lifelong Learning. *ICML*. (2014)

—— Cited on page 129.

MARÍA PÉREZ-ORTIZ, OMAR RIVASPLATA, JOHN SHAWE-TAYLOR, and CSABA SZEPESVÁRI. Tighter Risk Certificates for Neural Networks. *Journal of Machine Learning Research*. (2021)

—— Cited on pages 23, 160–162, 164, 184.

RAFAEL PINOT, LAURENT MEUNIER, ALEXANDRE ARAUJO, HISASHI KASHIMA, FLORIAN YGER, CÉDRIC GOUY-PAILLER, and JAMAL ATIF. Theoretical evidence for adversarial robustness through randomization. *NeurIPS*. (2019)

—— Cited on page 85.

RAFAEL PINOT, LAURENT MEUNIER, FLORIAN YGER, CÉDRIC GOUY-PAILLER, YANN CHEVALEYRE, and JAMAL ATIF. On the robustness of randomized classifiers to adversarial examples. *Machine Learning*. (2022)

—— Cited on page 85.

ROSS QUINLAN. Induction of Decision Trees. *Machine Learning*. (1986)

—— Cited on page 34.

ROSS QUINLAN. C4.5: Programs for Machine Learning. *Morgan Kaufmann*. (1993)

—— Cited on page 34.

MAXIM RAGINSKY, ALEXANDER RAKHLIN, and MATUS TELGARSKY. Non-convex learning via Stochastic Gradient Langevin Dynamics: a nonasymptotic analysis. *COLT*. (2017)

—— Cited on page 181.

LIVA RALAIVOLA, MARIE SZAFRANSKI, and GUILLAUME STEMPEL. Chromatic PAC-Bayes Bounds for Non-IID Data: Applications to Ranking and Stationary β -Mixing

Processes. *Journal of Machine Learning Research*. (2010)

—— Cited on pages 90, 250.

IEVGEN REDKO, EMILIE MORVANT, AMAURY HABRARD, MARC SEBBAN, and YOUNÈS BENNANI. Advances in Domain Adaptation Theory. *Elsevier*. (2019)

—— Cited on page 202.

IEVGEN REDKO, EMILIE MORVANT, AMAURY HABRARD, MARC SEBBAN, and YOUNÈS BENNANI. A survey on domain adaptation theory. *CoRR*. abs/2004.11829. (2020)

—— Cited on page 202.

DAVID REEB, ANDREAS DOERR, SEBASTIAN GERWINN, and BARBARA RAKITSCH. Learning Gaussian Processes by Minimizing PAC-Bayesian Generalization Bounds. *NeurIPS*. (2018)

—— Cited on pages 71, 111.

KUI REN, TIANHANG ZHENG, and XUE LIU. Adversarial Attacks and Defenses in Deep Learning. *Engineering*. (2020)

—— Cited on page 84.

RAMI AL-RFOU *et al.* Theano: A Python framework for fast computation of mathematical expressions. *CoRR*. abs/1605.02688. (2016)

—— Cited on page 35.

OMAR RIVASPLATA, ILJA KUZBORSKIJ, CSABA SZEPESVÁRI, and JOHN SHAWE-TAYLOR. PAC-Bayes Analysis Beyond the Usual Bounds. *Advances in Neural Information Processing System (NeurIPS)*. (2020)

—— Cited on pages 76, 146, 150, 155, 159, 162, 171, 173, 178, 201, 230, 316.

FRANK ROSENBLATT. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*. (1958)

—— Cited on page 35.

JEAN-FRANCIS ROY, FRANÇOIS LAVIOLETTE, and MARIO MARCHAND. From PAC-Bayes Bounds to Quadratic Programs for Majority Votes. *ICML*. (2011)

—— Cited on pages 60, 102, 106, 121, 122.

JEAN-FRANCIS ROY, MARIO MARCHAND, and FRANÇOIS LAVIOLETTE. A Column Generation Bound Minimization Approach with PAC-Bayesian Generalization Guar-

References

antees. *AISTATS*. (2016)

—— Cited on pages 102, 105, 106, 266.

STUART RUSSELL and PETER NORVIG. *Artificial Intelligence: A Modern Approach* (4th Edition). *Pearson*. (2020)

—— Cited on page 17.

HADI SALMAN, JERRY LI, ILYA RAZENSHTEYN, PENGCHUAN ZHANG, HUAN ZHANG, SÉBASTIEN BUBECK, and GREG YANG. Provably Robust Deep Learning via Adversarially Trained Smoothed Classifiers. *NeurIPS*. (2019)

—— Cited on page 85.

ROBERT SCHAPIRE, YOAV FREUND, PETER BARLETT, and WEE SUN LEE. Boosting the margin: A new explanation for the effectiveness of voting methods. *The annals of statistics*. (1998)

—— Cited on page 171.

ROBERT SCHAPIRE and YORAM SINGER. Improved Boosting Algorithms using Confidence-Rated Predictions. *COLT*. (1998)

—— Cited on page 57.

ROBERT SCHAPIRE and YORAM SINGER. Improved Boosting Algorithms Using Confidence-rated Predictions. *Machine Learning*. (1999)

—— Cited on page 57.

ROBERT SCHAPIRE and YORAM SINGER. BoosTexter: A Boosting-based System for Text Categorization. *Machine Learning*. (2000)

—— Cited on page 57.

EDWARD SCHEINERMAN and DANIEL ULLMAN. Fractional Graph Theory: A Rational Approach to the Theory of Graphs. *Courier Corporation*. (2011)

—— Cited on page 250.

MATTHIAS SEEGER. PAC-Bayesian Generalisation Error Bounds for Gaussian Process Classification. *Journal of Machine Learning Research*. (2002)

—— Cited on pages 52, 66, 67, 69, 70, 74, 77, 90, 102, 106, 107, 110, 133, 134, 138, 223, 227, 286.

VERA SHALAEVA, ALIREZA FAKHRIZADEH ESFAHANI, PASCAL GERMAIN, and MIHÁLY PETRECKZY. Improved PAC-Bayesian Bounds for Linear Regression. *AAAI*.

(2020)

—— Cited on page 65.

SHAI SHALEV-SHWARTZ and SHAI BEN-DAVID. Understanding Machine Learning - From Theory to Algorithms. *Cambridge University Press*. (2014)

—— Cited on page 17.

JOHN SHAWE-TAYLOR and ROBERT WILLIAMSON. A PAC Analysis of a Bayesian Estimator. *COLT*. (1997)

—— Cited on pages 19, 52, 56, 64.

RAVID SHWARTZ-ZIV and NAFTALI TISHBY. Opening the Black Box of Deep Neural Networks via Information. *CoRR*. abs/1703.00810. (2017)

—— Cited on page 204.

ADRIAN SMITH and GARETH ROBERTS. Bayesian computation via the Gibbs sampler and related Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society, Series B (Methodological)*. (1993)

—— Cited on page 185.

JOST TOBIAS SPRINGENBERG, ALEXEY DOSOVITSKIY, THOMAS BROX, and MARTIN RIEDMILLER. Striving for Simplicity: The All Convolutional Net. *ICLR*. (2015)

—— Cited on page 163.

CHRISTIAN SZEGEDY, WOJCIECH ZAREMBA, ILYA SUTSKEVER, JOAN BRUNA, DUMITRU ERHAN, IAN GOODFELLOW, and ROB FERGUS. Intriguing properties of neural networks. *ICLR*. (2014)

—— Cited on pages 82, 83.

NIKLAS THIEMANN, CHRISTIAN IGEL, OLIVIER WINTENBERGER, and YEVGENY SELDIN. A Strongly Quasiconvex PAC-Bayesian Bound. *ALT*. (2017)

—— Cited on pages 105, 135, 156, 157, 308.

NAFTALI TISHBY, FERNANDO PEREIRA, and WILLIAM BIALEK. The information bottleneck method. *CoRR*. physics/0004057. (2000)

—— Cited on page 204.

AAD VAN DER VAART and JON WELLNER. Weak convergence and empirical processes. *Springer series in statistics*. Springer. (1996)

—— Cited on page 51.

References

LESLIE VALIANT. A Theory of the Learnable. *Communications of the ACM*. (1984)

—— Cited on page 43.

VLADIMIR VAPNIK. Statistical Learning Theory. *A Wiley-Interscience publication*. Wiley. (1998)

—— Cited on pages 41, 46.

VLADIMIR VAPNIK and ALEXEY CHERVONENKIS. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Doklady Akademii Nauk USSR*. (1968)

—— Cited on pages 18, 19, 41, 43, 45.

VLADIMIR VAPNIK and ALEXEY CHERVONENKIS. On the Uniform Convergence of Relative Frequencies of Events to Their Probabilities. *Theory of Probability and its Applications*. (1971)

—— Cited on pages 18, 19, 41, 43, 45, 193.

VLADIMIR VAPNIK and ALEXEY CHERVONENKIS. Theory of Pattern Recognition. *Nauka, Moscow*. (1974)

—— Cited on pages 18, 19, 22, 41.

SERGIO VERDÚ. α -mutual information. *Information Theory and Applications Workshop*. (2015)

—— Cited on pages 319, 320, 326.

MAX WELLING and YEE WHYE TEH. Bayesian Learning via Stochastic Gradient Langevin Dynamics. *ICML*. (2011)

—— Cited on page 181.

YIHONG WU. Lecture notes on: Information-theoretic methods for high-dimensional statistics. (2020)

—— Cited on page 219.

HAN XIAO, KASHIF RASUL, and ROLAND VOLLGRAF. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *CoRR*. abs/1708.07747. (2017)

—— Cited on pages 96, 121, 139, 163, 186.

AOLIN XU and MAXIM RAGINSKY. Information-theoretic analysis of generalization capability of learning algorithms. *NIPS*. (2017)

—— Cited on page 173.

HUAN XU and SHIE MANNOR. Robustness and Generalization. *COLT*. (2010)

—— Cited on page 48.

HUAN XU and SHIE MANNOR. Robustness and generalization. *Machine Learning*. (2012)

—— Cited on pages 48, 50, 195.

DONG YIN, KANNAN RAMCHANDRAN, and PETER BARTLETT. Rademacher Complexity for Adversarially Robust Generalization. *ICML*. (2019)

—— Cited on page 85.

MALIK YOUNSI. Proof of a Combinatorial Conjecture Coming from the PAC-Bayesian Machine Learning Theory. *CoRR*. abs/1209.0824. (2012)

—— Cited on page 269.

MALIK YOUNSI and ALEXANDRE LACASSE. A combinatorial conjecture from PAC-Bayesian machine learning. *CoRR*. abs/2006.01387. (2020)

—— Cited on page 269.

VALENTINA ZANTEDESCHI, MARIA-IRINA NICOLAE, and AMBRISH RAWAT. Efficient Defenses Against Adversarial Attacks. *ACM Workshop on Artificial Intelligence and Security, AISec@CCS*. (2017)

—— Cited on pages 82, 85, 86.

WENDA ZHOU, VICTOR VEITCH, MORGANE AUSTERN, RYAN ADAMS, and PETER ORBANZ. Non-vacuous Generalization Bounds at the ImageNet Scale: a PAC-Bayesian Compression Approach. *ICLR*. (2019)

—— Cited on pages 23, 158, 161.

JI ZHU, HUI ZOU, SAHARON ROSSET, and TREVOR HASTIE. Multi-class AdaBoost. *Statistics and its Interface*. (2009)

—— Cited on page 57.

Abstract. In machine learning, a model is learned from data to solve a task automatically. In the supervised classification setting, the model aims to predict the label associated with an input. The model is learned using a limited number of examples, each consisting of an input and its associated label. However, the model's performance on the examples, computed by the empirical risk, does not necessarily reflect the performance on the task, which is represented by the true risk. Moreover, since it is not computable, the true risk is upper-bounded by a generalization bound that mainly depends on two quantities: the empirical risk and a complexity measure. One way to learn a model is to minimize a bound by a type of algorithm called self-bounding. PAC-Bayesian bounds are well suited to the derivation of this type of algorithm. In this context, the first contribution consists in developing self-bounding algorithms that minimize PAC-Bayesian bounds to learn majority votes. If these bounds are well adapted to majority votes, their use for other models becomes less natural. To overcome this difficulty, a second contribution focuses on the disintegrated PAC-Bayesian bounds that are natural for more general models. In this framework, we provide the first empirical study of these bounds. In a third contribution, we derive bounds that allow us to incorporate complexity measures defined by the user.

Keywords. Machine Learning, Generalization, PAC-Bayesian Bound, Disintegrated PAC-Bayesian Bound, Self-Bounding Algorithm, Majority Vote, Neural Network, Complexity Measure.

Résumé. En apprentissage automatique, un modèle est appris à partir de données pour résoudre une tâche de manière automatique. Dans le cadre de la classification supervisée, le modèle vise à prédire la classe associée à une entrée. Le modèle est appris à l'aide d'un nombre limité d'exemples, chacun étant constitué d'une entrée et de sa classe associée. Cependant, la performance du modèle sur les exemples, calculée par le risque empirique, ne reflète pas nécessairement la performance sur la tâche qui est représentée par le risque réel. De plus, n'étant pas calculable, le risque réel est majoré pour obtenir une borne en généralisation qui dépend principalement de deux quantités : le risque empirique et une mesure de complexité. Une façon d'apprendre un modèle est de minimiser une borne par un type d'algorithme appelé auto-certié (ou auto-limitatif). Les bornes PAC-Bayésiennes sont bien adaptées à la dérivation de ce type d'algorithmes. Dans ce contexte, la première contribution consiste à développer des algorithmes auto-certiés qui minimisent des bornes PAC-Bayésiennes pour apprendre des votes de majorité. Si ces bornes sont bien adaptées aux votes de majorité, leur utilisation pour d'autres modèles devient moins naturelle. Pour pallier cette difficulté, une seconde contribution se concentre sur les bornes PAC-Bayésiennes désintégrées qui sont naturelles pour des modèles plus généraux. Dans ce cadre, nous apportons la première étude empirique de ces bornes. Dans une troisième contribution, nous dérivons des bornes permettant d'incorporer des mesures de complexité pouvant être définies par l'utilisateur.

Mot-clés. Apprentissage Automatique, Généralisation, Borne PAC-Bayésienne, Borne PAC-Bayésienne Désintégrée, Algorithme Auto-certié, Algorithme Auto-limitatif, Vote de Majorité, Réseau de Neurones, Mesure de Complexité.