



HAL
open science

Random Matrix Theory for AI: From Theory to Practice

Mohamed El Amine Seddik

► **To cite this version:**

Mohamed El Amine Seddik. Random Matrix Theory for AI: From Theory to Practice. Signal and Image Processing. Université Paris-Saclay, 2020. English. NNT : 2020UPASG010 . tel-03125586

HAL Id: tel-03125586

<https://theses.hal.science/tel-03125586>

Submitted on 29 Jan 2021

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.

Random Matrix Theory for AI: From Theory to Practice

Thèse de doctorat de l'Université Paris-Saclay

École doctorale n° 580 Sciences et Technologies de
l'Information et de la Communication (STIC)
Spécialité de doctorat : Traitement du signal et des images
Unité de recherche : Université Paris-Saclay, CNRS, CentraleSupélec,
Laboratoire des signaux et systèmes, 91190, Gif-sur-Yvette, France.
Réfèrent : CentraleSupélec

**Thèse présentée et soutenue en visioconférence totale, le
03/11/2020, par**

Mohamed El Amine SEDDIK

Composition du jury:

Alexandre Gramfort Directeur de recherche, INRIA Paris	Président
Julie Delon Professeure, Université Paris Descartes	Rapporteuse & Examinatrice
Jamal Najim Directeur de recherche, Université Marne la Vallée	Rapporteur & Examineur
Mérouane Debbah Professeur, Huawei Labs & CentraleSupélec	Examineur
Florent Krzakala Professeur, EPFL	Examineur
Florent Chatelain Maître de conférences, Université de Grenoble-Alpes	Examineur
Romain Couillet Professeur, CentraleSupélec	Directeur
Mohamed Tamaazousti Ingénieur-chercheur, CEA Paris-Saclay	Co-encadrant

Titre: La théorie des matrices aléatoires pour l'IA : de la théorie à la pratique

Mots clés: Apprentissage automatique, théorie des matrices aléatoires, concentration de la mesure, réseaux de neurones, réseaux de neurones génératifs adverses

Résumé: De nos jours, l'IA repose en grande partie sur l'utilisation de données de grande taille et sur des méthodes d'apprentissage machine améliorées qui consistent à développer des algorithmes de classification et d'inférence en tirant parti de grands ensembles de données de grande taille. Ces grandes dimensions induisent de nombreux phénomènes contre-intuitifs, conduisant généralement à une mauvaise compréhension du comportement de nombreux algorithmes d'apprentissage machine souvent conçus avec des intuitions de petites dimensions de données. En tirant parti du cadre multidimensionnel (plutôt que d'en souffrir), la théorie des matrices aléatoires (RMT) est capable de prédire les performances de nombreux algorithmes non linéaires aussi complexes que certains réseaux de neurones aléatoires, ainsi que de nombreuses méthodes du noyau telles que les SVM, la classification semi-supervisée, l'analyse en composantes principales ou le regroupement spectral. Pour caractériser théoriquement les performances de ces algorithmes, le modèle de données sous-jacent est souvent un modèle de mélange gaussien (MMG) qui semble être une hypothèse forte étant donné la structure complexe des données réelles (par exemple, des images). En outre, la performance des algorithmes d'apprentissage automatique dépend du choix de la représentation des données (ou des caractéristiques) sur lesquelles ils sont appliqués. Encore une fois, considérer les représentations de données comme des vecteurs gaussiens semble être une hypothèse assez restrictive. S'appuyant sur la théorie des matrices aléatoires, cette thèse vise à aller au-delà de la simple hypothèse du MMG, en étudiant les outils classiques d'apprentissage machine sous l'hypothèse de vecteurs aléatoires concentrés qui généralisent

les vecteurs Gaussiens. Cette hypothèse est particulièrement motivée par l'observation que l'on peut utiliser des modèles génératifs (par exemple, les GAN) pour concevoir des structures de données complexes et réalistes telles que des images, grâce à des transformations Lipschitzienne de vecteurs gaussiens. Cela suggère notamment que l'hypothèse de concentration sur les données mentionnée ci-dessus est un modèle approprié pour les données réelles et qui est tout aussi mathématiquement accessible que les MMG. Par conséquent, nous démontrons à travers cette thèse, en nous appuyant sur les GANs, l'intérêt de considérer le cadre des vecteurs concentrés comme un modèle pour les données réelles. En particulier, nous étudions le comportement des matrices de Gram aléatoires qui apparaissent au cœur de divers modèles linéaires, des matrices à noyau qui apparaissent dans les méthodes à noyau et également des méthodes de classification qui reposent sur une solution implicite (par exemple, la couche de Softmax dans les réseaux de neurones), avec des données aléatoires supposées concentrées. En particulier, la compréhension du comportement de ces matrices/méthodes, pour des données concentrées, nous permet de caractériser les performances (sur des données réelles si nous les assimilons à des vecteurs concentrés) de nombreux algorithmes d'apprentissage machine, tels que le clustering spectral, les SVM, l'analyse en composantes principales et l'apprentissage par transfert. L'analyse de ces méthodes pour des données concentrées donne le résultat surprenant qu'elles ont asymptotiquement le même comportement que pour les données de MMG. Ce résultat suggère fortement l'aspect d'universalité des grands classificateurs d'apprentissage machine par rapport à la distribution sous-jacente des données.

Title: Random Matrix Theory for AI: From Theory to Practice

Keywords: Machine learning, random matrix theory, concentration of measure, neural networks, GANs

Abstract: AI nowadays relies largely on using large data and enhanced machine learning methods which consist in developing classification and inference algorithms leveraging large datasets of large sizes. These large dimensions induce many counter-intuitive phenomena, leading generally to a misunderstanding of the behavior of many machine learning algorithms often designed with small data dimension intuitions. By taking advantage of (rather than suffering from) the multidimensional setting, random matrix theory (RMT) is able to predict the performance of many non-linear algorithms as complex as some random neural networks as well as many kernel methods such as Support Vector Machines, semi-supervised classification, principal component analysis or spectral clustering. To characterize the performance of these algorithms theoretically, the underlying data model is often a Gaussian mixture model (GMM) which seems to be a strong assumption given the complex structure of real data (e.g., images). Furthermore, the performance of machine learning algorithms depends on the choice of data representation (or features) on which they are applied. Once again, considering data representations as Gaussian vectors seems to be quite a restrictive assumption. Relying on random matrix theory, this thesis aims at going beyond the simple GMM hypothesis, by studying classical machine learning tools under the hypothesis of Lipschitz-ally transformed Gaussian vectors also called concentrated random vectors, and which are more generic than Gaussian vectors. This hypothesis is particularly motivated by the observation that one can use generative models (e.g., GANs) to design complex and realistic data structures such as im-

ages, through Lipschitz-ally transformed Gaussian vectors. This notably suggests that making the aforementioned concentration assumption on data is a suitable model for real data and which is just as mathematically accessible as GMM models. Moreover, in terms of data representation, the concentration framework is compatible with one of the most widely used data representations in practice, namely deep neural nets (DNNs) representations, since they consist of a Lipschitz transformation of the input data (e.g., images). Therefore, we demonstrate through this thesis, leveraging on GANs, the interest in considering the framework of concentrated vectors as a model for real data. In particular, we study the behavior of random Gram matrices which appear at the core of various linear models, kernel matrices that appear in kernel methods, and also classification methods that rely on an implicit solution (e.g., Softmax layer in neural networks), with concentrated random inputs. Indeed, these methods are at the heart of many classifications, regression, and clustering machine learning algorithms. In particular, understanding the behavior of these matrices/methods, for concentrated data, allows us to characterize the performances (on real data if we assimilate them to concentrated vectors) of many machine learning algorithms, such as spectral clustering, SVMs, principal component analysis, and transfer learning. Analyzing these methods for concentrated data yields the surprising result that they have asymptotically the same behavior as for GMM data (with the same first and second-order statistics). This result strongly suggests the universality aspect of large machine learning classifiers w.r.t. the underlying data distribution.

Université Paris-Saclay

Espace Technologique / Immeuble Discovery

Route de l'Orme aux Merisiers RD 128 / 91190 Saint-Aubin, France

Random Matrix Theory for AI: From Theory to Practice

Mohamed El Amine SEDDIK

Under the supervision of
Prof. Romain COUILLET & Dr. Mohamed TAMAAZOUSTI

January 21, 2021

Acknowledgments

First of all, I would like to deeply thank my thesis director Prof. Romain Couillet as well as my co-supervisor Dr. Mohamed Tamaazousti for having constantly supported me throughout my thesis, I also thank them for their confidence, patience, motivation and for having guided me during these three unforgettable years. I had the immense pleasure of working with them both professionally and personally.

I would also like to express my sincere thanks for the two rapporteurs of my thesis: Prof. Julie Delon and Prof. Jamal Najim, so I would like to thank them for the time spent reading my manuscript and their very enriching comments.

I would also like to thank the other members of the jury for my defense: Dr. Alexandre Gramfort, Prof. Mérouane Debbah, Prof. Florent Krzakala and Dr. Florent Chatelain who were present despite the health crisis of Covid-19, I thank them for their questions and for their evaluation of my work.

I would like to express my gratitude to all the institutions (CEA and CentraleSupélec) and laboratories (LVML and L2S) that ensured the smooth running of my thesis. I would also like to thank all the people who were involved in my thesis. I am thinking in particular of Mrs. Odile Caminondo and Mrs. Virginie Valaire for their help on the administrative level within the CEA, and also to Mrs. Anne Batalie and Prof. Gilles Duc within CentraleSupélec. I would also like to thank all my colleagues in the SIALV department at the CEA and also my colleagues at GIPSA in Grenoble, where I spent very memorable moments.

I would like to make a special mention to my teammates during my thesis: Cosme, Youssef, Abdallah, John, Hafiz, Malik, Xiaoyi, Zhenyu, Lorenzo, Cyprien, Souheil, Alex, Alexandre, Tayeb, Charles, Hassane, Vincent, Harold and Vasily. I really enjoyed working with you and I was very pleased to have you as teammates. Finally, I would like to thank my family for their support, for believing in me and for supporting me throughout my years of study and even throughout my life. Special thanks goes out to my parents: Fatima and Mimoun, my grandparents: Malika and Chaib, my brothers: Moussa and Issam, my sister and her husband: Khadija and Mohamed.

Contents

1	Introduction	13
1.1	Illustrative High-dimensional Examples	15
1.1.1	Large Sample Covariance Matrices	15
1.1.2	Large Kernel Matrices	17
1.2	From GMM to Concentration through GAN	21
1.3	Some ML methods under Concentration	22
1.3.1	Behavior of Gram Matrices	23
1.3.2	Behavior of Kernel Matrices	24
1.3.3	Beyond Kernels to Neural Networks	24
1.3.4	Summary of Section 1.3	27
1.4	Outline and Contributions	27
2	Random Matrix Theory & Concentration of Measure Theory	29
2.1	Fundamental Random Matrix Theory Notions	30
2.1.1	The resolvent matrix notion	30
2.1.2	Stieltjes transform and spectral measure	30
2.1.3	Cauchy's integral and statistical inference	32
2.1.4	Deterministic and random equivalents	33
2.2	Fundamental Random Matrix Theory Results	34
2.2.1	Matrix identities and key lemmas	34
2.2.2	The Marčenko-Pastur Law	37
2.2.3	Random matrices of mixture models	38
2.3	Connections with machine learning through spiked models	39
2.3.1	Simple example of unsupervised learning	41
2.3.2	Simple example of supervised learning	45
2.4	Extensions with concentration of measure theory	48
2.4.1	The notion of concentrated vectors	48
2.4.2	Resolvent of the sample covariance matrix	50
3	Universality of Large Random Matrices	55
3.1	GAN Data are Concentrated Data	55
3.2	Random Gram Matrices of Concentrated Data	61
3.2.1	Motivation	61
3.2.2	Model and Main Results	62
3.2.2.1	Mixture of Concentrated Vectors	62
3.2.2.2	Behavior of the Gram matrix of concentrated vectors	63
3.2.2.3	Application to GAN-generated Images	65
3.2.3	Central Contribution	68

4	Random Kernel Matrices of Concentrated Data	71
4.1	Kernel Spectral Clustering	71
4.1.1	Motivation	72
4.1.2	Model and Main Results	72
4.1.2.1	Behavior of Large Kernel Matrices	76
4.1.2.2	Application to GAN-generated Images	87
4.1.3	Central Contribution and perspectives	89
4.2	Sparse Principal Component Analysis	89
4.2.1	Motivation	90
4.2.2	Model and Main Results	91
4.2.2.1	Random Matrix Equivalent	92
4.2.2.2	Application to sparse PCA	94
4.2.2.3	Experimental Validation	96
4.2.3	Central Contribution and Perspectives	99
5	Beyond Kernel Matrices, to Neural Networks	103
5.1	A random matrix analysis of Softmax layers	103
5.1.1	Motivation	104
5.1.2	Model setting: the Softmax classifier	105
5.1.3	Assumptions & Main Results	106
5.1.3.1	Concentration of the weights vector of the Softmax classifier	107
5.1.3.2	Experimental validation	112
5.1.4	Central Contribution	114
5.2	A random matrix analysis of Dropout layers	115
5.2.1	Motivation	115
5.2.2	Model and Main Results	116
5.2.2.1	Deterministic equivalent	119
5.2.2.2	Generalization Performance of α -Dropout	121
5.2.2.3	Training Performance of α -Dropout	122
5.2.3	Experiments	122
5.2.4	Central Contribution and Perspectives	123
6	Conclusions & Perspectives	125
6.1	Conclusions	125
6.2	Limitations and Perspectives	126
A	Synthèse de la thèse en Français	129
B	Practical Contributions	133
B.1	Generative Collaborative Networks for Super-resolution	133
B.1.1	Motivation	133
B.1.2	Proposed Methods	135
B.1.2.1	Proposed Framework	135
B.1.2.2	Existing Loss Functions	137
B.1.3	Experiments for Single Image Super-Resolution	137
B.1.3.1	Proposed Methods	137
B.1.3.2	Evaluation Metrics	138
B.1.3.3	Experiments	139
B.1.4	Central Contribution and Discussions	143
B.2	Neural Networks Compression	143

B.2.1	Motivation	143
B.2.2	Proposed Methods	144
B.2.2.1	Setting & Notations	144
B.2.2.2	Neural Nets PCA-based Distillation (Net-PCAD)	145
B.2.2.3	Neural Nets LDA-based Distillation (Net-LDAD)	145
B.2.3	Experiments	146
B.2.4	Central Contribution and Discussions	148
C	Proofs	157
C.1	Proofs of Chapter 3	157
C.1.1	Setting of the proof	157
C.1.2	Basic tools	158
C.1.3	Main body of the proof	159
C.2	Proofs of Chapter 4	162
C.2.1	Proofs of Section 4.1	162
C.2.2	Proofs of Section 4.2	165
C.2.2.1	Proof of Theorem 4.3	165
C.2.2.2	Proof of Theorem 4.4	167
C.3	Proofs of Chapter 5	168
C.3.1	Proofs of Section 5.1	168
C.3.2	Proofs of Section 5.2	172
C.3.2.1	Proof of Theorem 5.5	172
C.3.2.2	Proof of Theorem 5.6	173
C.3.2.3	Optimality	175
	References	177

List of Figures

1.1	Histogram of the sample covariance matrix	16
1.2	Behavior of kernel matrices in low-dimension versus high-dimension . . .	20
1.3	Examples of images generated by the BigGAN model	21
2.1	Spectrum of Gram matrices of spiked models	42
2.2	Simulated versus theoretical alignment of spiked models	44
2.3	Histogram of the decision function of linear classifiers	47
3.1	Standard GAN architecture	56
3.2	Deep learning representations of GAN-data are concentrated vectors . . .	58
3.3	Dynamics of spectral normalization	59
3.4	Examples of images generated by GANs.	65
3.5	Spectrum and dominant eigenspace of the Gram matrix of CNN represen- tations of GAN and real images.	66
3.6	Spectrum of the Gram matrix of CNN representations of GAN and real images.	67
3.7	Linear SVM performances on GAN data.	69
3.8	Linear SVM performances on GAN real.	70
4.1	Histogram of distances with Alexnet representations	73
4.2	Histogram of distances across different representation networks	73
4.3	Kernel matrix of Alexnet representations	74
4.4	Kernel matrix of CNN representations of GAN data	75
4.5	Kernel matrix of CNN representations of Real data	76
4.6	Spectrum of the kernel matrix of Alexnet representations	77
4.7	Spectrum of the kernel matrix for CNN representations of GAN data . . .	78
4.8	Spectrum of the kernel matrix for CNN representations of Real data . . .	79
4.9	Largest eigenvectors of the kernel matrix of Alexnet representations . . .	80
4.10	Largest eigenvectors of the kernel matrix for CNN representations of GAN data	81
4.11	Largest eigenvectors of the kernel matrix for CNN representations of Real data	82
4.12	Largest eigenspace of the kernel matrix of Alexnet representations	82
4.13	Largest eigenspace of the kernel matrix for GAN data	83
4.14	Largest eigenspace of the kernel matrix for real data	83
4.15	Kernel spectral clustering on Alexnet representations	84
4.16	Kernel spectral clustering on CNN representations of GAN images	85
4.17	Kernel spectral clustering on CNN representations of real images	86
4.18	Element-wise kernel function for Sparse PCA.	96

4.19	Spectrum of the sample covariance versus Sparse PCA	97
4.20	PCA recovery	98
4.21	PCA performances	99
4.22	Multiple PCAs model	100
4.23	Validation of the conditions $f'(0) = f''(0) = 0$	100
5.1	Training of the Softmax layer	105
5.2	Learned softmax weights for GAN data	113
5.3	Scores of the softmax classifier for GAN data	114
5.4	Learned softmax weights with real data	115
5.5	Scores of the softmax classifier with real data	116
5.6	Illustration of learning with α -Dropout	118
5.7	Histogram of the decision function of the α -Dropout model	122
5.8	Performances of the α -Dropout model	123
B.1	Difference between images from Imagenet and satellite images as data from different domains	134
B.2	Overview of the GCN framework	149
B.3	Used datasets for super-resolution	150
B.4	Super-resolution architectures	150
B.5	Qualitative super-resolution results	152
B.6	Super-resolution on satellite images	153
B.7	Qualitative super-resolution results	154
B.8	Learning curves of the student network	155
B.9	Learned Homoscedastic loss for the student network	155
B.10	Performances of the student network	156

Notations

Mathematical Symbols

\mathbb{R}	Set of real numbers.
\mathbb{C}	Set of complex numbers, we denote \mathbb{C}^+ the set $\{z \in \mathbb{C}, \Im[z] > 0\}$.
$\mathcal{M}_{p,n}$	Set of matrices of size $p \times n$.
\mathcal{M}_p	Set of squared matrices of size p .
\mathcal{D}_p	Set of diagonal matrices of size p .
$(\cdot)^\top$	Transpose operator.
$\text{tr}(\cdot)$	Trace operator.
$\text{diag}(\cdot)$	Diagonal operator, for $A \in \mathcal{M}_n$, $\text{diag}(A) \in \mathbb{R}^n$ is the vector with entries $\{A_{ii}\}_{i=1}^n$; for $\mathbf{a} \in \mathbb{R}^n$, $\text{diag}(\mathbf{a}) \in \mathcal{D}_n$ is the diagonal matrix taking \mathbf{a} as its diagonal.
$\ \cdot\ $	Operator (or spectral) norm of a matrix and Euclidean norm of a vector.
$\ \cdot\ _F$	Frobenius norm of a matrix, $\ A\ _F = \sqrt{\text{tr}(AA^\top)}$.
$\ \cdot\ _\infty$	Infinite norm of a matrix, $\ A\ _\infty = \max_{i,j} A_{ij} $.
$\text{dist}(\cdot)$	Distance between elements in a metric space.
$P(\cdot)$	Probability of an event with respect to the underlying probability measure space (Ω, \mathcal{F}, P) .
$\mathbb{E}[\cdot]$	Expectation operator, $\mathbb{E}[f] = \int f dP = \int_{\Omega} f(x)P(dx)$
$\text{Var}[\cdot]$	Variance operator, $\text{Var}[x] = \mathbb{E}[x^2] - (\mathbb{E}[x])^2$ if the first two moments of x exist
$\xrightarrow{a.s.}$	Almost surely convergence. We say a sequence $x_n \xrightarrow{a.s.} x$ if $P(\lim_{n \rightarrow \infty} x_n = x) = 1$
$Q(x)$	Q-function: the tail distribution of standard Gaussian $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty \exp(-t^2/2) dt$
$\mathcal{O}(1), o(1)$	A sequence u_n is bounded or converges to zero as $n \rightarrow \infty$, respectively.

Vectors and Matrices

I_n	Identity matrix of size n .
$\mathbf{1}_n$	(Column) vector of size n with all entries equal to one.
$\mathbf{0}_n$	(Column) vector of size n with all entries equal to zero.
$\mathbf{x}_i \in \mathbb{R}^p$	Input data/feature vector.
$\mathbf{z}_i \in \mathbb{R}^p$	Random vector having i.i.d. zero mean and unit variance entries.
X	Data/feature matrix having \mathbf{x}_i as column vectors.
Z	Random matrix having \mathbf{z}_i as column vectors such that $\mathbb{E}[\mathbf{z}_i] = \mathbf{0}_p$.
$\boldsymbol{\mu}$	Mean vector.
C	Covariance matrix.
K	Kernel matrix.
W	Weight matrix (in a neural network).
$\boldsymbol{\beta}$	Regression vector.
Q	Resolvent matrix, see Definition 3.

Chapter 1

Introduction

Contents

1.1 Illustrative High-dimensional Examples	15
1.1.1 Large Sample Covariance Matrices	15
1.1.2 Large Kernel Matrices	17
1.2 From GMM to Concentration though GAN	21
1.3 Some ML methods under Concentration	22
1.3.1 Behavior of Gram Matrices	23
1.3.2 Behavior of Kernel Matrices	24
1.3.3 Beyond Kernels to Neural Networks	24
1.3.4 Summary of Section 1.3	27
1.4 Outline and Contributions	27

Artificial intelligence (AI) is known as the set of theories and techniques used to create machines capable of simulating human intelligence. One of the most attractive sub-fields of AI is machine learning (ML) which aims at providing computer algorithms that “learn” automatically through experience in order to make future decisions without being explicitly programmed. Basically, ML algorithms leverage on building mathematical – very often parametric – models which will be optimized based on sample training data and then used afterwards to perform various AI tasks such as classification, regression, clustering etc.

Quite naturally, AI finds applications in various domains and hence one of the most important challenges of ML is to provide algorithms that can be applied to different kind of data (e.g., images, texts, graphs etc.). By construction, these data can be represented in different forms and therefore the performance of ML algorithms will rely largely on the chosen representation. This representation should ideally contain relevant information about the data in order to achieve learning with simple models and small amount of data. Historically, a huge amount of works were focused on the design of hand-crafted representations (or features) and then providing them to simple ML algorithms to resolve the desired tasks. But for most tasks and given the various types of data, these approaches are not easily scaleable to achieve effective AI.

Since the arrival of deep neural networks (DNNs), the ideas of developing hand-crafted features were immediately left out. Indeed, DNNs have surpassed most of the approaches by demonstrating their incredible ability to automatically learning relevant

representations from raw data in a wide range of applications, including computer vision, pattern recognition and natural language processing. Despite their success, a lot of questions are still unanswered regarding the theoretical foundations of DNNs and which are very crucial notably for their explainability. For instance, the full characterization of their learnt representations and/or parameters are still open problems.

One of the main aspects that made DNNs effective in practice is their being over-parametrized models. Indeed, it has been shown that deep architectures of these models surpass shallow ones when dealing with n p -dimensional data when both n and p are large, which is often the case in real life scenarios¹. Moreover, most effective DNNs happen to have a number of parameters N which is at least of order p or even much larger (e.g., LeNet-5 [LeC98] contains $N = 60000$ parameters).

In essence, these large dimensions induce many counter-intuitive phenomena that cause the intuitions of the small dimensions to collapse completely. For better understanding of these phenomena, we will provide subsequently in Section 1.1 some illustrative examples which reveal these counter-intuitive aspects. In the particular case when both $p, n \rightarrow \infty$ with $p/n \rightarrow 0 \in (0, \infty)$, random matrix theory (RMT) provides powerful tools to assess the performance of various ML algorithms by taking into account the effect of these dimensions. Indeed, RMT provides access to the internal mechanism of a large wide of ML methods thereby allowing a deeper understanding and systematic improvements of these methods. We shall refer the reader to the thesis of Z. Liao [Lia19] for applications of RMT to kernel methods, random shallow neural networks and neural networks dynamics; the thesis of X. Mai [Mai19] which addresses applications of RMT to semi-supervised learning and support vector machines.

The aforementioned works leverage largely on Gaussian assumptions² on the processed data. One of the major outcome of this thesis it to go beyond the Gaussian hypothesis to address the applicability of RMT to real data which are unlikely close to Gaussian vectors. In particular, working under the more generic statistical model of concentrated vectors [LC20], we provide justifications – leveraging on generative adversarial networks (GANs) – about the relevance of such model for realistic data modelling, and we further analyze, under the concentration assumption on the data, the behavior of large kernel matrices (which happen to be at the core of various ML algorithms) as well as some essential components of neural networks such as the last Softmax layer. A major result of the developed works in this thesis is the universality result stated as:

“Only first and second order statistics of concentrated data matter to describe the behavior of these methods”

thereby justifying the Gaussianity assumption on data as per the results of [Lia19, Mai19].

Let us now give some illustrative examples of the effect of high-dimensional data which leads to fundamentally different behaviors of ML algorithms comparing to low-dimensional data.

¹As an example, the MNIST dataset [LeC98] contains $n = 70000$ images of dimension $p = 28 \times 28 = 784$.

²Modelling the data as a k -class of Gaussian mixture model (GMM).

1.1 Illustrative High-dimensional Examples

In this section we highlight the effect of high dimensions on two classical examples: the first one describes the behavior of large sample covariance matrices presented in Subsection 1.1.1; and the second example describes the behavior of large kernel matrices presented in Subsection 1.1.2 which are at the core of a wide range of ML algorithms.

1.1.1 Large Sample Covariance Matrices

Covariance matrices are omnipresent in a wide range of signal processing and machine learning approaches, therefore understanding their behavior is of particular interest to achieve consistent estimations and inferences. We will present in this subsection a classical and yet fundamental result from random matrix theory which describes the counter-intuitive behavior of large covariance matrices.

Let $x_1, \dots, x_n \in \mathbb{R}^p$ some *i.i.d.* random data vectors such that $\mathbb{E}[x_i] = \mathbf{0}$ and $\mathbb{E}[x_i x_i^\top] = \mathbf{I}_p$. In the case where $x_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, the maximum likelihood estimator for the population covariance matrix (here \mathbf{I}_p) is the *sample covariance matrix* given by

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n x_i x_i^\top = \frac{1}{n} \mathbf{X} \mathbf{X}^\top \quad (1.1)$$

where $\mathbf{X} = [x_1, \dots, x_n] \in \mathcal{M}_{p,n}$. In the classical setting where p does not scale with the number of samples n , we know from the *strong law of large numbers* that, as $n \rightarrow \infty$

$$\hat{\Sigma} \xrightarrow{a.s.} \mathbf{I}_p$$

In particular, $\|\hat{\Sigma} - \mathbf{I}_p\| \xrightarrow{a.s.} 0$ in operator norm or any matrix norm since they are equivalent in finite dimension. However, it turns out that in the *random matrix theory* regime when both $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$,

$$\|\hat{\Sigma} - \mathbf{I}_p\| \not\xrightarrow{a.s.} 0$$

Indeed, in particular, in the under-sampling setting when $c > 1$ the sample covariance matrix converges entry-wise to the population matrix since

$$\max_{1 \leq i, j \leq p} |[\hat{\Sigma} - \mathbf{I}_p]_{ij}| \xrightarrow{a.s.} 0$$

while there is a clear *eigenvalue mismatch* between the two quantities which simply unfolds from

$$\begin{aligned} 0 &= \lambda_1(\hat{\Sigma}) = \dots = \lambda_{p-n}(\hat{\Sigma}) \leq \lambda_{p-n+1}(\hat{\Sigma}) \leq \dots \leq \lambda_p(\hat{\Sigma}) \\ 1 &= \lambda_1(\mathbf{I}_p) = \dots = \lambda_p(\mathbf{I}_p) \end{aligned}$$

therefore implies no convergence in spectral norm.

Figure 1.1 depicts the histogram of eigenvalues of $\hat{\Sigma}$ for $p = 2000$ and $n = 10000$, as can be seen, instead of having the eigenvalue 1 with multiplicity p one observes a spreading of the eigenvalues of $\hat{\Sigma}$ in the vicinity of 1. This spreading of eigenvalues is particularly described by the so-called *Marčenko-Pastur Law* [MP67] which predicts the limiting eigenvalues distribution of $\hat{\Sigma}$ in terms of the ratio $c = \lim_n p/n$ as follows

$$f_{MP}(x) = \frac{1}{2\pi c} \frac{\sqrt{(\lambda^+ - x)(x - \lambda^-)}}{x} + \max\left(1 - \frac{1}{c}, 0\right) \delta(x),^3 \quad (1.2)$$

³ δ denotes here the dirac function.

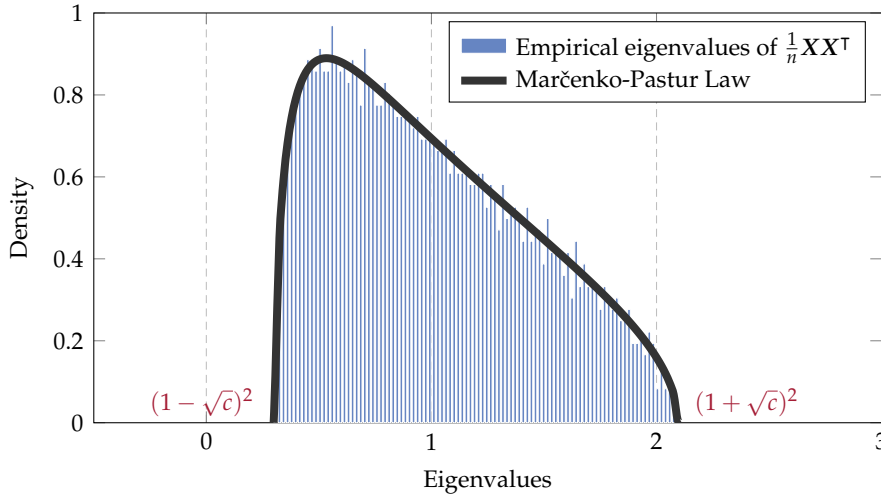


Figure 1.1: Histogram of the eigenvalues of $\hat{\Sigma} = \frac{1}{n}XX^T$ where X is a random matrix with *i.i.d.* $\mathcal{N}(0,1)$ entries. We considered the dimensions $p = 2000$, $n = 10000$ and we recall $c = p/n$.

where $\lambda^\pm = (1 \pm \sqrt{c})^2$ are respectively the right and left edges of the limiting distribution. Indeed, the theorem for Marčenko & Pastur [MP67] states that the random *empirical spectral distribution* (e.s.d.) μ of $\hat{\Sigma}$ defined in terms of its eigenvalues $\lambda_1, \dots, \lambda_p$ as

$$\mu(x) = \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i}(x) \quad (1.3)$$

converges, in the random matrix theory regime $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, to a *deterministic measure* μ_{MP}

$$\mu \xrightarrow{a.s.} \mu_{MP} \quad (1.4)$$

in distribution⁴. Essentially, μ_{MP} satisfies $\mu_{MP}(\{0\}) = \max(0, 1 - c^{-1})$ and on $(0, \infty)$ it has a continuous density function defined on the compact support $[(1 - \sqrt{c})^2, (1 + \sqrt{c})^2]$ by f_{MP} defined in equation 1.2. A detailed statement along with a proof of this fundamental result is provided in Subsection 2.2.2 of this manuscript.

Remark 1.1 (Universality of the Marčenko-Pastur Law). *The Marčenko-Pastur result holds for any random matrix X having random i.i.d. entries with zero mean, unit variance and finite fourth order moment whatever the underlying distribution of its entries.*

We therefore know from this result that the behavior of the sample covariance matrix is fundamentally different in high-dimension leading to counter-intuitive phenomena, which is highlighted here through the behavior of its spectrum. A fundamental question therefore occurs regarding *when data are of high-dimension?* It should be noted that high-dimensionality is relative to the number of observations, which unfolds – through the above example – from the observation that the limiting Marčenko-Pastur density function in equation 1.2 depends only on the ratio $c = \lim_n \frac{p}{n}$. Specifically, we know from this result that the eigenvalues of $\hat{\Sigma}$ spread from $(1 - \sqrt{c})^2$ to $(1 + \sqrt{c})^2$ instead of being concentrated around 1. As a result, the eigenvalues deviate on a range of diameter

$$(1 + \sqrt{c})^2 - (1 - \sqrt{c})^2 = 4\sqrt{c}$$

⁴In the sense that $\int \varphi(t)\mu(dt) \xrightarrow{a.s.} \int \varphi(t)\mu_{MP}(dt)$ for all bounded continuous test function φ .

For instance, even taking $n = 100p$ will result in a spread of eigenvalues around 1 with deviation $4\sqrt{c} = 0.4$ which is relatively large.

Remark 1.2 (Extensions to structured data models). *The Marčenko-Pastur result has been extended to some structured data models, namely by Bai and Silverstein in [BS10] for data of the form $\mathbf{x}_i = \Sigma^{1/2}\mathbf{z}_i$ for some positive semi-definite matrix Σ and \mathbf{z}_i random vectors with i.i.d. entries of zero mean, unit variance and finite fourth order moment. And also to a more convenient data model in the context of machine learning which is the Gaussian mixture model (GMM) by Benaych and Couillet in [BGC16]. The latter result will be presented in more details in Subsection 2.2.3.*

We have highlighted in this subsection the intriguing behavior of the sample covariance matrix when the underlying data are of high-dimension. In particular, we saw that in the random matrix theory regime when both $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, the sample covariance matrix $\hat{\Sigma}$ is no longer a consistent estimation of the population covariance matrix. In the next subsection, we will present another example of the effect of dimensionality on kernel methods, which are an essential component in modern machine learning.

1.1.2 Large Kernel Matrices

Kernel methods are a class of machine learning algorithms which were introduced in order to avoid transforming raw data into hand-crafted features and rather simply making the choice of a *kernel function* $\kappa(\mathbf{x}, \mathbf{y})$ instead, i.e., a *similarity function* between pairs of data samples \mathbf{x} and \mathbf{y} in their raw representation. One of the nice and very useful properties of (nonnegative definite) kernel methods is that they *implicitly* transform the input data into a high-dimensional feature space without requiring the computation of the data in that space, but rather computing the *inner similarity* between pairs of data, which is known as the *kernel trick*. However, the considered kernel function in practice are very often *non-linear* functions which makes the analysis of the resulting *random kernel matrix* non-trivial to tackle from the random matrix theory perspective.

One of the classical algorithms from kernel methods is the so-called *kernel spectral clustering* which specifically performs unsupervised classification of a sequence of data vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ into k distinct classes. It basically relies on the computation of large pair-wise kernel matrix $\mathbf{K} = \{\kappa(\mathbf{x}_i, \mathbf{x}_j)\}_{i,j=1}^n$ and performing a classical k -means clustering on the *dominant eigenspace* of the random kernel matrix \mathbf{K} .

In order to study the behavior of \mathbf{K} and so to determine the information encoded in its dominant eigenspace, a crucial step requires the *linearization* of the non-linearity $\kappa(\mathbf{x}, \mathbf{y})$. Under *asymptotically non-trivial* growth rate assumptions on the data statistics, which basically maintains a feasible yet not too easy clustering problem, it has been shown by Couillet and Benayach in [CBG⁺16] that for the similarity

$$\kappa(\mathbf{x}, \mathbf{y}) = f\left(\frac{1}{p}\|\mathbf{x} - \mathbf{y}\|^2\right)$$

for sufficiently smooth function f , the off-diagonal entries of \mathbf{K} tend – in the large dimensional regime when $p/n \rightarrow c$ as $p \rightarrow \infty$ – to a limiting constant independently of the data classes – *the between and within class vectors are “equidistant” in high-dimension*. This

intriguing high-dimensional behavior allows one to study \mathbf{K} through a Taylor expansion yielding to a *random matrix equivalent* of \mathbf{K} , thereby giving access to the characterization of functionals of \mathbf{K} and its (informative) eigenspace in the large dimensional setting. The aforementioned *non-trivial assumptions* notably permit an accurate approximation of \mathbf{K} by its random matrix equivalent in the large p, n limit.

Couillet and Benaych in [CBG⁺16] have notably analyzed the behavior of \mathbf{K} under the so-called *Gaussian mixture model* assumption on data, which is defined in the following.

Definition 1 (Gaussian Mixture Model (GMM)). *A sequence of data vectors $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ is said to form a k -class Gaussian mixture model of k classes $\mathcal{C}_1, \dots, \mathcal{C}_k$ with distinct means and covariances $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ respectively, if for $\mathbf{x}_i \in \mathcal{C}_\ell$*

$$\mathbf{x}_i = \boldsymbol{\mu}_\ell + \boldsymbol{\Sigma}_\ell^{\frac{1}{2}} \mathbf{z}_i \quad \text{with} \quad \mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$$

We further define the following quantities which shall be used subsequently

- (Data) $n_\ell = |\mathcal{C}_\ell|$ the cardinality of class \mathcal{C}_ℓ and $c_\ell = \lim_n \frac{n_\ell}{n} \in (0, 1)$.
- (Means) $\bar{\boldsymbol{\mu}} = \sum_{\ell=1}^k c_\ell \boldsymbol{\mu}_\ell$ and $\bar{\boldsymbol{\mu}}_\ell = \boldsymbol{\mu}_\ell - \bar{\boldsymbol{\mu}}$.
- (Covariances) $\bar{\boldsymbol{\Sigma}} = \sum_{\ell=1}^k c_\ell \boldsymbol{\Sigma}_\ell$ and $\bar{\boldsymbol{\Sigma}}_\ell = \boldsymbol{\Sigma}_\ell - \bar{\boldsymbol{\Sigma}}$.

Under the GMM model in Definition 1, we recall the growth rate assumptions made in [CBG⁺16], which are essentially, as $p \rightarrow \infty$,

- (Data) $p/n \rightarrow c \in (0, \infty)$ and $n_\ell/n \rightarrow (0, 1)$.
- (Means) $\limsup_p \max_\ell \|\bar{\boldsymbol{\mu}}_\ell\| < \infty$.
- (Covariances) $\limsup_p \max_\ell \|\boldsymbol{\Sigma}_\ell\| < \infty$ and $\limsup_p \max_\ell \frac{1}{\sqrt{p}} \text{tr} \bar{\boldsymbol{\Sigma}}_\ell < \infty$.

Remark 1.3 (On the non-trivial growth rates). *The above assumptions essentially ensure, as discussed in [CBG⁺16], that the classification is neither too easy nor too hard as the dimension grows to infinity. For instance, in the case of binary classification (i.e., $k = 2$), these assumptions imply $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\| = \mathcal{O}(1)$ and $\frac{1}{\sqrt{p}} \text{tr}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2) = \mathcal{O}(1)$ as $p \rightarrow \infty$. Thus if $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|$ and $\frac{1}{\sqrt{p}} \text{tr}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)$ increase with p , then a simple Bayesian analysis will reveal that the classification will become too easy even with a trivial algorithm, in contrast, if the two quantities vanish as p grows large, the classification becomes theoretically impossible whatever the chosen algorithm. In particular, $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\| = \mathcal{O}(1)$ unfolds from the following: Consider a Gaussian mixture of two classes such that $\mathcal{C}_\ell : \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_\ell, \mathbf{I}_p)$ for $\ell \in \{1, 2\}$ (i.e., $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \mathbf{I}_p$). In the case where $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ are perfectly known, one has the following optimal decision by the Neyman-Pearson test, specifically for some $\mathbf{x} \in \mathcal{C}_1$*

$$(\mathbf{x} - \boldsymbol{\mu}_2)^\top (\mathbf{x} - \boldsymbol{\mu}_2) - (\mathbf{x} - \boldsymbol{\mu}_1)^\top (\mathbf{x} - \boldsymbol{\mu}_1) \underset{\mathcal{C}_1}{\overset{\mathcal{C}_2}{\gtrless}} \log \frac{\det(\boldsymbol{\Sigma}_1)}{\det(\boldsymbol{\Sigma}_2)} = 0$$

Decomposing $\mathbf{x} = \boldsymbol{\mu}_1 + \mathbf{z}$ with $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$, the above test can be equivalently written as

$$g(\mathbf{z}) = \frac{1}{p} \|\Delta \boldsymbol{\mu}\|^2 + \frac{2}{p} (\Delta \boldsymbol{\mu})^\top \mathbf{z} \underset{\mathcal{C}_1}{\overset{\mathcal{C}_2}{\gtrless}} 0 \quad \text{where} \quad \Delta \boldsymbol{\mu} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$$

Moreover, since $g(\mathbf{z})$ is a sum of p independent random variables, we have by Lyapunov's central limit theorem 2.7, as $p \rightarrow \infty$

$$v_g^{-\frac{1}{2}}(g(\mathbf{z}) - m_g) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1) \quad \text{with} \quad m_g = \frac{1}{p} \|\Delta \boldsymbol{\mu}\|^2, \quad v_g = \frac{4}{p^2} \|\Delta \boldsymbol{\mu}\|^2$$

Therefore, the asymptotic (as $p \rightarrow \infty$) classification performance of $\mathbf{x} \in \mathcal{C}_1$ is non-trivial if and only if the two quantities m_g and $\sqrt{v_g}$ are of the same order, as such $\|\Delta \boldsymbol{\mu}\|$ must be $\mathcal{O}(1)$.

Under the above assumptions a fundamental high-dimensional behavior occurs which specifically states that the point-wise distance $\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2$ concentrates around the quantity

$$\tau = \frac{2}{p} \text{tr } \boldsymbol{\Sigma}$$

i.e., the between and within class data vectors are ‘‘equidistant’’ in high dimensions, essentially, one has with probability one that

$$\max_{1 \leq i \neq j \leq n} \left| \frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \tau \right| \xrightarrow{a.s.} 0$$

This result unfolds from the decomposition of the pair-wise distance $\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2$ into asymptotically controllable terms as follows. For some $\mathbf{x}_i \in \mathcal{C}_a$, we denote $\mathbf{z}_i = \frac{\mathbf{x}_i - \boldsymbol{\mu}_a}{\sqrt{p}}$ and let $\psi_i = \|\mathbf{z}_i\|^2 - \frac{1}{p} \text{tr } \boldsymbol{\Sigma}_a$. We therefore have the following decomposition by [CBG⁺16] for $\mathbf{x}_i \in \mathcal{C}_a$ and $\mathbf{x}_j \in \mathcal{C}_b$ with $a \neq b \in [k]$

$$\begin{aligned} \frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 &= \|\mathbf{z}_i - \mathbf{z}_j\|^2 + \frac{1}{p} \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\|^2 + \frac{2}{\sqrt{p}} (\boldsymbol{\mu}_a - \boldsymbol{\mu}_b)^\top (\mathbf{z}_i - \mathbf{z}_j) \\ &= \tau + \frac{1}{p} \text{tr } \bar{\boldsymbol{\Sigma}}_a + \frac{1}{p} \text{tr } \bar{\boldsymbol{\Sigma}}_b + \psi_i + \psi_j - 2\mathbf{z}_i^\top \mathbf{z}_j \\ &\quad + \frac{1}{p} \|\bar{\boldsymbol{\mu}}_a - \bar{\boldsymbol{\mu}}_b\|^2 + \frac{2}{\sqrt{p}} (\bar{\boldsymbol{\mu}}_a - \bar{\boldsymbol{\mu}}_b)^\top (\mathbf{z}_i - \mathbf{z}_j) \end{aligned}$$

From this decomposition, it turns out that, except for τ , all the other terms are at least of order $\mathcal{O}(p^{-\frac{1}{2}})$. Indeed, it is easily seen that $\psi_i = \mathcal{O}(p^{-\frac{1}{2}})$ and $\mathbf{z}_i^\top \mathbf{z}_j = \mathcal{O}(p^{-\frac{1}{2}})$ while $\frac{2}{\sqrt{p}} (\bar{\boldsymbol{\mu}}_a - \bar{\boldsymbol{\mu}}_b)^\top (\mathbf{z}_i - \mathbf{z}_j)$ is of order $\mathcal{O}(p^{-1})$.

This observation notably states that the entries of the kernel matrix \mathbf{K} will converge to the same value $f(\tau)$ at a first order approximation, specifically with probability one

$$\max_{1 \leq i \neq j \leq n} |\mathbf{K}_{ij} - f(\tau)| \xrightarrow{a.s.} 0$$

which suggests that \mathbf{K} is approximated at a first order by the rank one matrix $f(\tau) \mathbf{1}_n \mathbf{1}_n^\top$. However, it should be noted that the uniform convergence on the entries of \mathbf{K} does not imply at all a convergence in the spectral norm as we saw in the example of the sample covariance matrix from the previous subsection. Still, a finer Taylor expansion analysis and control of the vanishing matrices in spectral norm provides the actual mechanism and information encoded by the kernel matrix. Figure 1.2 visually confirms that the entries of \mathbf{K} converge to the same value $f(\tau)$ at a first approximation. We particularly considered in this figure two classes \mathcal{C}_1 and \mathcal{C}_2 such that $\boldsymbol{\mu}_\ell = (-1)^\ell \boldsymbol{\mu}$ for $\boldsymbol{\mu} = [2, \mathbf{0}_{p-1}]^\top \in \mathbb{R}^p$

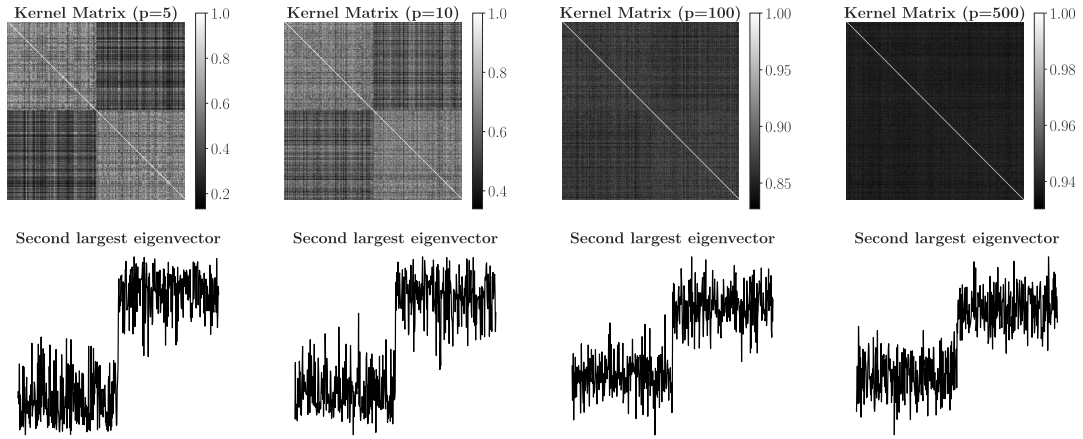


Figure 1.2: **(First line)** Kernel matrix with entries $K_{ij} = f(\frac{1}{p}\|x_i - x_j\|^2)$ and **(second line)** its corresponding second largest eigenvector which contains clustering information. We considered the kernel function $f(t) = \exp(-t)$ and two classes of means $-\mu$ and $+\mu$ and identity covariance. $\mu = [2, \mathbf{0}_{p-1}]^\top \in \mathbb{R}^p$, $n = 500$ and $p \in \{5, 10, 100, 500\}$.

and isotropic covariances. As we see from this figure, as p gets large the behavior of the kernel matrix is fundamentally different from the low-dimension setting, in which case one visualizes a block-structure of the kernel matrix. In contrast, looking at the second largest eigenvector of \mathbf{K} , we see that the class structure is preserved even in the high-dimensional setting. Indeed, as explained in [CBG⁺16], the class information appears in the second order approximation of \mathbf{K} since the first eigenvector with eigenvalue of order $\mathcal{O}(n)$ is proportional to $\mathbf{1}_n$, therefore not informative about the classes. Essentially, when $p, n \rightarrow \infty$, the kernel matrix \mathbf{K} is asymptotically well approximated by a so-called spiked random matrix of the form⁵

$$\mathbf{K} = \mathbf{P}_f + f'(\tau)\mathbf{Z}^\top\mathbf{Z} + f''(\tau)\mathbf{W} + o_p(1)$$

where \mathbf{P}_f is a low-rank informative matrix which contains information about the classes through the data statistics $\{\mu_\ell\}_{\ell=1}^k$ and $\{\Sigma_\ell\}_{\ell=1}^k$ and also depends on the kernel function f through its local first and second derivatives at τ . The remaining matrix terms are non-informative full-rank noise matrices which exhibit a spreading of eigenvalues in the spectrum of \mathbf{K} in the same way as we saw in the previous subsection though the example of the sample covariance matrix.

Remark 1.4 (On the analysis of the noise random matrices). *The analysis of the noise terms in the approximation of \mathbf{K} (e.g., the term $f'(\tau)\mathbf{Z}^\top\mathbf{Z}$) comes through generalizations of the Marčenko-Pastur result to structured data models as we discussed in the previous subsection in Remark 1.2.*

The main consequence of the above approximation is the access to the actual spectral behavior of \mathbf{K} , in particular, through the exact description of its dominant eigenvectors, thereby accessing the exact theoretical estimation of the performances of kernel spectral clustering as well as a wide range of kernel methods which rely on such kernel matrices

⁵The notation $\mathbf{A} = \mathbf{B} + o_p(1)$ means that $\|\mathbf{A} - \mathbf{B}\| = o(1)$ as $p \rightarrow \infty$ where $\|\cdot\|$ stands for the spectral norm.



Figure 1.3: Examples of images generated by the BigGAN model [BDS18].

K . Indeed, an aftermath of such analysis allows one to “tune” the optimal kernel function choice for the considered data. Moreover, based upon this fundamental result, Liao and Couillet gave the exact performance estimation of kernel LS-SVM in [LC17], while Mai and Couillet have analyzed and improved semi-supervised learning in [MC17] based on the analysis of K .

1.2 From GMM to Concentration through GAN

The starting point to obtain the results of the previous section and more fundamentally to analyze the behavior of ML algorithms is to design so-called *deterministic equivalents* which basically *encode* the behavior of large random matrices. We will recall in the next chapter fundamental random matrix theory results along with some applications to simple machine learning models such as spectral clustering with the Gram matrix and classification with a linear classifier.

So far, the considered assumption on data to design such deterministic equivalents is a GMM model (see Definition 1) as developed by Benaych and Couillet [BGC16], where we recalled their main result in Theorem 2.5. However, real data (e.g., images) are *unlikely close* to simple Gaussian vectors and therefore one needs a more realistic model to describe them. Following the well-known quote of R. Feynman: “*What I cannot create, I do not understand*”, it is fundamentally important to be able to create real data in order to fully understand their nature. To this end though, generative models are of particular interest. In particular, since the advent of Generative Adversarial Nets [GPAM⁺14], it is now possible to create neural network models that can generate complex data structures. Some examples of artificially generated images by the BigGAN model [BDS18] are depicted in Figure 1.3. In particular, GAN generated data are constructed by applying successive *Lipschitz* transformations to high-dimensional Gaussian random vectors. Schematically,

$$\text{Real data} \approx \text{GAN data} = \mathcal{F}_L \circ \dots \circ \mathcal{F}_1(z) \quad \text{with} \quad z \sim \mathcal{N}(\mathbf{0}, I_d)$$

where the \mathcal{F}_i 's are essentially Lipschitz transformations.

On the other hand, the fundamental *concentration of measure* phenomenon [Led05a] states that Lipschitz-ally transformed Gaussian vectors are *concentrated vectors*, i.e., they satisfy the following concentration property: A random vector $x \in E$ is said to be *concentrated* if for any 1-Lipschitz function $f : E \rightarrow \mathbb{R}$, for $q \geq 0$ there exists $C, \sigma > 0$ such

that

$$\forall t > 0, \quad \mathbb{P} \{|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| \geq t\} \leq Ce^{-(t/\sigma)^q}$$

Consequently, GAN data are concentrated random vectors by construction, and therefore this class of random vectors constitute a more appropriate and convincing statistical model for realistic data compared to simple Gaussian vectors. The main objective of this thesis is to exploit the framework of concentrated vectors developed by Louart and Couillet [LC18b] relying on the earlier works of El Karoui [EK⁺10b], which is particularly motivated by the fact that GAN data belong to this class of random vectors, in order to analyze the behavior of various ML methods. We particularly provide in the following section some key results from this thesis that concern mainly the study of large Gram and kernel matrices under the concentration assumption (see Subsections 1.3.1-1.3.2) and further applications related to neural networks through the study of the Softmax layer in Subsection 1.3.3.

1.3 Some ML methods under Concentration

In this section we summarize the main findings of this thesis which basically describe the behavior of standard ML methods under the concentration assumption on data. Precisely, we will assume throughout the following subsections that data are distributed as a *Mixture of Concentrated Vectors* as per the following definition.

Definition 2 (Mixture of Concentrated Vectors (MCV) [LC18b]). *Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ be a data matrix which is constituted of n random vectors distributed in k different classes $\mathcal{C}_1, \dots, \mathcal{C}_k$, such that the data classes are characterized by the moments, for $\mathbf{x}_i \in \mathcal{C}_\ell$*

$$\mathbb{E}[\mathbf{x}_i] = \boldsymbol{\mu}_\ell, \quad \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top$$

In particular, the data matrix \mathbf{X} satisfy a concentration assumption in the sense that, for any 1-Lipschitz function $f : \mathcal{M}_{p,n} \rightarrow \mathbb{R}$ with $\mathcal{M}_{p,n}$ enrolled by the Frobinuous norm $\|\cdot\|_F$, for $q > 0$ there exists $C, \sigma > 0$ ⁶ independent of p and n such that

$$\forall t > 0, \quad \mathbb{P} \{|f(\mathbf{X}) - \mathbb{E}f(\mathbf{X})| \geq t\} \leq Ce^{-(t/\sigma)^q}$$

Remark 1.5 (On the concentration of GAN random vectors). *We will see in Chapter 3 that GAN random vectors have notably the same concentration as standard Gaussian vectors which unfolds from the fact that GAN's generators networks have controlled Lipschitz norm; i.e., for Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ inputs (as commonly considered) whose observable diameter does not depend on the dimension d , the observable diameter of GAN's outputs does not increase with the data dimension. Moreover, the concentration of the data matrix \mathbf{X} implies the concentration of its columns vectors \mathbf{x}_i 's since $\mathbf{X} \mapsto \mathbf{X} \mathbf{j}_i = \mathbf{x}_i$ is 1-Lipschitz transformation, where $\mathbf{j}_i \in \mathbb{R}^n$ is the canonical vector defined as $(\mathbf{j}_i)_j = \delta_{i=j}$.*

In addition to the concentration assumption on data, we will further assume the classical random matrix theory regime as per the following assumption.

Assumption 1 (Growth rate). *As $p \rightarrow \infty$, assume*

1. $p/n \rightarrow c \in (0, \infty)$; $|\mathcal{C}_\ell|/n \rightarrow c_\ell \in (0, 1)$.
2. *The number of classes k is bounded.*

⁶ σ is called the *observable diameter*.

1.3.1 Behavior of Gram Matrices

The first contribution of this thesis concerns the analysis of the behavior of the Gram matrix defined as $\mathbf{G} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ under the concentration assumption in Definition 2 and the high-dimensional regime in Assumption 1. In particular, this contribution is presented in more details in Chapter 3 and is based on the following works:

- (C1) **MEA. Seddik**, C. Louart, M. Tamaazousti, R. Couillet, “Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures”, International Conference on Machine Learning (ICML’20), Online, 2020.
- (C1’) **MEA. Seddik**, M. Tamaazousti, R. Couillet, “Pourquoi les matrices aléatoires expliquent l’apprentissage ? Un argument d’universalité offert par les GANs”, Colloque francophone de traitement du signal et des images (Gretsi’19), Lille, France, 2019.

As will be recalled in Chapter 2, the spectral behavior of \mathbf{G} can be analyzed through its resolvent

$$\mathbf{R}(z) = (\mathbf{G} + z\mathbf{I}_n)^{-1}$$

The main result from this contribution is to provide a *deterministic equivalent* (see Definition 6) for $\mathbf{R}(z)$ which is given by the following Theorem.

Theorem 1.1 (Deterministic Equivalent for $\mathbf{R}(z)$). *Under the concentration model in Definition 2, the growth rate Assumptions 1 and further assume that for all $\ell \in [k]$, $\|\boldsymbol{\mu}_\ell\| = \mathcal{O}(\sqrt{p})$ ⁷. Then, $\mathbf{R}(z)$ concentrates with an observable diameter of order $\mathcal{O}(p^{-\frac{1}{2}})$. Furthermore,*

$$\|\mathbb{E}\mathbf{R}(z) - \tilde{\mathbf{R}}(z)\| = \mathcal{O}\left(\sqrt{\frac{\log(p)}{p}}\right), \quad \tilde{\mathbf{R}}(z) = \frac{1}{z} \text{diag} \left\{ \frac{\mathbf{I}_{n_\ell}}{1 + \delta_\ell^*(z)} \right\}_{\ell=1}^k + \frac{1}{pz} \mathbf{J} \Omega_z \mathbf{J}^\top$$

with $\Omega_z = \mathbf{M}^\top \bar{\mathbf{Q}}(z) \mathbf{M} \odot \text{diag} \left\{ \frac{\delta_\ell^*(z) - 1}{\delta_\ell^*(z) + 1} \right\}_{\ell=1}^k$ where $\mathbf{M} = [\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k]$, $\mathbf{J} = [j_1, \dots, j_k]$

and $\bar{\mathbf{Q}}(z) = \left(\frac{1}{ck} \sum_{\ell=1}^k \frac{\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top}{1 + \delta_\ell^*(z)} + z\mathbf{I}_p \right)^{-1}$ where $\delta^*(z) = [\delta_1^*(z), \dots, \delta_k^*(z)]^\top$ is the unique fixed point of the system of equations for each $\ell \in [k]$

$$\delta_\ell(z) = \frac{1}{p} \text{tr} \left((\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top) \left(\frac{1}{ck} \sum_{j=1}^k \frac{\boldsymbol{\Sigma}_j + \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top}{1 + \delta_j(z)} + z\mathbf{I}_p \right)^{-1} \right).$$

Theorem 1.1 along with the deterministic equivalent of the sample covariance matrix developed in [LC18b] particularly generalize the result of [BGC16] to the class of concentrated vectors. A remarkable outcome from this result is that the behavior of the Gram matrix \mathbf{G} depends strictly on the first and second statistical moments $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ of the data thereby providing a *universal behavior* of ML methods which rely on the Gram matrix regardless of the data distribution. This universality behavior has notably been verified (see Chapter 3) using GAN generated images as they satisfy the concentration assumption by design.

⁷This hypothesis is less restrictive in the sense that a p -dimensional vector with entries of order $\mathcal{O}(1)$ has an ℓ_2 -norm of order $\mathcal{O}(\sqrt{p})$. Furthermore, from a technical standpoint, this assumption also explains the contribution of the class-wise means in the $\delta_\ell(z)$ ’s subsequently.

1.3.2 Behavior of Kernel Matrices

The second contribution of this thesis concerns the analysis of large kernel matrices under the concentration hypothesis. The results of this part are presented in Section 4.1 and are particularly based on the following work:

- (C2) **MEA. Seddik**, M. Tamaazousti, R. Couillet, “Kernel Random Matrices of Large Concentrated Data: The Example of GAN-generated Images”, IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP’19), Brighton, United-Kingdom, 2019.

Specifically, we have analyzed in this work the behavior of large kernel matrices of the form

$$\mathbf{K} = \left\{ f \left(\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right) \right\}_{i,j=1}^n$$

where $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is supposed to be three times continuously-differentiable in the vicinity of $\tau = \frac{2}{p} \text{tr} \left(\sum_{\ell=1}^k c_\ell \boldsymbol{\Sigma}_\ell \right)$, and the data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ is assumed to be concentrated in the sense of Definition 2. Under these assumptions and further the non-trivial growth rate assumptions from Subsection 1.1.2, our first result states that *the between and within class data are “equidistant” in high-dimension independently from the classes*, specifically, for some $\delta > 0$ with probability $1 - \delta$

$$\max_{1 \leq i \neq j \leq n} \left\{ \left| \frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \tau \right| \right\} = \mathcal{O} \left(\frac{\log \left(\frac{p}{\sqrt{\delta}} \right)^{\frac{1}{q}}}{\sqrt{p}} \right)$$

Following the same approach as [EK⁺10b, CBG⁺16], the kernel matrix \mathbf{K} can therefore be Taylor expanded entry-wise leading to an approximation of it in spectral norm by a spiked random matrix model of the form⁸

$$\mathbf{K} = \mathbf{P}_f + f'(\tau) \mathbf{Z}^\top \mathbf{Z} + f''(\tau) \mathbf{W} + o_p(1)$$

where \mathbf{P}_f is a low-rank informative matrix which contains information about the classes though the data statistics $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ and also depends on the kernel function f through its local first and second derivatives at τ . The remaining matrix terms are non-informative full-rank noise matrices which exhibit a spreading of eigenvalues in the spectrum of \mathbf{K} , in particular, the behavior of the term $f'(\tau) \mathbf{Z}^\top \mathbf{Z}$ is described by the deterministic equivalent from the previous subsection. See Section 4.1 for more details. As for the Gram matrix, \mathbf{K} exhibits a *universal behavior* since its random matrix equivalent depends *strictly* on the classes statistics $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$.

1.3.3 Beyond Kernels to Neural Networks

Kernel matrices appear naturally as a backbone of kernel methods. For instance, it has been shown in [Lia19] that the MSE⁹ of random neural networks, the misclassification error rate of kernel ridge regression and the performance of kernel/random feature-based

⁸The notation $\mathbf{A} = \mathbf{B} + o_p(1)$ means that $\|\mathbf{A} - \mathbf{B}\| = o(1)$ as $p \rightarrow \infty$ where $\|\cdot\|$ stands for the spectral norm.

⁹Mean squared error.

spectral clustering methods, all depend *explicitly* on the eigenspectrum or on a certain functional of a particular random kernel/nonlinear Gram matrix. In this subsection, the aim is to go beyond kernel methods in order to analyze ML methods which have an *implicit* relationship with the input data, i.e., ML methods which are implicitly determined by (convex) optimization problems. In particular, relying on [MLC19] which studies the behavior of logistic regression under the RMT regime and using Gaussian assumptions on data, we push forward this study by assuming a k -class concentration model (see Definition 2) and therefore considering the more general Softmax classifier. Specifically, this subsection briefly presents our findings concerning the analysis of the Softmax classifier under the concentration hypothesis and the high-dimensional regime. This contribution is particularly detailed in Section 5.1 and is based on the following work:

(C4) MEA. Seddik, C.Louart, R. Couillet, M. Tamaazousti, “The Unexpected Deterministic and Universal Behavior of Large Softmax Classifiers”, AISTATS 2021.

Specifically, given a set of one-hot-vector¹⁰ labels $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^k$ corresponding to each data vector $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$, the Softmax classifier consists in minimizing the following objective function

$$\mathcal{L}(\mathbf{w}_1, \dots, \mathbf{w}_k) = -\frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{i\ell} \log p_{i\ell} + \frac{1}{2} \sum_{\ell=1}^k \lambda_\ell \|\boldsymbol{\mu}_\ell\|^2 \quad \text{with} \quad p_{i\ell} = \frac{\varphi(\mathbf{w}_\ell^\top \mathbf{x}_i)}{\sum_{j=1}^k \varphi(\mathbf{w}_j^\top \mathbf{x}_i)}$$

where $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top]^\top \in \mathbb{R}^{pk}$ stands for the class-weights vectors, λ_ℓ 's are class-wise regularization parameters and $\varphi : \mathbb{R} \rightarrow \mathbb{R}_+$. Cancelling the loss function gradient with respect to each weight vector \mathbf{w}_ℓ yields

$$\lambda_\ell \mathbf{w}_\ell = -\frac{1}{n} \sum_{i=1}^n \left(y_{i\ell} \psi(\mathbf{w}_\ell^\top \mathbf{x}_i) - \frac{\varphi(\mathbf{w}_\ell^\top \mathbf{x}_i)}{\sum_{j=1}^k \varphi(\mathbf{w}_j^\top \mathbf{x}_i)} \sum_{j=1}^k y_{ij} \psi(\mathbf{w}_j^\top \mathbf{x}_i) \right) \mathbf{x}_i$$

with $\psi = \varphi'/\varphi$. Our approach consists in writing the above expression as a contracting fixed point equation of \mathbf{W} . Specifically, for well chosen $\tilde{\mathbf{x}}_i$ and f_i , we have

$$\boldsymbol{\Lambda} \mathbf{W} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \quad \Rightarrow \quad \mathbf{W} = \Psi(\mathbf{W})$$

where $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n] \in \mathcal{M}_{kp, kn}$, $f_i(\tilde{\mathbf{X}}^\top \mathbf{W}) \in \mathbb{R}^{kn}$ and $\boldsymbol{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_k) \otimes \mathbf{I}_p \in \mathcal{M}_{kp}$. Therefore, Ψ is requested to be $(1 - \varepsilon)$ -Lipschitz for some $\varepsilon > 0$ under some assumptions on φ and the regularization parameters λ_ℓ . Due to randomness in the data matrix \mathbf{X} , Ψ is a contraction conditionally on some high-probable event \mathcal{A}_X (See Section 5.1 for more details).

Under the concentration hypothesis from Definition 2 and additional assumptions on φ and $(\lambda_\ell)_{\ell \in [k]}$, our first result states that there exists an event \mathcal{A}_X with $\mathbb{P}(\mathcal{A}_X) > 1 - Ce^{-cn}$ for some constants $C, c > 0$ such that the weights vector \mathbf{W} concentrates conditionally on \mathcal{A}_X with an observable diameter of order $\mathcal{O}(\sqrt{\log n/n})$, which we compactly write as

$$(\mathbf{W} \mid \mathcal{A}_X) \propto \mathcal{E}_q \left(\sqrt{\frac{\log n}{n}} \right)$$

¹⁰Defined as $y_{i\ell} = 1$ if $\mathbf{x}_i \in \mathcal{C}_\ell$ and $y_{i\ell} = 0$ otherwise.

The further characterization of the statistics of the weights vector \mathbf{W} requires the study of a resolvent of the form

$$\mathbf{Q} = \left(\mathbf{\Lambda} - \frac{1}{n} \tilde{\mathbf{X}} \mathbf{D} \tilde{\mathbf{X}}^\top \right)^{-1}$$

where \mathbf{D} is some block-diagonal matrix whose entries depend on the data matrix $\tilde{\mathbf{X}}$. The above resolvent is studied through a deterministic equivalent (see Section 5.1 for more details), which implies the following result.

Theorem 1.2 (Asymptotic statistics of the Softmax weights). *Under the concentration hypothesis from Definition 2 and supposing additional assumptions on φ and $(\lambda_\ell)_{\ell \in [k]}$. There exists a deterministic mapping*

$$\mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} = \mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \left(\{\boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell\}_{\ell=1}^k \right) : \mathbb{R}^{pk} \times \mathcal{M}_{pk} \rightarrow \mathbb{R}^{pk} \times \mathcal{M}_{pk}$$

depending only on the data statistics $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$, such that the equation

$$(\mathbf{m}, \mathbf{C}) = \mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{m}, \mathbf{C}) \quad \text{with} \quad \mathbf{m} \in \mathbb{R}^{pk}, \mathbf{C} \in \mathcal{M}_{pk}$$

admits a unique solution $(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W)$. Furthermore,

$$\|\mathbb{E}[\mathbf{W}] - \bar{\mathbf{m}}_W\| = \mathcal{O} \left(\sqrt{\frac{\log n}{n}} \right) \quad \text{and} \quad \|\mathbb{E}[\mathbf{W}\mathbf{W}^\top] - \bar{\mathbf{C}}_W\|_* = \mathcal{O} \left(\sqrt{\frac{\log n}{n}} \right)$$

Theorem 1.2 is fundamental in the sense that it states that the Softmax classifier retrieves information from the data *only* through the class-wise means and covariances, namely the statistics $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$, thereby highlighting *the universal character* of the Softmax classifier. In particular, our result generalizes the study of Mai and Liao [MLC19] of the logistic regression classifier studied in their work under a GMM assumption on data, while we extended their results to a k -class MCV (see Definition 2). Moreover, based on Theorem 1.2, the performances of the Softmax classifier become theoretically predictable as per the following corollary.

Corollary 1.1 (Generalization performance of the Softmax classifier). *For $\ell \in [k]$ let some new test data $\mathbf{x} \in \mathcal{C}_\ell$ and define $p_\ell(\mathbf{x}) = \varphi(\mathbf{w}_\ell^\top \mathbf{x}) / \sum_{j \in [k]} \varphi(\mathbf{w}_j^\top \mathbf{x})$. There exists $\bar{\kappa}^\ell \in \mathbb{R}^{k-1}$ and $\bar{\mathbf{K}}^\ell \in \mathcal{M}_{k-1}$ both depending only on $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ such that the test error*

$$E_t(\mathbf{x} \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\forall j \in [k] \setminus \{\ell\} : p_\ell(\mathbf{x}) \geq p_j(\mathbf{x}))$$

is asymptotically close to

$$\bar{E}_t(\mathbf{x} \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\mathbf{Z}_\ell \in \mathbb{R}_+^{k-1}) \quad \text{with} \quad \mathbf{Z}_\ell \sim \mathcal{N}(\bar{\kappa}^\ell, \bar{\mathbf{K}}^\ell)$$

In a nutshell, Corollary 1.1 affirms that the generalization error of the Softmax classifier is nothing but the cumulative distribution of a low-dimensional Gaussian vector, the mean and covariance of which depend strictly on the class wise means and covariances $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ of the input data. This notably demonstrates the *universality* property of the Softmax classifier regardless of the data distribution as long as it satisfies the concentration assumptions in Definition 2.

1.3.4 Summary of Section 1.3

The main core of this thesis concerns the analysis of the methods presented in the previous subsections under the concentration hypothesis on data as per Definition 2. A major outcome from these studies is that the behavior of the studied methods *solely* depend on the class-wise means and covariances $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$ of data, thereby highlighting the *universality* aspect of these methods. Our results notably support the validity of the GMM model in the high-dimensional regime as considered in [BGC16, CBG⁺16, Lia19, Mai19], and support the applicability of random matrix theory to the analysis of AI methods on realistic data as the surprising images generated by GANs. Note that the universality aspect does not suggest that real data are Gaussian vectors but rather that ML algorithms (at least for the studied examples) “see” the data only through its first and second order statistics.

1.4 Outline and Contributions

The rest of the manuscript is organized as follows: In Chapter 2 we provide some essential background and notions from random matrix theory along with some toy applications to ML through the so-called spiked models (see Section 2.3). We further provide in Section 2.4 some basic notions from the concentration of measure theory and recall some essential results from [LC18b] which will be used throughout the manuscript. In chapter 3 we present our first contribution concerning the analysis of large Gram matrices under the concentration assumption, from the works:

- (C1) MEA. Seddik, C. Louart, M. Tamaazousti, R. Couillet, “*Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures*”, International Conference on Machine Learning (ICML’20), Online, 2020.
- (C1’) MEA. Seddik, M. Tamaazousti, R. Couillet, “*Pourquoi les matrices aléatoires expliquent l’apprentissage ? Un argument d’universalité offert par les GANs*”, Colloque francophone de traitement du signal et des images (Gretsi’19), Lille, France, 2019.

In Chapter 4, we present our contributions to the analysis of kernel matrices. In particular, Section 4.1 presents our second contribution from the work:

- (C2) MEA. Seddik, M. Tamaazousti, R. Couillet, “*Kernel Random Matrices of Large Concentrated Data: The Example of GAN-generated Images*”, IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP’19), Brighton, United-Kingdom, 2019.

While Section 4.2 presents the analysis of kernel matrices of the type $f(\hat{\boldsymbol{\Sigma}})$ where $\hat{\boldsymbol{\Sigma}}$ stands for the sample covariance matrix, along with an application to the Sparse PCA problem. This work constitutes our third contribution from the work:

- (C3) MEA. Seddik, M. Tamaazousti, R. Couillet, “*A Kernel Random Matrix-Based Approach for Sparse PCA*”, International Conference on Learning Representations (ICLR’19), New Orleans, United-States, 2019.

In Chapter 5 we present applications to neural networks. In particular, Section 5.1 treats the study of the Softmax layer in neural networks which constitutes our fourth contribution from the work:

- (C4) **MEA. Seddik**, C.Louart, R. Couillet, M. Tamaazousti, “*The Unexpected Deterministic and Universal Behavior of Large Softmax Classifiers*”, AISTATS 2021.

Moreover, Section 5.2 presents an analysis of α -Dropout in a single hidden-layer network which basically consists in dropping out features at random and replacing them with some arbitrary value α . Our analysis exhibits a value of $\alpha \neq 0$ which maximizes the generalization accuracy of the studied network. This analysis is based on the following work:

- (C5) **MEA. Seddik**, R. Couillet, M. Tamaazousti, “*A Random Matrix Analysis of Learning with α -Dropout*”, The art of learning with missing values ICML workshop (ICML’20), Online, 2020.

In the appendix Chapter B, we present some practical contributions conducted during this thesis. In particular, the work:

- (C6) **MEA. Seddik**, M. Tamaazousti, J. Lin, “*Generative Collaborative Networks for Single Image SuperResolution*”, Neurocomputing’2019.

presents deep learning models for single-image super-resolution. The last contribution presents compression methods for dense neural networks, which is based on the work:

- (C7) **MEA. Seddik**, H.Essafi, A.Benzine, M.Tamaazousti, “*Lightweight Neural Networks from PCA & LDA Based Distilled Dense Neural Networks*”, International Conference on Image Processing (ICIP’20), Online, 2020.

Finally, Chapter 6 concludes the manuscript and discusses perspectives and the appendix Chapter C presents the proofs of the different theoretical results of this thesis.

Chapter 2

Random Matrix Theory & Concentration of Measure Theory

Contents

2.1 Fundamental Random Matrix Theory Notions	30
2.1.1 The resolvent matrix notion	30
2.1.2 Stieltjes transform and spectral measure	30
2.1.3 Cauchy's integral and statistical inference	32
2.1.4 Deterministic and random equivalents	33
2.2 Fundamental Random Matrix Theory Results	34
2.2.1 Matrix identities and key lemmas	34
2.2.2 The Marčenko-Pastur Law	37
2.2.3 Random matrices of mixture models	38
2.3 Connections with machine learning through spiked models	39
2.3.1 Simple example of unsupervised learning	41
2.3.2 Simple example of supervised learning	45
2.4 Extensions with concentration of measure theory	48
2.4.1 The notion of concentrated vectors	48
2.4.2 Resolvent of the sample covariance matrix	50

Random matrix theory, originally, aims to describe the eigenvalue distribution (also known as the spectral measure) of large random matrices. There exist different approaches in the literature to analyze these random matrices such as the moments method, or approaches involving tools from free probabilities, see [Tao12] and the references therein for an overview. In this manuscript we will consider the Stieltjes transform approach, which is often presented as the central notion of the theory [BS⁺98a].

Tackling machine learning applications of random matrix theory, the description of random matrices eigenvalues distribution is not the central interest, and the objects of interest are more fundamentally sub-spaces which are described by largest eigenvectors of random matrices (in the case of unsupervised learning) or quadratic forms involving random matrices (in the case of supervised learning). Examples (but not limited to) of machine learning methods which rely on eigenvectors of random matrices are: principal component analysis (PCA) [WEG87], spectral clustering [VLBB08], some semi-supervised learning approaches [AMGS12], Least-squared support vector machines [SV99]

and random neural networks also known as extreme learning machines [HZS06].

Therefore, a more general notion than the Stieltjes transform is of central interest for machine learning application of the theory. This notion, generally referred to as the *resolvent* of large random matrices and will be the central notion of this manuscript. The resolvent of a matrix allows the access and consistent estimation of complex functionals of this matrix (such as the aforementioned quadratic forms), and fundamentally describes its spectral measure, permits the location and description of its isolated eigenvalues and eigenvectors, and provides estimations for bilinear forms, therefore permits the estimation of the performances of various machine learning models. We refer the reader to the works [Mai19, Lia19] where applications of the theory to various machine learning algorithms are presented relying on Gaussian mixture models.

We begin this section by first introducing the fundamental notions from random matrix theory, then give fundamental random matrix theory results and at the end present extensions involving tools from the concentration of measure theory which will allow the theory, as we will present subsequently, to find applications using realistic data.

2.1 Fundamental Random Matrix Theory Notions

2.1.1 The resolvent matrix notion

We first introduce the notion of the resolvent matrix.

Definition 3 (Resolvent). *For a given symmetric matrix $\mathbf{M} \in \mathcal{M}_p$, the resolvent $\mathbf{Q}(z)$ of \mathbf{M} is defined, for $z \in \mathbb{C} \setminus \mathbb{R}^+$, as*

$$\mathbf{Q}(z) \equiv (\mathbf{M} - z\mathbf{I}_p)^{-1}. \quad (2.1)$$

We will systematically adopt the following notations for the resolvent matrix throughout this manuscript.

Notation 1. *Let $\mathbf{X} \in \mathcal{M}_{p,n}$ and $z \in \mathbb{C} \setminus \mathbb{R}^+$, the matrix \mathbf{M} in Definition 3, for which we will define a resolvent, will be either of the form $\mathbf{X}\mathbf{X}^\top$ or $\mathbf{X}^\top\mathbf{X}$, we will thus denote*

$$\mathbf{Q}(z) \equiv (\mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p)^{-1} \in \mathcal{M}_p, \quad \mathbf{R}(z) \equiv (\mathbf{X}^\top\mathbf{X} - z\mathbf{I}_n)^{-1} \in \mathcal{M}_n \quad (2.2)$$

Note that the dependence on z will be some times omitted if there is no ambiguity and we will simply write \mathbf{Q} or \mathbf{R} instead of $\mathbf{Q}(z)$ or $\mathbf{R}(z)$ respectively.

2.1.2 Stieltjes transform and spectral measure

The resolvent \mathbf{Q} is naturally related to the *empirical spectral measure* μ of \mathbf{M} , through the Stieltjes transform m_μ , which we will be defined subsequently.

Definition 4 (Empirical spectral measure). *For a given symmetric matrix $\mathbf{M} \in \mathcal{M}_p$, the empirical spectral measure μ associated to \mathbf{M} is defined through the normalized counting measure of the eigenvalues $\lambda_1, \dots, \lambda_p$ of \mathbf{M} , specifically,*

$$\mu(x) \equiv \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i}(x). \quad (2.3)$$

The spectral measure μ of \mathbf{M} is a (random if \mathbf{M} is random) probability measure, since for all $x \in \mathbb{R}$, $\mu(x) \geq 0$ and since $\int_{\mathbb{R}} \mu(x) dx = 1$. Therefore, we can define its associated Stieltjes transform as follows.

Definition 5. (Stieltjes Transform) Given some real probability measure μ with support $\mathcal{S}(\mu)$, the Stieltjes transform $q(z)$ is defined, for all $z \in \mathbb{C} \setminus \mathcal{S}(\mu)$, as

$$q(z) \equiv \int_{\mathbb{R}} \frac{d\mu(\lambda)}{\lambda - z}. \quad (2.4)$$

Note that the Stieltjes transform $q(z)$ is closely related to its associated resolvent matrix \mathbf{Q} through the following algebraic link.

$$q(z) = \frac{1}{p} \sum_{i=1}^p \int_{\mathbb{R}} \frac{\delta_{\lambda_i}(\lambda)}{\lambda - z} = \frac{1}{p} \sum_{i=1}^p \frac{1}{\lambda_i - z} = \frac{1}{p} \text{tr}(\mathbf{Q}(z)) \quad (2.5)$$

The Stieltjes transform q has various benefits and interesting properties, which are:

1. q is complex analytic on its definition domain $\mathbb{C} \setminus \mathcal{S}(\mu)$.
2. $q(z)$ is bounded for $z \in \mathbb{C} \setminus \mathcal{S}(\mu)$ as

$$|q(z)| \leq \frac{1}{\text{dist}(z, \mathcal{S}(\mu))}.$$

3. if $\Im(z) > 0$ then $\Im[q(z)] > 0$.
4. Since $q'(z) = \int_{\mathbb{R}} (t - z)^{-2} dt > 0$, m is an increasing function with $\lim_{x \rightarrow \pm\infty} q(x) = 0$ if $\mathcal{S}(\mu)$ is bounded.

The Stieltjes transform q as its name implies (transform), admits an inverse formula which notably provides the recovery of its underlying spectral measure μ , as per the following theorem.

Theorem 2.1 (Inverse formula of the Stieltjes transform). *Let a, b be some continuity points of the probability measure μ , therefore the segment $[a, b]$ is measurable with μ and we precisely have*

$$\mu([a, b]) = \frac{1}{\pi} \lim_{\epsilon \rightarrow 0} \int_a^b \Im[q(x + i\epsilon)] dx. \quad (2.6)$$

Moreover, if μ admits a density function f at some point x , i.e., $\mu(x)$ is differentiable in a neighborhood of x with $\lim_{\epsilon \rightarrow 0} \epsilon^{-1} \mu([x - \epsilon/2, x + \epsilon/2]) = f(x)$, then we have the inverse formula

$$f(x) = \frac{1}{\pi} \lim_{\epsilon \rightarrow 0} \Im[q(x + i\epsilon)]. \quad (2.7)$$

And finally, if μ has an isolated mass at some point x , then

$$\mu(\{x\}) = -\frac{1}{\pi} \lim_{\epsilon \rightarrow 0} \Im[q(x + i\epsilon)] \quad (2.8)$$

The resolvent matrix \mathbf{Q} can be viewed as a matrix-valued Stieltjes transform. Indeed, relying on [HLN⁺07], the definition of the Stieltjes transform can be extended to $p \times p$ matrix-valued positive measures $\mathbf{M}(d\lambda)$ in the sense that

$$\mu_v(d\lambda) \equiv v^\top \mathbf{M}(d\lambda) v, \quad (2.9)$$

is a positive real-valued measure for all $v \in \mathbb{R}^p$. Using the spectral decomposition $M = \mathbf{U} \text{diag}\{\lambda_i\}_{i=1}^p \mathbf{U}^\top$ where $\lambda_1, \dots, \lambda_p$ stand for the eigenvalues of M , the resolvent matrix $\mathbf{Q}(z)$ of M can be written as

$$\mathbf{Q}(z) = \int_{\mathbb{R}} \frac{M(d\lambda)}{\lambda - z} = \mathbf{U} \text{diag} \left\{ \frac{1}{\lambda_i - z} \right\}_{i=1}^p \mathbf{U}^\top \quad (2.10)$$

And therefore, $\mathbf{Q}(z)$ verifies similar properties as the real-valued Stieltjes transform, in particular, $\mathbf{Q}(z)$ is complex analytic and satisfies $\|\mathbf{Q}\| \leq \text{dist}(z, \mathcal{S}(\mu))^{-1}$.

2.1.3 Cauchy's integral and statistical inference

Since \mathbf{Q} is complex analytic, it benefits from complex analysis tools which allow us to infer statistical quantities involving complex functionals of M . In particular, we have the following theorem which provides an elegant bridge between the resolvent and Cauchy's integral theorem.

Theorem 2.2 (Cauchy's integral). *Given $\Gamma \subset \mathbb{C}$ some positively oriented contour and a complex function $\zeta(z)$ which is analytic in a region containing the contour Γ , we have*

- if Γ is surrounding some $z' \in \mathbb{C}$, then $\zeta(z') = -\frac{1}{2\pi i} \oint_{\Gamma} \frac{\zeta(z)}{z' - z} dz$;
- otherwise, $\frac{1}{2\pi i} \oint_{\Gamma} \frac{\zeta(z)}{z' - z} dz = 0$.

Using this result, one can immediately access *linear functionals of the eigenvalues* $\lambda_1, \dots, \lambda_p$ of M and using the Stieltjes transform or the resolvent matrix $\mathbf{Q}(z)$ by

$$\frac{1}{p} \sum_{i=1}^p \zeta(\lambda_i) = -\frac{1}{2\pi i} \oint_{\Gamma} \zeta(z) \frac{1}{p} \text{tr}(\mathbf{Q}(z)) dz = -\frac{1}{2\pi i} \oint_{\Gamma} \zeta(z) q(z) dz \quad (2.11)$$

for all ζ complex analytic in a compact neighborhood of $\mathcal{S}(\mu)$, the support of μ .

In addition to estimate functionals of the eigenvalues values of M , one can access its eigenvectors and eigenspaces with similar arguments. Indeed, through the spectral decomposition of $M = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ with $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p] \in \mathcal{M}_p$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$, we have

$$\mathbf{Q}(z) = \sum_{i=1}^p \frac{\mathbf{u}_i \mathbf{u}_i^\top}{\lambda_i - z}. \quad (2.12)$$

And if some eigenvalue λ_i is of multiplicity one, we can directly access its corresponding eigenvector \mathbf{u}_i through

$$\mathbf{u}_i \mathbf{u}_i^\top = -\frac{1}{2\pi i} \oint_{\Gamma_i} \mathbf{Q}(z) dz, \quad (2.13)$$

where Γ_{λ_i} a positive oriented contour surrounding the eigenvalue λ_i . We will be particularly interested in measuring the alignment between \mathbf{u}_i and some deterministic vector \mathbf{u} , so the evaluation of the quantity

$$|\mathbf{u}^\top \mathbf{u}_i|^2 = -\frac{1}{2\pi i} \oint_{\Gamma_{\lambda_i}} \mathbf{u}^\top \mathbf{Q}(z) \mathbf{u} dz \quad (2.14)$$

Note that the resolvent matrix \mathbf{Q} provides access to the quantities $\frac{1}{p} \sum_i \zeta(\lambda_i)$ and $|\mathbf{u}^\top \mathbf{u}_i|$ through scalar quantities and quadratic forms involving \mathbf{Q} , i.e., the quantities $\text{tr}(\mathbf{Q}(z))$ and $\mathbf{u}^\top \mathbf{Q}(z) \mathbf{u}$. The estimation of these quantities will be notably possible through the notion of *deterministic equivalent* [HLN⁺07] which will be used as a proxy for statistical inference. The next subsection provides a description for this notion.

2.1.4 Deterministic and random equivalents

Throughout this manuscript the matrix M is a *large dimensional random matrix*, which is of either the following forms: 1. a sample covariance matrix $\frac{1}{n}XX^\top$; 2. a Gram matrix $\frac{1}{p}X^\top X$; 3. a kernel covariance matrix $f\left(\frac{1}{n}XX^\top\right)$ where f is applied entry-wise; 4. or a kernel random of the form $\left\{f\left(\frac{1}{p}\|x_i - x_j\|^2\right)\right\}_{i=1}^n$, where $X = [x_1, \dots, x_n] \in \mathcal{M}_{p,n}$ is a random matrix and f a function with some regularity conditions. In the case of linear models, i.e., matrices of the form $M = \frac{1}{n}XX^\top$ and under some growth rate conditions, the associated resolvent matrix $Q(z)$ has a deterministic behavior [HLN⁺07] and therefore equivalent in some sense to a deterministic matrix $\bar{Q}(z)$. Whereas, for kernel models involving some non-linear function f , the internal mechanism of the matrix M is accessible through a random matrix equivalent [EK⁺10b, CBG⁺16] \tilde{M} which results from a linearization of M and will allow the study of M via a deterministic equivalent of the resolvent of \tilde{M} .

The notion of *deterministic equivalent* of a resolvent matrix Q is related to the existence of a *non-asymptotic* deterministic matrix having, in probability or almost surely, the same scalar observations as the random observations through Q . This character is notably a manifestation of the concentration of measure phenomenon which we will discuss at the end of this chapter. We therefore have the following definition for the notion of deterministic equivalents [HLN⁺07].

Definition 6 (Deterministic Equivalent). *A squared deterministic matrix $\bar{Q} \in \mathcal{M}_p$ is said to be a deterministic equivalent for the symmetric random matrix $Q \in \mathcal{M}_p$ if, for all deterministic matrix $A \in \mathcal{M}_p$ and vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^p$ of bounded norms (operator and Euclidean, respectively), we have, as $p \rightarrow \infty$, with some probability or almost surely*

$$\frac{1}{p} \operatorname{tr} A(Q - \bar{Q}) \rightarrow 0, \quad \mathbf{a}^\top (Q - \bar{Q}) \mathbf{b} \rightarrow 0 \quad (2.15)$$

Remark 2.1 (On the definition of deterministic equivalents). *Definition 6 can be compactly expressed thanks to the linear concentration notion introduced by [LC18b, Definition 2.1] that is particularly adapted with concentration of measure theory presented subsequently in Section 2.4. Specifically, a deterministic matrix $\bar{Q} \in \mathcal{M}_p$ is a deterministic equivalent of a random matrix $Q \in \mathcal{M}_p$ if for any bounded linear form $u : \mathcal{M}_p \rightarrow \mathbb{R}$ with a unit operator norm (i.e., for all $A \in \mathcal{M}_p$, $|u(A)| \leq \|A\|$) we have*

$$\forall t > 0 : \mathbb{P} \{ |u(Q) - u(\bar{Q})| \geq t \} \leq \alpha(t)$$

where $\alpha : \mathbb{R}_+ \rightarrow [0, 1]$ is any non-increasing and left continuous concentration function. Note that \mathcal{M}_p could be replaced with any normed space $(\mathcal{E}, \|\cdot\|)$ but for convenience we restrict the notion to \mathcal{M}_p enrolled with the operator norm $\|\cdot\|$.

Therefore, Definition 6 allows to access spectral information about the random matrix M . For instance, the Stieltjes transform of its (random) spectral measure and bilinear forms involving $Q(z)$, the resolvent of M . For instance, if a deterministic equivalent $\bar{Q}(z)$ for $Q(z)$ exists in the sense of Definition 6, we have that almost surely $\frac{1}{p} \operatorname{tr}(Q(z) - \bar{Q}(z)) \rightarrow 0$ which implies that the Stieltjes transform m of the spectral distribution μ converges to the deterministic quantity $\frac{1}{p} \operatorname{tr} \bar{Q}(z)$. Similarly, and exploiting Definition 6 using the bilinear form involving $Q(z)$, the (isolated) eigenvectors of M are thus accessible through the deterministic equivalent $Q(z)$.

Note that to exhibit such deterministic equivalent, it suffices to have a control of $\|\mathbb{E}Q - \bar{Q}\|$ which generally unfolds from a control of the variance of the quantities $\frac{1}{p} \text{tr}(AQ)$ and $\mathbf{a}^\top Q \mathbf{b}$. $\mathbb{E}Q$ is also a deterministic equivalent for Q but is not tractable in practice, however, we can generally find \bar{Q} such that $\|\mathbb{E}Q - \bar{Q}\| \rightarrow 0$ where \bar{Q} is computable through fixed-point equations or in some rare cases in close form, such as the classical examples of Marchenko-Pastur and the semi-circle laws.

We consider the following notations when a deterministic equivalent \bar{Q} exists for a random resolvent matrix Q , and when a random matrix equivalent \tilde{M} exists for a random matrix M .

Notation 2 (Deterministic equivalent). *Let $Q \in \mathcal{M}_p$ be a random matrix and $\bar{Q} \in \mathcal{M}_p$ a deterministic matrix, we write*

$$Q \leftrightarrow \bar{Q} \quad (2.16)$$

if \bar{Q} satisfies Definition 6.

Notation 3 (Random equivalent). *Let $M, \tilde{M} \in \mathcal{M}_p$ be two random matrices, we write*

$$M \rightsquigarrow \tilde{M} \quad (2.17)$$

if $\|\mathbb{E}[M - \tilde{M}]\| \rightarrow 0$.

Note that we will denote, throughout the manuscript, deterministic equivalents with an up bar symbol (i.e., \bar{Q}) and the random equivalents with an up tilde symbol (i.e., \tilde{M}).

2.2 Fundamental Random Matrix Theory Results

This section presents fundamental random matrix theory results, which provide explanation about some machine learning algorithms. We particularly provide the historical result of Marchenko-Pastur [MP67] which describes the spectral distribution of large sample covariance matrices. Since machine learning algorithms are generally applied to data described by distinct “clusters”, this section will also present the extension of the Marchenko-Pastur result to mixture models data [BGC16]. The nice thing about random matrix theory is that it relies on tools from different branches of mathematics such as complex analysis (as we previously saw), linear algebra and probability theory. We will first briefly present some essential matrix identities and statistical tools and then provide the aforementioned results with short sketches of proofs (for readability).

2.2.1 Matrix identities and key lemmas

As we discussed in the previous section, to exhibit a deterministic equivalent $\bar{Q}(z)$ of a given resolvent random matrix $Q(z)$, one needs to control the quantity $\|\mathbb{E}[Q(z) - \bar{Q}(z)]\|$. Since both $Q(z)$ and $\bar{Q}(z)$ are inverse of matrices, the following identity is of central interest.

Lemma 2.1 (Resolvent identity). *Let A and B be some invertible matrices, we have*

$$A^{-1} - B^{-1} = A^{-1}(B - A)B^{-1} \quad (2.18)$$

Proof. Simply through $A^{-1}(B - A)B^{-1} = A^{-1}BB^{-1} - A^{-1}AB^{-1} = A^{-1} - B^{-1}$. \square

As we have discussed, the random matrix \mathbf{M} can be of the form $\mathbf{X}\mathbf{X}^\top$ or $\mathbf{X}^\top\mathbf{X}$ for some random matrix $\mathbf{X} \in \mathcal{M}_{p,n}$, thus we have the following lemma which relates the resolvent matrices of $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{X}^\top\mathbf{X}$.

Lemma 2.2. *For $\mathbf{X} \in \mathcal{M}_{p,n}$, we have*

$$\mathbf{X}(\mathbf{X}^\top\mathbf{X} - z\mathbf{I}_n)^{-1} = (\mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p)^{-1}\mathbf{X}, \quad (2.19)$$

for $z \in \mathbb{C} \setminus \{0\}$ distinct from the eigenvalues of $\mathbf{X}\mathbf{X}^\top$.

We also have the following lemma, known as Sylvester's identity, which provides a link between the resolvents of $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{X}^\top\mathbf{X}$ through the determinant operator.

Lemma 2.3 (Sylvester's identity). *For $\mathbf{X} \in \mathcal{M}_{p,n}$, we have*

$$\det(\mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p) = \det(\mathbf{X}^\top\mathbf{X} - z\mathbf{I}_n)(-z)^{p-n}. \quad (2.20)$$

Sylvester's identity notably shows that $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{X}^\top\mathbf{X}$ share the same non-zero eigenvalues except $n - p$ zero eigenvalues of $\mathbf{X}^\top\mathbf{X}$ if $n \geq p$. We thus have the following result which relates the Stieltjes transforms of the spectral measures $\mu_{\mathbf{X}\mathbf{X}^\top}$ and $\mu_{\mathbf{X}^\top\mathbf{X}}$ of the random matrices $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{X}^\top\mathbf{X}$ respectively.

Lemma 2.4. *Let $\mathbf{X} \in \mathcal{M}_{p,n}$ with $n \geq p$, and $z \in \mathbb{C} \setminus \{0\}$ not an eigenvalue of $\mathbf{X}^\top\mathbf{X}$. Denote $\mathbf{Q}(z) \equiv (\mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p)^{-1}$ and $\mathbf{R}(z) \equiv (\mathbf{X}^\top\mathbf{X} - z\mathbf{I}_n)^{-1}$, we have*

$$\text{tr } \mathbf{Q}(z) = \text{tr } \mathbf{R}(z) + \frac{n-p}{z}. \quad (2.21)$$

In particular,

$$\mu_{\mathbf{X}\mathbf{X}^\top} = \frac{n}{p}\mu_{\mathbf{X}^\top\mathbf{X}} + \frac{n-p}{p}\delta_0. \quad (2.22)$$

Generally, the key approach to prove random matrix theory results relies on *leave-one-out* or *perturbation* approach to handle the dependencies. This approach consists in removing some i -th contribution from the resolvent \mathbf{Q} of \mathbf{M} to build a resolvent \mathbf{Q}_{-i} of \mathbf{M}_{-i} deprived of the i -th contribution. The main idea being that \mathbf{Q} and \mathbf{Q}_{-i} still have the same behavior but \mathbf{Q}_{-i} allows one to handle the dependencies. Technically, \mathbf{Q}_{-i} is built using matrix identities such as the Woodbury identity given in the following lemma.

Lemma 2.5 (Woodbury identity). *For $\mathbf{M} \in \mathcal{M}_p$ and $\mathbf{U}, \mathbf{V} \in \mathcal{M}_{p,k}$, such that \mathbf{M} and $\mathbf{M} + \mathbf{U}\mathbf{V}^\top$ are invertible, we have*

$$(\mathbf{M} + \mathbf{U}\mathbf{V}^\top)^{-1} = \mathbf{M}^{-1} - \mathbf{M}^{-1}\mathbf{U}(\mathbf{I}_k + \mathbf{V}^\top\mathbf{M}\mathbf{U})\mathbf{V}^\top\mathbf{M}^{-1} \quad (2.23)$$

The perturbation $\mathbf{U}\mathbf{V}^\top$ is generally of low rank k and the special case $k = 1$, \mathbf{U} and \mathbf{V} reduces to vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$ respectively. In this case, we simply have the Sherman-Morrison formula given in the following lemma.

Lemma 2.6 (Sherman-Morrison). *Let $\mathbf{M} \in \mathcal{M}_p$ invertible and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^p$ such that $1 + \mathbf{v}^\top\mathbf{M}^{-1}\mathbf{u} \neq 0$, we have*

$$(\mathbf{M} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{M}^{-1} - \frac{\mathbf{M}^{-1}\mathbf{u}\mathbf{v}^\top\mathbf{M}^{-1}}{1 + \mathbf{v}^\top\mathbf{M}^{-1}\mathbf{u}}, \quad (\mathbf{M} + \mathbf{u}\mathbf{v}^\top)^{-1}\mathbf{u} = \frac{\mathbf{M}^{-1}\mathbf{u}}{1 + \mathbf{v}^\top\mathbf{M}^{-1}\mathbf{u}}. \quad (2.24)$$

We give in the following an example where these identities help to handle the dependencies in random matrices.

Example 1. Suppose that $\mathbf{Q}(z) = (\frac{1}{n}\mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p)^{-1}$ where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ and \mathbf{x}_i 's are independent random vectors. Further denote $\mathbf{Q}_{-i}(z) = (\frac{1}{n}\mathbf{X}\mathbf{X}^\top - \frac{1}{n}\mathbf{x}_i\mathbf{x}_i^\top - z\mathbf{I}_p)^{-1}$ the resolvent $\mathbf{Q}(z)$ deprived from the i -th vector \mathbf{x}_i . Using Lemma 2.6 we therefore have the identities

$$\mathbf{Q}(z) = \mathbf{Q}_{-i}(z) - \frac{\mathbf{Q}_{-i}(z)\frac{1}{n}\mathbf{x}_i\mathbf{x}_i\mathbf{Q}_{-i}(z)}{1 + \frac{1}{n}\mathbf{x}_i^\top\mathbf{Q}_{-i}(z)\mathbf{x}_i}, \quad \mathbf{Q}(z)\mathbf{x}_i = \frac{\mathbf{Q}_{-i}(z)\mathbf{x}_i}{1 + \frac{1}{n}\mathbf{x}_i^\top\mathbf{Q}_{-i}(z)\mathbf{x}_i} \quad (2.25)$$

where $\mathbf{Q}_{-i}(z)$ is independent of \mathbf{x}_i by construction.

Lemma 2.6 will also be of particular interest when handling informative random matrix models as we will present in Section 2.3. A useful application of Lemma 2.6 is when replacing $\mathbf{M} \leftarrow \mathbf{M} - z\mathbf{I}_p$ for $z \in \mathbb{C}$, and letting $\mathbf{u} = \gamma\mathbf{u}$ where $\gamma \in \mathbb{R}$, which leads to the following rank-1 perturbation lemma for the resolvent of \mathbf{M} .

Lemma 2.7 (Perturbation lemma [SB95]). *Let $\mathbf{A}, \mathbf{M} \in \mathcal{M}_p$ some symmetric matrices, $\mathbf{u} \in \mathbb{R}^p$, $\gamma \in \mathbb{R}$ and $z < 0$, then*

$$\left| \operatorname{tr} \mathbf{A}(\mathbf{M} + \gamma\mathbf{u}\mathbf{u}^\top - z\mathbf{I}_p)^{-1} - \operatorname{tr} \mathbf{A}(\mathbf{M} - z\mathbf{I}_p)^{-1} \right| \leq \frac{\|\mathbf{A}\|}{|\Im(z)|} \quad (2.26)$$

Note that the bound in Lemma 2.7 does not depend on $\|\mathbf{u}\|$. In particular, denoting the (perturbed) resolvent $\mathbf{Q}_\gamma(z) \equiv (\mathbf{M} + \gamma\mathbf{u}\mathbf{u}^\top - z\mathbf{I}_p)^{-1}$, and letting $\mathbf{A} = \frac{1}{p}\mathbf{I}_p$, we obtain

$$\frac{1}{p} \operatorname{tr} \mathbf{Q}_\gamma(z) = \frac{1}{p} \operatorname{tr} \mathbf{Q}(z) + \mathcal{O}(p^{-1}), \quad (2.27)$$

which shows, since $\frac{1}{p} \operatorname{tr} \mathbf{Q}(z)$ is the Stieltjes transform of the empirical spectral measure of \mathbf{M} , that the spectral measure of $\mathbf{M} + \gamma\mathbf{u}\mathbf{u}^\top$ is asymptotically close to that of \mathbf{M} for any \mathbf{u} , in the large limit of p .

We further have the following fundamental trace Lemma which is at the core of the study of various random matrix models.

Lemma 2.8 (Trace Lemma [BS08]). *Let $\mathbf{x} \in \mathbb{R}^p$ a random vector with i.i.d. entries with zero mean, unit variance and finite $2k$ order moment for some $k > 1$. Let $\mathbf{Q} \in \mathcal{M}_p$ some deterministic or random (independent of \mathbf{x}) matrix, then there exists a constant $C > 0$ such that*

$$\mathbb{E} \left[\left| \frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \frac{1}{p} \operatorname{tr} \mathbf{Q} \right|^k \right] \leq C \frac{\|\mathbf{Q}\|^k}{p^{\frac{k}{2}}}$$

In particular, if $\limsup_p \|\mathbf{Q}\| < \infty$, and \mathbf{x} has entries with finite fourth-order moment, we have by Markov's inequality and Borel Cantelli Lemma,

$$\frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \frac{1}{p} \operatorname{tr} \mathbf{Q} \xrightarrow{a.s.} 0$$

Lemma 2.8 basically shows that the quadratic form $\frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ concentrates around its expectation at large p . This result will particularly be exploited in the following Subsection to prove the fundamental Marčenko-Pastur law.

2.2.2 The Marčenko-Pastur Law

Having set the principal tools and notions presented previously, we are now in place to state one of the fundamental random matrix theory results, namely the Marčenko-Pastur law.

Theorem 2.3 (Marčenko-Pastur law [MP67]). *Let $\mathbf{X} \in \mathcal{M}_{p,n}$ be a random matrix with i.i.d. entries having zero mean, unit variance and bound fourth order moment. Then, as $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, the empirical spectral distribution μ of $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$ satisfies*

$$\mu \xrightarrow{a.s.} \mu_{MP}$$

in the weak convergence sense, where μ_{MP} is a deterministic measure and particularly satisfies $\mu_{MP}(\{0\}) = \max\{0, 1 - c^{-1}\}$, moreover μ_{MP} has a continuous density function f_{MP} on the compact support $[(1 - \sqrt{c})^2, (1 + \sqrt{c})^2]$ defined as

$$f_{MP}(x) = \frac{1}{\sqrt{2\pi cx}} \sqrt{(x - \lambda_-)(\lambda_+ - x)}$$

where $\lambda_{\pm} = (1 \pm \sqrt{c})^2$.

Proof. The proof starts by writing the Stieltjes transform $m(z)$ of μ as

$$m(z) = \frac{1}{p} \operatorname{tr} \left(\frac{1}{n} \mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p \right)^{-1} = \frac{1}{p} \sum_{i=1}^p \left[\left(\frac{1}{n} \mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p \right)^{-1} \right]_{ii}$$

and writing

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{X}_{-1} \end{bmatrix}$$

therefore, for $\Im(z) > 0$,

$$\begin{aligned} \left[\left(\frac{1}{n} \mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p \right)^{-1} \right]_{11} &= \left[\begin{pmatrix} \frac{1}{n} \mathbf{x}_1^\top \mathbf{x}_1 - z & \frac{1}{n} \mathbf{x}_1^\top \mathbf{X}_{-1} \\ \frac{1}{n} \mathbf{X}_{-1} \mathbf{x}_1 & \frac{1}{n} \mathbf{X}_{-1} \mathbf{X}_{-1}^\top - z\mathbf{I}_{p-1} \end{pmatrix}^{-1} \right]_{11} \\ &= \frac{1}{-z - z \frac{1}{n} \mathbf{x}_1^\top \left(\frac{1}{n} \mathbf{X}_{-1}^\top \mathbf{X}_{-1} - z\mathbf{I}_{p-1} \right)^{-1} \mathbf{x}_1} \end{aligned}$$

where the last equality comes from standard block matrix inverse formula from linear algebra. Therefore, by Trace Lemma 2.8, as $p, n \rightarrow \infty$

$$\left[\left(\frac{1}{n} \mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p \right)^{-1} \right]_{11} = \frac{1}{-z - z \frac{1}{n} \operatorname{tr} \left(\frac{1}{n} \mathbf{X}_{-1}^\top \mathbf{X}_{-1} - z\mathbf{I}_{p-1} \right)^{-1}} \xrightarrow{a.s.} 0$$

Moreover, since

$$\mathbf{X}^\top \mathbf{X} = \mathbf{X}_{-1}^\top \mathbf{X}_{-1} + \mathbf{x}_1 \mathbf{x}_1^\top$$

we have by Rank-1 perturbation Lemma 2.7

$$\left[\left(\frac{1}{n} \mathbf{X}\mathbf{X}^\top - z\mathbf{I}_p \right)^{-1} \right]_{11} = \frac{1}{-z - z \frac{1}{n} \operatorname{tr} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} - z\mathbf{I}_n \right)^{-1}} \xrightarrow{a.s.} 0$$

and recalling Lemma 2.4, we have $\frac{1}{n} \text{tr} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} - z \mathbf{I}_n \right)^{-1} = \frac{1}{n} \text{tr} \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1} - \frac{n-p}{zn}$ we have

$$\left[\left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1} \right]_{11} - \frac{1}{1 - \frac{p}{n} - z - z \frac{1}{n} \text{tr} \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1}} \xrightarrow{a.s.} 0$$

Repeating the same procedure for the diagonal entries of $\left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1}$ up to p , and averaging, we have for $\Im(z) > 0$

$$m(z) - \frac{1}{1 - \frac{p}{n} - z - z \frac{p}{n} m(z)} \xrightarrow{a.s.} 0$$

Therefore, $m(z) \xrightarrow{a.s.} m_{MP}(z)$ which is a solution of the *fixed point equation* (with positive branch)

$$m_{MP}(z) = \frac{1}{1 - c - z - cz m_{MP}(z)}$$

thus

$$m_{MP}(z) = \frac{1-c}{2cz} - \frac{1}{2c} + \frac{\sqrt{(\lambda_+ - z)(z - \lambda_-)}}{2cz}$$

Finally, by the inverse Stieltjes Transform in Theorem 2.1, we have for $x \in [\lambda_-, \lambda_+]$

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\pi} \Im[m_{MP}(x + i\varepsilon)] = \frac{\sqrt{(\lambda_+ - z)(z - \lambda_-)}}{2\pi cx}$$

And for $x = 0$,

$$\lim_{\varepsilon \rightarrow 0} i\varepsilon \Im[m_{MP}(i\varepsilon)] = (1 - c^{-1}) \mathbf{1}_{c > 1}$$

which concludes the proof. \square

2.2.3 Random matrices of mixture models

From the signal processing or machine learning perspectives, one is interested in *inferring* and *identifying* patterns from the data. Naturally, this supposes that data are made by some *correlation* structures instead of pure noise. Therefore, generalizations of the Marčenko-Pastur result to *structured* data models is more relevant for signal processing and machine learning applications. A first generalization to the sample covariance matrix model has been proposed by Silverstein and Bai [SC95] which we recall in the following theorem.

Theorem 2.4 (Sample covariance matrix [SC95]). *Let $\Sigma \in \mathcal{M}_p$ be some nonnegative definite matrix with e.s.d. $\nu \rightarrow \bar{\nu}$ in the weak sense. Let $\mathbf{Z} \in \mathcal{M}_{p,n}$ a random matrix with i.i.d. entries having zero mean and unit variance. Consider the data matrix*

$$\mathbf{X} = \Sigma^{\frac{1}{2}} \mathbf{Z} \in \mathcal{M}_{p,n}$$

Then, as $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$, the e.s.d. μ of $\frac{1}{n} \mathbf{X} \mathbf{X}^\top \in \mathcal{M}_p$ satisfies

$$\mu \xrightarrow{a.s.} \bar{\mu}$$

in the weak convergence sense, and the Stieltjes transform $\bar{m}(z)$ of $\bar{\mu}$ is the unique solution with $\Im \bar{m}(z) > 0$ of the fixed point equation

$$\bar{m}(z) = \left(-z + c \int \frac{t}{1 + t\bar{m}(z)} \bar{\nu}(dt) \right)^{-1}$$

In particular, the limiting measure $\bar{\mu}$ is continuous on \mathbb{R}^+ and real analytic.

Moreover, in practice, one can face data with hierarchical structure and thus can be distributed in different clusters with different covariance profiles, i.e., a *mixture* model. Benaych and Couillet [BGC16] have notably extended Theorem 2.4 under the Gaussian mixture model (see Definition 1) which is more relevant from the machine learning point of view.

Theorem 2.5 (Gram matrices of Gaussian mixture models [BGC16]). *Let $\Sigma_1, \dots, \Sigma_k \in \mathcal{M}_p$ be some nonnegative definite matrices. Let a data matrix $\mathbf{X} = [\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(k)}] \in \mathcal{M}_{p,n}$ with $\mathbf{X}^{(\ell)} = [\mathbf{x}_1^{(\ell)}, \dots, \mathbf{x}_{n_\ell}^{(\ell)}] \in \mathcal{M}_{p,n_\ell}$ and $\mathbf{x}_i^{(\ell)} = \Sigma_\ell^{\frac{1}{2}} \mathbf{z}_i^{(\ell)}$ where $\mathbf{z}_i^{(\ell)}$ is a random Gaussian vector with i.i.d. zero mean and unit variance entries. Let the resolvents*

$$\mathbf{Q}(z) = \left(\frac{1}{p} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1}, \quad \mathbf{R}(z) = \left(\frac{1}{p} \mathbf{X}^\top \mathbf{X} - z \mathbf{I}_n \right)^{-1}$$

Then, as $n_1, \dots, n_k, p \rightarrow \infty$ with $n_\ell/n \rightarrow c_\ell \in (0, 1)$, $p/n \rightarrow c \in (0, \infty)$ and k being bounded, we have

$$\begin{aligned} \mathbf{Q}(z) &\leftrightarrow \bar{\mathbf{Q}}(z) = -\frac{1}{z} \left(\mathbf{I}_p + \sum_{\ell=1}^k \frac{c_\ell}{c} g_\ell(z) \Sigma_\ell \right)^{-1} \\ \mathbf{R}(z) &\leftrightarrow \bar{\mathbf{R}}(z) = \text{diag} \{ g_\ell(z) \mathbf{1}_{n_\ell} \}_{\ell=1}^k \end{aligned}$$

where $g_\ell(z)$ is the unique solution to the fixed point equation

$$g_\ell(z) = -\frac{1}{z} (1 + \bar{g}_\ell(z))^{-1}, \quad \bar{g}_\ell(z) = -\frac{1}{pz} \text{tr} \Sigma_\ell \left(\mathbf{I}_p + \sum_{a=1}^k \frac{c_a}{c} g_a(z) \Sigma_a \right)^{-1}$$

In essence, Theorem 2.5 provides deterministic equivalents for the resolvents $\mathbf{Q}(z)$ and $\mathbf{R}(z)$ which is more fundamental, as discussed in Subsection 2.1.3 than just describing the spectrum of the underlying random matrices. Indeed, such deterministic equivalents can then be used as proxies to infer the performances of the studied methods, an example of study concerns kernel methods as performed in [CBG⁺16]. In the next subsection, we provide two simple examples to illustrate how these deterministic can be exploited to study simple machine learning algorithms. In particular, we first provide a deterministic equivalent for a random matrix model of the form “*Information + Noise*” also referred to as *spiked models* in the random matrix theory community.

2.3 Connections with machine learning through spiked models

In this section, we consider a simple toy example model for the data, namely we consider data distributed in two classes \mathcal{C}_1 and \mathcal{C}_2 with opposite means and isotropic covariances.

Specifically, we consider a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ such that, for $\mathbf{x}_i \in \mathcal{C}_\ell$ for $\ell \in \{1, 2\}$

$$\mathbf{x}_i = (-1)^\ell \boldsymbol{\mu} + \mathbf{z}_i \quad \text{with} \quad \mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p) \quad (2.28)$$

where $\boldsymbol{\mu} \in \mathbb{R}^p$. Without loss of generality, supposing that the data are arranged in classes in the data matrix \mathbf{X} and denoting $\mathbf{y} = [-1, \dots, -1, +1, \dots, +1]^\top \in \mathbb{R}^n$ the vector of labels, in matrix form from the data matrix is given by

$$\mathbf{X} = \boldsymbol{\mu} \mathbf{y}^\top + \mathbf{Z} \quad \text{with} \quad \mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_n] \in \mathcal{M}_{p,n} \quad (2.29)$$

The model in equation 2.29 falls within the class of so-called spiked random matrix models which have been largely studied in random matrix theory [BGN11, BAP⁺05]. In particular, exploiting the result of Hachem *et al.* [HLN⁺07], one can design deterministic equivalents (in the sense of Definition 6) for the corresponding resolvents of the sample covariance and Gram matrices defined respectively as follows, for $z \in \mathbb{C} \setminus \mathbb{R}_+$

$$\mathbf{Q}(z) = \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1}, \quad \mathbf{R}(z) = \left(\frac{1}{n} \mathbf{X}^\top \mathbf{X} - z \mathbf{I}_n \right)^{-1} \quad (2.30)$$

We particularly need the following assumptions subsequently for the design of such deterministic equivalents. These assumptions particularly ensure that the above resolvents are of bounded spectral norm asymptotically when $p/n \rightarrow (0, \infty)$ as $n \rightarrow \infty$.

Assumption 2 (Growth rate). *As $n \rightarrow \infty$,*

1. $p/n \rightarrow c \in (0, \infty)$; 2. $|\mathcal{C}_\ell|/n \rightarrow c_\ell \in (0, 1)$; 3. $\|\boldsymbol{\mu}\| = \mathcal{O}(1)$.

Exploiting the result of [HLN⁺07], deterministic equivalents of the above resolvents are respectively given by

$$\begin{aligned} \mathbf{Q}(z) &\leftrightarrow \bar{\mathbf{Q}}(z) = \left(q^{-1}(z) \mathbf{I}_p - \frac{1}{n} z r(z) \|\mathbf{y}\|^2 \boldsymbol{\mu} \boldsymbol{\mu}^\top \right)^{-1} \\ \mathbf{R}(z) &\leftrightarrow \bar{\mathbf{R}}(z) = \left(r^{-1}(z) \mathbf{I}_n - \frac{1}{n} z q(z) \|\boldsymbol{\mu}\|^2 \mathbf{y} \mathbf{y}^\top \right)^{-1} \end{aligned}$$

where $(q(z), r(z))$ are the unique solution to the fixed point system of equations

$$q(z) = \frac{-1}{z(1 + \frac{1}{n} \text{tr} \bar{\mathbf{R}}(z))}, \quad r(z) = \frac{-1}{z(1 + \frac{1}{n} \text{tr} \bar{\mathbf{Q}}(z))}$$

Since $\|\mathbf{y}\|^2 = n$ and denoting $\bar{\mathbf{y}} = \frac{1}{\sqrt{n}} \mathbf{y}$, we obtain the following result using the Sherman-Morrison identity from Lemma 2.6.

Theorem 2.6 (Deterministic equivalents for $\mathbf{Q}(z)$ and $\mathbf{R}(z)$). *Under Assumption 2, deterministic equivalents for $\mathbf{Q}(z)$ and $\mathbf{R}(z)$ as defined in equation 2.30 are given by*

$$\begin{aligned} \mathbf{Q}(z) &\leftrightarrow \bar{\mathbf{Q}}(z) = q(z) \mathbf{I}_p - \frac{q^2(z)}{1 + (c + \|\boldsymbol{\mu}\|^2) q(z)} \boldsymbol{\mu} \boldsymbol{\mu}^\top \\ \mathbf{R}(z) &\leftrightarrow \bar{\mathbf{R}}(z) = r(z) \mathbf{I}_n - \frac{r^2(z) \|\boldsymbol{\mu}\|^2}{1 + (1 + \|\boldsymbol{\mu}\|^2) r(z)} \bar{\mathbf{y}} \bar{\mathbf{y}}^\top \end{aligned}$$

where

$$q(z) = \frac{1 - c - z + \sqrt{(1 - c - z)^2 - 4zc}}{2zc}, \quad r(z) = \frac{c - 1 - z + \sqrt{(c - 1 - z)^2 - 4z}}{2z}$$

We will provide in the following subsections two simple examples of how these deterministic equivalents can be exploited to infer the performances of ML algorithms. In particular, we first consider an example of unsupervised learning involving the Gram matrix $\frac{1}{n}\mathbf{X}^\top\mathbf{X}$ in Subsection 2.3.1 and then an example of supervised linear ridge regression involving the sample covariance matrix $\frac{1}{n}\mathbf{X}\mathbf{X}^\top$ in Subsection 2.3.2.

2.3.1 Simple example of unsupervised learning

As we discussed in the introduction, kernel spectral clustering is an unsupervised learning algorithm which aims at applying a clustering algorithm on the subspace corresponding to the largest eigenvalues of the kernel matrix. In this subsection, we consider the example of a linear kernel and we will illustrate with this example how Theorem 2.6 can be exploited in order to infer the performances of such algorithm. We therefore consider the kernel matrix defined as

$$\mathbf{K} = \frac{1}{n}\mathbf{X}^\top\mathbf{X} \quad (2.31)$$

From Theorem 2.6, the corresponding deterministic equivalent is

$$\mathbf{R}(z) \leftrightarrow \bar{\mathbf{R}}(z) = r(z)\mathbf{I}_n - \frac{r^2(z)\|\boldsymbol{\mu}\|^2}{1 + (1 + \|\boldsymbol{\mu}\|^2)r(z)}\bar{\mathbf{y}}\bar{\mathbf{y}}^\top$$

which namely shows that the behavior of \mathbf{K} involves a noise term $r(z)\mathbf{I}_n$ and a rank-1 informative term $\propto \bar{\mathbf{y}}\bar{\mathbf{y}}^\top$, in terms of spectrum \mathbf{K} corresponds to a *spiked* Marčenko-Pastur distribution, i.e., the spectrum of \mathbf{K} will contain a *bulk* of eigenvalues with limiting law corresponding to the Marčenko-Pastur law and an isolated (above some phase transition) eigenvalue corresponding to the information about the two classes in data. Moreover, the eigenvector of \mathbf{K} corresponding to the isolated eigenvalue will be aligned to the labels vector $\bar{\mathbf{y}}$ above some phase transition. Let us first determine the isolated eigenvalue in the spectrum of \mathbf{K} .

Isolated eigenvalue. We have the following proposition which provides the asymptotic largest eigenvalue of \mathbf{K} under the high-dimensional regime.

Proposition 2.1 (Isolated spike [BAP⁺05]). *Under Assumption 2, denoting λ_{\max} the largest eigenvalue of \mathbf{K} , then*

$$\lambda_{\max} \xrightarrow{a.s.} \bar{\lambda}_{\max} = \begin{cases} (1 + \|\boldsymbol{\mu}\|^2) \left(1 + \frac{c}{\|\boldsymbol{\mu}\|^2}\right) > (1 + \sqrt{c})^2 & \text{if } \|\boldsymbol{\mu}\|^2 \geq \sqrt{c} \\ (1 + \sqrt{c})^2 & \text{otherwise} \end{cases}$$

Proof. Given the above form of the deterministic equivalent $\bar{\mathbf{R}}(z)$, clearly a spike $\lambda > (1 + \sqrt{c})^2$ appears in the spectrum of \mathbf{K} when $\bar{\mathbf{R}}(z)$ gets singular. Indeed, given the expression of $r(z)$, it is defined for real values in the interval $((1 + \sqrt{c})^2, \infty)$, thus to find the isolated spike λ it suffices to solve the equation $1 + (1 + \|\boldsymbol{\mu}\|^2)r(z) = 0$ in z from which we obtain

$$\lambda = (1 + \|\boldsymbol{\mu}\|^2) \left(1 + \frac{c}{\|\boldsymbol{\mu}\|^2}\right)$$

Furthermore, to obtain the condition on $\|\boldsymbol{\mu}\|^2$ under which such a spike gets out of the right edge $\lambda > (1 + \sqrt{c})^2$, we need to express z evaluated at the spike λ in terms of $m = \|\boldsymbol{\mu}\|^2$ and then study the resulting function. Indeed, $r(z)$ satisfies

$$zr^2(z) + (1 + z - c)r(z) + 1 = 0$$

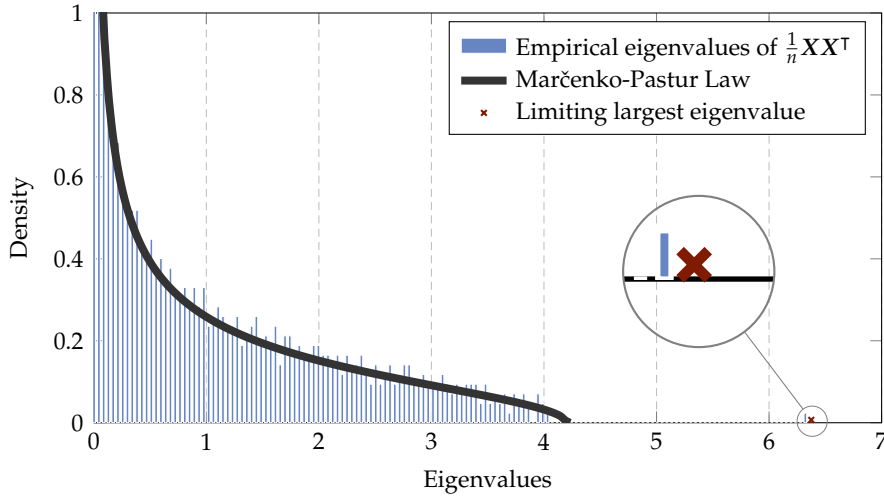


Figure 2.1: Histogram of the eigenvalues of the kernel matrix $K = \frac{1}{n}X^T X$ with X defined in equation 2.29. We considered the dimensions $p = 1100$, $n = 1000$ and $\boldsymbol{\mu} = [2, \mathbf{0}_{p-1}]^T \in \mathbb{R}^p$.

from which we can express the reciprocal function

$$z(r) = \frac{c-1}{1+r} - \frac{1}{(1+r)r}$$

Besides, the spike λ satisfies (where we recall $m = \|\boldsymbol{\mu}\|^2$)

$$r(\lambda) = \frac{-1}{1+m}, \quad 1+r(\lambda) = \frac{m}{1+m}, \quad (1+r(\lambda))r(\lambda) = \frac{-m}{(1+m)^2}$$

Therefore, the spike λ gets isolated from the main bulk at the following value

$$z_\lambda(m) = z(r(\lambda)) = \frac{(c-1)(1+m)}{m} + \frac{(1+m)^2}{m}$$

which defines a function *w.r.t.* m that is first decreasing and then increasing for large values of m . Consequently, the value of m which minimizes the above function corresponds to the phase transition at which the spike starts to get away from the bulk. Therefore, solving $z'_\lambda(m) = \frac{m^2-c}{m^2} = 0$ yields to $m = \sqrt{c}$ since $m \geq 0$. \square

Proposition 2.1 shows that the largest eigenvalue of K converges to a limiting value $\bar{\lambda}_{\max}$ which depend on the *signal strength* only above some condition $\|\boldsymbol{\mu}\|^2 \geq \sqrt{c}$, i.e., only when the classes are theoretically separable (for a large enough value of $\|\boldsymbol{\mu}\|^2$). Figure 2.1 depicts the spectrum of the random kernel matrix K , as theoretically predicted, the spectrum is made of a *bulk* of eigenvalues converging to the Marčenko-Pastur Law along with a largest eigenvalue – *spike* – which corresponds to the informative part of the model through the signal $\boldsymbol{\mu}$.

Remark 2.2 (On the fluctuations of the spike [BAP⁺05]). *It has been shown in [BAP⁺05] that the largest eigenvalue of K concentrates around its converging value, and the underlying converges in law has been determined depending on the condition $\|\boldsymbol{\mu}\|^2 \geq \sqrt{c}$. Indeed, when $\|\boldsymbol{\mu}\|^2 < \sqrt{c}$ it has been shown that $p^{\frac{2}{3}}(1+\sqrt{c})^{-\frac{4}{3}}c^{-\frac{1}{2}}(\lambda_{\max} - (1+\sqrt{c})^2) \xrightarrow{\mathcal{D}} \mathcal{T}$, where \mathcal{T} stands for the Tracy-Widom law. In contrast, when $\|\boldsymbol{\mu}\|^2 \geq \sqrt{c}$ then with $v = \left(\frac{(1+\|\boldsymbol{\mu}\|^2)^2}{c} - \frac{(1+\|\boldsymbol{\mu}\|^2)^2}{\|\boldsymbol{\mu}\|^4}\right) p^{\frac{1}{2}}$ we have $v(\lambda_{\max} - \bar{\lambda}_{\max}) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$.*

We will further show in the following that the corresponding eigenvector to the largest eigenvalue of K gets aligned with the true class labels vector $\bar{\mathbf{y}}$ for sufficiently separable classes, i.e., when $\|\boldsymbol{\mu}\|$ is large enough.

Largest eigenvector of K . Let $\hat{\mathbf{y}}$ be the largest eigenvector of K . In particular, \hat{y}_i stands for the score of the datum x_i to belong to class \mathcal{C}_1 if $\hat{y}_i < 0$ or class \mathcal{C}_2 otherwise. In order to estimate the performance of spectral clustering on the considered example, one is interested in evaluating the alignment $|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2$ between the estimated scores $\hat{\mathbf{y}}$ and the ground truth ones $\bar{\mathbf{y}}$. The deterministic equivalent $\bar{\mathbf{R}}(z)$ can therefore be exploited to evaluate the quantity $|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2$ as discussed in Subsection 2.1.3. Precisely,

$$|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2 = \frac{-1}{2\pi i} \oint_{\Gamma_{\lambda_{\max}}} \bar{\mathbf{y}}^\top \mathbf{R}(z) \bar{\mathbf{y}} dz = \frac{-1}{2\pi i} \oint_{\Gamma_{\bar{\lambda}_{\max}}} \bar{\mathbf{y}}^\top \bar{\mathbf{R}}(z) \bar{\mathbf{y}} dz + o_p(1) \quad (2.32)$$

where Γ_x stands for a sufficiently small positively-oriented complex contour surrounding x and $\bar{\lambda}_{\max} = (1 + \|\boldsymbol{\mu}\|^2) \left(1 + \frac{c}{\|\boldsymbol{\mu}\|^2}\right)$ is the converging value – as $p, n \rightarrow \infty$ – of λ_{\max} the largest eigenvalue of K . Note however that the above equality holds only when $\bar{\lambda}_{\max}$ is isolated, i.e., when it gets outside the bulk that is only valid if $\|\boldsymbol{\mu}\|^2 > \sqrt{c}$, as we saw in Proposition 2.1. Computing the above complex integral yields to the following result.

Proposition 2.2 (Largest eigenvector [Pau07]). *Under Assumption 2, denoting $\hat{\mathbf{y}}$ the eigenvector of K corresponding to the largest eigenvalue, then*

$$|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2 \xrightarrow{a.s.} \frac{1 - c\|\boldsymbol{\mu}\|^{-4}}{1 + \|\boldsymbol{\mu}\|^{-2}} \cdot \mathbf{1}_{\|\boldsymbol{\mu}\|^2 > \sqrt{c}}$$

Proof. The result comes simply through the evaluation of the complex integral

$$\begin{aligned} |\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2 &= \frac{-1}{2\pi i} \oint_{\Gamma_{\bar{\lambda}_{\max}}} \bar{\mathbf{y}}^\top \bar{\mathbf{R}}(z) \bar{\mathbf{y}} dz + o_p(1) \\ &= \frac{-1}{2\pi i} \oint_{\Gamma_{\bar{\lambda}_{\max}}} r(z) - \frac{r^2(z)\|\boldsymbol{\mu}\|^2}{1 + (1 + \|\boldsymbol{\mu}\|^2)r(z)} dz + o_p(1) \\ &= -\text{Res}_{\bar{\lambda}_{\max}} \left(\frac{r(z)(1 + r(z))}{1 + (1 + \|\boldsymbol{\mu}\|^2)r(z)} \right) + o_p(1) \\ &= -\lim_{z \rightarrow \bar{\lambda}_{\max}} \frac{(z - \bar{\lambda}_{\max})r(z)(1 + r(z))}{1 + (1 + \|\boldsymbol{\mu}\|^2)r(z)} + o_p(1) \end{aligned}$$

Using the Hospital rule, the above limit can therefore be evaluated as

$$|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2 = \frac{-r(\bar{\lambda}_{\max})(1 + r(\bar{\lambda}_{\max}))}{(1 + \|\boldsymbol{\mu}\|^2)r'(\bar{\lambda}_{\max})} + o_p(1)$$

It remains thus to evaluate r and r' at the limiting spike $\bar{\lambda}_{\max} = (1 + \|\boldsymbol{\mu}\|^2) \left(1 + \frac{c}{\|\boldsymbol{\mu}\|^2}\right)$. Indeed, we already know that $\bar{\lambda}_{\max}$ satisfies

$$1 + (1 + \|\boldsymbol{\mu}\|^2)r(\bar{\lambda}_{\max}) = 0 \quad \Rightarrow \quad r(\bar{\lambda}_{\max}) = \frac{-1}{1 + \|\boldsymbol{\mu}\|^2}$$

And since $r(z)$ satisfies

$$zr^2(z) + (1 + z - c)r(z) + 1 = 0$$

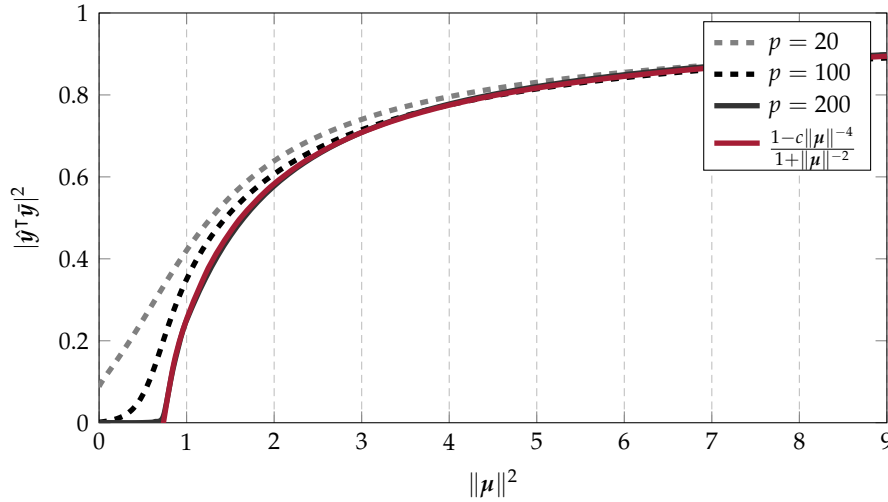


Figure 2.2: Simulated versus limiting values of $|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2$ for $\mathbf{K} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$ with $\mathbf{X} = \mathbf{Z} + \boldsymbol{\mu} \mathbf{y}^\top$. We considered $c = p/n = 1/2$, $\boldsymbol{\mu} = [\alpha, \mathbf{0}_{p-1}]^\top \in \mathbb{R}^p$ varying α .

Taking the derivative w.r.t. z and taking $z = \bar{\lambda}_{\max}$, we end up with

$$r'(\bar{\lambda}_{\max}) = \frac{-r(\bar{\lambda}_{\max})(1 + r(\bar{\lambda}_{\max}))}{2\bar{\lambda}_{\max}r(\bar{\lambda}_{\max}) + (1 + \bar{\lambda}_{\max} - c)}$$

□

Proposition 2.2 provides the asymptotic limit of the alignment $|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2$ between the true score labels $\bar{\mathbf{y}}$ and the estimated ones $\hat{\mathbf{y}}$ as the largest eigenvector of the kernel matrix \mathbf{K} , thereby providing the asymptotic performances¹ of spectral clustering when both p and n grow large. Figure 2.2 depicts simulated versus limiting values of $|\hat{\mathbf{y}}^\top \bar{\mathbf{y}}|^2$ which shows that as p grow large with the ration $c = \lim_p p/n$ being constant, the simulated alignment gets more and more closer to its asymptotic estimation.

Remark 2.3 (On other spiked models). *Similar results as Proposition 2.2 hold for a wide range of spiked models, mainly of one of the following forms:*

$$\begin{aligned} \mathbf{K} &= \frac{1}{n} \mathbf{Z}^\top \mathbf{Z} + \mathbf{P}, & \mathbf{K} &= \frac{1}{n} \mathbf{Z}^\top (\mathbf{I}_p + \mathbf{P}) \mathbf{Z} \\ \mathbf{K} &= \frac{1}{n} (\mathbf{Z} + \mathbf{P})^\top (\mathbf{Z} + \mathbf{P}), & \mathbf{K} &= \frac{1}{n} \mathbf{T} \mathbf{Z}^\top (\mathbf{I}_p + \mathbf{P}) \mathbf{Z} \mathbf{T} \end{aligned}$$

where \mathbf{T} and \mathbf{P} are deterministic matrices with \mathbf{P} of low rank, and $\mathbf{Z} \in \mathcal{M}_{p,n}$ a random matrix having zero mean and unit variance entries.

In the following subsection, we will see how to exploit Theorem 2.6 in order to express the performances of linear regression to separate the two classes defined by the data model in equation 2.29.

¹Let $\hat{\mathcal{C}}_i = \text{sign}(\hat{y}_i)$ be the estimated class \mathcal{C}_i of the datum x_i such that $\bar{y}_i \hat{y}_i > 0$. Then with probability one, the accuracy is given by $\frac{1}{n} \sum_{i=1}^n \delta_{\mathcal{C}_i = \hat{\mathcal{C}}_i} \xrightarrow{a.s.} Q\left(\sqrt{\frac{\zeta}{1-\zeta}}\right)$ with $\zeta = \frac{1-c\|\boldsymbol{\mu}\|^4}{1+\|\boldsymbol{\mu}\|^2}$ and $Q(x) = \frac{1}{2\pi} \int_x^\infty e^{-t^2/2} dt$.

2.3.2 Simple example of supervised learning

We consider in this subsection a linear classifier with an ℓ_2 regularization defined by the following optimization problem

$$\min_w \mathcal{E}(w) = \min_w \frac{1}{n} \|\mathbf{y} - \mathbf{X}^\top w\|^2 + \gamma \|w\|^2 \quad (2.33)$$

where w stands for the weights vector and $\gamma \in \mathbb{R}^+$ is the regularization parameter. The data matrix \mathbf{X} is defined in equation 2.29 and $\mathbf{y} \in \mathbb{R}^n$ stands for the vector of labels such that $y_i = (-1)^\ell$ if $x_i \in \mathcal{C}_\ell$ for $\ell \in \{1, 2\}$. The solution to the optimization problem in equation 2.33 is explicitly given by

$$w = \frac{1}{n} \mathbf{Q}(-\gamma) \mathbf{X} \mathbf{y} \quad \text{with} \quad \mathbf{Q}(z) = \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - z \mathbf{I}_p \right)^{-1} \quad (2.34)$$

where the resolvent $\mathbf{Q}(z)$ appears at the core of the classifier. In particular, one is interested on the behavior of the (hard) decision function for a new test datum $x \in \mathcal{C}_\ell$, for $\ell \in \{1, 2\}$, which is defined as

$$g(x) = x^\top w = \frac{1}{n} x^\top \mathbf{Q}(-\gamma) \mathbf{X} \mathbf{y} \stackrel{\mathcal{C}_1}{\leq} 0 \stackrel{\mathcal{C}_2}{\geq} 0 \quad (2.35)$$

Given the above form of the decision function, it is a simple sum of independent (not necessarily identically distributed) random variables. Therefore thanks to the Lyapunov's central limit theorem [Bil08], which we recall in the following theorem, $g(x)$ has a Gaussian approximation in the large limit when both $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$.

Theorem 2.7 (Lyapunov's CLT [Bil08]). *Let X_1, \dots, X_n be independent random variables of means $\mathbb{E}X_i = \mu_i$ such that $|X_i|$ have moments of order $2 + \varepsilon$ for some $\varepsilon > 0$ and let the Lyapunov condition $\lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\varepsilon}} \sum_{i=1}^n \mathbb{E}[|X_i - \mu_i|^{2+\varepsilon}] = 0$ hold, with $s_n = \text{Var}^{\frac{1}{2}}[\sum_{i=1}^n (X_i - \mu_i)]$. Then,*

$$\frac{1}{s_n} \sum_{i=1}^n (X_i - \mu_i) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$$

It suffices thus to estimate the first and second moments of the decision function $g(x)$ to fully describe its behavior and therefore infer the performances of the linear classifier. We start by estimating $m_\ell = \mathbb{E}g(x)$ for $x \in \mathcal{C}_\ell$ independent of the data matrix \mathbf{X} . In the following we will write $\mathbf{Q} = \mathbf{Q}(-\gamma)$ and $\bar{\mathbf{Q}} = \bar{\mathbf{Q}}(-\gamma)$ for simplicity.

$$m_\ell = \mathbb{E}g(x) = \mathbb{E} \left[\frac{1}{n} x^\top \mathbf{Q} \mathbf{X} \mathbf{y} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[y_i x^\top \mathbf{Q} x_i]$$

The above expectation requires to manage the statistical dependencies between \mathbf{Q} and x_i , which can be handled by Example 1, specifically through the identity

$$\mathbf{Q} x_i = \frac{\mathbf{Q}_{-i} x_i}{1 + \frac{1}{n} x_i^\top \mathbf{Q}_{-i} x_i} \quad (2.36)$$

Moreover, relying on the trace Lemma 2.8, the quadratic form $\frac{1}{n} x_i^\top \mathbf{Q}_{-i} x_i$ converges² to the quantity $\frac{1}{n} \text{tr}(\mathbb{E}[x_i x_i^\top] \bar{\mathbf{Q}})$ where $\bar{\mathbf{Q}}$ is the deterministic equivalent of \mathbf{Q} from Theorem 2.6.

²The resolvent \mathbf{Q}_{-i} has a bounded spectral norm, this will be highlighted in the next section.

Specifically, we have by the perturbation Lemma 2.7

$$\frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i - cq(-\gamma) \xrightarrow{a.s.} 0 \quad (2.37)$$

Consequently, the expectation of the decision function $g(\mathbf{x})$ can be estimated as

$$m_\ell - \frac{(-1)^\ell \boldsymbol{\mu}^\top \bar{\mathbf{Q}}(-\gamma) \boldsymbol{\mu}}{1 + cq(-\gamma)} \xrightarrow{a.s.} 0 \quad \Rightarrow \quad m_\ell - \frac{(-1)^\ell \|\boldsymbol{\mu}\|^2 q(-\gamma)}{1 + (c + \|\boldsymbol{\mu}\|^2) q(-\gamma)} \xrightarrow{a.s.} 0$$

Let us now turn into the estimation of the variance of $g(\mathbf{x})$. Let $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}_p + \boldsymbol{\mu}\boldsymbol{\mu}^\top$, we have

$$\begin{aligned} \mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] &= \frac{1}{n^2} \mathbb{E}[\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q} \mathbf{x} \mathbf{x}^\top \mathbf{Q} \mathbf{X} \mathbf{y}] = \frac{1}{n^2} \mathbb{E}[\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q} \boldsymbol{\Sigma} \mathbf{Q} \mathbf{X} \mathbf{y}] \\ &= \frac{1}{n^2} \sum_{i=1}^n y_i^2 \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \boldsymbol{\Sigma} \mathbf{Q} \mathbf{x}_i] + \frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \boldsymbol{\Sigma} \mathbf{Q} \mathbf{x}_j] \end{aligned}$$

Using again the identity in equation 2.36 and the above estimate of the quadratic form $\frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i$, we end up with

$$\mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] - \frac{1}{n} \frac{\text{tr}(\boldsymbol{\Sigma} \mathbb{E}[\mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-i}])}{(1 + cq(-\gamma))^2} - \frac{1}{n^2} \sum_{i \neq j} y_i y_j \frac{\mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-j} \mathbf{x}_j]}{(1 + cq(-\gamma))^2} \xrightarrow{a.s.} 0$$

where the term $\mathbb{E}[\mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-i}]$ is handled by the following identities

$$\begin{aligned} \eta(\mathbf{A}) &= \frac{1}{n} \text{tr}(\boldsymbol{\Sigma} \mathbb{E}[\mathbf{Q} \mathbf{A} \mathbf{Q}]) = \frac{(1 + cq(-\gamma)) \frac{1}{n} \text{tr}(\boldsymbol{\Sigma} \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}})}{((1 + cq(-\gamma)))^2 - \frac{1}{n} \text{tr}(\boldsymbol{\Sigma} \bar{\mathbf{Q}} \boldsymbol{\Sigma} \bar{\mathbf{Q}})} \\ \Delta(\mathbf{A}) &= \mathbb{E}[\mathbf{Q} \mathbf{A} \mathbf{Q}] = \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}} + \frac{\eta(\mathbf{A})}{1 + cq(-\gamma)} \bar{\mathbf{Q}} \boldsymbol{\Sigma} \bar{\mathbf{Q}} \end{aligned}$$

Which yield to

$$\frac{1}{n} \text{tr}(\boldsymbol{\Sigma} \mathbb{E}[\mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-i}]) - cq^2(-\gamma) \xrightarrow{a.s.} 0$$

And the remaining term $\mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-j} \mathbf{x}_j]$ develops thanks to Example 1, through the identity

$$\mathbf{Q} = \mathbf{Q}_{-i} - \frac{\mathbf{Q}_{-i} \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i}}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}$$

Indeed, we have for $i \neq j$

$$\begin{aligned} \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-i} \boldsymbol{\Sigma} \mathbf{Q}_{-j} \mathbf{x}_j] &= \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-ij} \boldsymbol{\Sigma} \mathbf{Q}_{-ji} \mathbf{x}_j] - \frac{1}{n} \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \boldsymbol{\Sigma} \mathbf{Q}_{-ji} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-ji} \mathbf{x}_j}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-ji} \mathbf{x}_i} \right] \\ &\quad - \frac{1}{n} \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_j \mathbf{x}_j^\top \mathbf{Q}_{-ij} \boldsymbol{\Sigma} \mathbf{Q}_{-ji} \mathbf{x}_j}{1 + \frac{1}{n} \mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{x}_j} \right] + \frac{1}{n^2} \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_j \mathbf{x}_j^\top \mathbf{Q}_{-ij} \boldsymbol{\Sigma} \mathbf{Q}_{-ji} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_j}{(1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_i)(1 + \frac{1}{n} \mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{x}_j)} \right] \end{aligned}$$

And recalling the convergence of the quadratic form in equation 4.16 and putting all the pieces together we end up with the following estimation of $\mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2]$

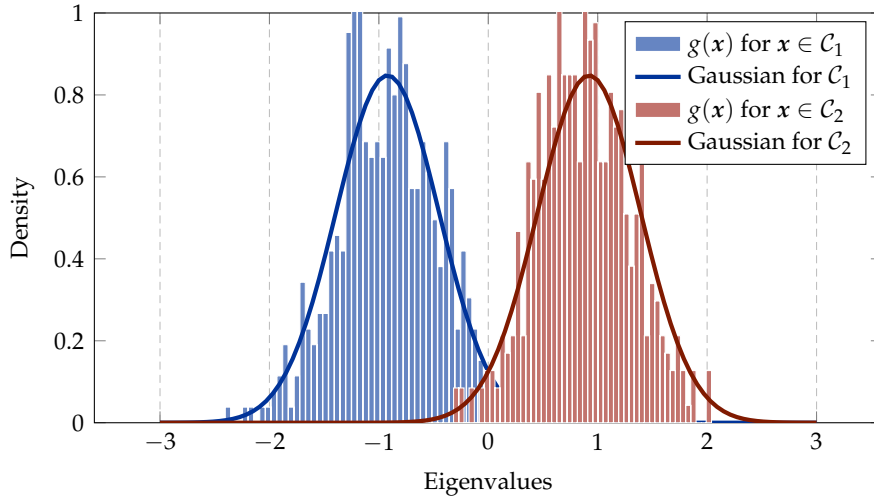


Figure 2.3: Histogram of the decision function $g(\mathbf{x})$ for new test data \mathbf{x} independent from the training set \mathbf{X} , versus its Gaussian approximation as per Theorem 2.8. We consider the parameters $p = 1100$, $n_1 = n_2 = 500$, $\gamma = 1 \cdot 10^{-1}$ and $\boldsymbol{\mu} = [3, \mathbf{0}_{p-1}]^\top$.

$$\mathbb{E} [(\mathbf{x}^\top \mathbf{w})^2] - \frac{cq^2(-\gamma)}{(1+cq(-\gamma))^2} - \frac{\boldsymbol{\mu}^\top \Delta(\boldsymbol{\Sigma}) \boldsymbol{\mu}}{(1+cq(-\gamma))^2} + \frac{2cq^2(-\gamma) \boldsymbol{\mu}^\top \bar{\mathbf{Q}}(-\gamma) \boldsymbol{\mu}}{(1+cq(-\gamma))^3} \xrightarrow{a.s.} 0$$

We therefore have the following result which provides a Gaussian approximation of the decision function $g(\mathbf{x})$.

Theorem 2.8 (Gaussian Approximation of $g(\mathbf{x})$). *Under Assumption 2, for $\mathbf{x} \in \mathcal{C}_\ell$ for $\ell \in \{1, 2\}$, we have*

$$(\nu - m_\ell^2)^{-\frac{1}{2}} (g(\mathbf{x}) - m_\ell) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$$

where

$$m_\ell = \frac{(-1)^\ell \|\boldsymbol{\mu}\|^2 q(-\gamma)}{1 + (c + \|\boldsymbol{\mu}\|^2) q(-\gamma)}$$

$$\nu = \frac{cq^2(-\gamma)}{(1+cq(-\gamma))^2} + \frac{\boldsymbol{\mu}^\top \Delta(\boldsymbol{\Sigma}) \boldsymbol{\mu}}{(1+cq(-\gamma))^2} - \frac{2cq^2(-\gamma) \boldsymbol{\mu}^\top \bar{\mathbf{Q}}(-\gamma) \boldsymbol{\mu}}{(1+cq(-\gamma))^3}$$

Theorem 2.8 states that the considered linear classifier is asymptotically equivalent to the thresholding of two monivariate Gaussian random variables. Figure 2.3 depicts simulated values of $g(\mathbf{x})$ along with its asymptotic theoretical Gaussian prediction where a perfect matching is observed. In particular, the generalization performance of the classifier can be expressed thanks to the Gaussian tail function $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-u^2/2} du$ as per the following Corollary to Theorem 2.8.

Corollary 2.1 (Generalization Performance of the linear classifier). *Under the setting and Assumptions of Theorem 2.8, for $\ell \in \{1, 2\}$, with probability one*

$$\mathbb{P} \left\{ (-1)^\ell g(\mathbf{x}) < 0 \mid \mathbf{x} \in \mathcal{C}_\ell \right\} - Q \left(\frac{m_\ell}{\sqrt{\nu - m_\ell^2}} \right) \rightarrow 0$$

The same analysis as the one presented in this subsection can be adopted to analyse a wide range of standard classifiers such as LS-SVM [LC17] and even some more sophisticated methods such as the logistic regression algorithm [MLC19]. The core of this analysis leverages the design of a deterministic equivalent for the resolvent \mathbf{Q} . In the next section, we will recall some recent results from Louart and Couillet [LC18b] which provide a systematic approach to design such deterministic equivalents leveraging on concentration assumptions on data, and which particularly generalizes the Gaussian assumption of the result of Benaych and Couillet [BGC16] recalled in Theorem 2.5. As discussed in the introduction and will be presented in more details in the next chapter, the concentration assumption is of particular interest in practice since it provides a realistic modeling for real data, in particular, for GAN data which satisfy the concentration assumption by construction.

2.4 Extensions with concentration of measure theory

The previous results are extensible to a richer class of random vectors, namely to the class of random concentrated vectors [Led05a]. As we discussed in the introduction, this class for random vectors are more appropriate for real data modelling since GANs data fall constructively within this class. In this section, we will briefly recall some essential concentration notions and properties, then we will provide their application to design a deterministic equivalent for the sample covariance matrix following the same approach developed in [LC18b].

2.4.1 The notion of concentrated vectors

Being the central tool of this manuscript, we introduce the notion of random concentrated vectors. Note that several and more advanced concentration notions have been recently developed in [LC18b] in order to specifically analyze the behavior of large sample covariance matrices, but for simplicity, we restrict ourselves in this manuscript to the sufficient so-called q -exponentially concentrated random vectors.

Definition 7 (Concentrated vector). *Given a set of indices \mathbb{S} , a sequence of normed vector spaces $(E_s, \|\cdot\|_s)_{s \in \mathbb{S}}$, a sequence of random vectors $\mathbf{Z}_s \in E_s$, a sequence of positive numbers σ_s , we say that \mathbf{Z}_s is q -exponentially concentrated with an observable diameter of order $\mathcal{O}(\sigma_s)$ if there exists two constants $C, c > 0$ such that for all sequence of 1-Lipschitz mappings $f_s : E_s \rightarrow \mathbb{R}$:*

$$\forall s \in \mathbb{S}, \forall t > 0 : \mathbb{P} \{ |f_s(\mathbf{Z}_s) - \mathbb{E}[f_s(\mathbf{Z}_s)]| \geq t \} \leq C e^{-c(t/\sigma_s)^q} \quad (2.38)$$

We note then $\mathbf{Z}_s \propto \mathcal{E}_q(\sigma_s)$; when $\sigma_s = \mathcal{O}(1)$, we write $\mathbf{Z}_s \propto \mathcal{E}_q$.

Therefore, concentrated vectors are defined through the concentration of any 1-Lipschitz real scalar “observation”. One of the most important examples of concentrated vectors are standard Gaussian vectors. Precisely, we have the following proposition. Note also that more examples such as uniform and Gamma distribution result in random vectors which satisfy the concentration property as per Definition 7, see [Led05a].

Proposition 2.3 (Gaussian vectors [Led05a]). *Let $d \in \mathbb{N}$ and $\mathbf{Z}_d \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$. Then \mathbf{Z}_d is a 2-exponentially concentrated vector independently on the dimension d , i.e. $\mathbf{Z}_d \propto \mathcal{E}_2$.*

But the richness of concentrated random vectors lies in their fundamental stability property through Lipschitz operations, which naturally generates wide families of concentrated random vectors.

Proposition 2.4 (Stability through Lipschitz transformations). *It is easily deduced from Definition 7 that given a sequence of positive numbers $L_s > 0$ and a sequence of L_s -Lipschitz transformations $\phi_s : (E_s, \|\cdot\|_s) \rightarrow (F_s, \|\cdot\|'_s)$,*

$$\text{if } \mathbf{Z}_s \propto \mathcal{E}_q(\sigma_s), \text{ then } \phi_s(\mathbf{Z}_s) \propto \mathcal{E}_q(L_s \sigma_s) \quad (2.39)$$

Proof. This comes simply by, for all $f_s : F_s \rightarrow \mathbb{R}$, 1-Lipschitz, $\frac{1}{L_s} f_s \circ \phi_s$ is 1-Lipschitz, and one can employ inequality 2.38 to $\frac{t}{L_s}$ which provides the result. \square

As we saw previously, we need to establish the concentration of quadratic forms of the form $\frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ with \mathbf{Q} deterministic or independent from \mathbf{x} . To this end, we will need the following property.

Proposition 2.5. *Let $X_n \in \mathbb{R}$ be a random variable depending on some integer $n \in \mathbb{N}$, satisfying $X_n \propto \mathcal{E}_q(\sigma_n)$ with $\sigma_n \rightarrow 0$ as $n \rightarrow \infty$ and $\limsup_n |\mathbb{E}X_n| < \infty$. The square of X_n remains concentrated and we specifically have*

$$X_n^2 - \mathbb{E}X_n^2 \propto \mathcal{E}_q(\sigma_n) + \mathcal{E}_{\frac{q}{2}}(\sigma_n^2)$$

Proof. From the algebraic identity $X_n^2 - (\mathbb{E}X_n)^2 = (X_n - \mathbb{E}X_n)^2 + 2(X_n - \mathbb{E}X_n)\mathbb{E}X_n$, one has

$$\mathbb{P} \left\{ |X_n^2 - (\mathbb{E}X_n)^2| \geq t \right\} \leq \mathbb{P} \left\{ |X_n - \mathbb{E}X_n| \geq \sqrt{\frac{t}{2}} \right\} + \mathbb{P} \left\{ |X_n - \mathbb{E}X_n| \geq \frac{t}{4|\mathbb{E}X_n|} \right\}$$

Applying the identities (where X'_n is an independent copy of X_n)

$$\mathbb{P} \left\{ |X_n^2 - m_{X_n^2}| \geq t \right\} \leq 2\mathbb{P} \left\{ |X_n^2 - (X'_n)^2| \geq t \right\} \leq 4\mathbb{P} \left\{ |X_n^2 - (\mathbb{E}X_n)^2| \geq t/2 \right\}$$

where $m_{X_n^2}$ is a median of X_n^2 , and the final result comes from the fact that the expectation and median of X_n^2 are asymptotically close to each other as $n \rightarrow \infty$. Indeed, we have

$$|\mathbb{E}X_n - m_{X_n}| \leq \mathbb{E}|X_n - m_{X_n}| = \mathcal{O} \left(\int_0^\infty e^{-c(t/\sigma_n)^q} dt \right) = \mathcal{O}(\sigma_n) \rightarrow 0$$

The same reasoning holds also for X_n^2 which concludes the proof. \square

The following Lemma exploits Proposition 2.5 to provide the concentration of the quadratic form $\frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x}$ for some concentrated vector $\mathbf{x} \in \mathbb{R}^p$ and bounded matrix $\mathbf{Q} \in \mathcal{M}_p$ (in spectral norm).

Lemma 2.9 (Trace Lemma with concentration). *Let $\mathbf{x} \in \mathbb{R}^p$ be a random vector such that $\mathbf{x} \propto \mathcal{E}_q$ and $\mathbf{Q} \in \mathcal{M}_p$ some deterministic (or random independent of \mathbf{x}) matrix with bounded spectral norm. Then,*

$$\frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \frac{1}{p} \text{tr}(\mathbb{E}[\mathbf{x}\mathbf{x}^\top] \mathbf{Q}) \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right) + \mathcal{E}_{\frac{q}{2}} \left(\frac{1}{p} \right)$$

In particular, there exists a constant $C > 0$ such that

$$\mathbb{E} \left[\left| \frac{1}{p} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \frac{1}{p} \text{tr}(\mathbb{E}[\mathbf{x}\mathbf{x}^\top] \mathbf{Q}) \right|^k \right] \leq C p^{-\frac{k}{2}}$$

Proof. The proof comes simply by rewriting $\frac{1}{p}\mathbf{x}^\top \mathbf{Q}\mathbf{x} = \|\frac{1}{\sqrt{p}}\mathbf{Q}^{\frac{1}{2}}\mathbf{x}\|^2$. Since, $\mathbf{x} \rightarrow \|\frac{1}{\sqrt{p}}\mathbf{Q}^{\frac{1}{2}}\mathbf{x}\|$ is $\mathcal{O}\left(\frac{1}{\sqrt{p}}\right)$ -Lipschitz we have $\|\frac{1}{\sqrt{p}}\mathbf{Q}^{\frac{1}{2}}\mathbf{x}\| \propto \mathcal{E}_q\left(\frac{1}{\sqrt{p}}\right)$, and the final concentration comes by Proposition 2.5. Moreover, the control of the k -th moment follows from the identity

$$\mathbb{E}|X|^k = \int_0^\infty k t^{k-1} \mathbb{P}\{|X| \geq t\} dt$$

for some random variable X . □

We will see in the next subsection how this result is exploited to design a deterministic equivalent for the sample covariance matrix.

2.4.2 Resolvent of the sample covariance matrix

In this subsection, we will consider a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ such that $\mathbf{X} \propto \mathcal{E}_q$ and we further suppose that the columns of \mathbf{X} are independent and have the same second order statistic $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \boldsymbol{\Sigma} \in \mathcal{M}_p$ with bounded spectral norm. As we saw previously, the analysis of the behavior of the sample covariance matrix $\hat{\boldsymbol{\Sigma}} = \frac{1}{n}\mathbf{X}\mathbf{X}^\top$ is related to the analysis of its resolvent, defined for $z \in \mathbb{C} \setminus \mathbb{R}_-$ ³ as

$$\mathbf{Q}(z) = (\hat{\boldsymbol{\Sigma}} + z\mathbf{I}_p)^{-1} = \left(\frac{1}{n}\mathbf{X}\mathbf{X}^\top + z\mathbf{I}_p\right)^{-1} \quad (2.40)$$

The resolvent matrix $\mathbf{Q}(z)$ notably satisfies several bounds which are of particular interest in the analysis of its behavior as we will see subsequently. Indeed, we recall the following Lemma from [LC18b].

Lemma 2.10 (The resolvent bounds). *The resolvent $\mathbf{Q}(z)$ satisfies the following bounds. For $z \in \mathbb{R}_+$, we have*

$$\mathbf{1.} \|\mathbf{Q}(z)\| \leq \frac{1}{z} \quad \mathbf{2.} \|\mathbf{Q}(z)\hat{\boldsymbol{\Sigma}}\| \leq 1 \quad \mathbf{3.} \|\mathbf{Q}(z)\mathbf{X}\| \leq \sqrt{\frac{n}{z}}$$

Proof. The upper bound for **1.** comes from the smallest eigenvalue of $\hat{\boldsymbol{\Sigma}} + z\mathbf{I}_p$ being larger than z . **2.** follows from the identity $\mathbf{Q}(z)\hat{\boldsymbol{\Sigma}} + z\mathbf{Q}(z) = \mathbf{I}_p$ and since $\mathbf{Q}(z)$ is symmetric positive definite. Finally, **3.** follows from combining **1.** and **2.**, giving the bound

$$\|\mathbf{Q}(z)\frac{1}{n}\mathbf{X}\mathbf{X}^\top\mathbf{Q}(z)\| \leq \frac{1}{z}$$

□

Now we provide the intuition behind the design of a deterministic equivalent for $\mathbf{Q}(z)$, we refer the reader to [LC18b] for a detailed proof. From a low-dimensional perspective (i.e., p fixed as $n \rightarrow \infty$), one would think that $\mathbf{Q}(z)$ would be close to $(\boldsymbol{\Sigma} + z\mathbf{I}_p)^{-1}$. This happens to be completely wrong in the random matrix theory regime, when both $p, n \rightarrow \infty$ with $p/n \rightarrow c \in (0, \infty)$. However, one can establish that $\mathbf{Q}(z)$ is equivalent to some matrix of the form $\bar{\mathbf{Q}}(z) = (\tilde{\boldsymbol{\Sigma}} + z\mathbf{I}_p)^{-1}$, in the sense of Definition 6,

³For convenience, we shall consider the notation of the resolvent with “+ z ” instead of “− z ”.

for some well-chosen deterministic matrix $\tilde{\Sigma} \in \mathcal{M}_p$. In order to find the explicit expression for $\tilde{\Sigma}$, we start by expressing the difference between $\tilde{Q}(z)$ and $\mathbb{E}[Q(z)]$, thanks to the resolvent identity in Lemma 2.1, as follows

$$\tilde{Q}(z) - \mathbb{E}[Q(z)] = \mathbb{E} \left[Q(z) \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - \tilde{\Sigma} \right) \tilde{Q}(z) \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} [Q(z) (x_i x_i^\top - \tilde{\Sigma}) \tilde{Q}(z)]$$

The next step now consists in handling the statistical dependencies between $Q(z)$ and x_i , this is made possible, as we saw in the previous Section, thanks to the Schur identities from Example 1. In particular, first using the identity

$$Q(z)x_i = \frac{Q_{-i}(z)x_i}{1 + \frac{1}{n}x_i^\top Q_{-i}(z)x_i}$$

We obtain

$$\begin{aligned} \tilde{Q}(z) - \mathbb{E}[Q(z)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(Q(z)x_i x_i^\top - Q(z)\tilde{\Sigma}) \tilde{Q}(z)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\left(\frac{Q_{-i}(z)x_i x_i^\top}{1 + \frac{1}{n}x_i^\top Q_{-i}(z)x_i} - Q(z)\tilde{\Sigma} \right) \tilde{Q}(z) \right] \end{aligned}$$

And thanks to the identity

$$Q(z) = Q_{-i}(z) - \frac{1}{n} \frac{Q_{-i}(z)x_i x_i^\top Q_{-i}(z)}{1 + \frac{1}{n}x_i^\top Q_{-i}(z)x_i}$$

We end up having

$$\begin{aligned} \tilde{Q}(z) - \mathbb{E}[Q(z)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[Q_{-i}(z) \left(\frac{x_i x_i^\top}{1 + \frac{1}{n}x_i^\top Q_{-i}(z)x_i} - \tilde{\Sigma} \right) \tilde{Q}(z) \right] \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [Q_{-i}(z)x_i x_i^\top Q_{-i}(z)\tilde{\Sigma}\tilde{Q}(z)] \end{aligned}$$

where the second line term happens to be of vanishing spectral norm due to the additional factor $1/n$ and since the matrix $\frac{1}{n} \sum_{i=1}^n \mathbb{E} [Q_{-i}(z)x_i x_i^\top Q_{-i}(z)\tilde{\Sigma}\tilde{Q}(z)]$ has a bounded spectral norm thanks to Lemma 2.10. Moreover, again by Lemma 2.10, since $Q_{-i}(z)$ has a bounded spectral norm and is independent from x_i , the quadratic form $\frac{1}{n}x_i^\top Q_{-i}(z)x_i$ converges to its expectation $\frac{1}{n} \text{tr}(\tilde{\Sigma}\mathbb{E}[Q_{-i}(z)])$ thanks to Lemma 2.9. Quite naturally, one would expect that the deterministic matrix $\tilde{\Sigma}$ should be

$$\tilde{\Sigma} = \frac{\Sigma}{1 + \frac{1}{n} \text{tr}(\Sigma\mathbb{E}[Q_{-i}(z)])}$$

so that the difference $\tilde{Q}(z) - \mathbb{E}[Q(z)]$ would be of vanishing spectral norm. Consequently, the term $\frac{1}{n} \text{tr}(\Sigma\mathbb{E}[Q_{-i}(z)])$ would naturally be close to $\frac{1}{n} \text{tr}(\Sigma\tilde{Q}(z))$, therefore defining the deterministic equivalent $\tilde{Q}(z)$ of the resolvent $Q(z)$ through an implicit fixed point equation rather than computing the expectation of $Q_{-i}(z)$ which is hard to compute explicitly in general cases. The deterministic equivalent $\tilde{Q}(z)$ is therefore given by

$$Q(z) \leftrightarrow \tilde{Q}(z) = \left(\frac{\Sigma}{1 + \delta(z)} + zI_p \right)^{-1}$$

where $\delta(z)$ is the unique⁴ solution to the fixed point equation

$$\delta(z) = \frac{1}{n} \operatorname{tr} \left(\boldsymbol{\Sigma} \left(\frac{\boldsymbol{\Sigma}}{1 + \delta(z)} + z \mathbf{I}_p \right)^{-1} \right)$$

Looking carefully to the above expression of $\bar{\mathbf{Q}}(z)$, we see that there is an additional term $\delta(z)$ comparing to what one would expect in a low-dimensional regime, i.e., the deterministic equivalent being $(\boldsymbol{\Sigma} + z \mathbf{I}_p)^{-1}$. Indeed, when p is fixed as $n \rightarrow \infty$, we particularly have $\delta(z) \approx \frac{1}{n} \operatorname{tr} (\boldsymbol{\Sigma} \mathbb{E}[\mathbf{Q}_{-i}(z)]) \rightarrow 0$, which recovers the low-dimensional case.

Remark 2.4 (On the concentration of $\frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i}(z) \mathbf{x}_i$). Note that Lemma 2.9 is not good enough to establish the concentration of $\frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i}(z) \mathbf{x}_i$. Indeed, $\mathbf{Q}_{-i}(z)$ being random demands to bound a probability involving \mathbf{x}_i . Specifically, one should decompose the calculus as

$$\begin{aligned} & \mathbb{P} \left\{ \left| \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i}(z) \mathbf{x}_i - \frac{1}{n} \operatorname{tr} (\boldsymbol{\Sigma} \mathbb{E}[\mathbf{Q}_{-i}(z)]) \right| \geq t \right\} \\ & \leq \mathbb{E} \left[\mathbb{P} \left\{ \left| \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i}(z) \mathbf{x}_i - \frac{1}{n} \operatorname{tr} (\boldsymbol{\Sigma} \mathbf{Q}_{-i}(z)) \right| \geq \frac{t}{2} \mid \mathbf{X}_{-i} \right\} \right] \\ & + \mathbb{P} \left\{ \left| \frac{1}{n} \operatorname{tr} \boldsymbol{\Sigma} (\mathbf{Q}_{-i}(z) - \mathbb{E}[\mathbf{Q}_{-i}(z)]) \right| \geq \frac{t}{2} \right\} \end{aligned}$$

where $\mathbf{X}_{-i} \in \mathcal{M}_{p, n-1}$ is the data matrix \mathbf{X} prevented from its i -th column \mathbf{x}_i . Therefore, the concentration of $\frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i}(z) \mathbf{x}_i$ follows after handling the above two terms (see [LC18b, Proposition 3.4]).

Under the more general case of a mixture of k -class concentrated vectors (see Definition 2), a deterministic equivalent of the sample covariance matrix is given by the following result.

Theorem 2.9 (Deterministic equivalent of the sample covariance matrix [LC18b]). Let a data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p, n}$ be satisfying Definition 2 and let Assumption 1 hold with $\|\boldsymbol{\mu}_\ell\| = \mathcal{O}(\sqrt{p})$ for all $\ell \in [k]$. Then the resolvent of the sample covariance matrix defined as $\mathbf{Q}(z) = (\frac{1}{n} \mathbf{X} \mathbf{X}^\top + z \mathbf{I}_p)^{-1}$ admits a deterministic equivalent $\bar{\mathbf{Q}}(z)$ given by

$$\mathbf{Q}(z) \leftrightarrow \bar{\mathbf{Q}}(z) = \left(\sum_{\ell=1}^k \frac{c_\ell (\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top)}{1 + \delta_\ell(z)} + z \mathbf{I}_p \right)^{-1}$$

where $(\delta_1(z), \dots, \delta_k(z))$ are the unique solution to system of fixed point equations defined for each $j \in [k]$ as

$$\delta_j(z) = \frac{1}{n} \left((\boldsymbol{\Sigma}_j + \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top) \left(\sum_{\ell=1}^k \frac{c_\ell (\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top)}{1 + \delta_\ell(z)} + z \mathbf{I}_p \right)^{-1} \right)$$

Theorem 2.9 generalizes the result of Benaych and Couillet [BGC16] (recalled in Theorem 2.5) to the more general class of concentrated vectors. However, the above expression of the deterministic equivalent involve strictly the class-wise means and covariances of the data namely the quantities $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$, thereby demonstrating the *universality* aspect of the behavior of the sample covariance matrix *w.r.t.* data distributions

⁴See in [LC18b, Proposition 3.8] for more details about the existence and uniqueness of the fixed point equation solution.

satisfying the concentration assumption in Definition 2. We will see in the next chapters, that this universality aspect goes far beyond the sample covariance matrix to more sophisticated non-linear ML methods such as kernel methods and the Softmax layer in neural networks.

Chapter 3

Universality of Large Random Matrices

This chapter is based on the following works:

- (C1) MEA. Seddik, C. Louart, M. Tamaazousti, R. Couillet, “*Random Matrix Theory Proves that Deep Learning Representations of GAN-data Behave as Gaussian Mixtures*”, International Conference on Machine Learning (ICML’20), Online, 2020.
- (C1’) MEA. Seddik, M. Tamaazousti, R. Couillet, “*Pourquoi les matrices aléatoires expliquent l’apprentissage ? Un argument d’universalité offert par les GANs*”, Colloque francophone de traitement du signal et des images (Gretsi’19), Lille, France, 2019.

Contents

3.1	GAN Data are Concentrated Data	55
3.2	Random Gram Matrices of Concentrated Data	61
3.2.1	Motivation	61
3.2.2	Model and Main Results	62
3.2.2.1	Mixture of Concentrated Vectors	62
3.2.2.2	Behavior of the Gram matrix of concentrated vectors	63
3.2.2.3	Application to GAN-generated Images	65
3.2.3	Central Contribution	68

This Chapter contains two main parts. The first part will highlight the importance of modeling real data as concentrated vectors, the second part will consist in studying the behavior of large Gram matrices with concentrated inputs.

3.1 GAN Data are Concentrated Data

Concentrated random vectors are particularly interesting from a practical standpoint for real data modeling. In fact, unlike simple Gaussian vectors, the former do not suffer from the constraint of having independent entries which is quite a restrictive assumption when modeling real data such as images or their non-linear features (*e.g.*, DL representations). The other modeling interest of concentrated vectors lies in their being already present in practice as alternatives to real data. Indeed, adversarial neural networks (GANs) have

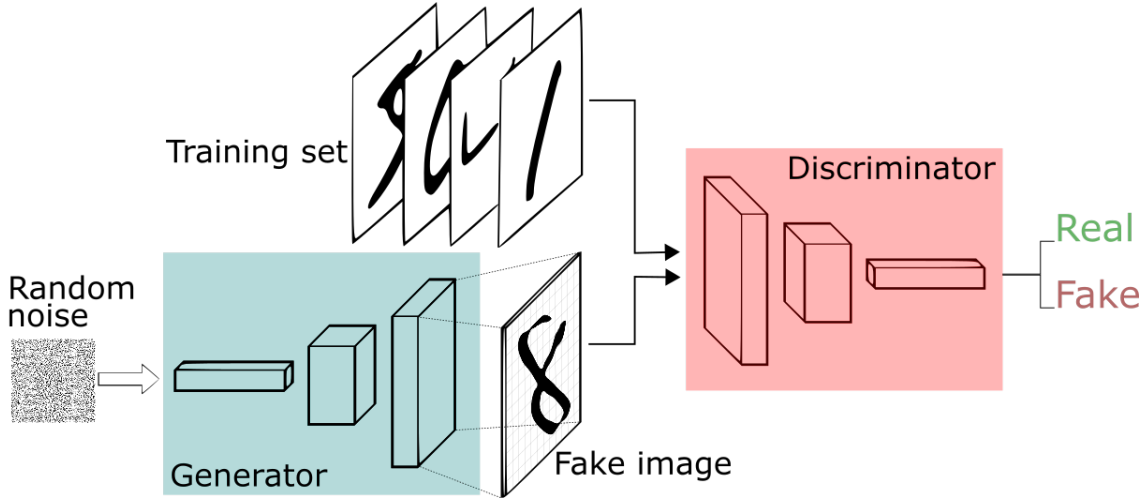


Figure 3.1: A generative adversarial model involves two networks; a Generator which transforms noise vectors to generated images and a Discriminator which seeks to identify real images from fake ones.

the ability nowadays to generate random *realistic* data (for instance realistic images) by applying successive Lipschitz operations to standard Gaussian vectors [GPAM⁺14].

A GAN architecture involves two networks, a generator model which maps random Gaussian noise to new plausible synthetic data and a discriminator model which classifies real data as real (from the dataset represented by p_{data}) or fake (for the generated data). The discriminator \mathcal{D} is updated directly through a binary classification problem, whereas the generator \mathcal{G} is updated through the discriminator based on the following Min-Max objective

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathbb{E}_{x \sim p_{\text{data}}} [\log \mathcal{D}(x)] + \mathbb{E}_{z \sim \mathcal{N}(0, I_d)} [\log (1 - \mathcal{D}(\mathcal{G}(z)))]$$

As such, the two models are trained alternatively in an adversarial manner, where the generator seeks to better deceive the discriminator and the former seeks to better identify the fake data [GPAM⁺14].

In particular, once both models are trained (when they reach a Nash equilibrium), DL representations of GAN-data –and GAN-data themselves– are schematically constructed in practice as follows:

$$\text{Real Data} \approx \text{GAN Data} = \mathcal{F}_N \circ \dots \circ \mathcal{F}_1(z), \quad (3.1)$$

where $z \sim \mathcal{N}(0, I_d)$, d stands for the input dimension of the generator model, N the number of layers, and the \mathcal{F}_i 's either Fully Connected Layers, Convolutional Layers, Pooling Layers, Up-sampling Layers and Activation Functions, Residual Layers or Batch Normalizations. All these operations happen to be *Lipschitz* applications. Precisely,

- **Fully Connected Layers and Convolutional Layers:** These are affine operations which can be expressed as

$$\mathcal{F}_i(x) = W_i x + b_i,$$

for W_i the weight matrix and b_i the bias vector. Here the Lipschitz constant is the operator norm (the largest singular value) of the weight matrix W_i , that is $\|\mathcal{F}_i\|_{\text{lip}} =$

$$\sup_{u \neq 0} \frac{\|W_i u\|_2}{\|u\|_2}.$$

- **Pooling Layers and Activation Functions:** Most commonly used activation functions and pooling operations are

$$\begin{aligned}\text{ReLU}(\mathbf{x}) &= \max(0, \mathbf{x}), \\ \text{MaxPooling}(\mathbf{x}) &= [\max(x_{S_1}), \dots, \max(x_{S_q})]^\top,\end{aligned}$$

where S_i 's are patches (*i.e.*, subsets of $[\dim(\mathbf{x})]$). These are at most 1-Lipschitz operations with respect to the Frobenius norm. Specifically, the maximum absolute sub-gradient of the ReLU activation function is 1, thus the ReLU operation has a Lipschitz constant of 1. Similarly, we can show that the Lipschitz constant of Max-Pooling layers is also 1.

- **Residual Connections:** Residual layers act the following way

$$\mathcal{F}_i(\mathbf{x}) = \mathbf{x} + \mathcal{F}_i^{(1)} \circ \dots \circ \mathcal{F}_i^{(\ell)}(\mathbf{x}),$$

where the $\mathcal{F}_i^{(j)}$'s are Fully Connected Layers or Convolutional Layers with Activation Functions, and which are Lipschitz operations. Thus \mathcal{F}_i is a Lipschitz operation with Lipschitz constant bounded by $1 + \prod_{j=1}^{\ell} \|\mathcal{F}_i^{(j)}\|_{lip}$.

- **Batch Normalization (BN) Layers:** They consist in statistically standardizing [IS15] the vectors of a small batch $\mathcal{B} = \{\mathbf{x}_i\}_{i=1}^b \subset \mathbb{R}^d$ as follows: for each $\mathbf{x}_k \in \mathcal{B}$

$$\mathcal{F}_i(\mathbf{x}_k) = \text{diag} \left(\frac{\mathbf{a}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) (\mathbf{x}_k - \mu_{\mathcal{B}} \mathbf{1}_d) + \mathbf{b}$$

where $\mu_{\mathcal{B}} = \frac{1}{db} \sum_{k=1}^b \sum_{i=1}^d [\mathbf{x}_k]_i$, $\sigma_{\mathcal{B}}^2 = \frac{1}{db} \sum_{k=1}^b \sum_{i=1}^d ([\mathbf{x}_k]_i - \mu_{\mathcal{B}})^2$, $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ are parameters to be learned and $\text{diag}(\mathbf{v})$ transforms a vector \mathbf{v} to a diagonal matrix with its diagonal entries being those of \mathbf{v} . Thus BN is a Lipschitz transformation with Lipschitz constant $\|\mathcal{F}_i\|_{lip} = \sup_i \left| \frac{\mathbf{a}_i}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right|$.

Therefore, as illustrated in Figure B.2, since standard Gaussian vectors are concentrated vectors as mentioned in Proposition 2.3 and since the notion of concentrated vectors is stable by Lipschitz transformations thanks to Proposition 2.4, GAN-data (and their DL representations) are concentrated vectors by design given the construction in Equation (3.1). Moreover, in order to generate data belonging to a specific class, Conditional GANs have been introduced [MO14]; once again data generated by these models are concentrated vectors as a consequence of the following Corollary.

Corollary 3.1. *Let $\mathcal{G}_1, \dots, \mathcal{G}_n : \mathbb{R}^d \rightarrow \mathbb{R}^p$ a set of n Lipschitz applications with Lipschitz constants $\|\mathcal{G}_i\|_{lip}$. Let $\mathcal{G} : \mathbb{R}^{d \times n} \rightarrow \mathbb{R}^{p \times n}$ be defined for each $\mathbf{X} \in \mathbb{R}^{d \times n}$ as*

$$\mathcal{G}(\mathbf{X}) = [\mathcal{G}_1(\mathbf{X}_{:,1}), \dots, \mathcal{G}_n(\mathbf{X}_{:,n})].$$

Then, for $\mathbf{Z} \in \mathcal{M}_{d,n}$

$$\mathbf{Z} \propto \mathcal{E}_q \quad \Rightarrow \quad \mathcal{G}(\mathbf{Z}) \propto \mathcal{E}_q \left(\sup_i \|\mathcal{G}_i\|_{lip} \right). \quad (3.2)$$

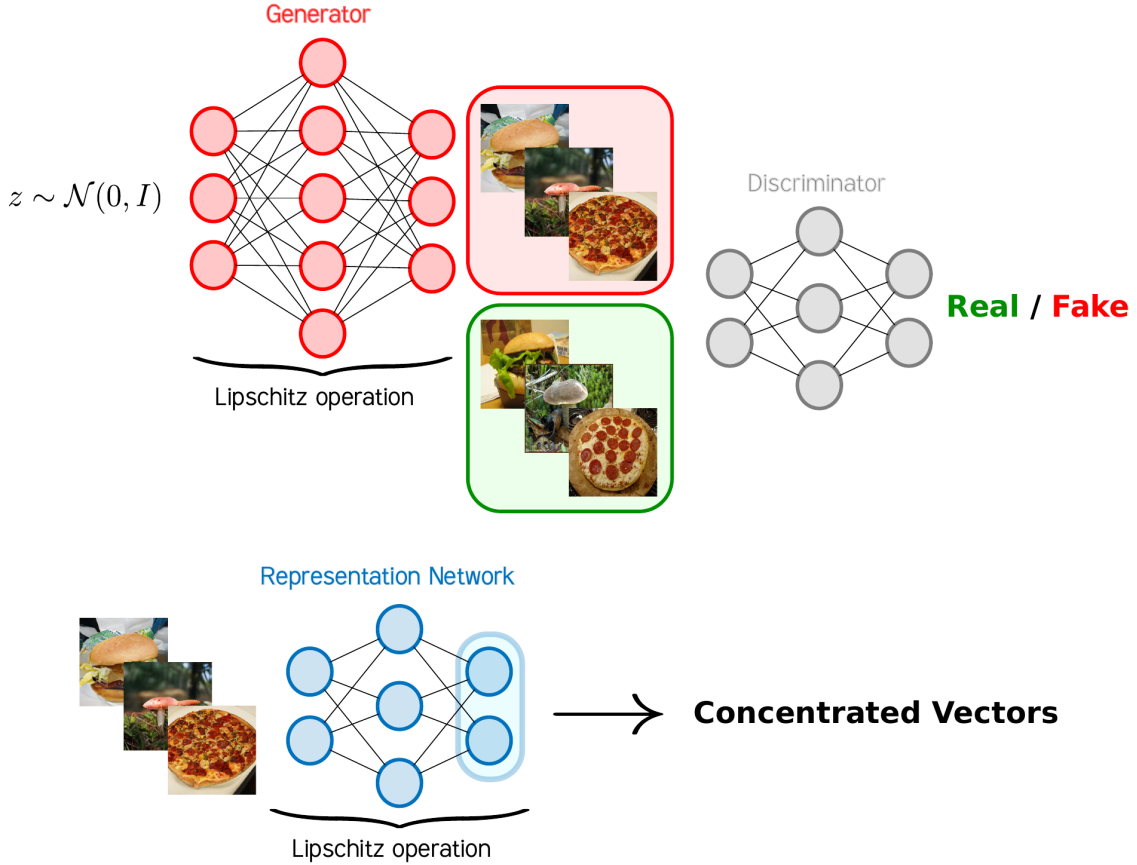


Figure 3.2: Deep learning representations of GAN-data are constructed by applying successive Lipschitz operations to Gaussian vectors, therefore they are *concentrated* vectors by design, since Gaussian vectors are concentrated and thanks to the Lipschitz stability in Proposition 2.4

Proof. This is a consequence of Proposition 2.4 since the map \mathcal{G} is $\sup_i \|\mathcal{G}_i\|_{lip}$ -Lipschitz with respect to (w.r.t.) the Frobenius norm. Indeed, for $\mathbf{X}, \mathbf{H} \in \mathbb{R}^{d \times n}$: $\|\mathcal{G}(\mathbf{X} + \mathbf{H}) - \mathcal{G}(\mathbf{X})\|_F^2 \leq \sum_{i=1}^n \|\mathcal{G}_i\|_{lip}^2 \cdot \|\mathbf{H}_{:,i}\|^2 \leq \sup_i \|\mathcal{G}_i\|_{lip}^2 \cdot \|\mathbf{H}\|_F^2$. \square

Indeed, a generator of a Conditional GAN model can be seen as a set of multiple generators where each generates data of a specific class conditionally on the class label (e.g., BigGAN model [BDS18]).

Yet, in order to ensure that the resulting Lipschitz constant of the combination of the above operations does not scale with the network or data size, so to maintain good concentration behaviors, a careful control of the learned network parameters is needed. This control happens to be already considered in practice in order to ensure the stability of GANs during the learning phase, notably to generate realistic and high-resolution images [RLNH17, BDS18]. The control of the Lipschitz constant of representation networks is also needed in practice in order to make them robust against adversarial examples [SZS⁺13, GAA⁺17]. This control is particularly ensured through spectral normalization of the affine layers [BDS18], such as Fully Connected Layers, Convolutional Layers and Batch Normalization. Indeed, spectral normalization [MKKY18] consists in applying the operation $\mathbf{W} \leftarrow \mathbf{W} / \sigma_1(\mathbf{W})$ to the affine layers at each backward iteration of the back-

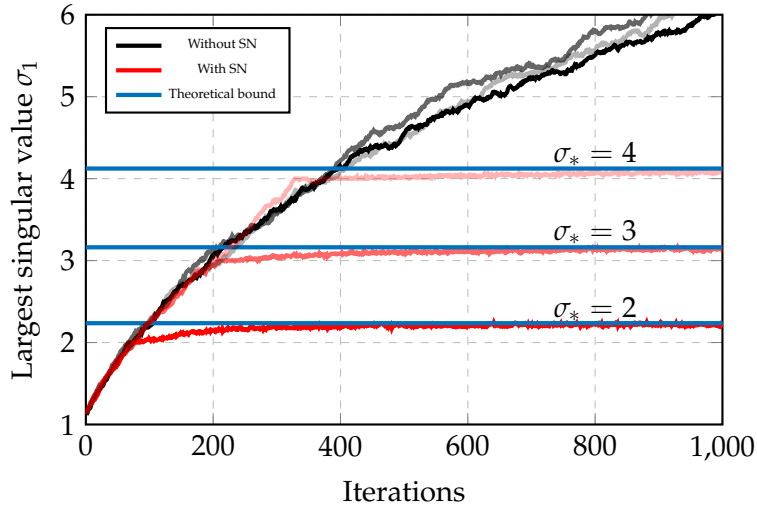


Figure 3.3: Behavior of the largest singular value of a weight matrix in terms of the iterations of a random walk (see proposition 3.1), without spectral normalization in **(black)** and with spectral normalization in **(red)**. The **(blue)** lines correspond to the theoretical bound $\sqrt{\sigma_*^2 + \eta^2 d_1 d_0}$ for different σ_* 's. We took $d_0 = d_1 = 100$ and $\eta = 1/d_0$.

propagation algorithm, where $\sigma_1(\mathbf{W})$ stands for the largest singular value of the weight matrix \mathbf{W} . [BDS18], have notably observed that, without spectral constraints, a subset of the generator layers grow throughout their GAN training and explode at collapse. They thus suggested the following spectral normalization –which happens to be less restrictive than the standard spectral normalization $\mathbf{W} \leftarrow \mathbf{W}/\sigma_1(\mathbf{W})$ [MKKY18]– to the affine layers:

$$\mathbf{W} \leftarrow \mathbf{W} - (\sigma_1(\mathbf{W}) - \sigma_*) \mathbf{u}_1(\mathbf{W})\mathbf{v}_1(\mathbf{W})^\top \quad (3.3)$$

where $\mathbf{u}_1(\mathbf{W})$ and $\mathbf{v}_1(\mathbf{W})$ denote respectively the left and right largest singular vectors of \mathbf{W} , and σ_* is an hyper-parameter fixed during training.

To get an insight about the influence of this operation and to ensure that it controls the Lipschitz constant of the generator, the following proposition provides the dynamics of a random walk in the space of parameters along with the spectral normalization in Equation (3.3). Indeed, since stochastic gradient descent (SGD) consists in estimating the gradient of the loss function on randomly selected batches of data, it can be assimilated to a random walk in the space of parameters [ASD18].

Proposition 3.1 (Lipschitz constant control). *Let $\sigma_* > 0$ and \mathcal{G} be a neural network composed of N affine layers, each one of input dimension d_{i-1} and output dimension d_i for $i \in [N]$, with 1-Lipschitz activation functions. Assume that the weights of \mathcal{G} at layer $i + 1$ are initialized as $\mathcal{U}([-\frac{1}{\sqrt{d_i}}, \frac{1}{\sqrt{d_i}}])$, and consider the following dynamics with learning rate η :*

$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \eta \mathbf{E}, \text{ with } \mathbf{E}_{i,j} \sim \mathcal{N}(0,1) \\ \mathbf{W} &\leftarrow \mathbf{W} - \max(0, \sigma_1(\mathbf{W}) - \sigma_*) \mathbf{u}_1(\mathbf{W})\mathbf{v}_1(\mathbf{W})^\top. \end{aligned} \quad (3.4)$$

Then, $\forall \varepsilon > 0$, the Lipschitz constant of \mathcal{G} is bounded at convergence with high probability as:

$$\|\mathcal{G}\|_{lip} \leq \prod_{i=1}^N \left(\varepsilon + \sqrt{\sigma_*^2 + \eta^2 d_i d_{i-1}} \right). \quad (3.5)$$

Proof. Since the Lipschitz constant of a composition of Lipschitz functions is bounded by the product of their Lipschitz constants, we consider the case $N = 1$ and a linear activation function. In this case, the Lipschitz constant corresponds to the largest singular value of the weight matrix. We consider the following notations for the proof

$$\begin{aligned} \bar{\mathbf{W}}_t &= \mathbf{W}_t - \eta E_t \text{ with } [E_t]_{i,j} \sim \mathcal{N}(0, 1) \\ \mathbf{W}_{t+1} &= \bar{\mathbf{W}}_t - \max(0, \bar{\sigma}_{1,t} - \sigma_*) \bar{\mathbf{u}}_{1,t} \bar{\sigma}_{1,t}^\top \end{aligned}$$

where $\bar{\sigma}_{1,t} = \sigma_1(\bar{\mathbf{W}}_t)$, $\bar{\mathbf{u}}_{1,t} = \mathbf{u}_1(\bar{\mathbf{W}}_t)$ and $\bar{\sigma}_{1,t} = v_1(\bar{\mathbf{W}}_t)$. The effect of spectral normalization is observed in the case where $\sigma_* > \bar{\sigma}_{1,t}$, otherwise the Lipschitz constant is bounded by σ_* . We therefore have

$$\|\bar{\mathbf{W}}_t\|_F^2 \leq \|\mathbf{W}_t\|_F^2 + \eta^2 d_1 d_0 \quad (3.6)$$

$$\|\mathbf{W}_{t+1}\|_F^2 = \|\bar{\mathbf{W}}_t\|_F^2 + \sigma_*^2 - \bar{\sigma}_{1,t}^2 \quad (3.7)$$

- If $\|\mathbf{W}_{t+1}\|_F \geq \|\mathbf{W}_t\|_F$, we have by equation 3.6 and equation 3.7

$$\|\bar{\mathbf{W}}_t\|_F^2 \leq \|\mathbf{W}_t\|_F^2 + \sigma_*^2 - \bar{\sigma}_{1,t}^2 + \eta^2 d_1 d_0 \Rightarrow \|\bar{\mathbf{W}}_t\| = \bar{\sigma}_{1,t} \leq \sqrt{\sigma_*^2 + \eta^2 d_1 d_0} = \delta$$

And since $\|\mathbf{W}_{t+1}\| \leq \|\bar{\mathbf{W}}_t\|$, we have $\|\mathbf{W}_{t+1}\| \leq \delta$.

- Otherwise, if there exists τ such that $\|\mathbf{W}_{\tau+1}\|_F < \|\mathbf{W}_\tau\|_F$, then for all $\varepsilon > 0$ there exists an iteration $\tau' \geq \tau$ such that $\|\mathbf{W}_{\tau'}\| \leq \delta + \varepsilon$. Indeed, otherwise we denote $\varepsilon_t = \|\mathbf{W}_t\|^2 - \delta^2$ and $\varepsilon_t > 0$ for all $t \geq \tau$. And if for all $t \geq \tau$, $\|\mathbf{W}_{t+1}\|_F \leq \|\mathbf{W}_t\|_F$, we have by equation 3.6 and equation 3.7

$$\|\mathbf{W}_t\|_F^2 - \|\mathbf{W}_{t+1}\|_F^2 \geq \|\bar{\mathbf{W}}_t\|^2 - \delta^2 \geq \|\mathbf{W}_{t+1}\|^2 - \delta^2 = \varepsilon_{t+1}$$

Integrating the above expression from τ to $T - 1 \geq \tau$, we end up with

$$\|\mathbf{W}_\tau\|_F^2 - \|\mathbf{W}_T\|_F^2 \geq \sum_{t=\tau}^{T-1} \varepsilon_t \Rightarrow 0 \leq \|\mathbf{W}_T\|_F^2 \leq \|\mathbf{W}_\tau\|_F^2 - \sum_{t=\tau}^{T-1} \varepsilon_t,$$

therefore, when $T \rightarrow \infty$, ε_t has to tend to 0 otherwise the right hand-side of the last inequality will tend to $-\infty$ which is absurd. □

Proposition 3.1 shows that the Lipschitz constant of a neural network is controlled when trained with the spectral normalization in Equation (3.3). In particular, recalling the notations in Proposition 3.1, in the limit where $d_i \rightarrow \infty$ with $\frac{d_i}{d_{i-1}} \rightarrow \gamma_i \in (0, \infty)$ for all $i \in [N]$ and choosing the learning rate $\eta = \mathcal{O}(d_0^{-1})$, the Lipschitz constant of \mathcal{G} is of order $\mathcal{O}(1)$ if it has finitely many layers N and σ_* is constant. Therefore, with this spectral normalization, it can be assumed that $\|\mathcal{G}\|_{lip} = \mathcal{O}(1)$ when dimensions grow. Figure 3.3 depicts the behavior of the Lipschitz constant of a linear layer with and without spectral normalization in the setting of Proposition 3.1, which confirms the obtained bound.

3.2 Random Gram Matrices of Concentrated Data

Now we turn to the analysis of the behavior of large Gram matrices assuming the data being concentrated vectors.

3.2.1 Motivation

As it represents the canonical form of similarity, the Gram matrix is at the core of various machine learning algorithms. Moreover, it provides a natural way to quantify the quality of a given representation. Indeed, the performance of machine learning methods depends strongly on the choice of the data representation (or features) on which they are applied. This data representation should ideally contain *relevant information* about the learning task in order to achieve learning with *simple* models and *small* amount of samples. In this sense, the simplest machine learning model is naturally a linear classifier which relies on the Gram matrix, and considering a small amount of data put us naturally in the random matrix theory regime where both p and n are large and comparable. In terms of representation learning, Deep neural networks [RHW⁺88] have particularly shown impressive results by automatically learning representations from raw data (e.g., images). However, due to the complex structure of deep learning models, the characterization of their hidden representations is still an open problem [B⁺09].

Specifically, quantifying what makes a given deep learning representation better than another is a fundamental question in the field of *Representation Learning* [BCV13]. Relying on [MBM11] a data representation is said to be *good* when it is possible to build *simple* models on top of it that are *accurate* for the given learning problem. Authors in [MBM11] have notably quantified the layer-wise evolution of the representation in deep networks by computing the principal components of the Gram matrix $\mathbf{G}_\ell = \{\phi_\ell(\mathbf{x}_i)^\top \phi_\ell(\mathbf{x}_j)\}_{i,j=1}^n$ at each layer for n input data $\mathbf{x}_1, \dots, \mathbf{x}_n$, where $\phi_\ell(\mathbf{x})$ is the representation of \mathbf{x} at layer ℓ of the given DL model, and the number of components controls the model simplicity. In their study, the impact of the representation at each layer is quantified through the prediction error of a linear classifier trained on the principal subspace of \mathbf{G}_ℓ .

Pursuing on this idea, given a certain representation model $\mathbf{x} \mapsto \phi(\mathbf{x})$, we aim in this study at theoretically analyzing the large dimensional behavior, and in particular the spectral information (i.e., eigenvalues and dominant eigenvectors), of the corresponding Gram matrix $\mathbf{G} = \{\phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)\}_{i,j=1}^n$ in order to determine the information encoded (i.e., the sufficient statistics) by the representation model on a set of real data $\mathbf{x}_1, \dots, \mathbf{x}_n$. Indeed, standard classification and regression algorithms –along with the last layer of a neural network [YKYR18]– retrieve the data information directly from functionals or the eigenspectrum of \mathbf{G} ¹. To this end, though, one needs a statistical model for the representations given the distribution of the raw data (e.g., images) which is generally unknown and not analytically tractable. Yet, as we have shown in the previous section, due to recent advances in generative models since the advent of Generative Adversarial Nets [GPAM⁺14], it is now possible to generate complex data structures by applying successive *Lipschitz* operations to Gaussian random vectors. In particular, GAN-data are used in practice as substitutes of real data for data augmentation [ASE17]. On the other hand, the fundamental concentration of measure phenomenon [Led05a] tells us

¹For instance, spectral clustering uses the dominant eigenvectors of \mathbf{G} , while support vector machines use functionals (quadratic forms) involving \mathbf{G} .

that Lipschitz-ally transformed Gaussian vectors satisfy a concentration property. Precisely, defining the class of *concentrated* vectors $x \in E$ through concentration inequalities of $f(x)$, for any real Lipschitz observation $f : E \rightarrow \mathbb{R}$, implies that deep learning representations of GAN-data fall within this class of random vectors, since the mapping $x \mapsto \phi(x)$ is Lipschitz. Thus, GAN-data are concentrated random vectors and thus a more appropriate statistical model of realistic data, as we demonstrated in the previous section.

Targeting classification applications by assuming a mixture of concentrated random vectors model (see Definition 2), this study describes the spectral behavior of Gram matrices G in the large n, p regime. Precisely, we show that these matrices have asymptotically (as $n, p \rightarrow \infty$ with $p/n \rightarrow c < \infty$) the same first-order behavior as for a Gaussian Mixture Model (GMM). As a result, by generating images using the BigGAN model [BDS18] and considering different commonly used deep representation models, we show that the spectral behavior of the Gram matrix computed on these representations is the same as on a GMM model with the same p -dimensional means and covariances. A surprising consequence is that, for GAN data, the aforementioned *sufficient statistics* to characterize the quality of a given representation network are only the *first* and *second* order statistics of the representations. This behavior is shown by simulations to extend beyond random GAN-data to real images from the Imagenet dataset [DDS⁺09].

3.2.2 Model and Main Results

3.2.2.1 Mixture of Concentrated Vectors

In this section, we assume data to be a mixture of concentrated random vectors with controlled $\mathcal{O}(1)$ Lipschitz constant (e.g., DL representations of GAN-data as we discussed in the previous section). Precisely, let x_1, \dots, x_n be a set of mutually independent random vectors in \mathbb{R}^p . We suppose that these vectors are distributed as one of k classes of distribution laws $\mathcal{L}_1, \dots, \mathcal{L}_k$ with distinct means $\{\mu_\ell\}_{\ell=1}^k$ and “covariances” $\{\Sigma_\ell\}_{\ell=1}^k$ defined respectively as

$$\mu_\ell = \mathbb{E}_{x_i \sim \mathcal{L}_\ell}[x_i], \quad \Sigma_\ell = \mathbb{E}_{x_i \sim \mathcal{L}_\ell}[x_i x_i^\top]. \quad (3.8)$$

For some $q > 0$, we consider a q -exponential concentration property on the laws \mathcal{L}_ℓ , in the sense that for any family of independent vectors y_1, \dots, y_s sampled from \mathcal{L}_ℓ , $[y_1, \dots, y_s] \in \mathcal{E}_q$ (see Definition 7). Without loss of generality, we arrange the x_i 's in a data matrix $X = [x_1, \dots, x_n]$ such that, for each $\ell \in [k]$

$$x_{1+\sum_{j=1}^{\ell-1} n_j}, \dots, x_{\sum_{j=1}^{\ell} n_j} \sim \mathcal{L}_\ell(\mu_\ell, \Sigma_\ell)$$

where n_ℓ stands for the number of x_i 's sampled from \mathcal{L}_ℓ . In particular, we have the concentration of the data matrix X as in the following assumption

Assumption 3 (Concentrated data). *We assume $X \propto \mathcal{E}_q$ for some $q > 1$.*

Such a data matrix X can be constructed through Lipschitz-ally transformed Gaussian vectors ($q = 2$), with controlled Lipschitz constant, thanks to Corollary 3.1. In particular, DL representations of GAN-data are constructed as such, as shown in Section 3.1. We further introduce the following notations that will be used subsequently.

$$M = [\mu_1, \dots, \mu_k] \in \mathbb{R}^{p \times k}, \quad J = [j_1, \dots, j_k] \in \mathbb{R}^{n \times k}, \quad Z = [z_1, \dots, z_n] \in \mathbb{R}^{p \times n},$$

where $\mathbf{j}_\ell \in \mathbb{R}^n$ stands for the canonical vector selecting the x_i 's of distribution \mathcal{L}_ℓ , defined by $(\mathbf{j}_\ell)_i = \mathbf{1}_{x_i \sim \mathcal{L}_\ell}$, and the \mathbf{z}_i 's are the centered versions of the x_i 's, i.e. $\mathbf{z}_i = x_i - \boldsymbol{\mu}_\ell$ for $x_i \sim \mathcal{L}_\ell(\boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)$.

3.2.2.2 Behavior of the Gram matrix of concentrated vectors

Now we study the behavior of the Gram matrix $\mathbf{G} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ in the large n, p limit and under the model of the previous section. Indeed, \mathbf{G} appears as a central component in many classification, regression and clustering methods. Precisely, a finer description of the behavior of \mathbf{G} provides access to the internal functioning and performance evaluation of a wide range of machine learning methods such as Least Squares SVMs [A⁺02], Semi-supervised Learning [CSZ09] and Spectral Clustering [NJW02]. Indeed, the performance evaluation of these methods has already been studied under GMM models in [LC17, MC17, CBG⁺16] through RMT. On the other hand, analyzing the spectral behavior of \mathbf{G} for DL representations quantifies their quality –through its principal subspace [MBM11]– as we have discussed in the introduction. In particular, the Gram matrix decomposes as

$$\mathbf{G} = \frac{1}{p} \mathbf{J} \mathbf{M}^\top \mathbf{M} \mathbf{J}^\top + \frac{1}{p} \mathbf{Z}^\top \mathbf{Z} + \frac{1}{p} (\mathbf{J} \mathbf{M}^\top \mathbf{Z} + \mathbf{Z}^\top \mathbf{M} \mathbf{J}^\top). \quad (3.9)$$

Intuitively \mathbf{G} decomposes as a low-rank informative matrix containing the class canonical vectors through \mathbf{J} and a noise term represented by the other matrices and essentially $\mathbf{Z}^\top \mathbf{Z}$. Given the form of this decomposition, RMT predicts –through an analysis of the spectrum of \mathbf{G} and under a GMM model [BGC16]– the existence of a threshold ξ function of the ratio p/n and the data statistics for which the dominant eigenvectors of \mathbf{G} contain information about the classes only when $\|\mathbf{M}^\top \mathbf{M}\| \geq \xi$ asymptotically (i.e., only when the means of the different classes are sufficiently distinct). See Subsection 2.3.1 for an illustrative example.

As we saw in Chapter 2, in order to characterize the spectral behavior (i.e., eigenvalues and leading eigenvectors) of \mathbf{G} under the concentration assumption in Assumption 3 on \mathbf{X} , we will be interested in studying the behavior of its resolvent defined as, for $z \in \mathbb{C} \setminus \mathbb{R}_-$

$$\mathbf{R}(z) = (\mathbf{G} + z \mathbf{I}_n)^{-1} \quad (3.10)$$

Practically speaking for spectral clustering with \mathbf{G} , since $\mathbf{R}(z)$ and \mathbf{G} share the same eigenvectors with associated eigenvalues $\frac{1}{\lambda_i + z}$ for $\mathbf{R}(z)$ with $\{\lambda_i\}_{i=1}^n$ the eigenvalues of \mathbf{G} , the projector matrix corresponding to the top m eigenvectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ of \mathbf{G} can be calculated through a Cauchy integral $\mathbf{U} \mathbf{U}^\top = \frac{-1}{2\pi i} \oint_\gamma \mathbf{R}(-z) dz$ where γ is an oriented complex contour surrounding the top m eigenvalues of \mathbf{G} . See Subsection 2.1.3 for more details.

To study the behavior of $\mathbf{R}(z)$, we look for a so-called *deterministic equivalent* [HLN⁺07] (see Definition 6) $\tilde{\mathbf{R}}(z)$ for $\mathbf{R}(z)$. In the following, we present our main result which gives such a deterministic equivalent under the concentration assumption on \mathbf{X} in Assumption 3 and under the following assumptions.

Assumption 4. As $p \rightarrow \infty$,

1. $p/n \rightarrow c \in (0, \infty)$,

2. The number of classes k is bounded,
3. $\|\boldsymbol{\mu}_\ell\| = \mathcal{O}(\sqrt{p})$ and $\mathbb{E}\|\mathbf{x}_i\| = \mathcal{O}(\sqrt{p})$.

Theorem 3.1 (Deterministic Equivalent for $\mathbf{R}(z)$). *Under the model described in Section 3.2.2.1 and Assumptions 3-4, we have $\mathbf{R}(z) \in \mathcal{E}_q(p^{-1/2})$. Furthermore,*

$$\mathbf{R}(z) \leftrightarrow \bar{\mathbf{R}}(z) = \frac{1}{z} \text{diag} \left\{ \frac{\mathbf{I}_{n_\ell}}{1 + \delta_\ell^*(z)} \right\}_{\ell=1}^k + \frac{1}{pz} \mathbf{J} \Omega_z \mathbf{J}^\top$$

Specifically,

$$\|\mathbb{E}\mathbf{R}(z) - \bar{\mathbf{R}}(z)\| = \mathcal{O} \left(\sqrt{\frac{\log(p)}{p}} \right)$$

with $\Omega_z = \mathbf{M}^\top \bar{\mathbf{Q}}(z) \mathbf{M} \odot \text{diag} \left\{ \frac{\delta_\ell^*(z)-1}{\delta_\ell^*(z)+1} \right\}_{\ell=1}^k$ and $\bar{\mathbf{Q}}(z) = \left(\frac{1}{ck} \sum_{\ell=1}^k \frac{\boldsymbol{\Sigma}_\ell}{1 + \delta_\ell^*(z)} + z \mathbf{I}_p \right)^{-1}$ where $\delta^*(z) = [\delta_1^*(z), \dots, \delta_k^*(z)]^\top$ is the unique fixed point of the system of equations for each $\ell \in [k]$

$$\delta_\ell(z) = \frac{1}{p} \text{tr} \left(\boldsymbol{\Sigma}_\ell \left(\frac{1}{ck} \sum_{j=1}^k \frac{\boldsymbol{\Sigma}_j}{1 + \delta_j(z)} + z \mathbf{I}_p \right)^{-1} \right)$$

Sketch of proof. The first step of the proof is to show the concentration of $\mathbf{R}(z)$. This comes from the fact that the application $\mathbf{X} \mapsto \mathbf{R}(z)$ is $2z^{-3/2}p^{-1/2}$ -Lipschitz w.r.t. the Frobenius norm, thus we have by Proposition 2.4 that $\mathbf{R}(z) \in \mathcal{E}_q(p^{-1/2})$.

The second step consists in estimating $\mathbb{E}\mathbf{R}(z)$ through a deterministic matrix $\bar{\mathbf{R}}(z)$. Indeed, $\mathbf{R}(z)$ can be expressed as a function of $\mathbf{Q}(z) = (\mathbf{X}\mathbf{X}^\top/p + z\mathbf{I}_p)^{-1}$ as $\mathbf{R}(z) = z^{-1}(\mathbf{I}_n - \mathbf{X}^\top \mathbf{Q}(z) \mathbf{X}/p)$, where the statistical dependency between \mathbf{X} and $\mathbf{Q}(z)$ is handled through Propositions C.1 and C.2 and finally exploiting the result of [LC19] which shows that $\mathbb{E}\mathbf{Q}(z)$ can be estimated through $\bar{\mathbf{Q}}(z)$ as per Theorem 2.9, we obtain the estimator $\bar{\mathbf{R}}(z)$ for $\mathbb{E}\mathbf{R}(z)$.

A more detailed proof is provided the appendix in Section C.1. □

Remark 3.1 (Equivalence between Theorem 3.1 and Theorem 2.5). *Note that Theorem 3.1 and Theorem 2.5 are equivalent when $\mathbf{M} = \mathbf{0}$ with the change of variable $c_j g_j(z) = \frac{1}{1 + \delta_j(z)}$. In particular, Theorem 2.5 supposes a Gaussian mixture model with zero means while we generalize this result to non-zero means which appear through the matrix of means \mathbf{M} .*

This result allows specifically to (i) describe the limiting eigenvalues distribution of \mathbf{G} , (ii) determine the spectral detectability threshold mentioned above (See Subsection 2.3.1 in the case of binary unsupervised clustering), (iii) evaluate the asymptotic ‘‘content’’ of the leading eigenvectors of \mathbf{G} (see again Subsection 2.3.1) and, much more fundamentally, (iv) infer the asymptotic performances of machine learning algorithms (See Subsection 2.3.2) that are based on simple functionals of \mathbf{G} (e.g., LS-SVM, spectral clustering etc.). Looking carefully at Theorem 3.1 we see that the spectral behavior of the Gram matrix \mathbf{G} computed on concentrated vectors only depends on the *first* and *second* order statistics of the laws \mathcal{L}_ℓ (their means $\boldsymbol{\mu}_\ell$ and ‘‘covariances’’ $\boldsymbol{\Sigma}_\ell$). This suggests the surprising result that \mathbf{G} has the same behavior as when the data follow a GMM model with the same means and covariances. The asymptotic spectral behavior of \mathbf{G} is therefore *universal*



Figure 3.4: **(Top)** GAN generated images using the BigGAN model [BDS18]. **(Bottom)** Real images selected from the Imagenet dataset [DDS⁺09]. We considered $n = 1500$ images from $k = 3$ classes which are {mushroom, pizza, hamburger}.

with respect to the data distribution laws which satisfy the aforementioned concentration properties (for instance DL representations of GAN-data). We illustrate this universality result in the next section by considering data as CNN representations of GAN generated images.

3.2.2.3 Application to GAN-generated Images

This section presents experiments that confirm the result of Theorem 3.1. In particular, we compare, in the first part, the eigenvalues distribution and the largest eigenvectors of the Gram matrix computed on deep learning representations with those of the Gram matrix computed on Gaussian data with the same first and second order moments. In the second part of this section, we evaluate the performance of a linear SVM model on the principal subspace of the Gram matrix (computed on the representations or on the corresponding Gaussian data) by varying the number of components in the same vein as the work of montavon2011kernel. In the following, all representation networks are standard convolutional neural networks pre-trained on the Imagenet dataset deng2009imagenet, in particular, we used pre-trained models of the Pytorch deep learning framework.

Spectrum and Dominant Eigenspace of the Gram Matrix. We consider $n = 1500$ data $x_1, \dots, x_n \in \mathbb{R}^p$ as CNN representations –across popular CNN architectures of different sizes p – of GAN-generated images using the generator of the Big-GAN model [BDS18]. We further use real images from the Imagenet dataset [DDS⁺09] for comparison. In particular, we empirically compare the spectrum of the Gram matrix of this data with the

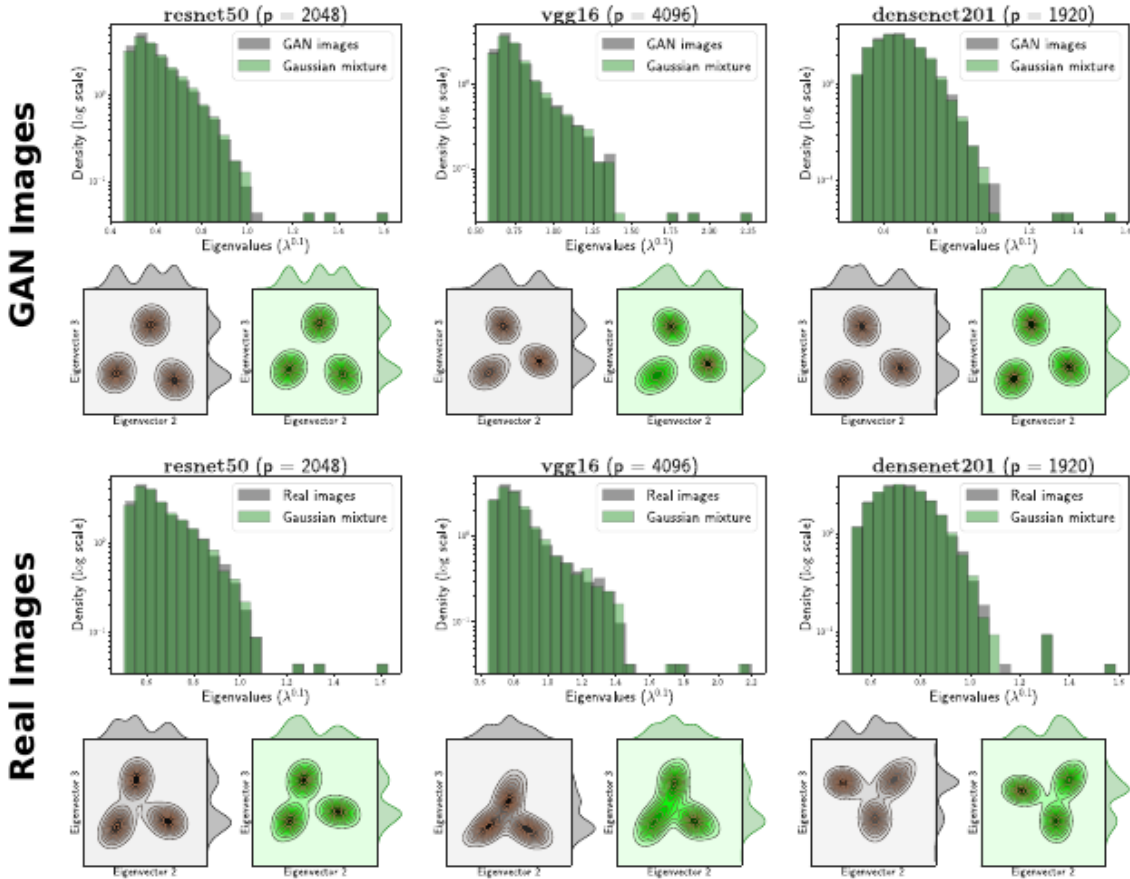


Figure 3.5: **(Top)** Spectrum and leading eigenspace of the Gram matrix for CNN representations of GAN generated images using the BigGAN model [BDS18]. **(Bottom)** Spectrum and leading eigenspace of the Gram matrix for CNN representations of real images selected from the Imagenet dataset [DDS⁺09]. Columns correspond to the three representation networks which are *resnet50*, *vgg16* and *densenet201*. We used $n = 1500$ images and considered $k = 3$ classes as depicted in Figure 3.4.

Gram matrix of a GMM model with the same means and covariances. We also consider the leading 2-dimensional eigenspace of the Gram matrix which contains clustering information as detailed in the previous section. Figure 3.4 depicts some images generated using the Big-GAN model (Top) and the corresponding real class images from the Imagenet dataset (Bottom). The Big-GAN model is visually able to generate highly realistic images which are by construction concentrated vectors, as discussed in Section 3.1 and therefore satisfy the assumptions of Theorem 3.1.

Figure 3.5 depicts the spectrum and leading 2D eigenspace of the Gram matrix computed on CNN representations of GAN generated and real images (in gray), and the corresponding GMM model with same first and second order statistics (in green). The Gram matrix is seen to follow the same spectral behavior for GAN-data as for the GMM model which is a natural consequence of the universality result of Theorem 3.1 with respect to the data distribution. Besides, and perhaps no longer surprisingly, we further observe that the spectral properties of G for real data (here CNN representations of real images) are conclusively matched by their Gaussian counterpart.

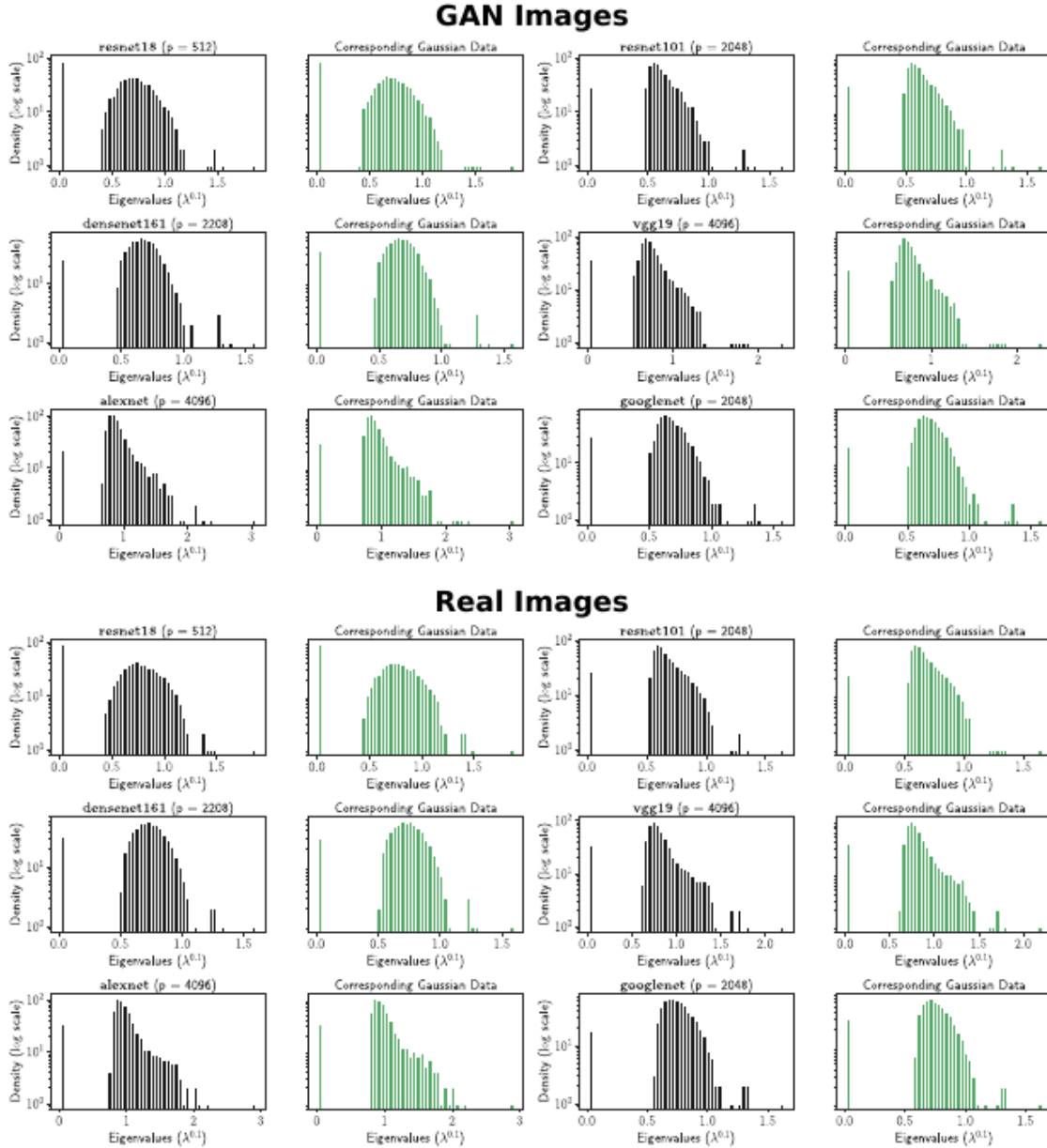


Figure 3.6: Spectrum of the Gram matrix for CNN representations in **(black)** and the corresponding Gaussian data **(in green)** for GAN generated images using the BigGAN model [BDS18] **(Top)** and for real images randomly selected from the Imagenet dataset [DDS⁺09] **(Bottom)**. The considered representation network are *resnet18*, *resnet101*, *densenet161*, *vgg19*, *alexnet* and *googlenet*. We used $n = 600$ images selected among $k = 6$ classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (100 images per class).

Figure 3.6 shows more results about the Gram matrix spectrum of the representations **(in black)** and the corresponding Gaussian data **(in green)**, by considering more representation networks and using $k = 6$ classes for both GAN images and real images, which confirms the result of Theorem 3.1. This both theoretically and empirically confirms that the proposed random matrix framework is fully compliant with the theoretical analysis

of real machine learning datasets. As a consequence, recalling the work of [MBM11], the *quality* of a given representation is quantified through the prediction accuracy of a linear classifier trained on the principal Gram matrix eigenvectors of the representations computed on a set of samples. Given our result in Theorem 3.1, and the fact that the top m eigenvectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ of \mathbf{G} are related to the resolvent matrix $\mathbf{R}(z)$ through the Cauchy integral $\mathbf{U}\mathbf{U}^\top = \frac{-1}{2\pi i} \oint_\gamma \mathbf{R}(-z)dz$ where γ is an oriented complex contour surrounding the top m eigenvalues of \mathbf{G} , we should expect that the prediction accuracy of a linear classifier trained on the principal eigenvectors of \mathbf{G} be the same for the representations themselves as for the corresponding Gaussian data with the same first and second order moments. Therefore, the purpose of the following section is to show simulations which confirm this result.

Linear SVM Performance on the Dominant k -dimensional Eigenspace of \mathbf{G} . Now we compare the performance of a linear SVM model trained on the dominant Gram matrix's k -dimensional eigenspace of the representations versus the corresponding Gaussian data with the same first and second order moments. Experiments were made in the following settings:

- **Data types:** We do the experiments for both GAN generated images using the BigGAN model [BDS18] and for real images randomly selected for the Imagenet dataset [DDS⁺09]. In both cases we consider $n = 6000$ images.
- **Classes:** We consider $k = 6$ classes which are: hamburger, mushroom, pizza, strawberry, coffee and daisy.
- **Representation networks:** We consider 9 representation networks pre-trained on the Imagenet dataset [DDS⁺09] which are: *vgg16* ($p = 4096$), *vgg19* ($p = 4096$), *resnet18* ($p = 512$), *resnet50* ($p = 2048$), *resnet101* ($p = 2048$), *densenet161* ($p = 2208$), *densenet201* ($p = 1920$), *alexnet* ($p = 4096$) and *googlenet* ($p = 2048$).

Figure 3.7 depicts the train and test accuracy of a linear SVM trained on the top k eigenvectors of \mathbf{G} , for the representations (of GAN generated images) and the corresponding Gaussian data, for different values of k . As we can notice, the performance of the SVM model on the representations matches its performance on the corresponding Gaussian data with the same first and second order statistics as predicted by Theorem 3.1. This matching seems to extend beyond GAN images (which are concentrated vectors) to real images as depicted in Figure 3.8. As a consequence, our results suggest that the *quality* of a given representation network can be quantified through their first two statistical moments.

3.2.3 Central Contribution

The central contributions parts of this work are to highlight the universality aspects of large random matrices on real data, by studying the behavior of the Gram matrix for a large class of random vectors, the so-called *concentrated vectors*, which are much richer than Gaussian vectors. The concentration assumption is particularly motivated by the fact that realistic data can be generated using GANs by Lipschitz transformations of Gaussian vectors, which fall within the class of concentrated vectors. Particularly, by generating images across GANs, we have shown in this work that the spectral behavior

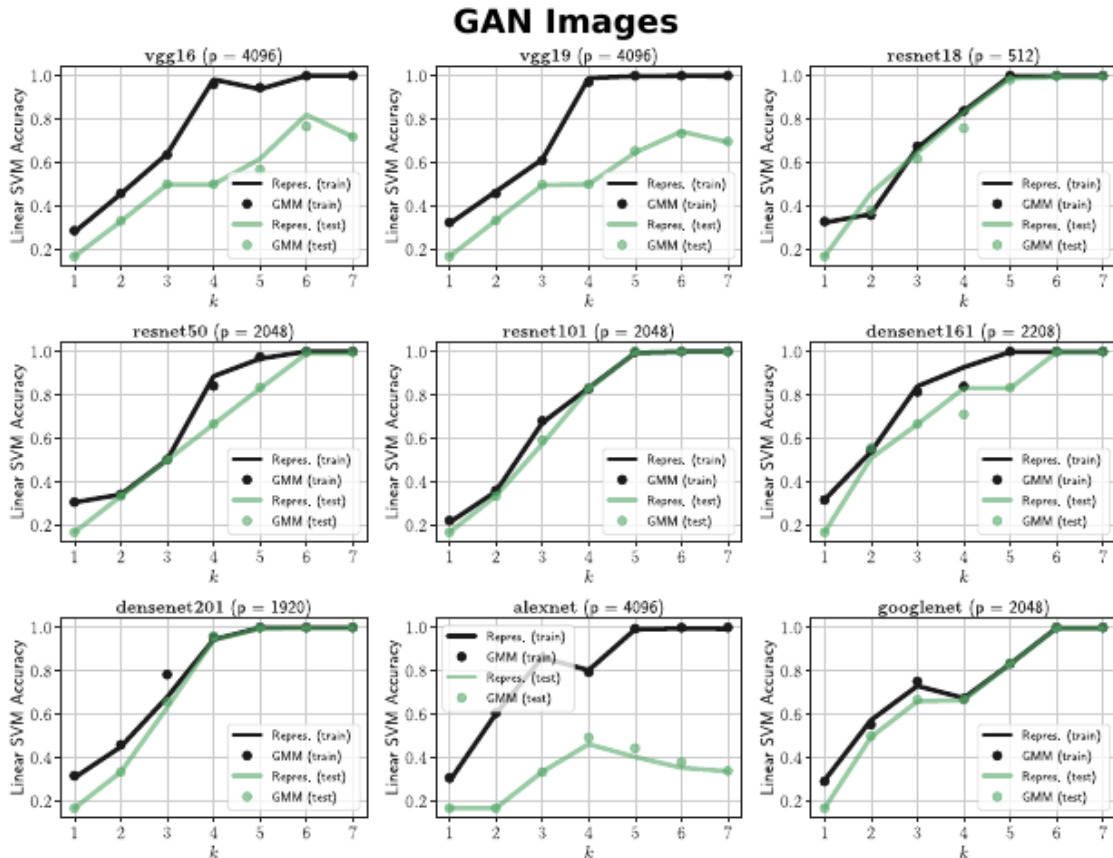


Figure 3.7: **Train** and **test** accuracy of a linear SVM model trained on the top k eigenvectors of the Gram matrix, computed on the representations of GAN generated images. We generated with the BigGAN model [BDS18] $n = 6000$ images belonging to the 6 classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (1000 images per class). We considered 9 representation networks which are *vgg16*, *vgg19*, *resnet18*, *resnet50*, *resnet101*, *densenet161*, *densenet201*, *alexnet* and *googlenet*. Lines represent the performance of the SVM model on the representations themselves whereas dots represent the performance of the SVM model on Gaussian data with the same first and second order moments. We used a train vs test split of 2/3 and 1/3 respectively.

of the Gram matrix is the same on the generated data as on a Gaussian mixture model, thereby making the estimation of machine learning algorithms (which are based on the Gram matrix) predictable through random matrix theory for real data if we can assimilate them to GAN generated data. In the next chapter we will go beyond the simple Gram matrix to the more general setting of kernel matrices, which constructively lead to *non-linear* ML methods, and therefore additional notions shall be introduced in order to break the non-linearity.

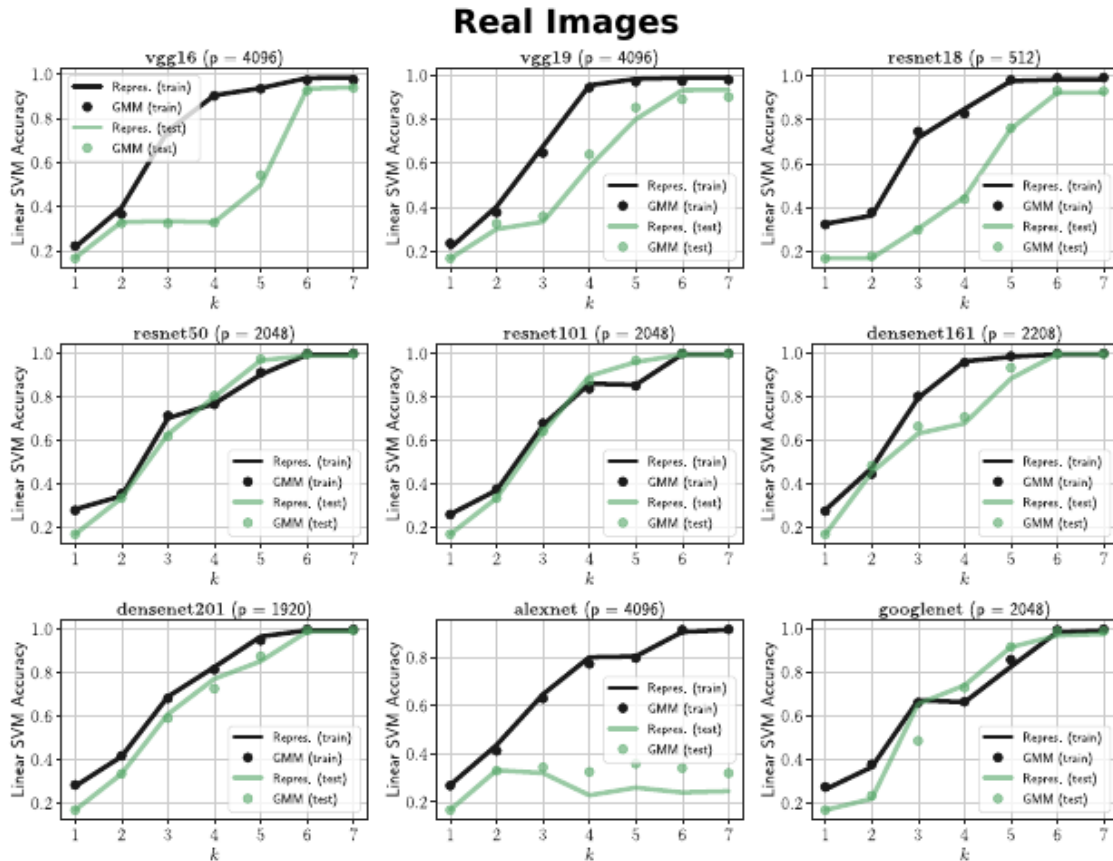


Figure 3.8: **Train** and **test** accuracy of a linear SVM model trained on the top k eigenvectors of the Gram matrix, computed on the representations of Real images. We randomly sampled $n = 6000$ images belonging to the 6 classes {hamburger, mushroom, pizza, strawberry, coffee, daisy} (1000 images per class) of the Imagenet dataset [DDS⁺09]. We considered 9 representation networks *vgg16*, *vgg19*, *resnet18*, *resnet50*, *resnet101*, *densenet161*, *densenet201*, *alexnet* and *googlenet*. Lines represent the performance of the SVM model on the representations themselves whereas dots represent the performance of the SVM model on Gaussian data with the same first and second order moments. We used a train vs test split of 2/3 and 1/3 respectively.

Chapter 4

Random Kernel Matrices of Concentrated Data

Contents

4.1 Kernel Spectral Clustering	71
4.1.1 Motivation	72
4.1.2 Model and Main Results	72
4.1.2.1 Behavior of Large Kernel Matrices	76
4.1.2.2 Application to GAN-generated Images	87
4.1.3 Central Contribution and perspectives	89
4.2 Sparse Principal Component Analysis	89
4.2.1 Motivation	90
4.2.2 Model and Main Results	91
4.2.2.1 Random Matrix Equivalent	92
4.2.2.2 Application to sparse PCA	94
4.2.2.3 Experimental Validation	96
4.2.3 Central Contribution and Perspectives	99

This chapter is composed of two main parts. The first part analyzes the behavior of kernel matrix of the for $K = \{\frac{1}{p}\|x_i - x_j\|^2\}$ where the x_i 's are supposed to be concentrated. The second part analyzes the behavior of kernel matrices of the form $f(\hat{\Sigma})$ where Σ stands for the sample covariance matrix.

4.1 Kernel Spectral Clustering

This section is based on the following work:

- (C2) MEA. Seddik, M. Tamaazousti, R. Couillet, “Kernel Random Matrices of Large Concentrated Data: The Example of GAN-generated Images”, IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'19), Brighton, United-Kingdom, 2019.

Kernel matrices generalize the inner product similarity of Gram matrices by introducing a non-linear function in order to capture non-linearities from data. the objective of this section is to analyze the behavior of these matrices under the concentration hypothesis and supposing the high-dimensional RMT regime.

4.1.1 Motivation

Gram matrices fall within the larger class of kernel matrices since they correspond to the particular case of inner-product similarity. In order to retrieve non-linear structures from data, the used methods for achieving classical classification or regression tasks rely on non-linear approaches including neural networks [KSH12, LC18c] and algorithms that are based on kernel methods, such as kernel-based support vector machines [LC17], semi-supervised classification [MC17], kernel-based PCA [STC19b] and spectral clustering [AKC18, CBG⁺16].

Due to their non-linear design, these methods are particularly difficult to analyze theoretically. For practical large and numerous data, the study of kernel-based methods relies on the characterization of kernel matrices $\mathbf{K} \in \mathbb{R}^{n \times n}$ in the large dimensional regime (*i.e.*, $p/n \rightarrow c_0$ as $n \rightarrow \infty$). Under *asymptotically non-trivial* growth rate assumptions on the data statistics (*i.e.*, maintaining a feasible get not too easy problem), the entries $K_{ij} = f(\mathbf{x}_i^\top \mathbf{x}_j / p)$ or $K_{ij} = f(\|\mathbf{x}_i - \mathbf{x}_j\|^2 / p)$ of \mathbf{K} tend to a limiting constant independently of the data classes – *the between and within class vectors are “equidistant” in high-dimension*. This observation allows one to study \mathbf{K} through a Taylor expansion, thereby giving access to the characterization of functionals of \mathbf{K} or its (informative) eigenvalue-eigenvector pairs in the large dimensional regime.

Indeed, such an analysis was initiated in [EK⁺10b] where it has been shown that \mathbf{K} has a linear behavior in the large p, n asymptotics. Under a k -class Gaussian mixture model, it has been shown in [CBG⁺16] that the normalized Laplacian matrix associated with \mathbf{K} behaves asymptotically as a so-called spiked random matrix, where some of the isolated eigenvalues and eigenvectors contain clustering information. In particular, the authors in [CBG⁺16] demonstrated that the obtained theoretical model agrees with empirical results using the popular MNIST dataset [LeC98], thereby suggesting a sort of *universality* of spectral clustering regarding the underlying data distribution.

The aim of this study is to confirm this observation by relaxing the Gaussianity assumption to the class of concentrated vectors leveraging on the observation that GAN data fall within this class of random vectors as presented in Section 3.1. In this study, we analyze the kernel matrix \mathbf{K} under a k -class *concentration* mixture model [LC19]. Precisely, we prove that \mathbf{K} behaves (up to centering) asymptotically as a spiked random matrix in the large p large n regime, thereby generalizing the results of [CBG⁺16] to a broader class of distributions. We particularly confirm our theoretical findings by considering the input data as popular CNN representations of images generated by the BigGAN model [BDS18], where the latter is trained to fit the manifold distribution of the well-known Imagenet dataset. We further consider real images for comparison.

4.1.2 Model and Main Results

We consider the same concentration assumptions as for the analysis of the Gram matrix in the previous chapter, but we use slightly different notations and assumptions. Consider n independent random vectors $x_1, \dots, x_n \in \mathbb{R}^p$ distributed in k classes represented by k distributions $\mathcal{L}_1, \dots, \mathcal{L}_k$ supposedly all distinct. We consider the hypothesis of q -exponential concentration, meaning that there exists $q \geq 2$ such that for all $s \in \mathbb{N}$, any $\ell \in [k]$ and any family of independent vectors $\mathbf{y}_1, \dots, \mathbf{y}_s$ following the distribution \mathcal{L}_ℓ , we

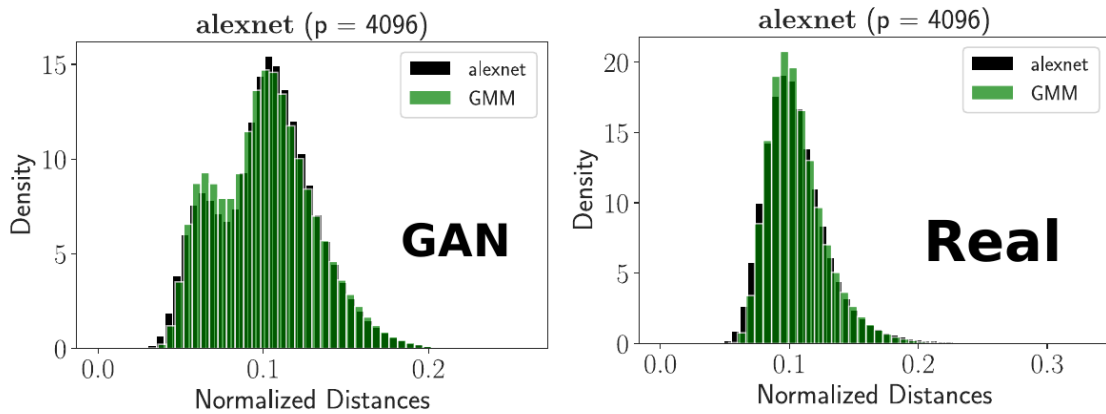


Figure 4.1: Histogram of distances with Alexnet representations for both GAN images (on the left) and real images (on the right). We can notice a perfect match between the representations and the corresponding GMM data. Note that, for real images, the pairwise distances *concentrate* around a constant quantity as per Lemma 4.1, this comes from the fact that real images are hardly separable compared to GANs ones.

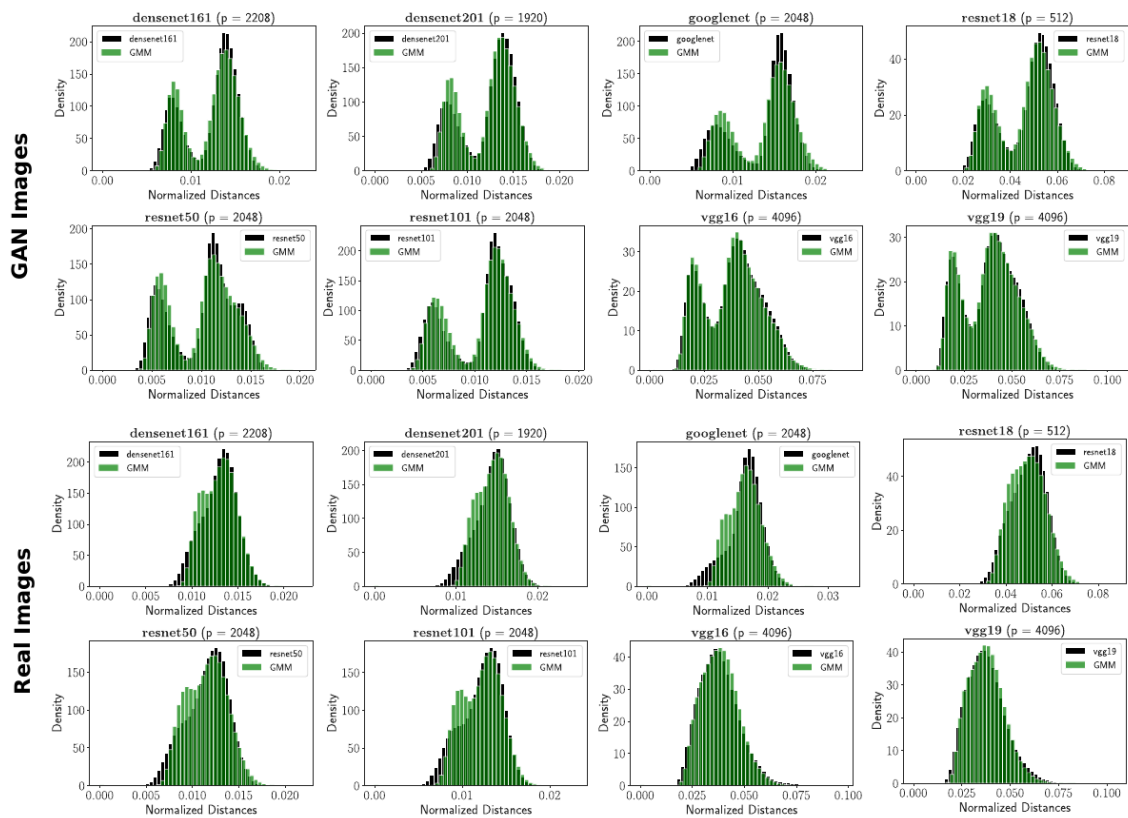


Figure 4.2: Histogram of distances across different representation networks for both GAN images (on top) and real images (on bottom). We can notice a perfect match between the representations and the corresponding GMM data. Note that, for real images, the pairwise distances *concentrate* around a constant quantity as per Lemma 4.1, this comes from the fact that real images are hardly separable compared to GANs ones.

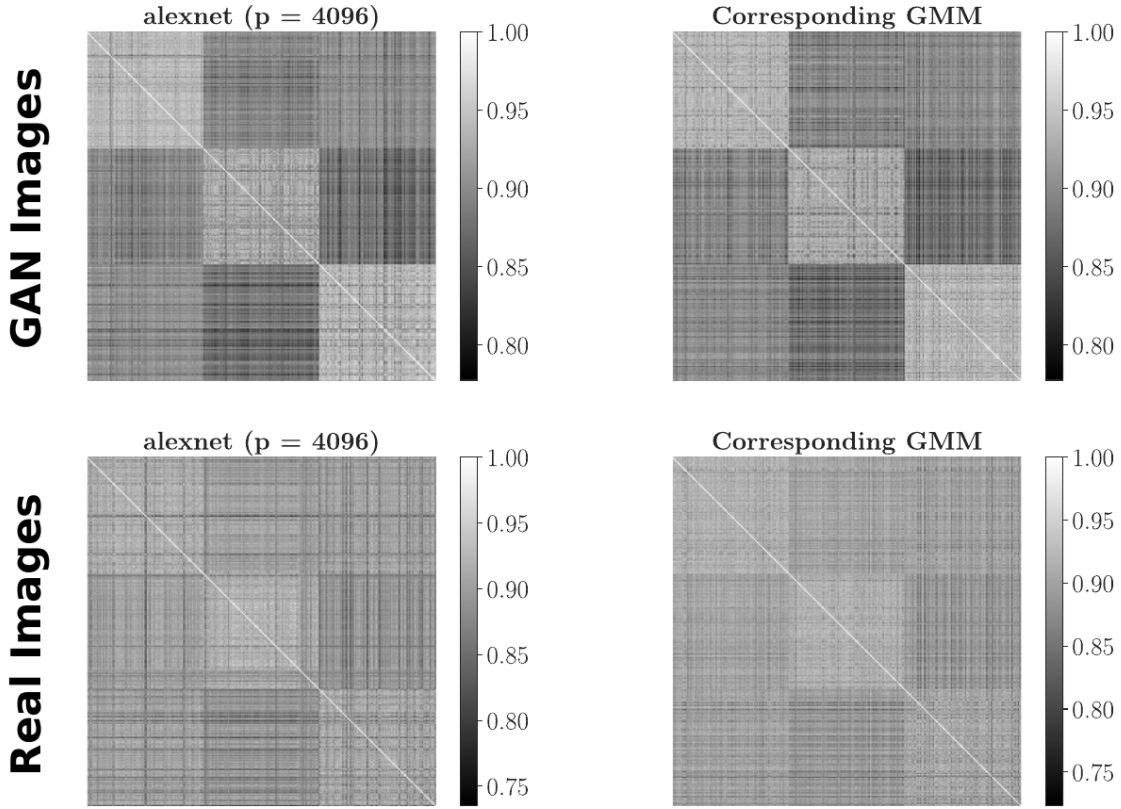


Figure 4.3: **(Left)** Kernel matrix of Alexnet representations for GAN images **(top)** and Real images **(bottom)** with the corresponding GMM data **(right)**. We can notice a perfect matching between the kernel matrix of the representations and the corresponding kernel with GMM data. We further note (at least for Real images) that all the entries of the kernel matrix tend to the same value $f(\tau_p)$ as a first order approximation. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

have the concentration $[y_1, \dots, y_s] \propto \mathcal{E}_q$ (see Definition 7). In particular, we consider the concentration of $X \equiv [x_1, \dots, x_n] \in \mathcal{M}_{p,n}$ as per the following assumption

Assumption 5 (Data concentration). $X \propto \mathcal{E}_q$.

For $\ell \in [k]$, we denote by μ_ℓ the mean of the distribution \mathcal{L}_ℓ , Σ_ℓ denotes its covariance matrix, defined respectively as

$$\mu_\ell \equiv \mathbb{E}_{x \sim \mathcal{L}_\ell}[x], \quad \Sigma_\ell \equiv \mathbb{E}_{x \sim \mathcal{L}_\ell}[xx^\top] - \mu_\ell \mu_\ell^\top, \quad (4.1)$$

and n_ℓ stands for the number of vectors among the x_i 's following \mathcal{L}_ℓ . Let $\mu \in \mathbb{R}^p$ and $\Sigma \in \mathbb{R}^{p \times p}$ be respectively defined as

$$\mu \equiv \sum_{\ell=1}^k \frac{n_\ell}{n} \mu_\ell, \quad \Sigma \equiv \sum_{\ell=1}^k \frac{n_\ell}{n} \Sigma_\ell \quad (4.2)$$

We further denote $\bar{\mu}_\ell \equiv \mu_\ell - \mu$ and $\bar{\Sigma}_\ell \equiv \Sigma_\ell - \Sigma$.

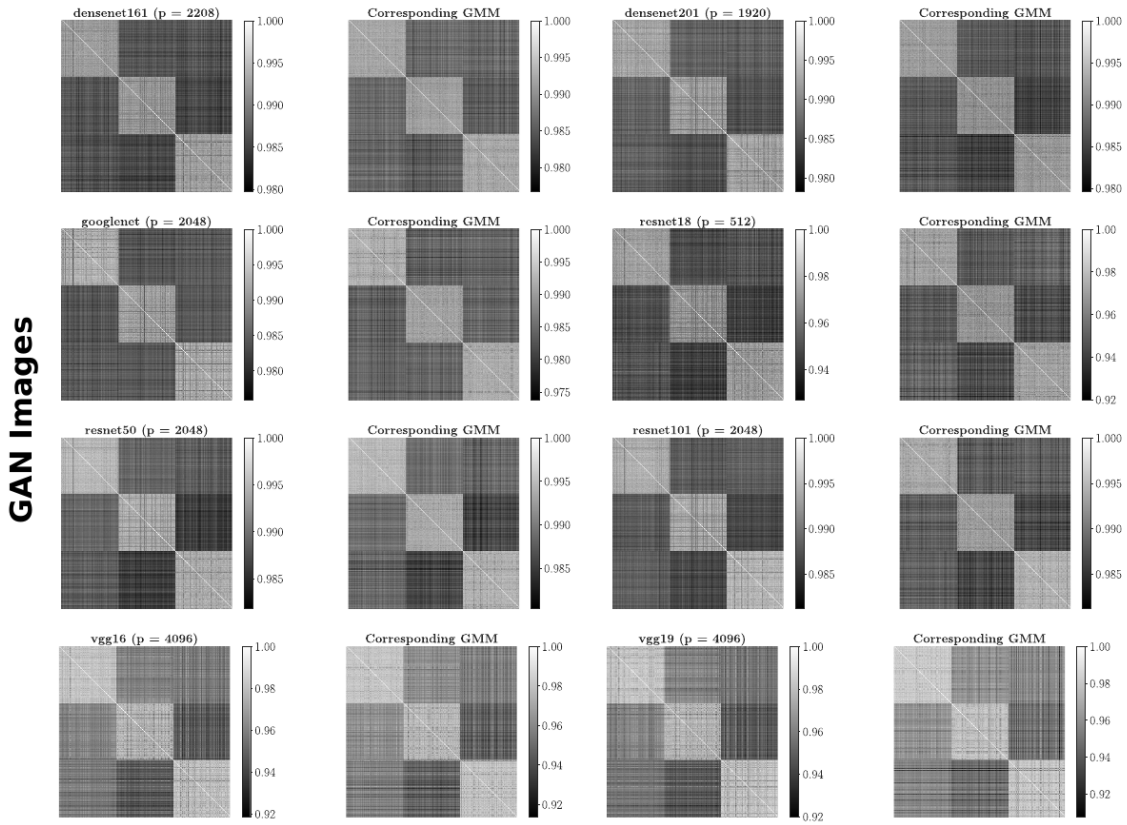


Figure 4.4: Kernel matrix of CNN representations of GAN data using different representations networks. The kernel matrix is defined as $K_{ij} = f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)$ where the \mathbf{x}_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

We shall consider the following set of assumptions on the data statistics and the kernel function in the large dimensional regime, meaning that both p and n grow at controlled joint rate. These assumptions notably guarantee the non-triviality of spectral clustering under the considered regime as we have shown in the introduction (see Subsection 1.3.2).

Assumption 6 (Growth rate). *As $p \rightarrow \infty$, consider the following conditions:*

- (Data) $\frac{p}{n} \rightarrow c_0 \in (0, \infty)$, $\frac{m_\ell}{n} \rightarrow c_0 \in (0, 1)$.
- (Means) $\limsup_p \max_\ell \|\bar{\boldsymbol{\mu}}_\ell\| < \infty$ and $\limsup_p \max_i \frac{1}{\sqrt{p}} \mathbb{E} \|\mathbf{x}_i\| < \infty$.
- (Covariances) $\limsup_p \max_\ell \|\bar{\boldsymbol{\Sigma}}_\ell\| < \infty$, $\limsup_p \max_{a,b} \frac{1}{\sqrt{p}} \text{tr} \bar{\boldsymbol{\Sigma}}_\ell < \infty$,
 $\limsup_p \max_{a,b} \frac{1}{p} \text{tr} \bar{\boldsymbol{\Sigma}}_a \bar{\boldsymbol{\Sigma}}_b < \infty$.

Assumption 7 (Kernel function). *Let $\tau \equiv \frac{2}{p} \text{tr} \boldsymbol{\Sigma}$ and let $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ be a three-times continuously differentiable function in a neighborhood of the values taken by τ and such that $\liminf_n f(\tau) > 0$.*

Without loss of generality, for each $\ell \in [k]$, we arrange the \mathbf{x}_i 's in \mathbf{X} as

$$\mathbf{x}_{1+\sum_{j=1}^{\ell-1} n_j}, \dots, \mathbf{x}_{\sum_{j=1}^{\ell} n_j} \sim \mathcal{L}_\ell(\boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}_\ell)$$

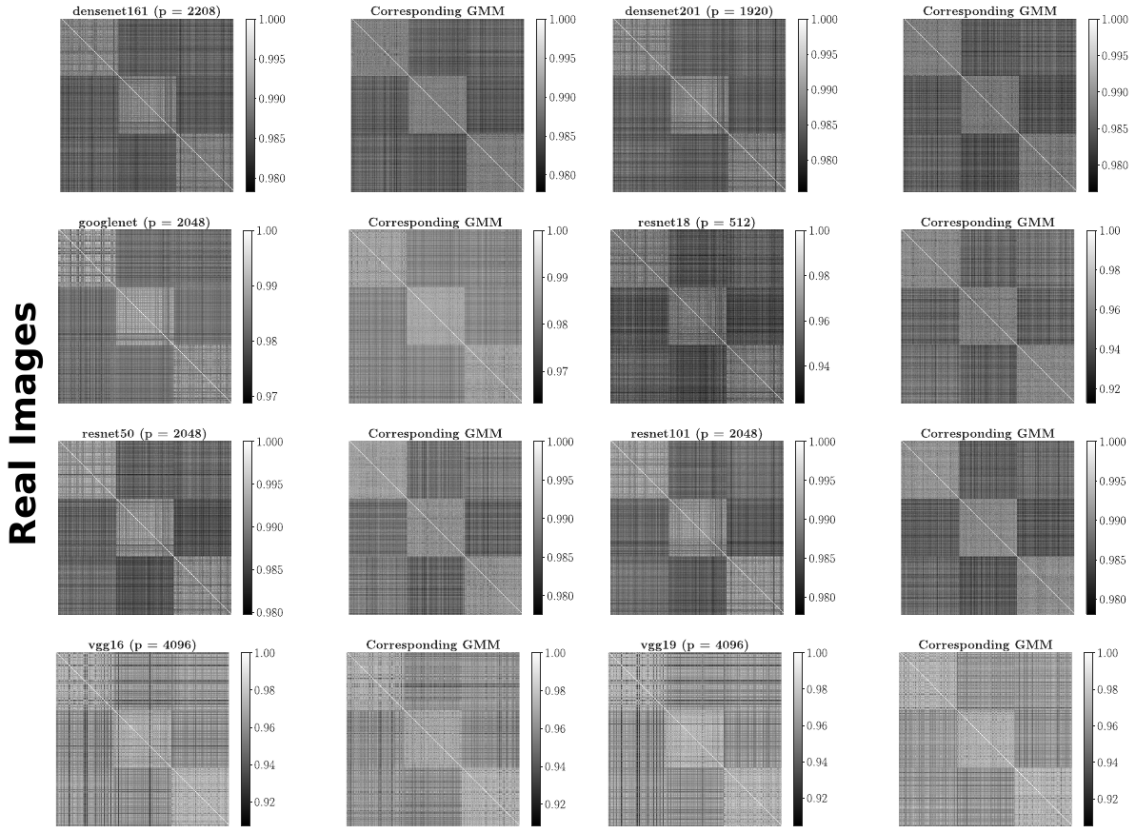


Figure 4.5: Kernel matrix of CNN representations of Real data using different representations networks. The kernel matrix is defined as $K_{ij} = f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)$ where the \mathbf{x}_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

and define the kernel matrix K as the translation-invariant random matrix

$$K \equiv \left\{ f \left(\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 \right) \right\}_{i,j=1}^n \quad (4.3)$$

4.1.2.1 Behavior of Large Kernel Matrices

Between and Within Class Data are “equidistant” in High-dimension. The first key and fundamental result states that the between and within class vectors are “equidistant” in the high-dimensional regime. Namely, we have the following lemma under the q -exponential concentration hypothesis and Assumption 6.

Lemma 4.1. Denote $\tau \equiv \frac{2}{p} \text{tr} \Sigma$ and let Assumption 6 hold. Then for any $\delta > 0$, we have with probability at least $1 - \delta$

$$\max_{1 \leq i \neq j \leq n} \left\{ \left| \frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2 - \tau \right| \right\} = \mathcal{O} \left(\frac{\log(\frac{p}{\sqrt{\delta}})^{1/q}}{\sqrt{p}} \right) \quad (4.4)$$

Proof. For $\mathbf{x}_i \sim \mathcal{L}_a(\boldsymbol{\mu}_a, \Sigma_a)$, denote

$$\mathbf{z}_i \equiv \frac{\mathbf{x}_i - \boldsymbol{\mu}_a}{\sqrt{p}} \quad \text{and} \quad \psi_i \equiv \|\mathbf{z}_i\|^2 - \mathbb{E} [\|\mathbf{z}_i\|^2] = \|\mathbf{z}_i\|^2 - \frac{1}{p} \text{tr} \Sigma_a$$

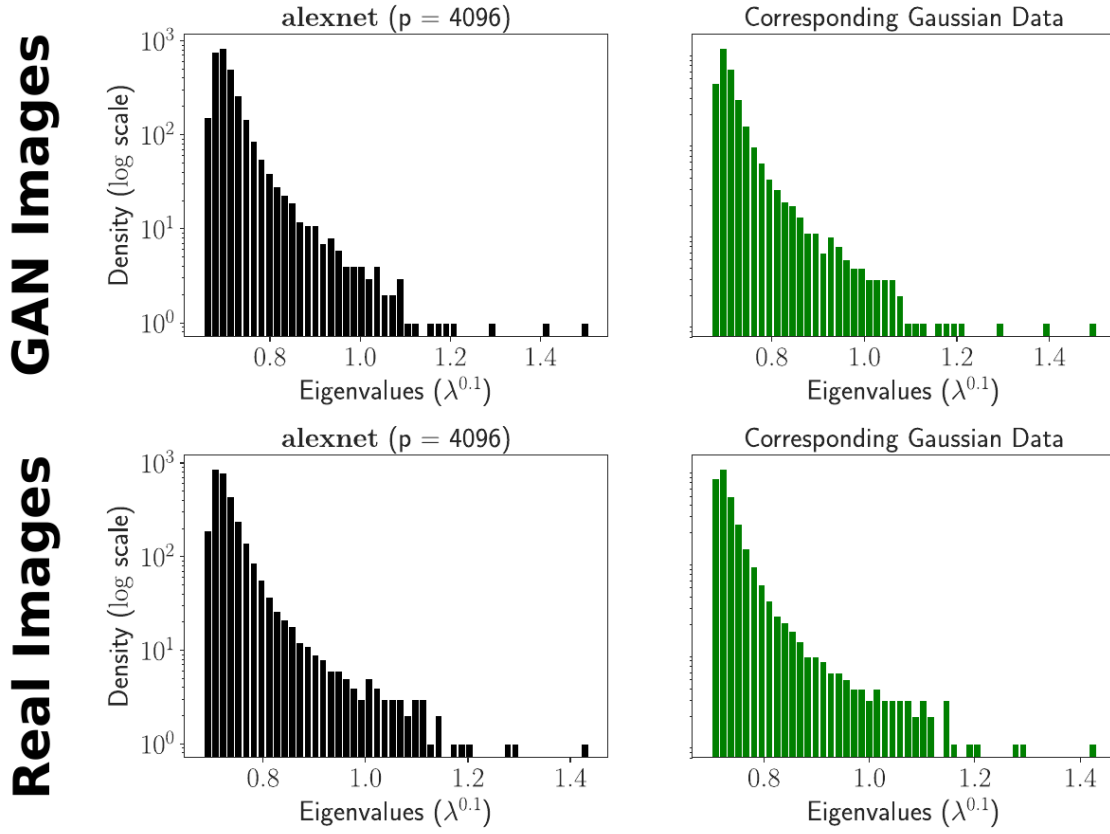


Figure 4.6: **(Left)** Spectrum of the kernel matrix of Alexnet representations for GAN images **(top)** and Real images **(bottom)** with the corresponding GMM data **(right)**. We can notice a perfect matching between the spectrum of the kernel matrix of the representations and the corresponding spectrum with GMM data. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

Let $a \neq b \in [k]$, $x_i \sim \mathcal{L}_a(\mu_a, \Sigma_a)$ and $x_j \sim \mathcal{L}_b(\mu_b, \Sigma_b)$, we decompose the normalized euclidean distance between x_i and x_j as

$$\begin{aligned} \frac{1}{p} \|x_i - x_j\|^2 &= \|z_i - z_j\|^2 + \frac{1}{p} \|\mu_a - \mu_b\|^2 + \frac{2}{\sqrt{p}} (\mu_a - \mu_b)' (z_i - z_j) \\ &= \tau + \frac{1}{p} \text{tr } \bar{\Sigma}_a + \frac{1}{p} \text{tr } \bar{\Sigma}_b + \psi_i + \psi_j - 2z_i' z_j \\ &\quad + \frac{1}{p} \|\bar{\mu}_a - \bar{\mu}_b\|^2 + \frac{2}{\sqrt{p}} (\bar{\mu}_a - \bar{\mu}_b)' (z_i - z_j) \end{aligned}$$

First, we show the asymptotic concentration of each stochastic term in the previous decomposition. Recalling the q -exponentially concentration assumption, each x_i is q -exponentially concentrated (since it can be seen as a 1-Lipschitz transformation of X), precisely $x_i \propto \mathcal{E}_q$. Thus, since $x \mapsto \frac{x - \mu_a}{\sqrt{p}}$ is $\frac{1}{\sqrt{p}}$ -Lipschitz, we obtain

$$z_i = \frac{x_i - \mu_a}{\sqrt{p}} \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right)$$

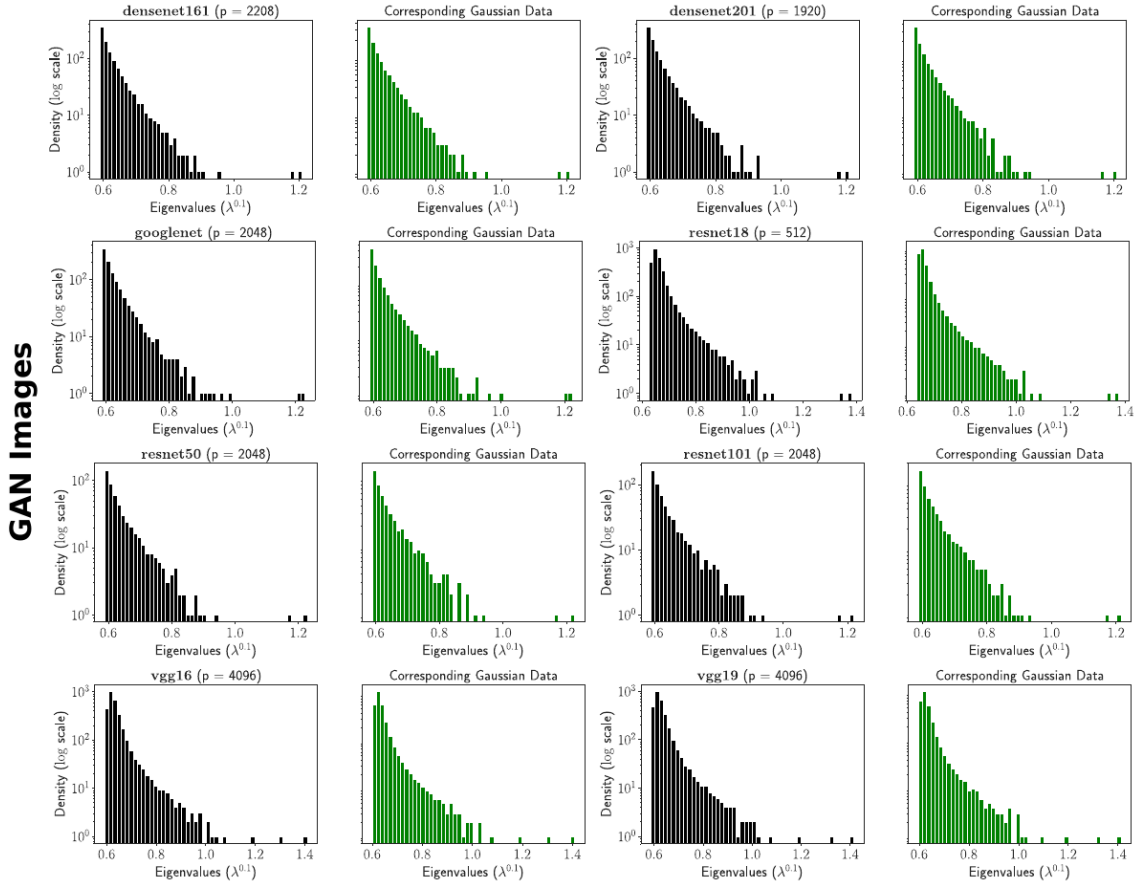


Figure 4.7: Spectrum of the kernel matrix for CNN representations of GAN data using different representations networks. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

We have by proposition 2.5 in [LC18b] the q -exponentially concentration of the concatenation $[z_i, z_j]$ since z_i and z_j are independent, namely

$$[z_i, z_j] \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right)$$

Moreover, since $[x, y] \mapsto x \pm y$ is 2-Lipschitz, we get

$$z_i \pm z_j \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right) \quad (4.5)$$

Now, since $z \rightarrow \frac{2}{\sqrt{p}}(\bar{\mu}_a - \bar{\mu}_b)'z$ is $(\frac{4}{\sqrt{p}} \limsup_n \max_\ell \|\bar{\mu}_\ell\|)$ -Lipschitz and recalling from Assumptions 6 that $\limsup_n \max_\ell \|\bar{\mu}_\ell\| < \infty$, we obtain

$$\frac{2}{\sqrt{p}}(\bar{\mu}_a - \bar{\mu}_b)'(z_i - z_j) \propto \mathcal{E}_q \left(\frac{1}{p} \right) \quad (4.6)$$

By the polarization identity $z_i'z_j = \frac{1}{4} [\|z_i + z_j\|^2 - \|z_i - z_j\|^2]$, recalling the concentration in equation 4.5 and applying Proposition 2.5 (involving the assumption $\frac{1}{\sqrt{n}}\mathbb{E}\|x_i\| < \infty$)

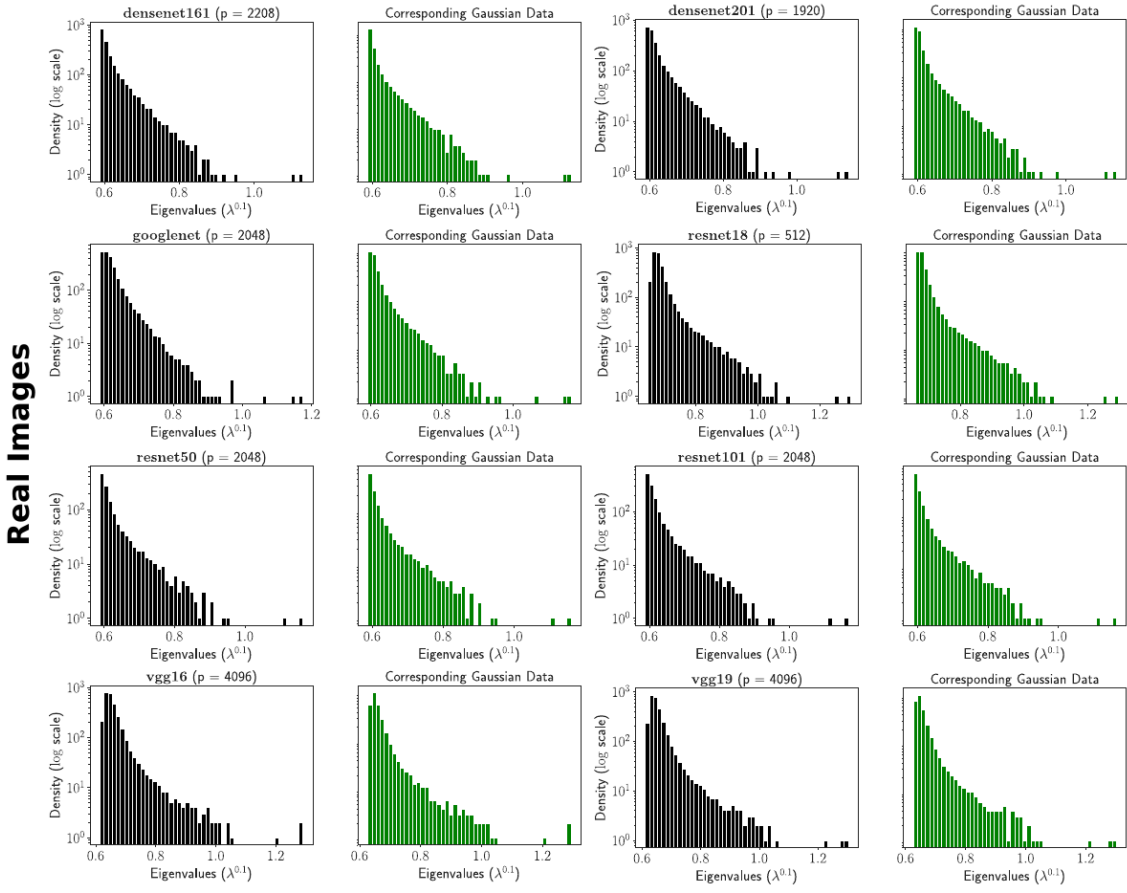


Figure 4.8: Spectrum of the kernel matrix for CNN representations of Real data using different representations networks. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

and Proposition 2.4, we have

$$\|z_i \pm z_j\|^2 \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right) + \mathcal{E}_{\frac{q}{2}} \left(\frac{1}{p} \right)$$

Once again, by Proposition 2.4, we obtain

$$z_i' z_j \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right) + \mathcal{E}_{\frac{q}{2}} \left(\frac{1}{p} \right) \quad (4.7)$$

Similarly, we also have

$$\psi_i \propto \mathcal{E}_q \left(\frac{1}{\sqrt{p}} \right) + \mathcal{E}_{\frac{q}{2}} \left(\frac{1}{p} \right)$$

Now that we have established the concentration of each stochastic term, we determine the order of their maximum values. By equation 4.6, there exists two absolute constants C and σ such that

$$\mathbb{P} \left\{ \left| \frac{2}{\sqrt{p}} (\bar{\mu}_a - \bar{\mu}_b)' (z_i - z_j) \right| \geq t \right\} \leq C e^{-(pt/\sigma)^q},$$

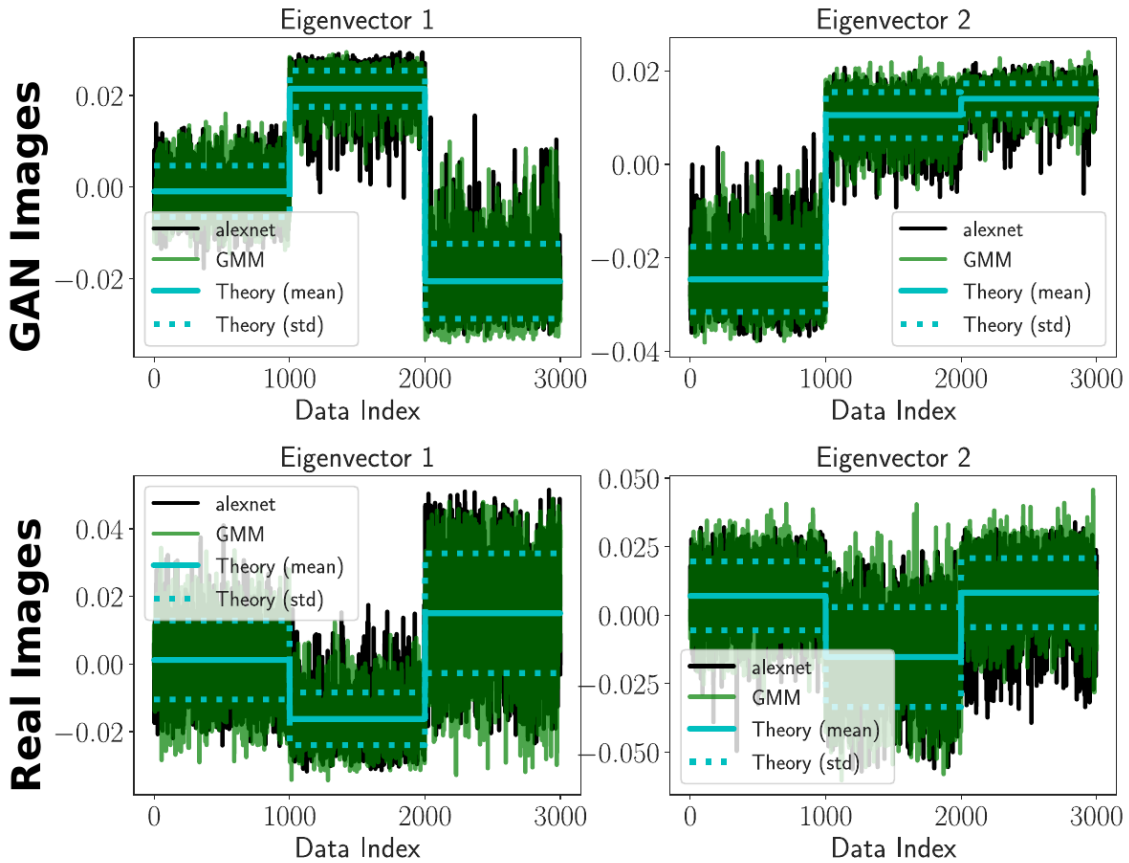


Figure 4.9: Two largest eigenvectors of the kernel matrix of Alexnet representations for GAN images (**top**) and real images (**bottom**) with the corresponding largest eigenvectors with GMM data and the theoretical predictions. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

by the union bound, we have in particular the concentration of the maximum over $i, j \in [n]$

$$\mathbb{P} \left\{ \max_{i,j} \left| \frac{2}{\sqrt{p}} (\bar{\mu}_a - \bar{\mu}_b)' (z_i - z_j) \right| \geq t \right\} \leq C n^2 e^{-(pt/\sigma)^q},$$

which implies that, for $\delta > 0$, we have with probability at least $1 - \delta$

$$\max_{i,j} \left| \frac{2}{\sqrt{p}} (\bar{\mu}_a - \bar{\mu}_b)' (z_i - z_j) \right| = \mathcal{O} \left(\frac{\log(p/\sqrt{\delta})^{1/q}}{p} \right) \quad (4.8)$$

Recalling the concentration of the dot product $z_i' z_j$ in equation 4.7, there exists two absolute constants C and σ such that

$$\mathbb{P} \{ |z_i' z_j| \geq t \} \leq C \left(e^{-(p/\sigma)^{q/2} t^{q/2}} + e^{-(p/\sigma)^{q/2} t^q} \right),$$

thus for $t \leq 1$ (only the q -concentration domains for small deviations) and by the union bound over $i, j \in [n]$

$$\mathbb{P} \left\{ \max_{i,j} |z_i' z_j| \geq t \right\} \leq 2C n^2 e^{-(p/\sigma)^{q/2} t^q},$$

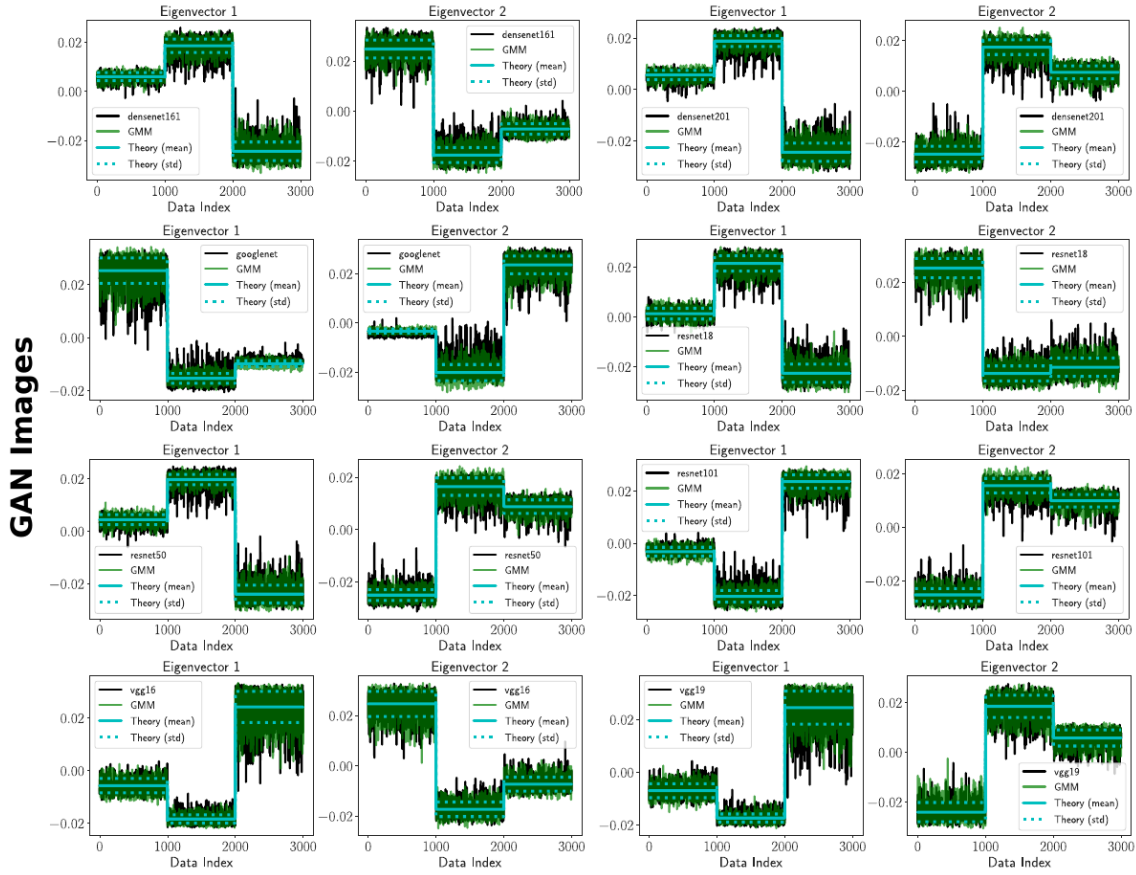


Figure 4.10: Two largest eigenvectors of the kernel matrix for CNN representations for GAN images with the corresponding largest eigenvectors with GMM data and the theoretical predictions. The kernel matrix is defined as $K_{ij} = f(\|\mathbf{x}_i - \mathbf{x}_j\|^2/p)$ where the \mathbf{x}_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

hence, for $\delta > 0$, the following holds with probability at least $1 - \delta$

$$\max_{i,j} |z'_i z_j| = \mathcal{O}\left(\frac{\log(p/\sqrt{\delta})^{1/q}}{\sqrt{p}}\right) \quad (4.9)$$

equivalently we have the concentration of the remaining stochastic term

$$\max_{i,j} |\psi_i + \psi_j| = \mathcal{O}\left(\frac{\log(p/\sqrt{\delta})^{1/q}}{\sqrt{p}}\right) \quad (4.10)$$

Recalling the considered setting, we further have

$$\frac{1}{p} \text{tr} \bar{\Sigma}_\ell = \mathcal{O}(p^{-1/2}), \quad \frac{1}{p} \|\bar{\mu}_a - \bar{\mu}_b\|^2 = \mathcal{O}(p^{-1}).$$

Combining the different orders yields the final result. \square

Random Matrix Equivalent for K . From the observation of Lemma 4.1, all the off-diagonal entries of the kernel matrix K tend to the same quantity $f(\tau)$ asymptotically.

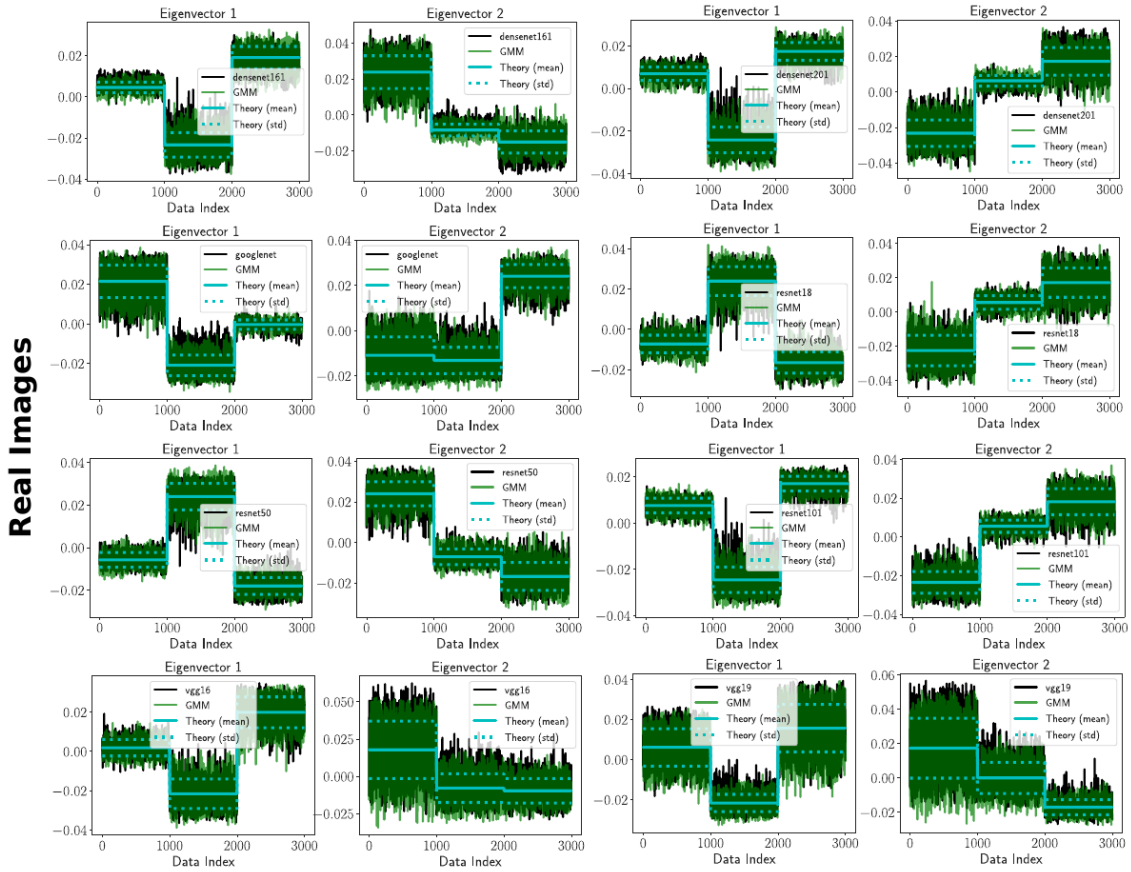


Figure 4.11: Two largest eigenvectors of the kernel matrix for CNN representations for Real images with the corresponding largest eigenvectors with GMM data and the theoretical predictions. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

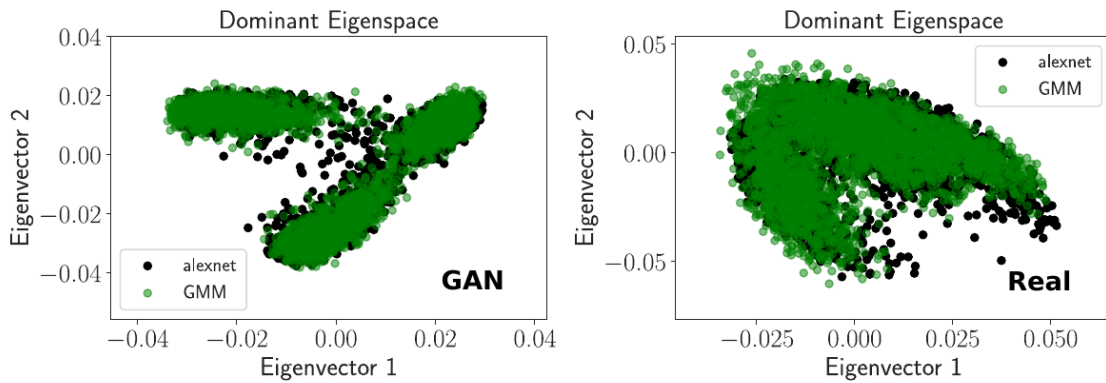


Figure 4.12: Largest eigenspace of the kernel matrix of Alexnet representations for GAN images (**left**) and real images (**right**) and the corresponding GMM data. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

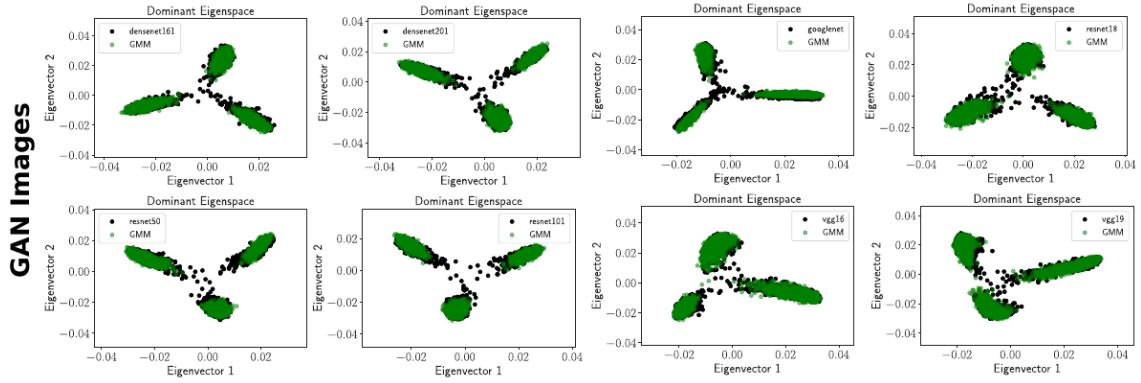


Figure 4.13: Largest eigenspace of the kernel matrix for CNN representations of GAN images and the corresponding GMM data. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

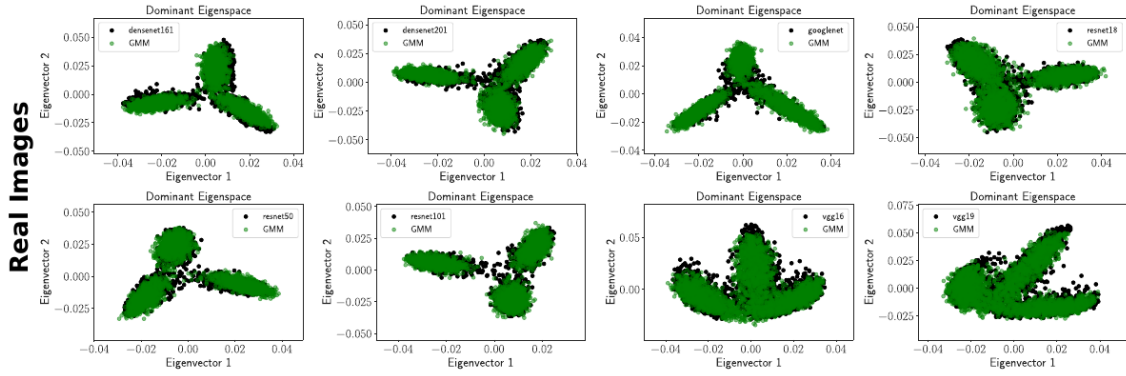


Figure 4.14: Largest eigenspace of the kernel matrix for CNN representations of real images and the corresponding GMM data. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

Therefore, K can be Taylor expanded entry-wise and we show in the following that it asymptotically has (up to centering) the same behavior as a spiked random matrix.

Before introducing this asymptotic equivalent and for subsequent use, we introduce the following quantities

$$\mathbf{M} = [\bar{\mu}_1, \dots, \bar{\mu}_k] \in \mathcal{M}_{p,k}, \quad \mathbf{t} = \left\{ \frac{\text{tr} \bar{\Sigma}_\ell}{\sqrt{p}} \right\}_{\ell=1}^k \in \mathbb{R}^k, \quad \mathbf{J} = [j_1, \dots, j_k] \in \mathcal{M}_{n,k}$$

$$\mathbf{T} = \left\{ \frac{\text{tr} \bar{\Sigma}_a \bar{\Sigma}_b}{p} \right\}_{a,b=1}^k \in \mathcal{M}_k, \quad \mathbf{Z} = [z_1, \dots, z_n] \in \mathcal{M}_{p,n}, \quad \mathbf{P} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$$

where $j_\ell \in \mathbb{R}^n$ stands for the canonical vector of the class represented by \mathcal{L}_ℓ , defined by $(j_\ell)_i = \delta_{x_i \sim \mathcal{L}_\ell}$. The vectors z_i are defined as $z_i \equiv (x_i - \bar{\mu}_\ell) / \sqrt{p}$ for each $\ell \in [k]$. We will further denote the matrix \mathbf{Z}_ℓ , the columns of which are the z_i 's in class corresponding to \mathcal{L}_ℓ .

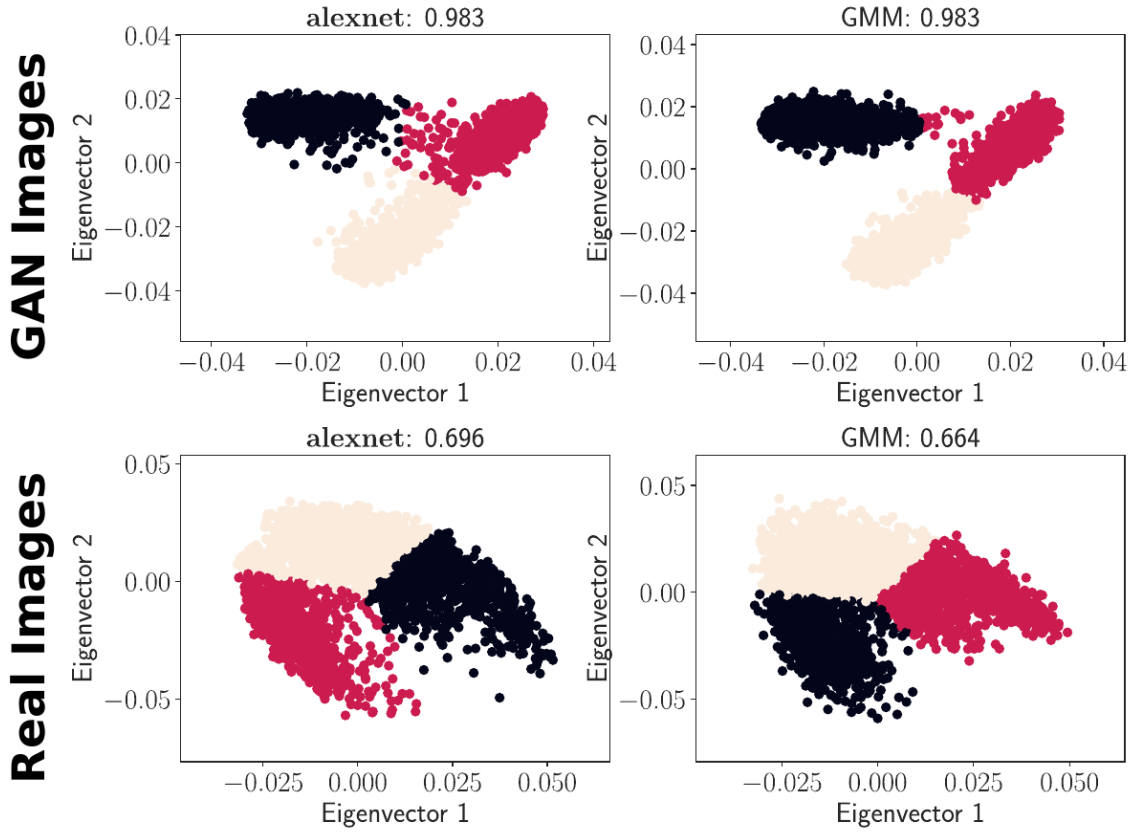


Figure 4.15: Kernel spectral clustering on Alexnet representations using GAN images (**top**) and real images (**bottom**). The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

The coming result states that there exists a random matrix \tilde{K} such that $PKP \rightsquigarrow \tilde{K}$, i.e., PKP admits as a random equivalent \tilde{K} asymptotically, where \tilde{K} has a tractable behavior from the random matrix theory standpoint.

Theorem 4.1 (Asymptotic Random Matrix Equivalent). *Let Assumptions 5, 6 and 7 hold and let \tilde{K} be defined as*

$$\begin{aligned} \tilde{K} &= -2f'(\tau) [PZ^\top ZP + UAU^\top] + F(\tau)P, \quad F(\tau) = (f(0) - f(\tau) + \tau f'(\tau)) \\ A &= \begin{bmatrix} A_{11} & I_k & -\frac{f''(\tau)}{2f'(\tau)} \mathbf{t} \\ I_k & \mathbf{0}_{k \times k} & \mathbf{0}_{k \times 1} \\ -\frac{f''(\tau)}{2f'(\tau)} \mathbf{t}^\top & \mathbf{0}_{1 \times k} & -\frac{f''(\tau)}{2f'(\tau)} \end{bmatrix}, \quad \psi_i = \|z_i\|^2 - \mathbb{E}\|z_i\|^2 = \|z_i\|^2 - \frac{1}{p} \text{tr} \Sigma_\ell \\ A_{11} &= M^\top M - \Xi - \frac{f''(\tau)}{2f'(\tau)} [\mathbf{t}\mathbf{t}^\top + 2T], \quad \Phi = Z^\top M - \{Z_\ell^\top \bar{\mu}_\ell \mathbf{1}_k^\top\}_{\ell=1}^k \\ U &= \left[\frac{J - \mathbf{1}_n \mathbf{c}^\top}{\sqrt{p}}, P\Phi, P\psi \right], \quad \Xi = \left\{ \frac{\|\bar{\mu}_a\|^2 + \|\bar{\mu}_b\|^2}{2} \right\}_{a,b=1}^k \end{aligned}$$

Then

$$PKP \rightsquigarrow \tilde{K}$$

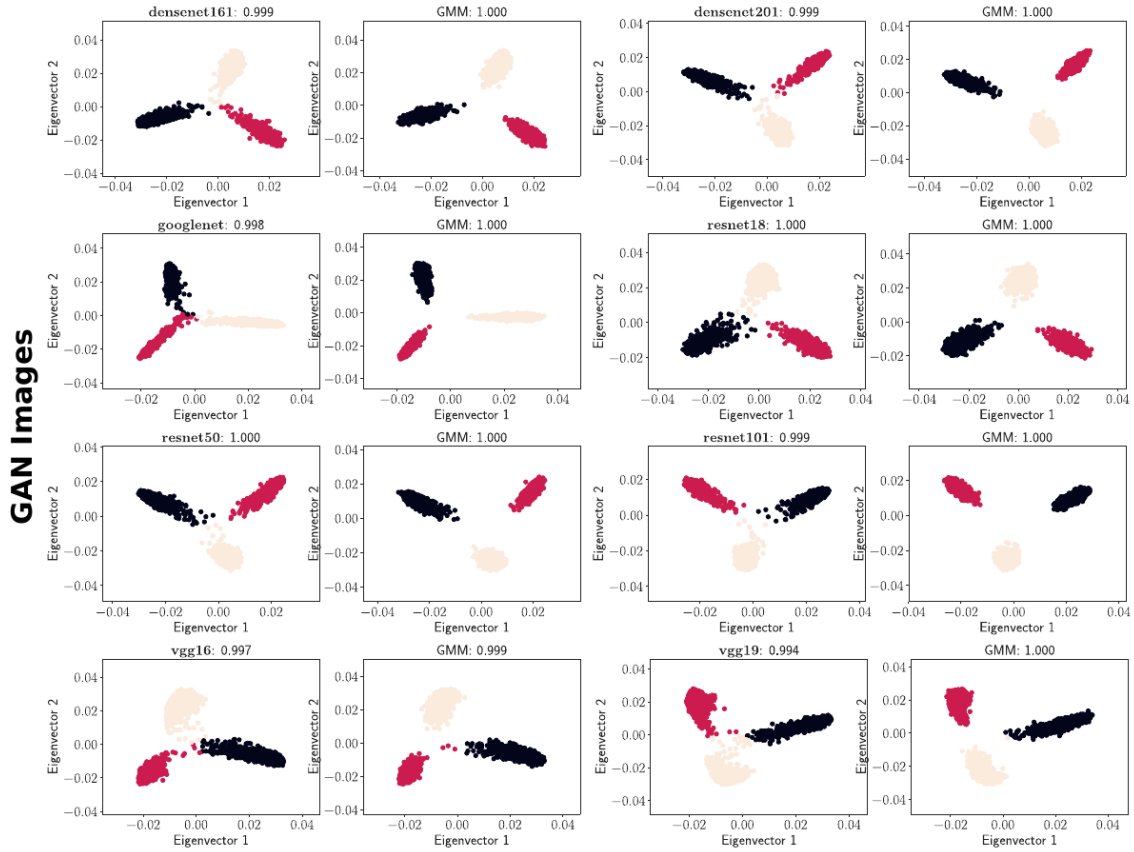


Figure 4.16: Kernel spectral clustering on CNN representations using GAN images. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

Specifically, for $\delta > 0$, there exists $C_\delta > 0$ such that for all $\gamma > 0$

$$\|PKP - \hat{K}\| \leq C_\delta p^{-1/2+\gamma} \log(p)^\gamma \text{ with probability at least } 1 - \delta.$$

Proof. See Subsection C.2.1 in the Appendix. \square

Theorem 4.1 shows that, up to centering by P , the kernel matrix K has asymptotically the same behavior as \hat{K} . In particular, the obtained approximation in operator norm implies that PKP and \hat{K} share the same eigenvalues (by Weyl's inequality [EI98, Thm 4.1]) and same *isolated* eigenvectors asymptotically. Therefore, the asymptotic spectral properties of K (i.e., the classification performance of algorithms involving K) may be studied through its equivalent \hat{K} .

Indeed, note that \hat{K} is made of a sum of a random matrix $PZ^\top ZP$ and a maximum $(k-1)$ -rank matrix containing linear combinations of the class-wise canonical vectors j_ℓ weighted by the inner-products between class means $M^\top M$ and class covariance-products and traces (through t and T). The matrix \hat{K} can then be identified as a *spiked random matrix model* [BGN12]. Note however that, unlike the standard spiked random matrices, the low-rank part of \hat{K} depends statistically on the noise part and the latter is a mixture between random matrices made of concentrated vectors. In particular, the spectrum of \hat{K} is composed of a *bulk* along with up to $k-1$ isolated eigenvalues, and the associated

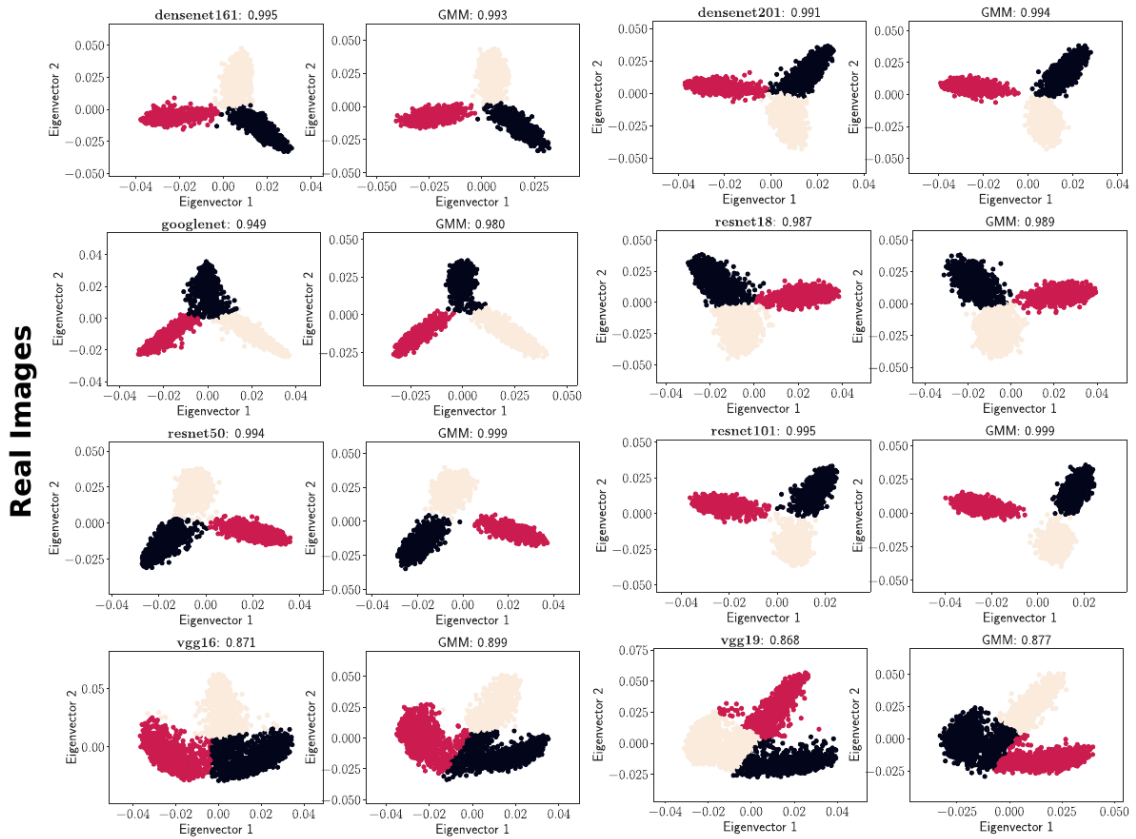


Figure 4.17: Kernel spectral clustering on CNN representations using real images. The kernel matrix is defined as $K_{ij} = f(\|x_i - x_j\|^2/p)$ where the x_i 's stand for the representations (or the corresponding GMM data) and we have chosen $f(t) = \exp(-t)$.

eigenvectors are aligned with the eigenvectors in \mathbf{U} , therefore with linear combinations of the class canonical vectors j_1, \dots, j_k . Consequently, characterizing the asymptotic performance of spectral clustering relies on the characterization of the isolated eigenvectors of \tilde{K} . In fact, these eigenvectors are *informative* if their associated eigenvalues are far away from the main eigenvalue *bulk*. The latter is due to the presence of the random Gram matrix $PZ^\top ZP$ for which the spectrum is described by Theorem 3.1 for concentrated vectors.

As for the Gram matrix, we see from Theorem 3.1 and Theorem 4.1 that the spectral behavior of the kernel matrix K depends only on the first and second order statistics of the laws \mathcal{L}_ℓ , namely their means μ_ℓ and covariances Σ_ℓ . Typically, K have the same behavior when data is described by a GMM model with the same means and covariances. The asymptotic spectral behavior of K is therefore *universal* w.r.t. the data distribution laws which fall within the introduced class of concentrated vectors. This notably explains the observations in [CBG⁺16] in which the obtained theoretical model, under GMM assumptions, fit with empirical results using the MNIST dataset [LeC98].

In the following, we provide the conditions under which the *informative* eigenvalues become visible in the spectrum of \tilde{K} which provide access to the performances of kernel spectral clustering. Having Theorem 3.1 which describes the noise term $PZ^\top ZP$, we can determine the conditions under which the spikes can be visible outside the main *bulk* of $PZ^\top ZP$, and the result concerning the isolated eigenvectors. We however need to

introduce the following resolvent from Theorem 2.9

$$\mathbf{Q}_\delta \equiv \left(\sum_{\ell=1}^k \frac{c_\ell}{k} \frac{\boldsymbol{\Sigma}_\ell}{1 + \delta_\ell(z)} + z\mathbf{I}_p \right)^{-1} \quad (4.11)$$

where $\delta_\ell(z)$ is the unique solution of the fixed point equation $\delta_\ell(z) = \frac{1}{n} \text{tr}(\boldsymbol{\Sigma}_\ell \mathbf{Q}_\delta)$. We further need the following technical assumption on the class-wise covariances to ensure that $\mathbf{PZ}^\top \mathbf{ZP}$ does not produce *non-informative* isolated eigenvalues.

Assumption 8 (Spikes control). Denote $\lambda_1^\ell, \dots, \lambda_p^\ell$ the eigenvalues of $\boldsymbol{\Sigma}_\ell$, for each $\ell \in [k]$. As $n \rightarrow \infty$, $\frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i^\ell} \xrightarrow{\mathcal{D}} \rho_\ell$ with support \mathcal{S}_ℓ , and $\max_{1 \leq i \leq p} \text{dist}(\lambda_i^\ell, \mathcal{S}_\ell) \rightarrow 0$.

Now we can state the theorem that ensures the presence of *informative* eigenvalues in the spectrum of $\tilde{\mathbf{K}}$, and gives the characterization of the corresponding isolated eigenvectors, which results from standard random matrix techniques [BGN12].

Theorem 4.2 (Isolated eigenvalues). Suppose that Assumptions 5-6-8 hold and let $\boldsymbol{\Lambda}_z$ be the matrix defined for $z \in \mathbb{C}_+$ as

$$\boldsymbol{\Lambda}_z = \alpha_\tau(z) \mathbf{I}_k - \left[\alpha_\tau(z) \left(z\mathbf{M}^\top \mathbf{Q}_\delta \mathbf{M} + \Xi + \frac{f''(\tau)}{f'(\tau)} \mathbf{T} \right) + \frac{f''(\tau)}{2f'(\tau)} \mathbf{t}\mathbf{t}^\top \right] \Gamma_z$$

where

$$\begin{aligned} \alpha_\tau(z) &\equiv 1 + \frac{f''(\tau)}{2zc_0pf'(\tau)} \sum_{\ell=1}^k \eta_\ell(z) \text{tr} \boldsymbol{\Sigma}_\ell^2 \\ \Gamma_z &\equiv \frac{1}{z} \left\{ \frac{c_a \eta_a(z) c_b \eta_b(z)}{\sum_{\ell=1}^k c_\ell \eta_\ell(z)} \right\}_{a,b=1}^k - \frac{1}{z} \mathcal{D} \{ c_\ell \eta_\ell(z) \}_{\ell=1}^k + \frac{2}{z} (\mathbf{I}_k - \mathcal{D}(\eta))^\top \\ \eta_\ell(z) &\equiv \frac{1}{1 + \delta_\ell(z)} \end{aligned}$$

Let λ^* be at an infinitesimal distance from the bulk support \mathcal{S} of $\mathbf{PZ}^\top \mathbf{ZP}$, such that $\alpha_\tau(\lambda^*) \neq 0$ and $\boldsymbol{\Lambda}_{\lambda^*}$ has a zero eigenvalue of multiplicity m^* . Then \mathbf{PKP} produces m^* spikes asymptotically close to

$$\rho^* \equiv -2f'(\tau)\lambda^* + F(\tau)$$

Furthermore, the eigenspace projector corresponding to the (asymptotically converging to ρ^*) isolated eigenvalues of \mathbf{PKP} has a non-vanishing projection onto $\text{span}(\mathbf{j}_1, \dots, \mathbf{j}_k)$.

Theorem 4.2 gives the conditions under which the spikes can be observed in the spectrum of $\tilde{\mathbf{K}}$, and states that the corresponding eigenvectors are aligned to some extent to the class canonical vectors $\mathbf{j}_1, \dots, \mathbf{j}_k$, which is important for spectral clustering. We refer the reader to [CBG⁺16] for the performances characterization of kernel spectral clustering in the k -class Gaussian mixture model case as the behavior of \mathbf{K} is universal.

4.1.2.2 Application to GAN-generated Images

The objective of this subsection is to present simulations to validate our findings concerning the behavior of large kernel matrices. To highlight this aspect, we evaluate the kernel matrix on $\mathbf{x}_1, \dots, \mathbf{x}_n$ being CNN representations of GAN-generated images and

we further use real images for comparison. Specifically, we consider the same setting as the previous chapter, i.e., GAN images are generated by the BigGAN model [BDS18] and then represented using several representation models pretrained on the Imagenet dataset [DDS⁺09]. Moreover, real images are samples from the Imagenet dataset [DDS⁺09].

GAN architecture & real images: We consider as a setting of our experiments the BigGAN model [BDS18] which takes as input a Gaussian noise vector of dimension $d = 140$ and a one-hot-vector to generate a specific class image of dimension $256 \times 256 \times 3$. We particularly consider $k = 3$ classes which are $\{pizza, mushroom, hamburger\}$. Examples of the generated images are shown in Figure 3.4. Real images are sampled from the same classes from the Imagenet dataset [DDS⁺09]. In both cases, we use $n = 3000$ images (i.e., $n_\ell = 1000$ for each class \mathcal{C}_ℓ) in all our experiments.

CNN representation: We consider popular CNN representations pretrained on the Imagenet dataset [DDS⁺09]. Specifically, we consider 9 representation networks which are: *vgg16* ($p = 4096$), *vgg19* ($p = 4096$), *resnet18* ($p = 512$), *resnet50* ($p = 2048$), *resnet101* ($p = 2048$), *densenet161* ($p = 2208$), *densenet201* ($p = 1920$), *alexnet* ($p = 4096$) and *googlenet* ($p = 2048$).

GMM data: For both GAN and real data (i.e., CNN representations of GAN and real images respectively) we build the corresponding GMM data of class \mathcal{C}_ℓ as $\mathbf{x}'_i = \hat{\boldsymbol{\mu}}_\ell + \hat{\boldsymbol{\Sigma}}_\ell^{1/2} \mathbf{z}_i$ with $\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ where $\hat{\boldsymbol{\mu}}_\ell$ and $\hat{\boldsymbol{\Sigma}}_\ell$ are respectively the empirical estimates of $\boldsymbol{\mu}_\ell$ and $\boldsymbol{\Sigma}_\ell$.

The key point in the analysis of the kernel matrix \mathbf{K} is the behavior of the quantity $\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2$. In Figure 4.1, we have depicted the histogram of this quantity for both GAN and real images using the *alexnet* ($p = 4096$) representation. As we can notice, in both cases the histogram of $\frac{1}{p} \|\mathbf{x}_i - \mathbf{x}_j\|^2$ matches its Gaussian counterpart $\frac{1}{p} \|\mathbf{x}'_i - \mathbf{x}'_j\|^2$. Moreover, for real data, this quantity seems to converge to the same value as per Lemma 4.1. We do not observe the same behavior for GAN data which may be the consequence of the fact that GANs have low entropy, i.e., GAN data are easily separable compared to real data. Figure 4.2 depicts the same histogram using other representations from which we obtain the same conclusions.

As a consequence, the entries of the kernel matrix \mathbf{K} converge to the same value, as we can visually observe from Figure 4.3 (at least for real images). Still we observe the same behavior of the kernel matrix for GAN data/real data versus their GMM counterparts. Figure 4.4 shows more examples using the different CNN representation for GAN images, and Figure 4.5 depicts other representations of real images (where we observe that the entries \mathbf{K} converge to the same quantity as p grows).

More interestingly, we see from Figure 4.6 that the spectrum of \mathbf{K} when computed on the representations of both GAN and real images, matches perfectly the spectrum of its GMM counterpart. A result which stays valid across all the representation networks for GAN images in Figure 4.7 as well as real images 4.8. As proved by Theorem 4.1, the kernel matrix \mathbf{K} behaves as spiked random matrix model which can be also observed from these figures. In particular, the matching between the representations themselves and the corresponding GAN data demonstrates the *universal behavior* of the kernel matrix \mathbf{K} . This universal behavior goes beyond the spectrum of \mathbf{K} and is also observed through

the largest eigenvectors of K which contain information about the classes as per Theorem 4.2. In particular, we can see from Figure 4.9 that the largest eigenvectors computed on the *alexnet* representations are perfectly aligned with their GMM counterparts, and also for the other representations using GAN images (see Figure 4.10) and real images (see Figure 4.11). The corresponding eigenspaces are also depicted in Figure 4.12 for the *alexnet* representation network and the other representation networks in Figure 4.13 and Figure 4.14.

Since kernel spectral clustering consists in applying the k-means algorithm on the largest subspace of the kernel matrix, we depicted in Figure 4.15 the result of k-means for *alexnet* representations and the corresponding GMM data along with the obtained accuracies. As we can see, the behavior of k-means on the representations is almost the same as on the corresponding GMM data which confirms the *universal aspect* of K . This observation holds for the remaining considered representation networks for both GAN images in Figure 4.16 and real images 4.17. Note that the accuracies obtained with GAN data are almost equal to 100% which is a consequence of the fact that GAN data are easily separable, while this is not true for real data as depicted in Figure 4.17.

4.1.3 Central Contribution and perspectives

We have studied in this work large kernel matrices for a wide class of random inputs, *i.e.*, concentrated data, which are more generic than Gaussian mixtures as we have shown in the first part of this manuscript. Our study has notably shown the universality aspect of kernel spectral clustering (based on the kernel matrix K) by highlighting the fact that the asymptotic behavior of kernel matrices, for concentrated inputs, depend on the first and second moments statistics, thereby match the behavior of a GMM model as empirically observed in [CBG⁺16]. Moreover, since the performance of spectral clustering is predictable for a GMM model [CBG⁺16], it is thereby predictable for concentrated data as well and thus for complex data as deep learning representations of the surprising realistic images generated by GANs. In the next section, we will provide the analysis of a different form of kernel matrices which fundamentally differ from this section by the fact that their entries concentrate around different values. Thereby, making the conclusions considerably different, we further apply our result to the problem of Sparse PCA which yield to an effective PCs recovery method.

4.2 Sparse Principal Component Analysis

This section is based on the following work:

- (C3) MEA. Seddik, M. Tamaazousti, R. Couillet, “A Kernel Random Matrix-Based Approach for Sparse PCA”, International Conference on Learning Representations (ICLR’19), New Orleans, United-States, 2019.

This section presents the analysis of kernel matrices of the form $f(\hat{\Sigma})$ where $\hat{\Sigma}$ is the sample covariance matrix. The behavior of these matrices is fundamentally different from those of the previous section, indeed, the entries of $f(\hat{\Sigma})$ concentrate around different values which makes the analysis a bit different. In particular, we apply our analysis to the Sparse PCA problem which results in competitive method to the existing methods in the state-of-the-art.

4.2.1 Motivation

Principal component analysis (PCA) is extensively used in data analysis and machine learning applications. It is a dimension reduction technique that aims to project a given dataset onto principal subspaces spanned by the leading eigenvectors of the sample covariance matrix [WEG87], which represent the principal modes of variance. Basically, the statistical interpretation of PCA lies in the fact that most of the variance in the data is captured by these modes. Consequently, PCA reduces the dimension of the feature space while keeping most of the information in the data. It is well-known [And63] that PCA performs efficiently in the traditional data setting where the number of features is small and the number of samples is large.

Consider a data matrix $\mathbf{X} \in \mathcal{M}_{p,n}$ consisting of n centered samples, each sample having p features. The standard PCA method requires the computation of the sample covariance matrix $\hat{\Sigma} = \mathbf{X}\mathbf{X}^\top/n$ and estimates the first principal components $\mathbf{u}_1, \mathbf{u}_2, \dots$ (i.e., the successive dominant eigenvectors of $\Sigma = \mathbb{E}[\mathbf{X}\mathbf{X}^\top/n]$) by the ordered eigenvectors $\hat{\mathbf{u}}_1, \hat{\mathbf{u}}_2, \dots$ of $\hat{\Sigma}$. Authors in [JL09] demonstrated that, in the high dimensional regime where $n, p \rightarrow \infty$ with $p/n \rightarrow c > 0$, the principal component $\hat{\mathbf{u}}_1$ estimated by standard PCA is inconsistent. Essentially, if $p/n \rightarrow 0$ then $\|\hat{\mathbf{u}}_1 - \mathbf{u}_1\|_2 \rightarrow 0$ in the high-dimensional asymptotic regime (see Subsection 2.3.1 for more details). As discussed in the background Chapter 2 of this manuscript, this phenomenon is well investigated within the field of random matrix theory for covariance models of the form $\hat{\Sigma} = \frac{1}{n}\Sigma^{1/2}\mathbf{Z}\mathbf{Z}^\top\Sigma^{1/2}$, where Σ is a positive semi-definite matrix and \mathbf{Z} is a $p \times n$ matrix with random *i.i.d.* entries. Essentially relying on the so-called spiked models (see Subsection 2.3) of random matrix theory, where Σ is a low-rank perturbation of the identity matrix, namely $\Sigma = \mathbf{I}_p + \sum_{i=1}^k \omega_i \mathbf{u}_i \mathbf{u}_i^\top$ with k fixed with respect to p, n . Authors in [BAP⁺05] and [Pau07] notably exhibited a phase transition phenomenon: as $p/n \rightarrow c$, if $\omega_i < \sqrt{c}$ the estimated principal component $\hat{\mathbf{u}}_i$ using standard PCA is (almost surely) asymptotically orthogonal to the true principal component \mathbf{u}_i (i.e., $\hat{\mathbf{u}}_i^\top \mathbf{u}_i \rightarrow 0$); on the other hand, if $\omega_i > \sqrt{c}$, $\liminf_n |\hat{\mathbf{u}}_i^\top \mathbf{u}_i| > 0$. This phase transition phenomenon has attracted recently much attention within the random matrix community [BGN11, CDMF09, FP07, KY13].

The inconsistency of standard PCA in high dimensions motivated the idea to look for more structural information on the principal components. In particular, considering that the principal components are sparse in an appropriate basis (e.g., in the wavelet domain), a large body of works have emerged and proposed improved PCA approaches that account for sparsity. One of the most consistent sparse PCA methods in the literature is the covariance thresholding (CT) algorithm [KNV15]. Based on the intuition that the small entries of the empirical covariance matrix $\hat{\Sigma}$ induce noise in its principal components, this method consists in applying the popular *soft-thresholding* function (with threshold $\tau > 0$);

$$\text{soft}(\cdot; \tau) : t \mapsto \text{sign}(t) \cdot (|t| - \tau)_+, \quad (4.12)$$

entry-wise to the empirical covariance matrix $\hat{\Sigma}$ and performing PCA on the resulting matrix. Authors in [DM14, DM16] have theoretically demonstrated that the covariance thresholding algorithm recovers the sought-for principal components with high probability under controlled growth rates between p, n and the sparsity level. In this work, we particularly show that the soft-thresholding method in fact falls within a broader class of kernel-based¹ PCA algorithms that are particularly suited to sparse PCA recov-

¹Note that *kernel-based* terminology is used to highlight that our work falls within the framework of

ery. This method consists in considering the matrix $f(\hat{\Sigma})$ instead of $\hat{\Sigma}$ where f is a function applied entry-wise to Σ . By imposing some constraints on f , most importantly that $f'(0) = f''(0) = 0$, we show that sparse PCA can be performed with provably high accuracy for sufficiently large n, p .

There has been a wide range of approaches to tackle the sparse PCA problem. Mainly, three classes of approaches emerge in the literature. Most popular techniques are optimization based algorithms [dGJL05, MWA06, ZS07, ZHT06, WGR⁺09], where the idea is to see the problem of sparse PCA through an optimization perspective, and to propose methods to solve the latter by either considering a different formulation – e.g. semi-definite programming (SDP) or convex relaxations – or adding penalties to the original optimization problem such as a LASSO regularization. The second class of approaches covers matrix decomposition-based techniques [APD14, PDK13, SH08], where sparse principal components are extracted through solving a low rank matrix approximation problem based on Singular Value Decomposition. Finally, most consistent sparse PCA methods adopt thresholding-based approaches: initial heuristics used factor rotation techniques and thresholding of eigenvectors to obtain sparsity [CJ95]. Based on the well-known power method, [YZ13] introduced an efficient sparse PCA approximation to obtain the exact level of required sparsity, by truncating to zero the principal components iteratively except for their largest entries. A step further, under a spiked covariance model (see Section 4.2.2), [M⁺13] proposed a very efficient iterative thresholding approach for estimating principal subspaces in the sparse setting. Similarly, assuming a single-spike model, [KNV15] proved that, when the sparsity level $s \geq \Omega(\sqrt{n})$, a standard SDP approach cannot recover consistently the sparse spike; in particular, the authors presented empirical results suggesting that for $s = \mathcal{O}(\sqrt{n})$, recovery is possible by a simple covariance thresholding algorithm. More recently, [DM16] analyzed and theoretically proved, under a spiked model, that indeed the covariance thresholding algorithm [KNV15] succeeds with high probability under controlled growth rates between p, n and s .

In this work, while restricting ourselves to a setting where p and n grow at a controlled joint rate, we provide an elementary argument, based on a matrix-wise Taylor expansion controlled through a concentration of measure approach, that generalizes the CT method to a large family of kernel-based methods, by means of a kernel random matrix approach [EK⁺10b, EK10a]. Concretely, we study kernel random matrices of the form $f(\mathbf{X}\mathbf{X}^\top/n)$ where $\mathbf{X} = \Sigma^{1/2}\mathbf{Z}$ and \mathbf{Z} is a random matrix with concentrated columns. [EK⁺10b] studied kernel matrices of the form $f(\mathbf{X}^\top\mathbf{X}/n)$ (i.e., the so-called inner-product kernel matrices), which is equivalent to the case $\Sigma = \mathbf{I}_p$ when considering the form $f(\mathbf{X}\mathbf{X}^\top/n)$. In particular, we elaborate from El Karoui’s study by Taylor expanding $f(\mathbf{X}\mathbf{X}^\top/n)$ in the vicinity of Σ entry-wise and controlling the resulting matrices via concentration arguments.

4.2.2 Model and Main Results

Consider a data matrix $\mathbf{X} \in \mathcal{M}_{p,n}$ defined as

$$\mathbf{X} \equiv \Sigma^{1/2}\mathbf{Z} = (\mathbf{I}_p + \mathbf{P})^{1/2}\mathbf{Z}, \quad (4.13)$$

kernel random matrices and should not be confused with the standard kernel PCA.

where $\mathbf{Z} \in \mathcal{M}_{p,n}$ some random matrix, $\mathbf{P} = \sum_{i=1}^k \omega_i \mathbf{u}_i \mathbf{u}_i^\top$ and $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \mathcal{M}_{p,k}$ is isometric. Here, k refers to the number of principal components (or eigenvectors) $\mathbf{u}_1, \dots, \mathbf{u}_k \in \mathbb{R}^p$ to be evaluated, with $\omega_1 > \dots > \omega_k > 0$ the corresponding eigenvalues respectively. We define the quantity $\beta_p \equiv \max_i \|[\boldsymbol{\Sigma}^{1/2}]_{\cdot,i}\|$.

We consider the following concentration and growth rate assumptions.

Assumption 9 (Concentrated data). $\mathbf{Z} \propto \mathcal{E}_2$ with $\mathbb{E}\mathbf{Z} = \mathbf{0}$ and $\mathbb{E}[\mathbf{Z}\mathbf{Z}^\top/n] = \mathbf{I}_p$.

Assumption 10 (Growth rate). There exists $B > 0$ independent of p, n such that $\max_{ij} |[\boldsymbol{\Sigma}_p]_{ij}| < B$. Besides, there exists $\epsilon > 0$ such that $\beta_p \leq B' n^{\frac{1}{4}-\epsilon}$ for all p, n and for some absolute constant $B' > 0$.

Under these assumptions, our main technical result is as follows:

4.2.2.1 Random Matrix Equivalent

Theorem 4.3 (Asymptotic Equivalent). For f a three-times continuously differentiable function, define the matrices \mathbf{F} and $\tilde{\mathbf{F}}$ respectively by²

$$\mathbf{F} \equiv \left\{ f \left(\left[\frac{1}{n} \mathbf{X}\mathbf{X}^\top \right]_{ij} \right) \right\}_{i,j=1}^p, \quad \tilde{\mathbf{F}} \equiv f(\boldsymbol{\Sigma}) + \sum_{k=1}^2 \frac{f^{(k)}(\boldsymbol{\Sigma})}{k!} \odot \left[\boldsymbol{\Sigma}^{1/2} \left(\frac{1}{n} \mathbf{Z}\mathbf{Z}^\top - \mathbf{I}_p \right) \boldsymbol{\Sigma}^{1/2} \right]^{\odot k}.$$

Then

$$\mathbf{F} \rightsquigarrow \tilde{\mathbf{F}}$$

Specifically, for $\eta > 0$, there exists an absolute constant $C_\eta > 0$ such that with probability at least $1 - \eta$

$$\|\mathbf{F} - \tilde{\mathbf{F}}\| \leq C_\eta \frac{\beta_p^6 p}{n^{3/2} \sqrt{\eta}}. \quad (4.14)$$

For a general smooth function f , the kernel random matrix $f(\hat{\boldsymbol{\Sigma}})$ is particularly difficult to analyze through the usual tools of random matrix theory, such as the moment or Stieltjes transform-based methods [Tao12]. Rather than directly analyzing such a kernel random matrix, Theorem 4.3 gives an asymptotic *random* equivalent to it, in operator norm, that has mainly two properties. First, the approximation matrix $\tilde{\mathbf{F}}$ contains “simple” objects that have already been analyzed in random matrix theory – in particular, the term $(\mathbf{Z}\mathbf{Z}^\top/n - \mathbf{I}_p)$ in the expression of $\tilde{\mathbf{F}}$. Second, the approximation in operator norm implies (by Weyl’s inequality [EI98, Theorem 4.1]) that, when $\|\mathbf{F} - \tilde{\mathbf{F}}\|_{op} \rightarrow 0$, \mathbf{F} and $\tilde{\mathbf{F}}$ have the same eigenvalues and same “isolated” eigenvectors asymptotically (see Corollary 4.2 subsequently).

Before sketching a proof for Theorem 4.3, we need the following technical Lemma which establishes the concentration of the entries of the noise term $(\mathbf{Z}\mathbf{Z}^\top/n - \mathbf{I}_p)$.

Lemma 4.2 (A Concentration Result). For all $i, j \in [p]$, the bilinear form

$$g_{ij}(\mathbf{Z}) \equiv [\boldsymbol{\Sigma}^{1/2}]_{i,\cdot} \left(\frac{1}{n} \mathbf{Z}\mathbf{Z}^\top - \mathbf{I}_p \right) [\boldsymbol{\Sigma}^{1/2}]_{\cdot,j} \quad (4.15)$$

² f and $f^{(k)}$ are applied entry-wise and $\odot k$ stands for the element-wise k -th power.

satisfies

$$g_{ij}(\mathbf{Z}) \propto \mathcal{E}_1 \left(\frac{\beta_p^2}{n} \right) + \mathcal{E}_2 \left(\frac{\beta_p^2}{\sqrt{n}} \right). \quad (4.16)$$

Proof. Denoting by \mathbf{v}_i the i -th column vector of the matrix $\mathbf{\Sigma}^{1/2}$, we have by the polarization identity, for all \mathbf{M} Hermitian,

$$\mathbf{v}_i^\top \mathbf{M} \mathbf{v}_j = \frac{1}{4} [(\mathbf{v}_i + \mathbf{v}_j)^\top \mathbf{M} (\mathbf{v}_i + \mathbf{v}_j) - (\mathbf{v}_i - \mathbf{v}_j)^\top \mathbf{M} (\mathbf{v}_i - \mathbf{v}_j)].$$

It thus suffices to prove the result for the quadratic form

$$g(\mathbf{Z}) = \mathbf{v}^\top \left(\frac{1}{n} \mathbf{Z} \mathbf{Z}^\top - \mathbf{I}_p \right) \mathbf{v}$$

where $\mathbf{v} \in \mathbb{R}^p$. Noticing that $\mathbf{v}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{v} = \|\mathbf{v}^\top \mathbf{Z}\|^2$ and $\mathbb{E} \left[\frac{1}{n} \mathbf{v}^\top \mathbf{Z} \mathbf{Z}^\top \mathbf{v} \right] = \mathbf{v}^\top \mathbf{v}$, we need to prove the concentration of the random variable $\|\mathbf{v}^\top \mathbf{Z}\|^2$. In fact, since $\mathbf{v}^\top \mathbf{Z}$ is a Concentrated vector, by Definition, $\|\mathbf{v}^\top \mathbf{Z}\| \propto \mathcal{E}_2(\|\mathbf{v}\|)$ since $M \mapsto \mathbf{v}^\top M$ and $u \mapsto \|u\|$ are respectively $\|\mathbf{v}\|$ -Lipschitz and 1-Lipschitz functions. We get the final result by Proposition 2.5. \square

Now we turn to the main ingredients to prove Theorem 4.3.

Sketch of Proof of Theorem 4.3. The main idea of the proof relies on the following intuition: for large n , the entries of $\mathbf{Z} \mathbf{Z}^\top / n - \mathbf{I}_p$ and its successive Hadamard products tend to zero at controllable rate. The concentration of measure framework then allows for the control of non-linear functions of the entries of $\mathbf{Z} \mathbf{Z}^\top / n - \mathbf{I}_p$ thanks to the previous Lemma. A Taylor expansion of F around $f(\mathbf{\Sigma})$ then leads to controlling the operator norm of

$$f^{(3)}(\zeta^n) \odot [\mathbf{\Sigma}_p^{1/2} (\mathbf{Z} \mathbf{Z}^\top / n - \mathbf{I}_p) \mathbf{\Sigma}^{1/2}]^{\odot 3}$$

for ζ^n a matrix with entries in the set $[[\mathbf{X} \mathbf{X}^\top / n]_{ij}, [\mathbf{\Sigma}]_{ij}]$ (or $[[\mathbf{\Sigma}]_{ij}, [\mathbf{X} \mathbf{X}^\top / n]_{ij}]$). This follows precisely from exploiting Lemma 4.2 twice, to control the fluctuations of the entries of both ζ^n (by the conditions on $\max_{ij} |[\mathbf{\Sigma}]_{ij}|$ and β_p) and $[\mathbf{\Sigma}^{1/2} (\mathbf{Z} \mathbf{Z}^\top / n - \mathbf{I}_p) \mathbf{\Sigma}^{1/2}]^{\odot 3}$, with the bound provided in the theorem statement, thereby completing the proof. A detailed proof of Theorem 4.3 is provided in the Appendix. \square

In order to simplify our arguments concerning the analysis of Sparse PCA, we make the additional assumption that both p and n are large and comparable. We further suppose that the deterministic quantities ω_i 's and β_p are bounded as $n \rightarrow \infty$. Specifically, we make the following set of assumptions

Assumption 11 (RMT Growth rate). *As $n \rightarrow \infty$,*

A1 $p/n \rightarrow c \in (0, \infty)$,

A2 $\limsup_n \max_i \omega_i < \infty$; specifically $\limsup_n \beta_p < \infty$.

Under this setting, we have the following important corollary to Theorem 4.3.

Corollary 4.1. *Define the matrices F and \tilde{F} as in Theorem 4.3 and let Assumptions A1 and A2 hold. Then, for $\eta > 0$*

$$F = \tilde{F} + \mathcal{O}_\eta(n^{-\frac{1}{2}}), \quad (4.17)$$

where $X = \mathcal{O}_\eta^m(n^{-\alpha})$ stands for the fact that $\mathbb{P} \left\{ \|X\| \geq C n^{-\alpha} \eta^{-\frac{1}{2m}} \right\} \leq \eta$ for some absolute constant $C > 0$ and non-negative integer m .

As a consequence of Corollary 4.1, we have, by the $\sin(\Theta)$ theorem of [DK70], the corollary below concerning the eigenvectors of the matrices F and \tilde{F} .

Corollary 4.2. *Let v_1, \dots, v_k and $\tilde{v}_1, \dots, \tilde{v}_k$ denote respectively the k principal eigenvectors of F and \tilde{F} . Denote by $\Delta_i = \omega_i - \omega_{i+1}$ for $i \in [k-1]$. Then for $\eta > 0$, we have*

$$\max_{i \in [k]} \min_{y \in \{+1, -1\}} \Delta_i^2 \|v_i - y \tilde{v}_i\|^2 = \mathcal{O}_\eta(n^{-1}). \quad (4.18)$$

4.2.2.2 Application to sparse PCA

Before presenting our results concerning the behavior of kernel matrices of the form $f(\hat{\Sigma})$ in the context of sparse PCA, we introduce in following a particular notion of sparsity for large matrices that will be used subsequently.

When considering a large-dimensional random matrix setting, the notion of sparsity for such matrices is particularly attached to the choice of the matrix norm.³ [EK08] introduced a definition (ε -sparsity) for sparsity of matrices that is compatible with spectral analysis, and specifically adapted to the operator norm. The ε -sparsity definition requires some notions from graph theory that we present in the following:

- To each $p \times p$ symmetric matrix M , we define its corresponding adjacency matrix as $\mathcal{A}(M) = \{1_{M_{ij} \neq 0}\}_{i,j=1}^p$, which corresponds to a graph \mathcal{G}_p with p vertices.
- A walk is said to be closed on this graph if it starts and finishes at the same vertex and the number of edges traversed by a walk defines the length of this walk. Denote $\mathcal{C}_p(k)$ the set of closed walks of length k on \mathcal{G}_p .

We these notions, we are now in place to introduce the definition of the ε -sparse notion.

Definition 8 (ε -sparse matrices [EK08, Definition 1]). *A sequence of covariance matrices $\{\Sigma_p\}_{p=1}^\infty$ is said to be ε -sparse if the sequence of their associated graphs $\{\mathcal{G}_p\}_{p=1}^\infty$ satisfies, for all $k \in 2\mathbb{N}$,*

$$|\mathcal{C}_p(k)| \leq C_k p^{\varepsilon(k-1)+1}$$

where $\varepsilon \in [0, 1]$, $C_k > 0$ independent of p and $|\mathcal{S}|$ denotes the cardinality of the set \mathcal{S} .

The ε -sparsity is both useful and convenient to our study for the following reasons: 1) it is adapted to the analysis of the operator norm of large sparse matrices (as we give concentration results on the operator norm); 2) it is also more general than other sparsity notions such as in [BL08]. In the latter, the authors developed a natural permutation-invariant notion of sparsity which is more specific than Definition 8 as pointed out in the introduction of their article. Furthermore, note that both sparsity notions (Definition 8 and the one in [BL08]) provide equivalent bounds for $\varepsilon < \frac{1}{2}$ and when considering the large dimensional $p \sim n$ setting (see subsection 2.4 in [BL08]); this is precisely the setting considered in Theorem 4.4 introduced subsequently (cf. $\mu > 0$).

³Considering the identity matrix (which is a sparse matrix), $\|I_p\| = 1$ while $\|I_p\|_F = \sqrt{p} \rightarrow \infty$.

Remark 4.1. As Definition 8 is based on a graph defined by its corresponding adjacency matrix, we have the following property: given an ε -sparse matrix \mathbf{M} and a function f such that $f(0) = 0$ and $f(x) \neq_{x \neq 0} 0$, the matrix $f(\mathbf{M})$, resulting from the application of f entry-wise to \mathbf{M} , remains ε -sparse; this is simply a consequence of $\mathcal{A}(\mathbf{M}) = \mathcal{A}(f(\mathbf{M}))$.

To get an insight on our coming results, consider the scenario where \mathbf{U} contains finitely many non-zero entries. In this case, the perturbation matrix \mathbf{P} in equation 4.13 contains finitely many non-zero entries (say s) on each line and a simple enumeration shows that $|\mathcal{C}_p(k)| \leq p s^{k-1}$, thus \mathbf{P} is 0-sparse in the sense of Definition 8. Similarly, \mathbf{I}_p is 0-sparse and by the additive stability⁴ of the ε -sparsity notion, $\mathbf{\Sigma}$ remains 0-sparse. More generally, if we assume that it exists $\varepsilon \in [0, \frac{1}{2})$ such that the population covariance matrix $\mathbf{\Sigma}$ is ε -sparse, we have the following set of consequences. By Corollary 4.1, choosing f in such a way that $f'(0) = f''(0) = 0$ ensures that the terms $f'(\mathbf{\Sigma}) \odot \dots$ and $f''(\mathbf{\Sigma}) \odot \dots$ vanish in the expression of $\tilde{\mathbf{F}}$. Indeed, for $k \in \{1, 2\}$

- (i) Only finitely entries of $f^{(k)}(\mathbf{\Sigma})$ do not vanish, precisely by Remark 4.1, since

$$\mathcal{A}(f^{(k)}(\mathbf{\Sigma})) = \mathcal{A}(\mathbf{\Sigma}),$$

the matrix⁵ $f^{(k)}(\mathbf{\Sigma})$ is also (almost) ε -sparse.

- (ii) The matrix $\mathbf{F}^{(k)} = [\mathbf{\Sigma}^{1/2} (\frac{1}{n} \mathbf{Z} \mathbf{Z}^\top - \mathbf{I}_p) \mathbf{\Sigma}^{1/2}]^{\odot k}$ has entries of order $\mathcal{O}(n^{-k/2})$. As a result,⁶ we have for $\eta > 0$ and for all $m > 0$, $\max_{i,j} |F_{ij}^{(k)}| = \mathcal{O}_\eta^m \left(n^{-\frac{k}{2} + \frac{1}{m}} \right)$.

Since in addition the operator norm of $\mathbf{\Sigma}^{1/2} (\mathbf{Z} \mathbf{Z}^\top / n - \mathbf{I}_p) \mathbf{\Sigma}^{1/2}$ is typically of order $\mathcal{O}(1)$ (see e.g., [BS⁺98a]), it is then easily seen that, for each $k \geq 1$, the operator norm of the Hadamard product $f^{(k)}(\mathbf{\Sigma}) \odot \mathbf{F}^{(k)}$ vanishes (see Lemma C.7 in the Appendix). In particular, note that the non-zero entries of $\mathbf{\Sigma}$ are controlled through the maximum entry of $\mathbf{F}^{(k)}$ which is vanishing asymptotically, as mentioned in item (ii) above. On the opposite $f(\mathbf{\Sigma})$ does not vanish since it has entries bounded away from zero (as long of course as $f \neq 0$). We precisely have the following result.

Theorem 4.4 (Sparse PCA). *Let $\mu > 0$ and suppose $\mathbf{\Sigma}$ is a $\frac{1}{2+\mu}$ -sparse matrix. For f a three-times continuously differentiable function such that $f'(0) = f''(0) = 0$ and for $\eta > 0$, we have for all $\varepsilon \in (0, \frac{\mu}{2(3+2\mu)})$*

$$\mathbf{F} = f(\mathbf{\Sigma}) + \mathcal{O}_\eta^{[1/\varepsilon]} \left(n^{\frac{-\mu}{2(2+\mu)} + \varepsilon \left(2 - \frac{1}{2+\mu} \right)} \right). \quad (4.19)$$

Proof. See Section C.2.2. □

Remark 4.2. Theorem 4.4 gives a general result concerning the estimation of ε -sparse covariance matrices (more precisely, element-wise functionals of sparse covariance matrices). In particular, the spiked model in equation 4.13 with \mathbf{U} sparse corresponds to the particular case when $\mu \rightarrow \infty$; in this case, for $\eta > 0$ and for all $\varepsilon \in (0, 1/4)$, $\mathbf{F} = f(\mathbf{\Sigma}) + \mathcal{O}_\eta^{[1/\varepsilon]} (n^{-\frac{1}{2} + 2\varepsilon})$.

One may then perform a PCA on \mathbf{F} for some function f with $f'(0) = f''(0) = 0$ (we denote by these two conditions in the following). But, while $\mathbf{\Sigma} = \mathbf{I}_p + \mathbf{P}$ is a low rank perturbation of the identity (therefore having only k eigenvalues strictly greater than 1 with

⁴See Fact .1 in [EK08].

⁵Given $\mathbf{M} \in \mathcal{M}_p$, its corresponding adjacency matrix is defined as $\mathcal{A}(\mathbf{M}) = \{\mathbf{1}_{M_{ij} \neq 0}\}_{i,j=1}^p$.

⁶See proof of Lemma C.7.

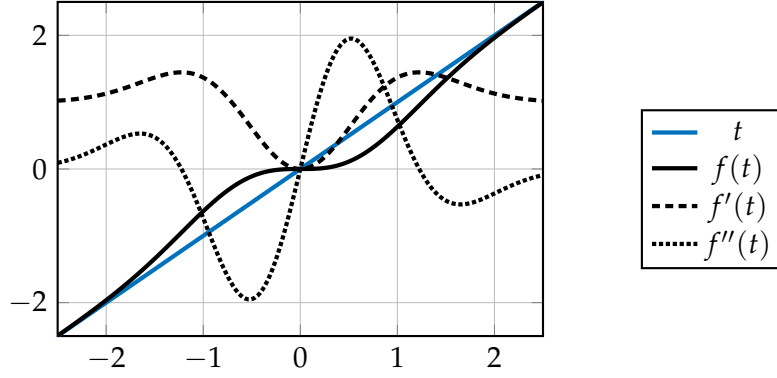


Figure 4.18: The function $f(t) = t(1 - e^{-t^2})$ along with its derivatives.

k being fixed w.r.t. p), $f(\Sigma)$ is likely more complex and not a mere low rank deformation of the identity. Now, if Σ has all its non-zero entries greater than a certain threshold τ , an appropriate choice for f that avoids the deformation of $I_p + P$ is such that $f(t) = t$ for all $|t| > \tau$.

Such a convenient choice is

$$f(t) = t(1 - e^{-at^2}), \quad (4.20)$$

for some $a > 0$, which is depicted along with its derivatives in Figure 4.18 for $a = 1$. This function notably satisfies

$$\begin{aligned} f'(t) &= 1 + e^{-at^2}(2at^2 - 1) \Rightarrow f'(0) = 0, \\ f''(t) &= -2ate^{-at^2}(2at^2 - 3) \Rightarrow f''(0) = 0. \end{aligned}$$

Note that a compromise in the choice of a must be made that both maintains a close approximation of the identity by f on a large range and rather small values of f'' in the vicinity of zero. Interestingly, it can be verified that the extrema of f' are independent of a but are found at $\pm\sqrt{\frac{3}{2a}}$ and thus smaller values of a create sharper f' in the vicinity of zero. Similarly, the extrema of f'' are found at $\pm\sqrt{\frac{3\pm\sqrt{6}}{2a}} \propto 1/\sqrt{a}$, and precisely given by the four values $2\sqrt{3a(3\pm\sqrt{6})}e^{-\frac{1}{2}(3\pm\sqrt{6})} \propto \sqrt{a}$. Thus smaller a induce larger maxima for f'' but no sharper slope.

4.2.2.3 Experimental Validation

In this section, we provide some experiments in the context of sparse PCA. Precise setting given in caption of Figure 4.19. The spectrum of the sample covariance matrix (in gray) is quite different from that of the population covariance Σ . One instead observes a “bulk” of eigenvalues spread in the vicinity of 1. The limiting measure of these eigenvalues is given by the Marčenko-Pastur density function

$$f_{MP}(x) = \frac{1}{2\pi c} \frac{\sqrt{(\lambda^+ - x)(x - \lambda^-)}}{x} + \max\left(1 - \frac{1}{c}, 0\right) \delta(x), \quad (4.21)$$

⁷ δ denotes here the dirac function.

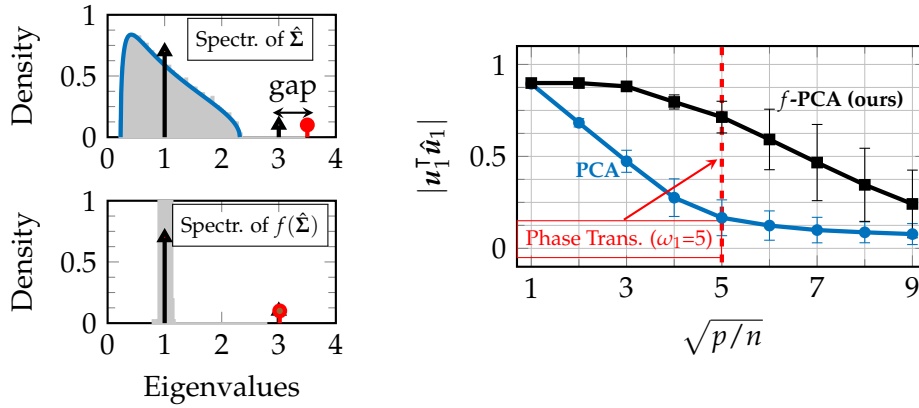


Figure 4.19: **(Left)** Spectrum of $\hat{\Sigma}$ (up) and $f(\hat{\Sigma})$ (bottom) for $p = 2048$ and $n = 7500$. Limiting Marčenko-Pastur density [MP67] in blue versus spectrum of Σ in black, with $\omega_1 = 2$; estimated largest eigenvalue in red. **(Right)** Alignment between estimated PC and GT (the “Three Peak” example of [JL09] in the “Symmlet 8” wavelet basis), in terms of $\sqrt{p/n}$. We considered $\omega_1 = 5$ and thus the phase transition for standard PCA occurs at $\sqrt{p/n} = 5$, thereby suggesting another phase transition for f -PCA. Curves obtained from 500 realizations of Z with random *i.i.d.* $\mathcal{N}(0, 1)$ entries.

which is represented by the blue curve in Figure 4.19 (left). Furthermore, one observes a gap between the true spike and the estimated spike (in red) through the sample covariance matrix. This phenomenon is well-understood in random matrix theory. In particular, the extreme eigenvalue in our setting converges almost surely to the quantity $(1 + \omega_1)(1 + \frac{c}{\omega_1})$, where we recall that $c = \lim_n p/n$.

However, thanks to sparsity, the spectrum of $F = f(\hat{\Sigma})$ closely matches that of Σ , as suggested by Theorem 4.4. In particular, the extreme eigenvalue, which corresponds to the principal component, is consistently estimated. Figure 4.19 (right) depicts the alignment between the estimated principal component and ground truth, by standard PCA (in blue) and our method (in black), in terms of $\sqrt{p/n}$. Our method retrieves the principal component even when the spike is not visible in the spectrum of $\hat{\Sigma}$; namely beyond the phase transition $\sqrt{p/n} \geq \omega_1$. In fact, the standard PCA result is too noisy compared with the one when considering $f(\hat{\Sigma})$, as depicted in Figure 4.20.

Higher Rank Case. In this paragraph, we provide further experiments by considering a rank three case and by using the “Three Peak”, “Piece Poly” and “Step New” signals of Johnstone et al. [JL09], in the “Symmlet 8” wavelet basis, as principal components. We compare the estimated PCs by our method with the kernel function in equation 4.20 to the estimated ones through standard PCA and the CT method [DM16]. As shown in Figure 4.22, the proposed method retrieves consistently the principal components compared to a standard PCA. In particular, we obtain results that are similar to the ones obtained by the CT method while generalizing it to the class of smooth functions with $f'(0) = f''(0) = 0$.

Other choices of the kernel function f . In this paragraph, we consider functions of the form $f(t) = \alpha t^3 + \beta t^2 + \gamma t$ where $\alpha, \beta, \gamma \in \mathbb{R}$ are some parameters to fix in order to allow or not the conditions $f'(0) = f''(0) = 0$. In particular, we set different parameters

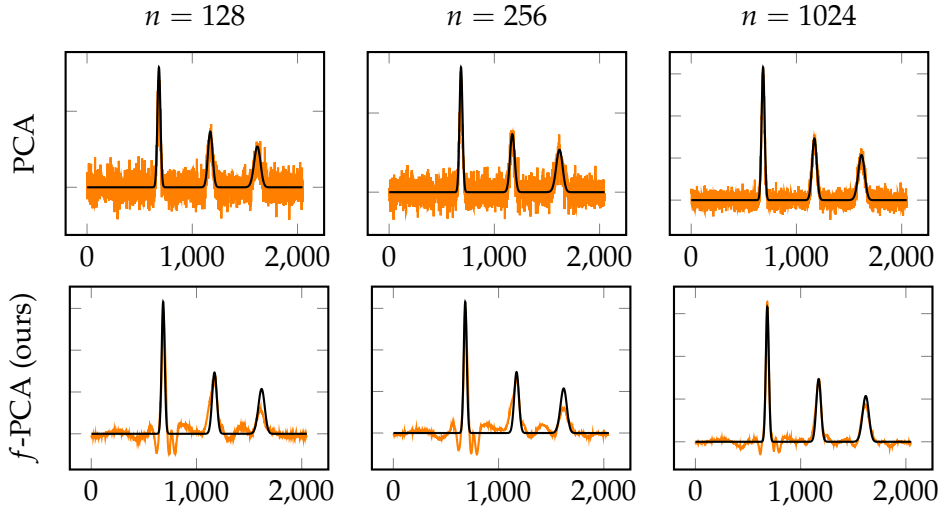


Figure 4.20: Principal component recovery (in orange) by standard PCA (**up**) and our method (**down**) for the “Three Peak” example of johnstone2009consistency. The signal is sparse in the “Symmlet 8” wavelet basis. We use $p = 2048$, $\omega_1 = 5$ for the strength of the spike and different values of n .

choices for α, β and γ in order to validate these conditions. Figure 4.23 depicts different PC recovery using the f -PCA method with the considered class of functions. As we can observe from this figure, the “cleanest” signal recovery is obtained when $\alpha \neq 0, \beta = 0, \gamma = 0$ (i.e., when $f'(0) = f''(0) = 0$) thereby validating our theoretical conditions on the kernel function f for a consistent sparse PCA recovery. Note that these conditions are necessary but not sufficient in the sense that f has to be linear for large values of t (In particular, this is the case for the function f given by equation 4.20). In fact, the outcome provided by f -PCA for $f(t) = \alpha t^3$ with $\alpha \neq 0$ is not optimal as the obtained signal is deformed (due to the unverified linearity condition), compared to the GT one.

Complexity and Performance of f -PCA. In terms of complexity, as our method consists in computing the sparse eigenvectors of a $p \times p$ matrix which can be done by power method, the complexity of estimating the principal component is about $\mathcal{O}(ps)$ where s is the sparsity level. And regarding the performance *w.r.t.* state-of-the-art methods, Figure 4.21 depicts the performances of standard PCA, different state-of-the-art sparse PCA methods and our method, in terms of total projections score (left) and total projections error (right), for different values of the amplitudes ω_i 's. We refer, in this figure, to standard PCA as PCA, TpowPCA for the method in [YZ13], ITSPCA for the method in [M⁺13], CT refers to the method in [DM16] and finally we refer to our method as f -PCA. The total projections score \mathcal{S} and error \mathcal{E} are given respectively by

$$\mathcal{S} = \frac{1}{k} \sum_{i=1}^k (\mathbf{u}_i^\top \hat{\mathbf{u}}_i)^2, \quad \mathcal{E} = \|\mathbf{U}\mathbf{U}^\top - \hat{\mathbf{U}}\hat{\mathbf{U}}^\top\|_F, \quad (4.22)$$

where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_k]$ are the ground truth principal components and $\hat{\mathbf{U}} = [\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_k]$ are the estimated ones.

As suggested theoretically and verified experimentally, our proposed method strongly attenuates the “noise component” of the sample covariance matrix and thus consistently

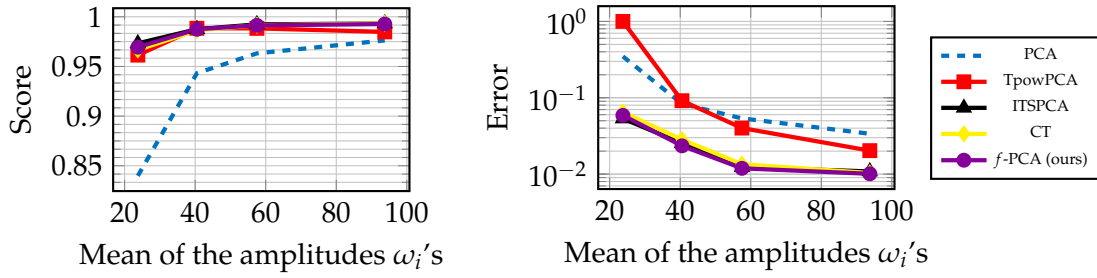


Figure 4.21: Performances of standard PCA, different state-of-the-art sparse PCA methods and our method in term of total projections score (**left**) and total projections error (**right**) for different values of the amplitudes ω_i . The PCs u_i , for $i \in [4]$ are the “Three Peak”, “Piece Poly”, “Step New” and “Sing” signals of johnstone2009consistency. We use $p = 2048$ and $n = 1024$. The soft-parameters a and τ (respectively for our method and CT) are selected by cross-validation using a validation set of size n . The selected parameters are $a = 20$ and $\tau = 0.1$.

estimates the principal components. In particular, in term of total projections score, PCA is the most inconsistent. In general, ITSPCA, CT and our method give equivalent results. The same holds when considering the total projections error as a metric, except that TpowPCA performs inconsistently, compared to PCA, for small values of amplitudes due to the initialization step from the PCA eigenvectors.

Discussion. The mostly used concurrent methods to PCA in a sparse context are iterative truncated power methods (such as the TPower [YZ13] algorithm or the ITSPCA [M⁺13] approach). These algorithms, despite great observed performances, as compared to standard PCA, suffer from two limitations. First, they are usually initialized from the PCA eigenvectors themselves and may not converge to good estimates. For weak signals, PCA is so impacted by noise that the mentioned initialization limitation may lead to non convergent or dramatically erroneous outcomes of the method. The proposed approach deals precisely with this limitation by strongly attenuating the “noise component” of the sample covariance matrix. In particular, our approach gives equivalent results to the CT method while generalizing it to the class of smooth functions f such that $f'(0) = f''(0) = 0$, in the considered regime. The second limitation concerns the choice of the hyper-parameters; in fact, TPower and ITSPCA need to set up an arbitrary deterministic threshold value that maintains at each iteration step only most powerful components. The proposed method as well as CT need also to set up a “soft” parameter (a and τ respectively). But, on the basis of [CS13, KC17], we believe that our present investigation can be extended to the *asymptotically non-trivial* setting where $\omega_i = \mathcal{O}(1/\sqrt{p})$ (in which case the dominant eigenmodes scale at a similar rate with residual noise); this setting may likely allow to exhibit and estimate optimal hyper-parameter choices. Notably, this setting has already been used in [TAKC18] in a different context, for hyper-parameters estimation.

4.2.3 Central Contribution and Perspectives

We tackled in this work the problem of sparse PCA in the large dimensional setting through a random matrix perspective, thereby generalizing recent ideas to a broader

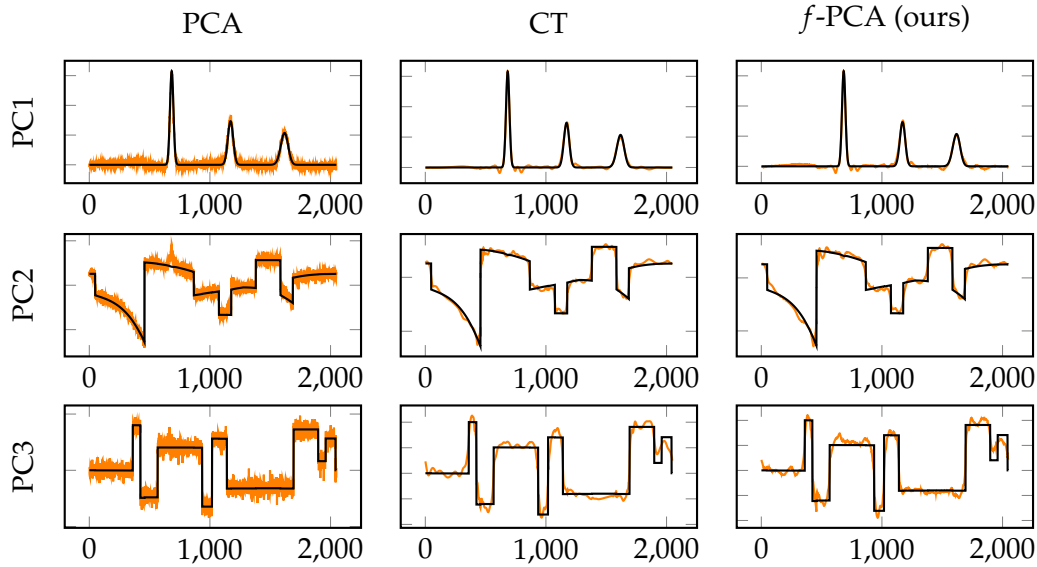


Figure 4.22: Multiple principal components ($k = 3$) recovery (in orange) with standard PCA (left), the CT method (middle) and our method (right) where the PCs are considered to be the “Three Peak”, “Piece Poly” and “Step New” signals of johnstone2009consistency, in the “Symmlet 8” wavelet basis. We use $p = 2048$, $n = 1024$ and the spikes strengths are set respectively as $\omega_1 = 100$, $\omega_2 = 75$ and $\omega_3 = 50$. In particular, we note the similarity between the results obtained by our method and CT.

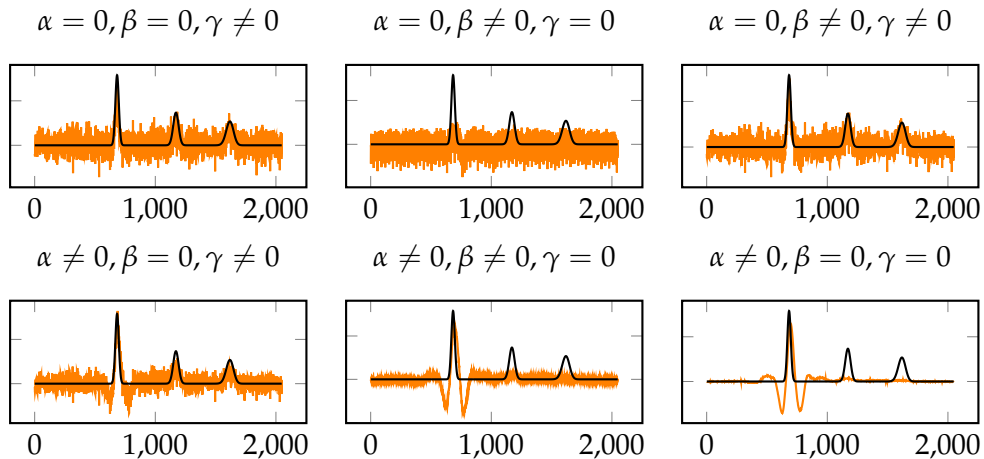


Figure 4.23: PC recovery (in orange) by f -PCA with the function $f(t) = \alpha t^3 + \beta t^2 + \gamma t$ for different values of the parameters $(\alpha, \beta, \gamma) \in \mathbb{R}^3$. We consider the “Three Peak” example of johnstone2009consistency which is sparse in the “Symmlet 8” wavelet basis. We use $p = 2048$, $n = 256$ and $\omega_1 = 5$. In particular, we notice that the “cleanest” signal is obtained when $\alpha \neq 0, \beta = 0, \gamma = 0$ which validate our theoretical conditions $f'(0) = f''(0) = 0$.

kernel-based method. Our analysis of this problem has yielded insights into how the principal components can be consistently estimated. Namely, given a spiked covariance model $\hat{\Sigma}$ and a smooth function f , we gave in this work sufficient conditions on f to consistently estimate the principal components through the matrix $f(\hat{\Sigma})$. Our result is based on the concentration assumption which generalizes the Gaussian case. However, the proposed method needs to set up a “soft” parameter (a). Based on [CS13, KC17], the present work can be extended to the *asymptotically non-trivial* setting where $\omega_i = \mathcal{O}(1/\sqrt{p})$ (in which case the dominant eigenmodes scale at a similar rate with residual noise); which may likely allow to exhibit and estimate optimal hyper-parameter choices.

Chapter 5

Beyond Kernel Matrices, to Neural Networks

Contents

5.1 A random matrix analysis of Softmax layers	103
5.1.1 Motivation	104
5.1.2 Model setting: the Softmax classifier	105
5.1.3 Assumptions & Main Results	106
5.1.3.1 Concentration of the weights vector of the Softmax classifier	107
5.1.3.2 Experimental validation	112
5.1.4 Central Contribution	114
5.2 A random matrix analysis of Dropout layers	115
5.2.1 Motivation	115
5.2.2 Model and Main Results	116
5.2.2.1 Deterministic equivalent	119
5.2.2.2 Generalization Performance of α -Dropout	121
5.2.2.3 Training Performance of α -Dropout	122
5.2.3 Experiments	122
5.2.4 Central Contribution and Perspectives	123

In the previous chapter, we have mainly presented kernel methods which explicitly depend on the input data, even when coupled with a classifier such as LS-SVM [LC17]. This Chapter will present methods that go beyond kernel matrices and which appear as sub-components of neural networks. Specifically, this chapter is composed of two main sections. The first section provides an analysis of the Softmax layer in neural networks, while the second section addresses questions about the Dropout operation in neural networks.

5.1 A random matrix analysis of Softmax layers

This section is based on the following work:

- (C4) MEA. Seddik, C.Louart, R. Couillet, M. Tamaazousti, “*The Unexpected Deterministic and Universal Behavior of Large Softmax Classifiers*”, AISTATS 2021.

5.1.1 Motivation

The intricate nature of deep neural network training leaves little insight on the specific information encoded into the inter-layer connectivity weights of a fully trained network, thereby so far not allowing for particularly useful interpretation and control of their performances [YKYR18].

At the very source of these difficulties are the multiple non-linearities and the implicit optimization scheme involved in the network design: the activation functions in the intermediate layers as well as the soft or hard final decision layer [LWL⁺17]. For lack of a tractable comprehensive analysis, literature studies have mostly focused on individual components which, when isolated, become tractable. For instance, the effect of non-linearities in a single-hidden layer network was analysed in [PW17, LC18c], the learning dynamics in elementary network designs in [SMG13, dCPS⁺18] and the overall understanding of the geometry of the loss surface in a largely approximated version of a deep neural net in [PB17, CHM⁺15].

These works are however restricted to the analysis either of the intermediate layers of practical neural nets, or oversimplify the network to an extent that makes the results rather impractical. The present study instead focuses on the training of the weights of the last decision layer, by specifically studying the widely used Softmax component in neural networks classifiers. The Softmax classifier has the property, which we will see to be of importance here, to be optimal for Gaussian mixture inputs with equal covariance [YW19]. Specifically, assuming the feature representations of the data fixed at the penultimate layer of the network, and modelling these features as *concentrated random vectors* [Led05b] (which is a natural assumption as concentrated random vectors enjoy the property to be stable through Lipschitz maps, and thus through the action of intermediate neural network layers [SLTC20]), in this work we have studied the statistical behavior of this last layer once trained (see Figure 5.1).

From a technical standpoint, as the Softmax classifier training corresponds to a (possibly non-convex) optimization problem, our analysis of the Softmax weights is performed by first expressing the optimization problem as a contracting fixed point equation, and then showing that the assumed *concentration properties* of the data features naturally transfer to the solution of the fixed-point equation, and thus to the Softmax weights. This has the major consequence that, as $n, p \rightarrow \infty$, the Softmax weights tend to have a deterministic behavior which we express explicitly as a function of the data statistics and the Softmax parameters.

Our most fundamental findings may be summarized as follows: **1.** the above deterministic behavior exhibits a surprising *universality* of the Softmax classifier, in the sense that the large dimensional statistics of the weights solely depend on the statistical means and covariances of the input data features; **2.** this suggests in turn that, quite counter-intuitively, at least as far as the last Softmax classification layer is concerned, no further discriminative feature of the data is extracted and, possibly most outstandingly, *the Softmax layer treats the input data as if they were Gaussian random vectors*; this, in passing, supports the Gaussianity assumption on the data representations commonly considered in the literature [HRU⁺17, PRU⁺18]; **3.** combined to the aforementioned optimality of the Softmax classifier on Gaussian mixture models with strongly discriminative class-wise means, this compellingly supports an overall classification optimality of the Softmax

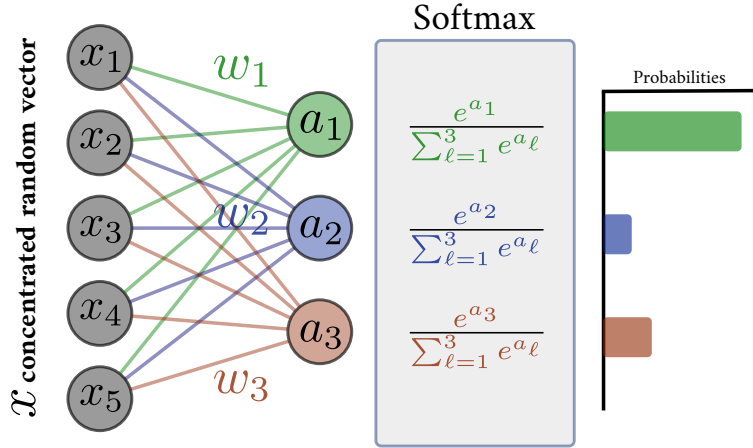


Figure 5.1: Illustration of the Softmax classifier with *concentrated random vectors* [Led05b, LC18b] (belonging to some space \mathcal{X}) as input data, i.e., satisfying the property $\mathbb{P}(|\varphi(x) - \mathbb{E}\varphi(x)| > t) \leq C e^{-(t/\sigma)^q}$, for all 1-Lipschitz $\varphi : \mathcal{X} \rightarrow \mathbb{R}$ (see Definition 7). GAN data as well as their deep network-based representations are practical examples of such random vectors [SLTC20].

classifier on large dimensional representations of real data. A similar behavior was already pointed out, yet not well understood, by the authors in [MVPC13, GCM18]; 4. Our findings are supported both theoretically and practically by considering the input data features as CNN-representations of images generated by the BigGAN model [BDS18].

5.1.2 Model setting: the Softmax classifier

Let $(x_1, y_1), \dots, (x_n, y_n)$ be a set of n labeled data associated to one of k classes $\mathcal{C}_1, \dots, \mathcal{C}_k$, where $x_i \in \mathbb{R}^p$ and $y_i \in \mathbb{R}^k$ is one-hot encoded vectors such that $y_{i\ell} = 1$ if $x_i \in \mathcal{C}_\ell$. The x_i 's are assumed to be the input of an ℓ_2 -regularized Softmax classifier with regularization parameters $(\lambda_\ell)_{\ell \in [k]} \in \mathbb{R}^+$, which aims to determine the class-wise weight vectors $w_1, \dots, w_k \in \mathbb{R}^p$ minimizing the loss¹, for some real-valued function $\phi : \mathbb{R} \rightarrow \mathbb{R}$:

$$\mathcal{L}(w_1, \dots, w_k) = -\frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{i\ell} \log p_{i\ell} + \frac{1}{2} \sum_{\ell=1}^k \lambda_\ell \|w_\ell\|^2 \quad \text{with} \quad p_{i\ell} = \frac{\phi(w_\ell^\top x_i)}{\sum_{j=1}^k \phi(w_j^\top x_i)}$$

In particular, the classical Softmax classifier corresponds to the case where $\phi(t) = e^t$ [GP17]. Cancelling the loss function gradient with respect to each weight vector w_ℓ yields

$$\lambda_\ell w_\ell = -\frac{1}{n} \sum_{i=1}^n \left(y_{i\ell} \psi(w_\ell^\top x_i) - \frac{\phi(w_\ell^\top x_i)}{\sum_{j=1}^k \phi(w_j^\top x_i)} \sum_{j=1}^k y_{ij} \psi(w_j^\top x_i) \right) x_i, \quad \ell \in [k], \quad (5.1)$$

where $\psi \equiv \phi' / \phi$. Denoting $a_\ell = w_\ell^\top x$ for some data vector x with corresponding one-hot label vector y . The expression of equation 5.1 can be obtained using the chain rule

¹Biases are not introduced in the present formulation as their effect is known to be negligible in practice [KXR⁺19] and would only decrease the readability and accessibility of our results.

through the following derivations.

$$\frac{\partial \mathcal{L}}{\partial a_j} = - \sum_{\ell} y_{\ell} \frac{1}{p_{\ell}} \frac{\partial p_{\ell}}{\partial a_j} = -y_j \psi(a_j) + p_j \sum_{\ell} y_{\ell} \psi(a_{\ell})$$

where we used the fact that $\frac{\partial p_i}{\partial a_i} = \frac{\phi'(a_i)}{\phi(a_i)} p_i (1 - p_i)$ if $i = j$ and $\frac{\partial p_i}{\partial a_j} = -\frac{\phi'(a_j)}{\phi(a_j)} p_i p_j$ otherwise.

Under the concentration assumptions on the data matrix $\mathbf{X} \equiv [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ and on ψ , and assuming p, n large, we subsequently show that the vector $\mathbf{W} \equiv [\mathbf{w}_1^{\top}, \dots, \mathbf{w}_k^{\top}]^{\top} \in \mathbb{R}^{pk}$ has a well defined behavior, which in turn allows us to accurately predict the performances of the Softmax classifier.

5.1.3 Assumptions & Main Results

We first characterize the data classes: if $y_{i\ell} = 1$, then $\mathbf{x}_i \in \mathbb{R}^p$ is a random vector with

$$\mathbb{E}[\mathbf{x}_i] \equiv \boldsymbol{\mu}_{\ell}, \quad \mathbb{E}_{\mathbf{x}_i}[\mathbf{x}_i \mathbf{x}_i^{\top}] - \boldsymbol{\mu}_{\ell} \boldsymbol{\mu}_{\ell}^{\top} \equiv \boldsymbol{\Sigma}_{\ell}.$$

The vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are further assumed to be independent and are such that $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathcal{M}_{p,n}$ satisfies a *concentration* property as stated by the following assumption.

Assumption 12 (Concentrated data). *Letting $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, $\mathbf{X} \propto \mathcal{E}_2$.*

In the terms of Definition 7, Assumption 12 holds here for $s = (p, n)$ with $\mathbb{S} = \mathbb{N}^2$. In order to be able to transfer the concentration of \mathbf{X} to the Softmax weights $\mathbf{w}_1, \dots, \mathbf{w}_k$, a further condition is needed: the number of data n must scale with the data dimension p , i.e., $\mathbb{S} = \{(p, n) \in \mathbb{N}^2, \kappa p \leq n \leq Kp\}$ for some $K > \kappa > 0$.² This is summarized by the request:

Assumption 13 (Growth rate). *$n = O(p)$ and $p = O(n)$.*

Concentrated vectors satisfy a host of interesting properties (the reader being referred to [Led05b] for a detailed account and to [LC18c] for their application to random matrix asymptotics, closer to the present work). We merely stress here one of these properties, of central importance to the present work, and which fundamentally justifies the appearance of Gaussian-like behaviors in large neural networks, even when the neural network input is far from Gaussian [KGC18, NBA⁺18].

Theorem 5.1 (CLT for concentrated vectors [Kla07, FGP07]). *Let $\mathbf{X} \in \mathbb{R}^p$ be a random vector with $\mathbb{E}[\mathbf{X}] = 0$ and $\mathbb{E}[\mathbf{X}\mathbf{X}^{\top}] = I_p$, and σ be the uniform measure on the sphere $\mathcal{S}^{p-1} \subset \mathbb{R}^p$ of radius 1. Then, if $\mathbf{X} \propto \mathcal{E}_2$, there exists two constants $C, c > 0$ and a set $\Theta \subset \mathcal{S}^{p-1}$ such that $\sigma(\Theta) \geq 1 - \sqrt{p} C e^{-c\sqrt{p}}$ and $\forall \boldsymbol{\theta} \in \Theta$:*

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(\boldsymbol{\theta}^{\top} \mathbf{X} \geq t) - G(t)| \leq p^{-1/4}$$

for G the cumulative distribution function of an $\mathcal{N}(0, 1)$ random variable.

²Formally, in the present setting, it is sufficient that $p \leq \frac{1}{\kappa} n$. However, to obtain simpler expressions, it is convenient to assume, in addition, that $n \leq Kp$.

5.1.3.1 Concentration of the weights vector of the Softmax classifier

Having set the assumptions and main technical tools, we turn now to the characterization of the the statistical behavior of the Softmax classifier weights $\mathbf{W} = [\mathbf{w}_1^\top, \dots, \mathbf{w}_k^\top]^\top \in \mathbb{R}^{pk}$, as a result, accesses the asymptotic performances of the classifier. To this end, our approach is to first write the implicit defining equation 5.1 of \mathbf{W} under the formal form $\mathbf{W} = \Psi(\mathbf{W})$, for $\Psi : \mathbb{R}^{pk} \rightarrow \mathbb{R}^{pk}$ to characterize, and then to *transfer* the concentration of \mathbf{X} (Assumption 12) to a concentration of \mathbf{W} .

For the concentration of \mathbf{X} to propagate into \mathbf{W} defined through the formal form $\mathbf{W} = \Psi(\mathbf{W})$, Ψ is required to have contracting properties, which in turn will enforce structural conditions on the operator ϕ and on the regularizers $(\lambda_\ell)_{\ell \in [k]}$. Specifically Ψ is requested to be $(1 - \varepsilon)$ -Lipschitz (for some $\varepsilon > 0$) so to ensure, thanks to the Banach fixed point theorem, the existence and uniqueness of $\mathbf{W} \in \mathbb{R}^{pk}$. However, being a *random map* depending on \mathbf{X} , Ψ is only contracting under the (asymptotically highly probable) event \mathcal{A}_X that the norm of \mathbf{X} is not too large. Indeed, we can rewrite equation 5.1 as

$$\Lambda \mathbf{W} = \frac{1}{n} \tilde{\mathbf{X}} f(\tilde{\mathbf{X}}^\top \mathbf{W}) \quad (5.2)$$

where $\Lambda = \text{diag}(\Lambda_1, \dots, \Lambda_k) \otimes \mathbf{I}_p$, $f(\tilde{\mathbf{X}}^\top \mathbf{W}) \in \mathbb{R}^{kn}$ concatenates the elements $f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \in \mathbb{R}^k$, $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{kp, kn}$, where we introduced $\tilde{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i & & \\ & \ddots & \\ & & \mathbf{x}_i \end{pmatrix} \in \mathcal{M}_{pk, k}$ $\forall i \in [n]$, and the functions:

$$\begin{aligned} f_i : \mathbb{R}^k &\longrightarrow \mathbb{R}^k \\ \mathbf{v} &\longmapsto \left[\frac{\phi(\mathbf{v}_\ell)}{\sum_{j=1}^k \phi(\mathbf{v}_j)} \sum_{j=1}^k y_{ij} \psi(\mathbf{v}_j) - y_{i\ell} \psi(\mathbf{v}_\ell) \right]_{1 \leq \ell \leq k}. \end{aligned}$$

Rewriting (5.5) as $\mathbf{W} = \Psi(\mathbf{W})$, one sees that Ψ is contracting if $\|d\Psi|_{\mathbf{w}}\| \leq 1 - \varepsilon$ for some $\varepsilon > 0$. Now, since $\|d\Psi|_{\mathbf{w}}\| \leq \frac{1}{n} \|\Lambda^{-1}\| \|\tilde{\mathbf{X}}\|^2 \|df|_{\tilde{\mathbf{X}}^\top \mathbf{w}}\|$ we first need $\|df\|_\infty = \sup_{\mathbf{v} \in \mathbb{R}^k} \|df|_{\mathbf{v}}\|$ to be bounded, that is ensured by:

Assumption 14 (Regularity). $\|\frac{\phi'}{\phi}\|_\infty \leq \infty$ and $\|\frac{\phi''}{\phi}\|_\infty \leq \infty$.

Besides, we also need to be able to bound $\frac{1}{n} \|\tilde{\mathbf{X}}\|^2 = \frac{1}{n} \|\mathbf{X}\|^2$. The spectral norm being lower than the Frobenius norm (involved in Assumption 1 giving the concentration of $\tilde{\mathbf{X}}$), it is a 1-Lipschitz observation of $\tilde{\mathbf{X}}$ thus there exists two constants $C, c > 0$ (independent of p, n), such that:

$$\mathbb{P} \left(\left| \frac{1}{\sqrt{n}} \|\tilde{\mathbf{X}}\| - \frac{1}{\sqrt{n}} \mathbb{E}[\|\tilde{\mathbf{X}}\|] \right| \geq t \right) \leq C e^{-cnt^2}. \quad (5.3)$$

This concentration inequality proves that the random variable $\frac{1}{\sqrt{n}} \|\tilde{\mathbf{X}}\|$ is almost deterministic and equal to $\frac{1}{\sqrt{n}} \mathbb{E}[\|\tilde{\mathbf{X}}\|]$. Bounding this last quantity necessities the following result:

$$|\mathbb{E}[\|\tilde{\mathbf{X}}\|] - \|\mathbb{E}[\tilde{\mathbf{X}}]\| | = O(\sqrt{p+n}), \quad (5.4)$$

It suffices to have a bound on $\|\mathbb{E}[\tilde{\mathbf{X}}]\|$ which comes from the following assumption:

Assumption 15. $\sup_{1 \leq \ell \leq k} \|\boldsymbol{\mu}_\ell\| \leq O(1)$

This assumption implies that $\|\mathbb{E}[\tilde{\mathbf{X}}]\| = O(\sqrt{n})$, and we can deduce from equation 5.4 that $\mathbb{E}[\|\tilde{\mathbf{X}}\|] = O(\sqrt{n})$ (recall from Assumption 2 that $p = O(n)$).

Remark 5.1. For classification problems, if $\forall a, b \in [k], a \neq b, \|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\| \gg 1$ then the classification becomes trivial in the large dimensional regime, a reasonable assumption then is $\|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\| \leq O(1)$. However, in cases where $\sup_{1 \leq \ell \leq k} \|\boldsymbol{\mu}_\ell\|$ is of order $O(\sqrt{p})$, it is still possible to work with the data matrix $\tilde{\mathbf{X}} - \frac{1}{n} \tilde{\mathbf{X}} \mathbf{1}_n \mathbf{1}_n^\top$ that satisfies Assumption 15 but has the drawback that the columns are then dependent. However, this is not a big issue since this dependence is very small and it can be managed thanks to some technical consideration that we want to avoid here.

Now, assuming

Assumption 16. $\frac{1}{n} \mathbb{E}[\|\mathbf{X}\|^2] \|df\|_\infty \|\Lambda^{-1}\| < 1$.

and noting $\varepsilon = \frac{1}{2} - \frac{1}{2} \|\Lambda^{-1}\| \|f'\|_\infty \mathbb{E}[\|\tilde{\mathbf{X}}\|/\sqrt{n}]^2$, it can be deduced from equation 5.3 that the event:

$$\mathcal{A}_X = \left\{ \frac{1}{n} \left| \|\tilde{\mathbf{X}}\|^2 - \mathbb{E}[\|\tilde{\mathbf{X}}\|^2] \right| \leq \frac{\varepsilon}{2 \|\Lambda^{-1}\| \|df\|_\infty} \right\}$$

has a very high probability to happen (bigger than $1 - Ce^{-cn}$ for two constant $C, c > 0$) and it satisfies $\mathcal{A}_X \subset \{\|d\Psi\|_\infty \leq 1 - \varepsilon\}$. As a consequence, our fixed point \mathbf{W} is uniquely determined under the event \mathcal{A}_X that appears in the following theorem which provides the concentration of the weights vector \mathbf{W} as stated by the following Theorem.

Theorem 5.2 (Concentration of \mathbf{W} [?]). *Under the previous Assumptions, there exist two constants $C, c > 0$ and an event \mathcal{A}_X with $\mathbb{P}(\mathcal{A}_X) > 1 - Ce^{-cn}$ such that³*

$$(\mathbf{W} \mid \mathcal{A}_X) \propto \mathcal{E}_2 \left(\sqrt{\log n/n} \right).$$

Remark 5.2. The concentration of the random vector \mathbf{W} is far from being trivial because \mathbf{W} is not explicitly written as a Lipschitz transformation of \mathbf{X} and additional tools are necessary to prove the concentration of \mathbf{W} . The complete proof will be provided in an extended version of this paper.

Since their observable diameter ($\sqrt{\log n/n}$) vanishes at large n , it therefore entails from Theorem 5.2 that the random weights vector \mathbf{W} tend to be deterministic as p, n grow large. In particular, \mathbf{W} concentrates around its expectation which can be estimated in terms of the data statistics as we will see subsequently. Indeed, the subsequent result further characterizes the first and second order statistics of \mathbf{W} at large p, n .

Theorem 5.3 (Asymptotic statistics of \mathbf{W}). *Define the statistics*

$$\mathbf{m}_W \equiv \mathbb{E}[\mathbf{W}], \quad \mathbf{C}_W \equiv [\mathbf{W}\mathbf{W}^\top] - \mathbf{m}_W \mathbf{m}_W^\top.$$

Then, under Assumptions 12 and 13 and additional assumptions on ϕ and $(\lambda_\ell)_{\ell \in [k]}$ (Assumptions 14, 15 and 16), there exists a deterministic mapping $\mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} = \mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k)$:

³Formally, the random vector \mathbf{W} is a measurable mapping $\Omega \rightarrow \mathbb{R}^{pk}$, where $(\Omega, \mathcal{F}, \mathbb{P})$ is the underlying probability space. If $\mathbb{P}(\mathcal{A}) > 0$, for $\mathcal{A} \in \mathcal{F}$, the random vector $(\mathbf{W} \mid \mathcal{A})$ is the measurable mapping $\mathcal{A} \rightarrow \mathbb{R}^{pk}$ such that, $\forall \omega \in \mathcal{A}, (\mathbf{W} \mid \mathcal{A})(\omega) = \mathbf{W}(\omega)$. The statistics of $(\mathbf{W} \mid \mathcal{A})$ are then computed in the probability space $(\mathcal{A}, \mathcal{F} \wedge \mathcal{A}, \mathbb{P}_{\mathcal{A}})$, where $\mathcal{F} \wedge \mathcal{A} = \{B \cap \mathcal{A}, B \in \mathcal{F}\}$ and $\forall B \in \mathcal{F}, \mathbb{P}_{\mathcal{A}}(B) = \mathbb{P}(B)/\mathbb{P}(\mathcal{A})$.

$\mathbb{R}^{pk} \times \mathcal{M}_{pk} \longrightarrow \mathbb{R}^{pk} \times \mathcal{M}_{pk}$ depending only on the statistics $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k$ of \mathbf{X} , such that the equation

$$(\mathbf{m}, \mathbf{C}) = \mathcal{F}_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}(\mathbf{m}, \mathbf{C}) \quad \text{with} \quad \mathbf{m} \in \mathbb{R}^{pk}, \mathbf{C} \in \mathcal{M}_{pk}$$

admits a unique solution $(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W)$. Besides,

$$\|\bar{\mathbf{m}}_W - \mathbf{m}_W\| \leq O\left(\sqrt{\log n/n}\right) \quad \text{and} \quad \|\bar{\mathbf{C}}_W - \mathbf{C}_W\|_* \leq O\left(\sqrt{\log n/n}\right).$$

The central outcome of Theorem 5.3 is that, under the data concentration Assumption 12, the behavior of the Softmax classifier *only depends on the class-wise means and covariances of the input data*. This arises as a direct consequence of the Lipschitz character of the Softmax classifier which preserves concentration (by the stability result of Remark ??), and of the presence of a *projection of the parameter vectors \mathbf{w}_ℓ onto the concentrated data \mathbf{x}_i* at the core of the optimization formulation: according to Theorem C.1, these projections induce an asymptotic Gaussian behavior with mean and variance depending *only* on the first statistics of the data and the weights vector \mathbf{W} .

We now provide the main ingredients to obtain the result of Theorem 5.3, which mainly unfolds from two essential steps: (i) the control of the statistical dependencies between \mathbf{W} and \mathbf{X} , and (ii) the estimation of the statistics \mathbf{m}_W and \mathbf{C}_W of the weight vector \mathbf{W} . We start by reformulating (5.1) in the compact and convenient form as we saw previously

$$\boldsymbol{\Lambda} \mathbf{W} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \quad (5.5)$$

(i) Control of the statistical dependencies. Applying the expectation operator both sides to (5.5), the main technical difficulty arises from the evaluation of $\mathbb{E}[\tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W})]$ due to the elaborate dependencies between the weight vector \mathbf{W} and the data $\tilde{\mathbf{x}}_i$. Note that $\tilde{\mathbf{x}}_i^\top \mathbf{W}$ *a priori* has no reason of being Gaussian (even in the limit) and the performances of the Softmax classifier may depend on high order statistics of \mathbf{X} . These statistical dependencies are dealt with by introducing a mapping $\mathbf{W}_{-i} : [0, 1] \rightarrow \mathbb{R}^{pk}$, defined for $i \in [n]$, as the unique solution to:

$$\forall t \in [0, 1] : \boldsymbol{\Lambda} \mathbf{W}_{-i}(t) = \frac{1}{n} \sum_{j \neq i} \tilde{\mathbf{x}}_j f_j(\tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i}(t)) + \frac{1}{n} t \tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}(t)). \quad (5.6)$$

This mapping can be seen as a path between the weights vector $\mathbf{W} = \mathbf{W}_{-i}(1)$ of the Softmax classifier and $\mathbf{W}_{-i}(0)$ which is completely independent of $\tilde{\mathbf{x}}_i$ and which will be simply denoted \mathbf{W}_{-i} .

Using the inverse function theorem, the mapping $t \mapsto \mathbf{W}_{-i}(t)$ is differentiable, we then deduce the following central close form formula:

$$\mathbf{W}'_{-i}(t) = \frac{1}{n} \mathbf{Q}_{-i}(t) \tilde{\mathbf{x}}_i \chi'_i(t) \quad \text{with} \quad \mathbf{Q}_{-i}(t) \equiv \left(\boldsymbol{\Lambda} - \frac{1}{n} \tilde{\mathbf{X}}_{-i} \mathbf{D}^{(i)}(t) \tilde{\mathbf{X}}_{-i}^\top \right)^{-1} \in \mathcal{M}_{kp}, \quad (5.7)$$

where $\chi_i(t) \equiv t f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}(t))$, $\mathbf{D}_j^{(i)}(t) \equiv d f_j \tilde{\mathbf{x}}_j^\top \mathbf{W}_{-i}(t) \in \mathcal{M}_k$, $\mathbf{D}^{(i)}(t) \in \mathcal{M}_{kn}$ is a block-diagonal matrix with block-diagonal matrices $\mathbf{D}_j^{(i)}(t) \in \mathcal{M}_k$ for $j \in [n]$, and finally

$\tilde{\mathbf{X}}_{-i} \equiv (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{i-1}, 0, \tilde{\mathbf{x}}_{i+1}, \dots, \tilde{\mathbf{x}}_n) \in \mathcal{M}_{pk, kn}$ ($\tilde{\mathbf{X}}_{-i}$ is independent of $\tilde{\mathbf{x}}_i$).

Relying on concentration of measure arguments [LC20], the random vector $\mathbf{Q}_{-i}(t)\tilde{\mathbf{x}}_i$ is almost constant in terms of t and thus almost equal to $\mathbf{Q}_{-i}(0)\tilde{\mathbf{x}}_i$. Moreover, the fact that $\mathbf{Q}_{-i}(0)$ (also simply denoted \mathbf{Q}_{-i}) is independent of $\tilde{\mathbf{x}}_i$ allows us to integrate the identity (5.7) to obtain the core result of the analysis of the Softmax weights relating \mathbf{W} and \mathbf{W}_{-i} .

Theorem 5.4. *Under the event \mathcal{A}_X , we have $\mathbf{W} - \mathbf{W}_{-i} + \frac{1}{n}\mathbf{Q}_{-i}\tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \propto \mathcal{E}_2\left(\frac{1}{n}\right)$. In particular, we have the bound:*

$$\mathbb{E} \left[\left\| \mathbf{W} - \mathbf{W}_{-i} + \frac{1}{n}\mathbf{Q}_{-i}\tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \right\| \right] \leq O\left(\frac{\sqrt{\log n}}{n}\right). \quad (5.8)$$

Specifically, we have the following concentration inequality, for some constants $C, c > 0$:

$$\forall t > 0 : \mathbb{P} \left(\left\| \tilde{\mathbf{x}}_i^\top \mathbf{W} - \tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i} + \frac{1}{n}\tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i}\tilde{\mathbf{x}}_i f_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}) \right\| \geq t \mid \mathcal{A}_X \right) \leq Ce^{-cnt^2/\log n}. \quad (5.9)$$

(ii) Estimation of the mean and covariance of the Softmax weights. By breaking the statistical dependencies of the problem through \mathbf{W}_{-i} , we may now access and estimate the statistics \mathbf{m}_W and \mathbf{C}_W . This precisely comes from a deterministic approximation of the quadratic form $\frac{1}{n}\tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i}\tilde{\mathbf{x}}_i$ in (5.9) in the large n, p limit, a result inspired by [LC20]:

Proposition 5.1. *For any $\ell \in [k]$, let n_ℓ be the number of columns of \mathbf{X} in class \mathcal{C}_ℓ and, for any block diagonal matrix $\Delta = \text{diag}(\Delta_\ell)_{1 \leq \ell \leq k} \in \mathcal{M}_{k^2}$ ($\forall \ell \in [k], \Delta_\ell \in \mathcal{M}_k$), then*

$$\mathbf{Q}(t) \leftrightarrow \bar{\mathbf{Q}}(\Delta) \equiv \left(\Delta - \sum_{a=1}^k \frac{n_a}{n} \Gamma_a(\Delta_a) \otimes \mathbf{S}_a \right)^{-1} = \begin{pmatrix} \bar{\mathbf{Q}}_{1,1}(\Delta) & \dots & \bar{\mathbf{Q}}_{1,k}(\Delta) \\ \vdots & \ddots & \vdots \\ \bar{\mathbf{Q}}_{k,1}(\Delta) & \dots & \bar{\mathbf{Q}}_{k,k}(\Delta) \end{pmatrix} \in \mathcal{M}_{kp}$$

where $\Gamma_\ell(\Delta_a) = \mathbb{E} \left[(\mathbf{I}_k - \mathbf{D}_j^{-i}(0)\Delta_a)^{-1} \mathbf{D}_j^{-i}(0) \right]$ for \mathbf{x}_j in class \mathcal{C}_ℓ and $\mathbf{S}_\ell = \mathbb{E}[\mathbf{x}_j \mathbf{x}_j^\top]$. Specifically, the fixed point equation

$$\Delta_\ell = \left[\frac{1}{n} \text{Tr}(\mathbf{S}_\ell \bar{\mathbf{Q}}_{a,b}(\Delta)) \right]_{1 \leq a, b \leq k}$$

admits a unique solution $\Delta \in \mathcal{M}_{k^2}$ that satisfies, for any \mathbf{x}_i is in class $\ell \in [k]$,

$$\forall t > 0 : \mathbb{P} \left(\left\| \frac{1}{n}\tilde{\mathbf{x}}_i^\top \mathbf{Q}_{-i}\tilde{\mathbf{x}}_i - \Delta_\ell \right\| \geq t \mid \mathcal{A}_X \right) \leq Ce^{-cnt^2/\log n} \quad \text{for some constants } C, c > 0.$$

From this result, using the identity (5.9), we then obtain an estimation for $\tilde{\mathbf{x}}_i^\top \mathbf{W}$:

Proposition 5.2. *For any $\mathbf{v} \in \mathbb{R}^k$, there exists a unique point $g_i(\mathbf{v}) \in \mathbb{R}^k$ satisfying:*

$$g_i(\mathbf{v}) = \mathbf{v} - \Delta_i f_i(g_i(\mathbf{v})),$$

and, for some constants $C, c > 0$,

$$\mathbb{P} \left(\left\| \tilde{\mathbf{x}}_i^\top \mathbf{W} - g_i(\tilde{\mathbf{x}}_i^\top \mathbf{W}_{-i}) \right\| \geq t \mid \mathcal{A}_X \right) \leq Ce^{-cnt^2/\log n}.$$

Therefore, letting $h_i = f_i \circ g_i$, by Hölder's inequality [Fin92], since \tilde{x}_i is concentrated,

$$\begin{aligned} \left\| \mathbf{m}_W - \frac{1}{n} \sum_{i=1}^n \Lambda^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i})] \right\| &= O\left(\sqrt{\frac{\log n}{n}}\right) \\ \left\| \mathbf{C}_W - \frac{1}{n^2} \sum_{i=1}^n \Lambda^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i}) h_i(\tilde{x}_i^\top \mathbf{W}_{-i})^\top \tilde{x}_i^\top] \Lambda^{-1} \right\|_* &= O\left(\sqrt{\frac{\log n}{n}}\right) \end{aligned}$$

Knowing from Theorem C.1 that $\tilde{x}_i^\top \mathbf{W}_{-i}$ is asymptotically Gaussian, $\mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i})]$ and $\mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i}) h_i(\tilde{x}_i^\top \mathbf{W}_{-i})^\top \tilde{x}_i^\top]$ can be explicitly evaluated (for instance using Stein's Lemma [LN08]), and only depend on the statistical means and covariances of $\tilde{x}_1, \dots, \tilde{x}_n$ and of \mathbf{W}_{-i} (which has the same statistics as \mathbf{W}). Their exact expressions are provided in the appendix in Section C.3.1. Finally, let us introduce the $2k$ functions $m_1, \dots, m_k : \mathbb{R}^{kp} \times \mathcal{M}_{kp} \rightarrow \mathbb{R}^{kp}$ and $c_1, \dots, c_k : \mathbb{R}^{kp} \times \mathcal{M}_{kp} \rightarrow \mathcal{M}_{kp}$ defined, $\forall i \in [n]$, by

$$\begin{aligned} m_{k(i)}(\mathbf{m}_W, \mathbf{C}_W) &= \Lambda^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i})] \\ c_{k(i)}(\mathbf{m}_W, \mathbf{C}_W) &= \frac{1}{n} \Lambda^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i}) h_i(\tilde{x}_i^\top \mathbf{W}_{-i})^\top \tilde{x}_i^\top] \Lambda^{-1}, \end{aligned}$$

where $k(i)$ denotes the class of \tilde{x}_i and \mathbf{m}_W and \mathbf{C}_W are respectively the mean and covariance of \mathbf{W} . The mappings $(m_\ell)_{1 \leq \ell \leq k}$ and $(c_\ell)_{1 \leq \ell \leq k}$ are uniquely determined by the means $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ and the covariances $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k$. In particular, the deterministic pair $(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W)$, defined as the unique solution of

$$\bar{\mathbf{m}}_W = \sum_{\ell=1}^k \frac{n_\ell}{n} m_\ell(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W) \quad \text{and} \quad \bar{\mathbf{C}}_W = \sum_{\ell=1}^k \frac{n_\ell}{n} c_\ell(\bar{\mathbf{m}}_W, \bar{\mathbf{C}}_W), \quad (5.10)$$

is a good approximation for $(\mathbf{m}_W, \mathbf{C}_W)$ as stated in Theorem 5.3.

Once the Softmax classifier is trained, the probability for a new datum \mathbf{x} to belong to class $\ell \in [k]$ is explicitly given by $p_\ell(\mathbf{x}) = \phi(\mathbf{w}_\ell^\top \mathbf{x}) / \sum_{j \in [k]} \phi(\mathbf{w}_j^\top \mathbf{x})$. As a consequence of Theorem C.1, $\mathbf{w}_\ell^\top \mathbf{x}$ has a high probability to be Gaussian (since \mathbf{x} is concentrated and \mathbf{w}_ℓ has a deterministic behavior). The performances of the Softmax classifier are therefore theoretically tractable.

Corollary 5.1 (Generalization performance of the Softmax classifier). *For $\ell \in [k]$, there exists $\bar{\kappa}^\ell \in \mathbb{R}^{k-1}$ and $\bar{\mathbf{K}}^\ell \in \mathcal{M}_{k-1}$ both depending only on $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_k$ such that the misclassification error $E_t(\mathbf{x} \in \mathcal{C}_\ell)$ of a new datum \mathbf{x} belonging to class $\ell \in [k]$ defined as $E_t(\mathbf{x} \in \mathcal{C}_\ell) \equiv 1 - \mathbb{P}(\forall j \in [k] \setminus \{\ell\} : p_\ell(\mathbf{x}) \geq p_j(\mathbf{x}))$ is asymptotically approximated as*

$$E_t(\mathbf{x} \in \mathcal{C}_\ell) - 1 + \mathbb{P}(\mathbf{Z}_\ell \in \mathbb{R}_+^{k-1}) \xrightarrow{a.s.} 0 \quad \text{with} \quad \mathbf{Z}_\ell \sim \mathcal{N}(\bar{\kappa}^\ell, \bar{\mathbf{K}}^\ell). \quad (5.11)$$

In essence, Corollary 5.1 states that the generalization performance of the Softmax classifier reduces to the cumulative distribution of a low-dimensional Gaussian vector, the mean and covariance of which only depend on the class-wise means and covariances of the input data. This demonstrates the remarkable *universality* property of the Softmax classifier with respect to the data distribution, which we recall is only requested to satisfy a very loose concentration behavior (Assumption 12). The exact expressions of $\bar{\kappa}^\ell$ and $\bar{\mathbf{K}}^\ell$, along with a justification of the corollary, are given subsequently.

Let us decompose the matrix $S_W = \mathbb{E}[\mathbf{W}\mathbf{W}^T]$ followingly:

$$S_W = \left(\begin{array}{c|c|c} S_W^{1,1} & \dots & S_W^{1,k} \\ \vdots & & \vdots \\ \hline S_W^{k,1} & \dots & S_W^{k,k} \end{array} \right) \in \mathcal{M}_{kp},$$

where for all $a, b \in [k]$, $S_W^{a,b} \in \mathcal{M}_p$, so that we can introduce the low-dimensional random vector $z \sim \mathcal{N}(\ell, \mathbf{K}^\ell)$ with:

$$\ell \equiv \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \quad \text{and} \quad \mathbf{K}^\ell \equiv (\text{Tr}((\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top) \mathbf{S}_W^{a,b}))_{1 \leq a, b \leq k} - \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \mathbf{m}_W^\top \boldsymbol{\mu}_\ell. \quad (5.12)$$

The expected misclassification error $E_t(x \in \mathcal{C}_\ell)$ on a test data x belonging to class $\ell \in [k]$ expresses followingly:

$$E_t(x \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\forall j \in [k] \setminus \{\ell\} : p_\ell(x) \geq p_j(x)) \quad \text{where} \quad p_j(x) = \frac{\phi(\mathbf{w}_j^\top x)}{\sum_{h=1}^k \phi(\mathbf{w}_h^\top x)}.$$

Since x and w are both concentrated and independent Theorem C.1 allows us to assume that for all $j \in [k]$, $\tilde{x}^\top \mathbf{W} = (\mathbf{w}_j^\top x)_{1 \leq j \leq k} \sim \mathcal{N}(\ell, \mathbf{K}^\ell)$ (the objects ℓ and \mathbf{K}^ℓ were introduced in equation C.10). To simplify the problem, we are going to make the additional hypothesis that ϕ is increasing since in that case

$$\forall j \in [k] \setminus \{\ell\} : p_\ell(x) \geq p_j(x) \iff \mathbf{w}_\ell^\top x \mathbf{1}_k - \tilde{x}^\top \mathbf{W} \in \mathbb{R}_+^k$$

Since the ℓ^{th} coordinate of $\tilde{x}^\top \mathbf{W}$ is, by definition, equal to $\mathbf{w}_\ell^\top x$, only the $k-1$ other are interesting. Let us then introduce the Gaussian vector $\mathbf{Z}_\ell \in \mathbb{R}^{k-1}$ defined for all $j \in [k] \setminus \ell$ as: $[\mathbf{Z}_\ell]_j = (\mathbf{w}_\ell - \mathbf{w}_j)^\top x$. Such a vector \mathbf{Z}_ℓ has then the mean $\bar{\equiv}^\ell \mathbf{P}^\ell$ and the covariance $\bar{\mathbf{K}} \equiv \mathbf{P} \mathbf{K}^\ell \mathbf{P}^\top$ with:

$$\mathbf{P} = \mathbf{1}_{k-1} \mathbf{e}_\ell^\top - \mathbf{I}_k^{-\ell} \in \mathcal{M}_{k-1,k}$$

where $\mathbf{1}_{k-1} \in \mathbb{R}^{k-1}$ is a vector full of 1, \mathbf{e}_ℓ is the ℓ^{th} vector of the canonical basis of \mathbb{R}^k (full of zeros with a one in the ℓ^{th} coordinate) and $\mathbf{I}_k^{-\ell}$ is the identity matrix of \mathcal{M}_k deprived of the ℓ^{th} row. We can then express

$$E_t(x \in \mathcal{C}_\ell) = 1 - \mathbb{P}(\mathbf{Z}_\ell \in \mathbb{R}_+^{k-1})$$

5.1.3.2 Experimental validation

This section provides an experimental setup to support our theoretical findings. The input data $\mathbf{X} = [x_1, \dots, x_n]$ are independent Resnet15⁴ representations of size $p = 512$ [SIVA17] of images generated by the BigGAN generative adversarial network model [BDS18]: as such, being the composition of two neural networks (BigGAN and Resnet15) applied to random standard Gaussian noise (as per the BigGAN model), \mathbf{X} is concentrated by construction and satisfies Assumption 12, as requested (see [SLTC20] for a detailed analysis of the Lipschitz properties of these networks). Under this setting, Figure B.9 depicts the learned Softmax weights against their expected large n, p asymptotics as per Theorem 5.3. Despite the finite p, n setting of the simulation, a perfect match is observed between the learned weights and the theoretical predictions. Further experiments are performed on *real images* from the ImageNet dataset [DDS⁺09], which again show a perfect match between theory and practice, thereby strongly suggesting that the

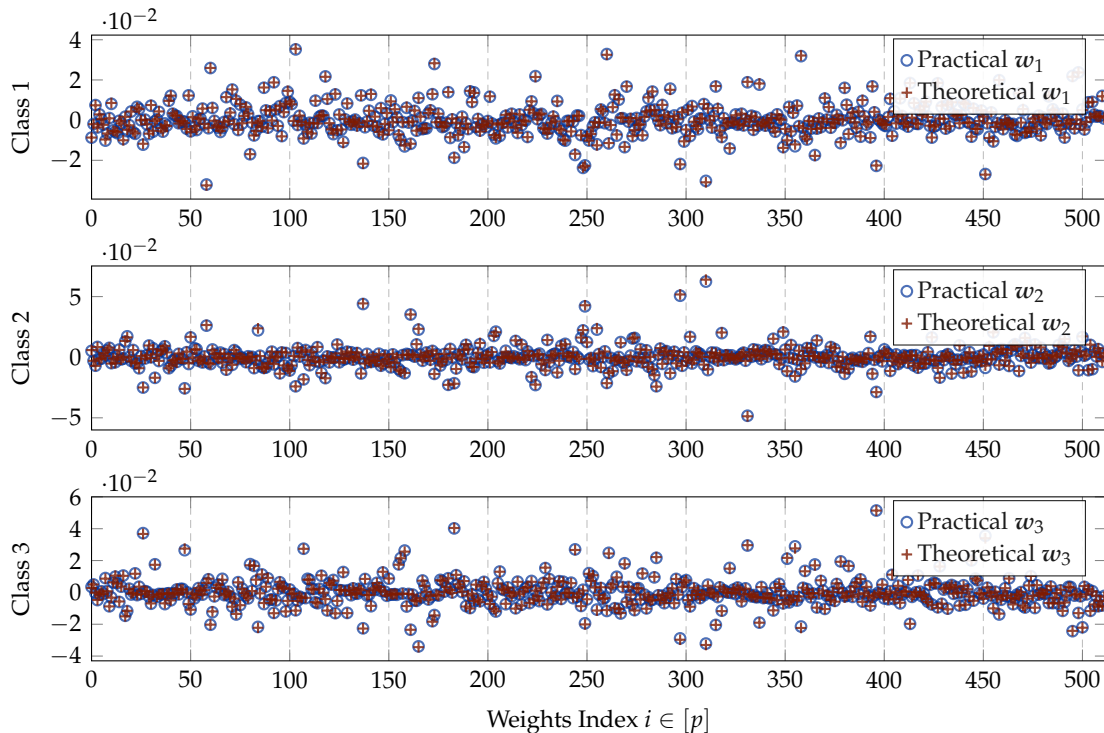


Figure 5.2: Learned weights (in blue circles) versus their theoretical estimations (in red crosses) as per Theorem 5.3. The used data are Resnet18 [SIVA17] representations ($p = 512$) of images generated by the BigGAN model [BDS18] which are concentrated vectors [SLTC20]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

conclusions of Theorem 5.3 extend to real data.

Figure 5.5 next displays the class-wise scores of a practical Softmax classifier on an independent test set against their estimated statistics according to Corollary 5.1. An almost perfect match is again observed between empirical values and theoretical statistics. Further experiments report similar outputs for real ImageNet data. We importantly stress that, as per Corollary 5.1, the theoretical estimates were obtained using *only the empirical class-wise means and covariances* of the input data. Figure 5.5 thus confirms the theoretically predicted universality of the Softmax classifier. A Python implementation is attached for reproducibility of these experiments.

We provide further experiments using real images from the Imagenet dataset [DDS⁺09]. Figure B.9 depicts the learned Softmax weights against their expected large p, n asymptotics as predicted by Theorem 3.2. As for GAN generated images, we observe a perfect match between the learned weights and the theoretical predictions. An almost perfect match is also observed for the scores (between the practical scores and their theoretical counterparts) as depicted in Figure 5.5, thereby strongly suggesting that the conclusions of Theorem 5.3 extend to real data.

⁴We used its Pytorch implementation [PGM⁺19] pre-trained on the Imagenet dataset [DDS⁺09].

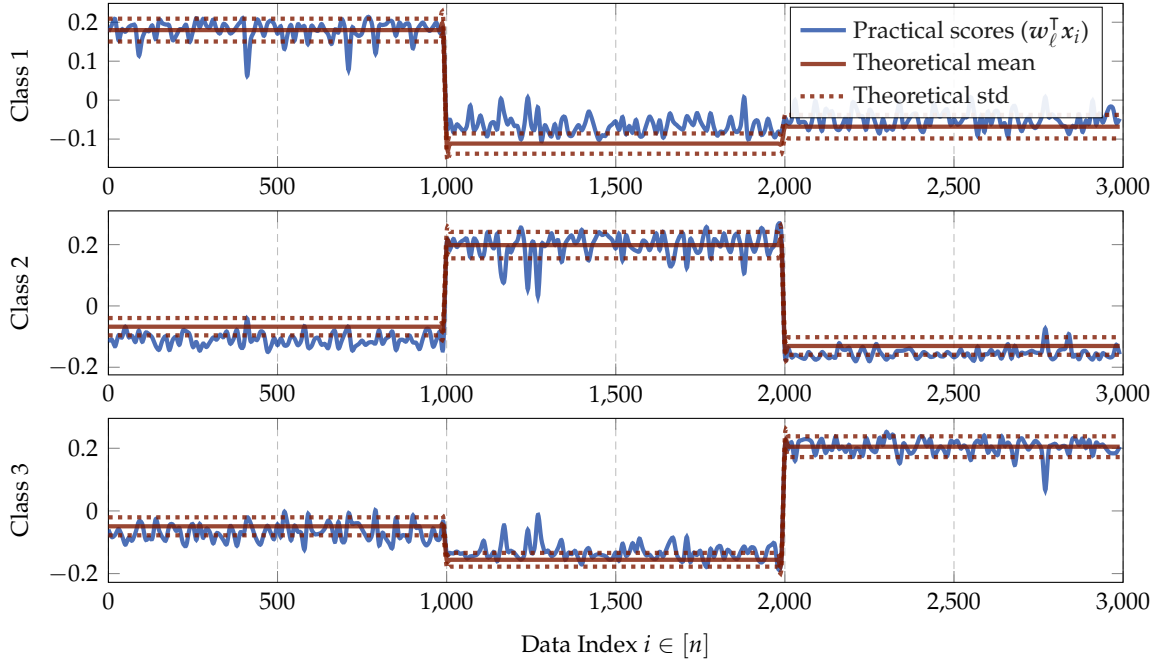


Figure 5.3: Scores (in blue) versus their theoretical estimations (in red) as per Corollary 5.1, with the theoretical means (through $\bar{\kappa}^\ell$) and standard deviations (through \bar{K}^ℓ), on a test set independent from the training set. The used data are Resnet18 [SIVA17] representations ($p = 512$) of images generated by the BigGAN model [BDS18] which are concentrated vectors [SLTC20]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

5.1.4 Central Contribution

As a consequence of Corollary 5.1, we have demonstrated that, even though the Softmax classifier has a non-linear nature, a property supposedly useful to extract “deep” non-linear features, for n, p rather large, the input data are in fact treated as if they were distributed as a mere Gaussian mixture model. This large dimensional universality phenomenon fundamentally revisits the conventional insights acquired along the years on non-linear classification methods.

As an aftermath, the Softmax classifier being optimal for Gaussian mixture inputs with common covariance, our study is strongly suggestive of the optimality of Softmax as the last layer of a deep neural network classifier.

Yet, the present study assumes a clear-cut separation between a back-end network training isolated from the front-end Softmax layer (as depicted in Figure 5.1). A thorough validation of the equivalence between full network training and this divided approach is a necessary final step to confirm the claimed optimality and anticipate the performances of Softmax classification.

On the other hand, our analysis leverages on the contracting fixed point approach by adding regularization parameters to the the weight vectors. Another approach would be to consider different regularization techniques such as the Dropout operation in neural

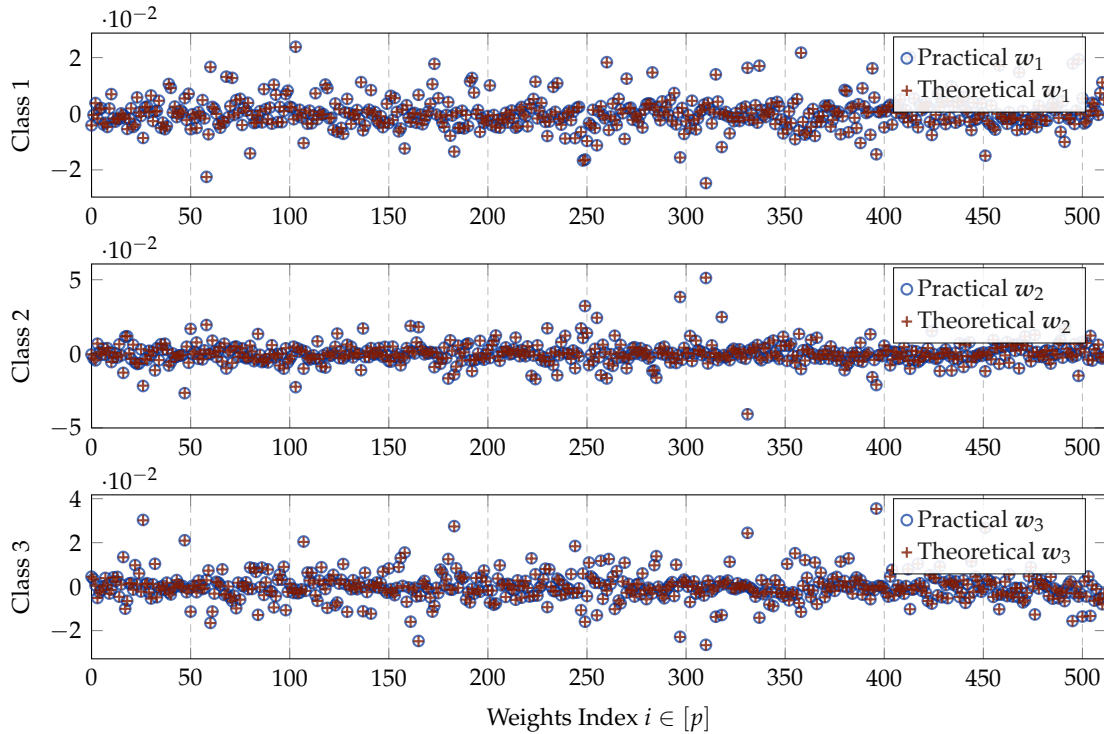


Figure 5.4: Learned weights (in blue circles) versus their theoretical estimations (in red crosses) as per Theorem 3.2 of the Main Paper. The used data are Resnet18 [SIVA17] representations ($p = 512$) of real images from the Imagenet dataset [DDS⁺09]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

networks. In the next section, we will provide an analysis of such layer on a one hidden-layer network.

5.2 A random matrix analysis of Dropout layers

This section is based on the following work:

- (C5) MEA. Seddik, R. Couillet, M. Tamaazousti, “A Random Matrix Analysis of Learning with α -Dropout”, The art of learning with missing values ICML workshop (ICML’20), Online, 2020.

5.2.1 Motivation

Many practical datasets contain samples with missing features which impair the behavior of machine learning models. Improperly handling these missing values results in biased predictions. While various imputation techniques exist in the literature, such as imputation of the global mean, the simplest is *zero imputation*, by which the missing features are simply replaced by zeros. Neural networks have been notably shown to be affected when trained on zero-imputed data [HLM15, ŠST⁺18, YLK⁺19].

In neural networks, zero imputation can be seen as applying a Dropout [SHK⁺14] operation to the input data features, or equivalently as applying a binary mask entry-wise

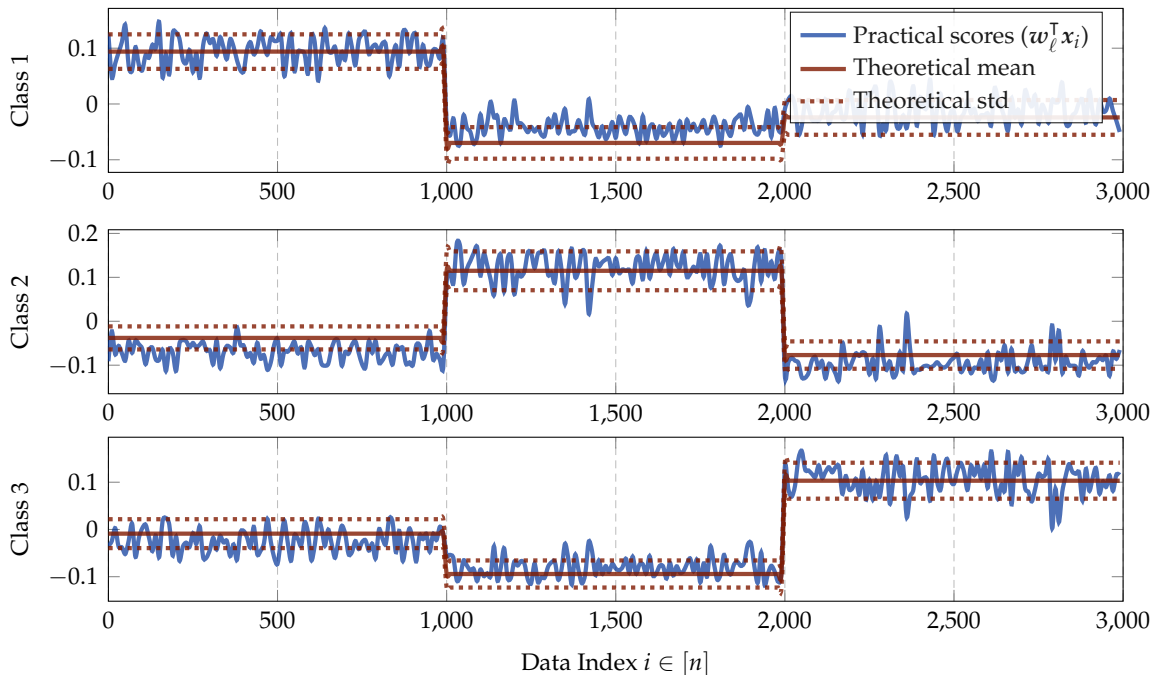


Figure 5.5: Scores (in blue) versus their theoretical estimations (in red) as per Corollary 3.3 of the Main Paper, with the theoretical means (through $\bar{\kappa}^\ell$) and standard deviations (through \bar{K}^ℓ), on a test set independent from the training set. The used data are Resnet18 [SIVA17] representations ($p = 512$) of real images from the Imagenet dataset [DDS⁺09]. We considered $k = 3$ classes which are: *hamburger*, *mushroom*, *pizza*. $n = 3000$ and the regularization constants $\lambda_1, \lambda_2, \lambda_3 = 1.5$. The data are normalized such that their norm is $0.1 \cdot \sqrt{p}$ to ensure \mathcal{A}_X .

to the data. The Dropout operation is commonly used as a regularization technique applied to certain hidden layers of a neural network during its training phase. However, since zero imputation is known to alter the behavior of neural networks [YLK⁺19], the Dropout operation must result in the same deleterious effects. Dropping features with other values than zero may thus improve the Dropout in neural networks and mitigate the effects of zero imputation [WWL13, SB16].

To prove and quantify the benefits of a α -Dropout approach, in this contribution we study a one hidden layer neural network with α -Dropout, i.e., in which the missing or dropped features are replaced by a fixed real value α . Training only the output layer, the network (sometimes referred to as an extreme learning machine [HZS06]) reduces to a ridge-regression classifier learnt on α -imputed data. Specifically, under the instrumental, yet instructive, setting of a network trained on a set of n data samples of p -dimensional features (or equivalently p neurons) distributed in two classes, we retrieve the exact generalization performance when both p and n grow large. A major outcome of our study is the identification of the optimal value of α which maximizes the generalization performances of the classifier.

5.2.2 Model and Main Results

Let the training data $d_1, \dots, d_n \in \mathbb{R}^q$ be independent vectors drawn from two distinct distribution classes \mathcal{C}_1 and \mathcal{C}_2 of respective cardinality n_1 and n_2 (and we denote $n =$

$n_1 + n_2$). We suppose the \mathbf{d}_i 's pass through a first (fixed) random neural network layer with Lipschitz activation $\sigma : \mathbb{R}^q \rightarrow \mathbb{R}^p$ in such a way that $\sigma(\mathbf{d}_i)$ is a *concentrated random vector* [LC18c]. This random projection is then followed by a random α -Dropout, i.e., entries of the feature vector $\sigma(\mathbf{d}_i)$ are dropped uniformly at random and replaced by some fixed value $\alpha \in \mathbb{R}$. Letting $\boldsymbol{\mu} \in \mathbb{R}^p$, we further suppose for simplicity of exposition that for $\mathbf{d}_i \in \mathcal{C}_a$,

$$\mathbb{E}[\sigma(\mathbf{d}_i)] = (-1)^a \boldsymbol{\mu}, \quad \mathbb{E}[\sigma(\mathbf{d}_i)\sigma(\mathbf{d}_i)^\top] = \mathbf{I}_p + \boldsymbol{\mu}\boldsymbol{\mu}^\top.$$

Remark 5.3. *As we saw in Chapter 3, this assumption could be largely relaxed but simplifies the interpretation of our results.*

Overall, after the α -Dropout layer, the feature vector $\tilde{\mathbf{x}}_i \in \mathcal{C}_a$ may thus be written

$$\tilde{\mathbf{x}}_i = ((-1)^a \boldsymbol{\mu} + \mathbf{z}_i) \odot \mathbf{b}_i + \alpha (\mathbf{1}_p - \mathbf{b}_i), \quad (5.13)$$

for $a \in \{1, 2\}$, where $\boldsymbol{\mu} \in \mathbb{R}^p$, \mathbf{z}_i is a concentrated random vector with zero mean and identity covariance, and \mathbf{b}_i is a random binary mask vector with *i.i.d.* entries $b_{ij} \sim \text{Ber}(\varepsilon)$. That is, features are discarded with an average dropout rate ε , as performed in the classical Dropout procedure in neural networks [SHK⁺14].

The model equation 5.13 thus describes a single hidden layer network with α -Dropout (dropped features are replaced by α) applied to a two-class mixture of concentrated random vectors of mean $(-1)^a \boldsymbol{\mu}$ for \mathbf{d}_i in class \mathcal{C}_a and isotropic covariance. As shown in Chapter 3, from a random matrix perspective, the asymptotic performance of the neural network under study is strictly equivalent to that of features $\tilde{\mathbf{x}}_i$ modelled as in equation 5.13 but with $\mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I}_p)$ thanks to the universality result, an assumption we will make from now on.

In a matrix form, the training features $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n] \in \mathcal{M}_{p,n}$ can be compactly written

$$\tilde{\mathbf{X}} = \mathbf{B}_\varepsilon \odot (\mathbf{Z} + \boldsymbol{\mu}\mathbf{y}^\top) + \alpha (\mathbf{1}_p \mathbf{1}_n^\top - \mathbf{B}_\varepsilon), \quad (5.14)$$

where \mathbf{Z} has *i.i.d.* $\mathcal{N}(0, 1)$ entries, $[\mathbf{B}_\varepsilon]_{ij} \sim \text{Ber}(\varepsilon)$ and $\mathbf{y} \in \mathbb{R}^n$ stands for the vector of class labels with $y_i = -1$ for $\tilde{\mathbf{x}}_i \in \mathcal{C}_1$ and $y_j = 1$ for $\tilde{\mathbf{x}}_j \in \mathcal{C}_2$.

For reasons that will be clarified latter, as depicted in Figure 5.6, we shall consider the standardized⁵ data matrix $\mathbf{X} \equiv \frac{\tilde{\mathbf{X}}\mathbf{P}_n}{\sqrt{\varepsilon + \alpha^2\varepsilon(1-\varepsilon)}}$, with $\mathbf{P}_n = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top$, i.e.,

$$\mathbf{X} = \frac{(\mathbf{B}_\varepsilon \odot (\mathbf{Z} + \boldsymbol{\mu}\mathbf{y}^\top)) \mathbf{P}_n + \alpha \mathbf{B}_\varepsilon \mathbf{P}_n}{\sqrt{\varepsilon + \alpha^2\varepsilon(1-\varepsilon)}}. \quad (5.15)$$

Under the features data model in equation 5.15, we aim in the following to study the generalization performance of a (fully connected) linear layer applied to the features x_i 's which is thus equivalent to optimizing (with an ℓ_2 regularization term)

$$\mathcal{E}(\mathbf{w}) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 + \gamma \|\mathbf{w}\|^2, \quad (5.16)$$

⁵Centring by the empirical mean and dividing by the standard deviation $\sqrt{\varepsilon + \alpha^2\varepsilon(1-\varepsilon)}$ as in batch normalization layers [IS15].

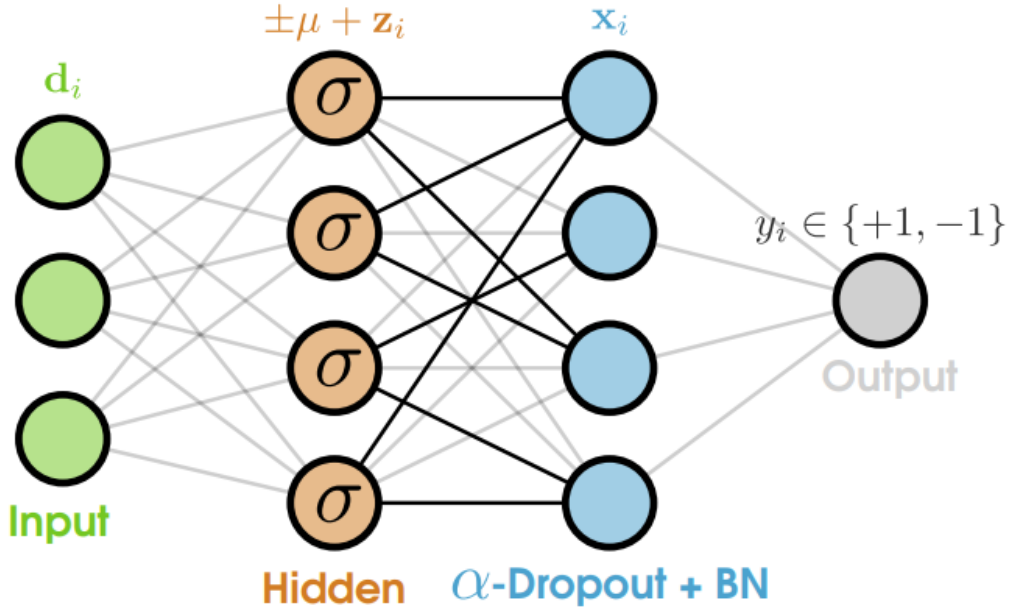


Figure 5.6: Illustration of the model under study in equation 5.15. It consists in a one-hidden-layer network followed by an α -Dropout layer and a Batch normalization (BN) layer where the output is given by a linear layer.

the solution of which is explicitly given by, for $z \in \mathbb{C} \setminus \mathbb{R}^-$

$$w = \frac{1}{n} Q(\gamma) X y, \quad Q(z) \equiv \left(\frac{1}{n} X X^T + z I_p \right)^{-1}. \quad (5.17)$$

The associated (hard) decision function for a new datum feature vector $x \in \mathcal{C}_a$, for $a \in \{1, 2\}$, then reads

$$g(x) \equiv x^T w = \frac{1}{n} x^T Q(\gamma) X y \underset{\mathcal{C}_2}{\overset{\mathcal{C}_1}{\approx}} 0. \quad (5.18)$$

The model in equation 5.15 coupled with the ridge loss function in equation 5.16 is that of an extreme learning machine trained with α -Dropout through the random matrix B_ε .

In mathematical terms, studying the generalization performance under a large dimensional network regime consists in studying the statistical behavior of the *resolvent matrix* $Q(z)$ defined in equation 5.17. The main technicality precisely arises from the unconventional presence of the matrix B_ε in the model. In the following, we derive a deterministic equivalent for $Q(z)$ which is the basic technical ingredient for the further analysis, as a function of α and ε , of the network generalization performance. Let us first present some technical growth rate assumptions before deriving a deterministic equivalent for $Q(z)$.

Assumption 17 (Growth rate). As $n \rightarrow \infty$,

1. $q/n \rightarrow r \in (0, \infty)$ and $p/n \rightarrow c \in (0, \infty)$;
2. For $a \in \{1, 2\}$, $\frac{n_a}{n} \rightarrow c_a \in (0, 1)$;

3. $\|\boldsymbol{\mu}\| = \mathcal{O}(1)$.

Remark 5.4. Assumption 17 along with the standardization in equation 5.15 ensure particularly that the operator norm of $\mathbf{Q}(z)$ is bounded asymptotically.

5.2.2.1 Deterministic equivalent

Now we provide the key technical tool to express the performances of the network under investigation, i.e., a deterministic equivalent for $\mathbf{Q}(z)$. Let $\mathbf{A} = \frac{\mathbf{P}(\mathbf{B}_\varepsilon \odot (\boldsymbol{\mu}\boldsymbol{\mu}^\top))}{\sqrt{\varepsilon + \alpha^2 \varepsilon(1-\varepsilon)}}$ be the signal part of the the model equation 5.15, from [HLN⁺07], a deterministic equivalent of $\mathbf{Q}(z)$ is given by

$$\bar{\mathbf{Q}}(z) \equiv \left(q^{-1}(z) + \frac{1}{1+cq(z)} \frac{1}{n} \mathbb{E}[\mathbf{A}\mathbf{A}^\top] \right)^{-1} \quad \text{with} \quad q(z) \equiv \frac{1+cq(z)}{1+z(1+cq(z))}$$

And by Assumption 17 ($\|\boldsymbol{\mu}\| = \mathcal{O}(1)$), we have

$$\begin{aligned} \frac{1}{n} \mathbb{E}[\mathbf{A}\mathbf{A}^\top] &= \frac{\varepsilon}{1+\alpha^2(1-\varepsilon)} \boldsymbol{\mu}\boldsymbol{\mu}^\top + \frac{1-\varepsilon}{1+\alpha^2(1-\varepsilon)} \text{diag}(\boldsymbol{\mu}^{\odot 2}) + \mathcal{O}_{\|\cdot\|}(n^{-\frac{1}{2}}) \\ &= a\boldsymbol{\mu}\boldsymbol{\mu}^\top + b \text{diag}(\boldsymbol{\mu}^{\odot 2}) + \mathcal{O}_{\|\cdot\|}(p^{-\frac{1}{2}}) \end{aligned}$$

Therefore, by Lemma 2.6, we have (we denote $r(z) = \frac{1}{1+cq(z)}$)

$$\begin{aligned} \bar{\mathbf{Q}}(z) &= \left(q^{-1} \mathbf{I}_p + br(z) \text{diag}(\boldsymbol{\mu}^{\odot 2}) + ar(z) \boldsymbol{\mu}\boldsymbol{\mu}^\top \right)^{-1} \\ &= \underbrace{\left(q^{-1} \mathbf{I}_p + br(z) \text{diag}(\boldsymbol{\mu}^{\odot 2}) \right)^{-1}}_{\mathcal{D}_z} - \frac{ar(z) \mathcal{D}_z \boldsymbol{\mu}\boldsymbol{\mu}^\top \mathcal{D}_z}{1+ar(z) \boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}} \end{aligned}$$

where $q(z) = \frac{1+cq(z)}{1+z(1+cq(z))}$ is a second order equation in $q(z)$. Therefore, we have

Proposition 5.3. Under Assumption 17,

$$\mathbf{Q}(z) \leftrightarrow \bar{\mathbf{Q}}(z) \equiv \mathcal{D}_z - \frac{\frac{\varepsilon}{1+\alpha^2(1-\varepsilon)} \mathcal{D}_z \boldsymbol{\mu}\boldsymbol{\mu}^\top \mathcal{D}_z}{1+cq(z) + \frac{\varepsilon}{1+\alpha^2(1-\varepsilon)} \boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}}$$

where $\mathcal{D}_z \equiv q(z) \text{diag} \left\{ \frac{1+cq(z)}{1+cq(z) + \frac{(1-\varepsilon)q(z)}{1+\alpha^2(1-\varepsilon)} \mu_i^2} \right\}_{i=1}^p$, and $q(z)$ is given by

$$q(z) \equiv \frac{c-z-1 + \sqrt{(c-z-1)^2 + 4zc}}{2zc}.$$

Proposition 5.3 shows that the deterministic equivalent $\bar{\mathbf{Q}}(z)$ involves two terms: a diagonal matrix \mathcal{D}_z (describing the noise part of the data model) and an informative scaled rank-1 matrix $\mathcal{D}_z \boldsymbol{\mu}\boldsymbol{\mu}^\top \mathcal{D}_z$. We see through the expression of $\bar{\mathbf{Q}}(z)$ that the informative term is linked to the noise term (through \mathcal{D}_z) if $\varepsilon \neq 1$, and for small values of ε or equivalently large values of α the “energy” of the informative term is transferred to the noise term which will result in a poor classification accuracy on the train set, still we will subsequently see that for a fixed value of ε , there exists a value of α which will provide optimal classification rates on the test set. We will next use the property that $\mathbf{a}^\top \mathbf{Q}(z) \mathbf{b} \simeq \mathbf{a}^\top \bar{\mathbf{Q}}(z) \mathbf{b}$

for all large n, p and deterministic bounded vectors \mathbf{a}, \mathbf{b} , to exploit $\bar{\mathbf{Q}}(z)$ as a proxy for the performance analysis (which is precisely related to a bilinear form on $\mathbf{Q}(z)$) of the α -Dropout neural network.

Further, let us introduce the following quantities which will be used subsequently. First, we have under Assumption 17, the statistics of the feature vector \mathbf{x}_i , for $\mathbf{x}_i \in \mathcal{C}_a$, are:

$$\begin{aligned} \mathbf{m}_a &\equiv \mathbb{E}[\mathbf{x}_i] = (-1)^a \sqrt{\frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)}} \boldsymbol{\mu}, \\ \mathbf{C}_\varepsilon &\equiv \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \mathbf{I}_p + \frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)} \boldsymbol{\mu} \boldsymbol{\mu}^\top \\ &\quad + \frac{1 - \varepsilon}{1 + \alpha^2(1 - \varepsilon)} \left(\text{diag}(\boldsymbol{\mu}^{\odot 2} + 2\alpha\boldsymbol{\mu}) - \frac{\alpha^2}{p} \mathbf{1}_p \mathbf{1}_p^\top \right). \end{aligned} \quad (5.19)$$

The above statistics are derived as follows. Let $\mathbf{x} = \mathbf{b} \odot (\mathbf{z} + \boldsymbol{\mu}) + \alpha \mathbf{P} \mathbf{b}$. Denote $\mathbb{E}[\mathbf{P} \mathbf{b}] = \mathbb{E}[\mathbf{b} - \frac{1}{p} \sum_{i=1}^n b_i \mathbf{1}_p] = \mathbf{0}$, we have $\mathbb{E}[\mathbf{b} \odot (\mathbf{z} + \boldsymbol{\mu})] = \mathbf{b} \odot \boldsymbol{\mu}$ and $\mathbb{E}[\mathbf{b} \mathbf{b}^\top] = \varepsilon \mathbf{1}_p \mathbf{1}_p^\top + (1 - \varepsilon) \mathbf{I}_p \equiv \mathbf{C}_b$. Therefore,

$$\begin{aligned} \mathbb{E}[\alpha \mathbf{P} \mathbf{b} (\alpha \mathbf{P} \mathbf{b})^\top] &= \alpha^2 \mathbf{P} \mathbf{C}_b \mathbf{P}^\top \\ \mathbb{E}[\mathbf{b} \odot (\mathbf{z} + \boldsymbol{\mu}) (\mathbf{b} \odot (\mathbf{z} + \boldsymbol{\mu}))^\top] &= \mathbf{C}_b \odot (\mathbf{I}_p + \boldsymbol{\mu} \boldsymbol{\mu}^\top) \\ \mathbb{E}[\mathbf{b} \odot (\mathbf{z} + \boldsymbol{\mu}) (\alpha \mathbf{P} \mathbf{b})^\top] &= \alpha \text{diag}(\boldsymbol{\mu}) \mathbf{C}_b \mathbf{P} \end{aligned}$$

Hence,

$$\mathbb{E}[\mathbf{x} \mathbf{x}^\top] = \mathbf{C}_b \odot (\mathbf{I}_p + \boldsymbol{\mu} \boldsymbol{\mu}^\top) + \alpha \text{diag}(\boldsymbol{\mu}) \mathbf{C}_b \mathbf{P} + \alpha \mathbf{P} \mathbf{C}_b \text{diag}(\boldsymbol{\mu}) + \alpha^2 \mathbf{P} \mathbf{C}_b \mathbf{P}^\top$$

Since $\|\boldsymbol{\mu}\| = \mathcal{O}(1)$, as in Assumption 17, we have

$$\begin{aligned} \mathbf{C}_b \odot (\mathbf{I}_p + \boldsymbol{\mu} \boldsymbol{\mu}^\top) &= \varepsilon \mathbf{I}_p + (1 - \varepsilon) \mathbf{I}_p + \varepsilon \boldsymbol{\mu} \boldsymbol{\mu}^\top + (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}^{\odot 2}) = \mathbf{I}_p + \varepsilon \boldsymbol{\mu} \boldsymbol{\mu}^\top + (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}^{\odot 2}) \\ \text{diag}(\boldsymbol{\mu}) \mathbf{C}_b \mathbf{P} &= (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}) \mathbf{P} = (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}) - \frac{1 - \varepsilon}{p} \boldsymbol{\mu} \mathbf{1}_p^\top = (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}) + \mathcal{O}_{\|\cdot\|}(p^{-\frac{1}{2}}) \\ \mathbf{P} \mathbf{C}_b \mathbf{P} &= (1 - \varepsilon) \mathbf{P} \end{aligned}$$

Thus

$$\begin{aligned} \mathbb{E}[\mathbf{x} \mathbf{x}^\top] &= (1 + \alpha^2(1 - \varepsilon)) \mathbf{I}_p + \varepsilon \boldsymbol{\mu} \boldsymbol{\mu}^\top - \frac{\alpha^2(1 - \varepsilon)}{p} \mathbf{1}_p \mathbf{1}_p^\top + (1 - \varepsilon) \text{diag}(\boldsymbol{\mu}^{\odot 2}) \\ &\quad + 2\alpha(1 - \varepsilon) \text{diag}(\boldsymbol{\mu}) + \mathcal{O}_{\|\cdot\|}(p^{-\frac{1}{2}}) \end{aligned}$$

Therefore,

$$\mathbf{m}_a \equiv \mathbb{E}[\mathbf{x}_i] = (-1)^a \sqrt{\frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)}} \boldsymbol{\mu} \quad (5.20)$$

$$\mathbf{C}_\varepsilon \equiv \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top] = \mathbf{I}_p + \frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)} \boldsymbol{\mu} \boldsymbol{\mu}^\top + \frac{1 - \varepsilon}{1 + \alpha^2(1 - \varepsilon)} (\text{diag}(\boldsymbol{\mu}^{\odot 2}) + 2\alpha \text{diag}(\boldsymbol{\mu})) \quad (5.21)$$

$$- \frac{\alpha^2(1 - \varepsilon)}{p(1 + \alpha^2(1 - \varepsilon))} \mathbf{1}_p \mathbf{1}_p^\top + \mathcal{O}_{\|\cdot\|}(p^{-\frac{1}{2}}). \quad (5.22)$$

We will also need in the following the quantity

$$\delta(z) \equiv \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \bar{\mathbf{Q}}(z)). \quad (5.23)$$

5.2.2.2 Generalization Performance of α -Dropout

The generalization performance of the classifier relates to misclassification errors

$$\mathbb{P}(g(\mathbf{x}) > 0 \mid \mathbf{x} \in \mathcal{C}_1), \quad \mathbb{P}(g(\mathbf{x}) < 0 \mid \mathbf{x} \in \mathcal{C}_2)$$

where $g(\cdot)$ is the decision function previously defined in equation 5.18.

Since the Dropout is deactivated at inference time, the statistics of \mathbf{x} correspond to the setting where $\varepsilon = 1$, and thus

$$\mathbb{E}[\mathbf{x}] = (-1)^a \boldsymbol{\mu}, \quad \mathbf{C}_1 = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}_p + \boldsymbol{\mu}\boldsymbol{\mu}^\top.$$

Further, define the following quantities which shall be used subsequently

$$\begin{aligned} \eta(\mathbf{A}) &\equiv \frac{(1 + \delta(\gamma)) \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \bar{\mathbf{Q}}(\gamma) \mathbf{A} \bar{\mathbf{Q}}(\gamma))}{(1 + \delta(\gamma))^2 - \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \bar{\mathbf{Q}}(\gamma) \mathbf{C}_\varepsilon \bar{\mathbf{Q}}(\gamma))}, \\ \Delta(\mathbf{A}) &\equiv \bar{\mathbf{Q}}(\gamma) \left(\mathbf{A} + \frac{\eta(\mathbf{A})}{1 + \delta(\gamma)} \mathbf{C}_\varepsilon \right) \bar{\mathbf{Q}}(\gamma). \end{aligned}$$

By Lyapunov's central limit theorem [Bil08], the decision function has the following Gaussian approximation as $n \rightarrow \infty$.

Theorem 5.5 (Gaussian Approximation of $g(\mathbf{x})$). *Under Assumption 17, for $\mathbf{x} \in \mathcal{C}_a$ with $a \in \{1, 2\}$,*

$$v^{-\frac{1}{2}} (g(\mathbf{x}) - m_a) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1)$$

where

$$\begin{aligned} m_a &\equiv (-1)^a \sqrt{\frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)}} \frac{\boldsymbol{\mu}^\top \bar{\mathbf{Q}}(\gamma) \boldsymbol{\mu}}{1 + \delta(\gamma)} \\ v &\equiv \frac{1}{(1 + \delta(\gamma))^2} \left(\eta(\mathbf{C}_1) + \frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)} \times \left[\boldsymbol{\mu}^\top (\Delta(\mathbf{C}_1) - \bar{\mathbf{Q}}(\gamma)) \boldsymbol{\mu} - \frac{2\eta(\mathbf{C}_1) \boldsymbol{\mu}^\top \bar{\mathbf{Q}}(\gamma) \boldsymbol{\mu}}{1 + \delta(\gamma)} \right] \right). \end{aligned}$$

Proof. See Subsection C.3.2.1. □

In a nutshell, Theorem 5.5 states that the one hidden layer network classifier with α -Dropout is asymptotically equivalent to the thresholding of two monivariate Gaussian random variables, the means and variances of which depend on $\boldsymbol{\mu}$, \mathbf{C}_ε and the parameters α and ε . As such, we have the corresponding (asymptotic) classification errors:

Corollary 5.2 (Generalization Performance of α -Dropout). *Under the setting of Theorem 5.5, for $a \in \{1, 2\}$, with probability one*

$$\mathbb{P}((-1)^a g(\mathbf{x}) < 0 \mid \mathbf{x} \in \mathcal{C}_a) - Q\left(\frac{m_a}{\sqrt{v}}\right) \rightarrow 0$$

with $Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-u^2/2} du$ the Gaussian tail function.

Corollary 5.2 can therefore be exploited to find the optimal value of α^* which minimizes the test misclassification error, since $Q'(x) < 0$ the optimal value α^* satisfies the equation $\frac{1}{m_a} \frac{\partial m_a}{\partial \alpha} = \frac{1}{\sqrt{v}} \frac{\partial \sqrt{v}}{\partial \alpha}$ which can be solved numerically.

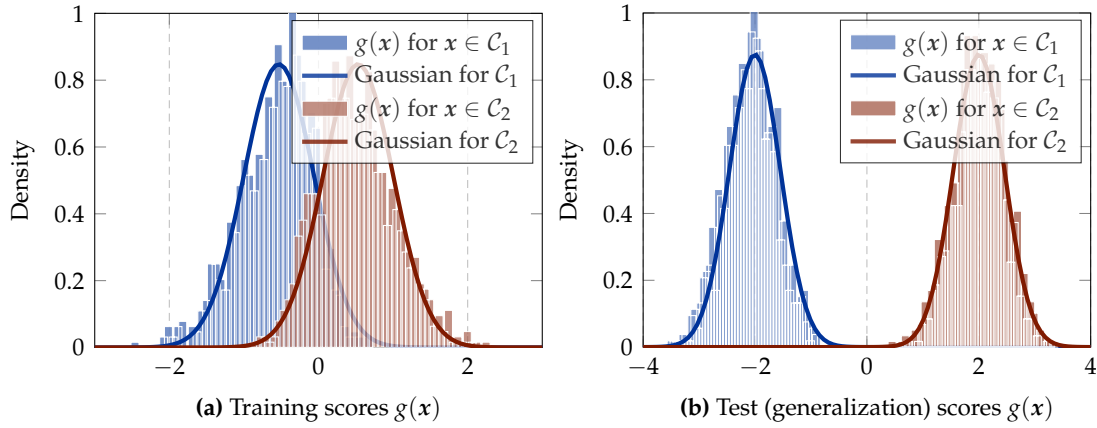


Figure 5.7: Histogram of the decision function $g(x)$ when applied to the training data **(a)** and test data **(b)**. The curves represent the Gaussian approximations as per Theorem 5.5 and Theorem 5.6 for test and training data respectively. We used the parameters $\boldsymbol{\mu} = \frac{5 \cdot \mathbf{u}}{\|\mathbf{u}\|}$ with $\mathbf{u} = [10, 10, -10, -10, \mathbf{v}]$ where $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{p-4})$, $p = 125$, $n_1 = n_2 = 1000$, $\gamma = 1 \cdot 10^{-2}$, $\varepsilon = 0.25$ and $\alpha = 2$.

5.2.2.3 Training Performance of α -Dropout

It is instructive to compare the generalization versus training performances of the network classifier with α -Dropout. Following similar arguments as in [LC18c], the central limit argument of the previous section also holds for $g(x)$ with $x \in \mathcal{C}_a$ taken from the training set \mathbf{X} .

Theorem 5.6 (Training performance of α -Dropout). *Under Assumption 17, for $x \in \mathcal{C}_a$ with $a \in \{1, 2\}$ a column of \mathbf{X} , with probability one,*

$$\mathbb{P}((-1)^a g(x) < 0 \mid x \in \mathcal{C}_a) - Q\left(\frac{\bar{m}_a}{\sqrt{\bar{v} - \bar{m}_a^2}}\right) \rightarrow 0$$

where

$$\begin{aligned} \bar{m}_a &\equiv \frac{\delta(\gamma)}{1 + \delta(\gamma)} + \frac{(-1)^a \varepsilon}{1 + \alpha^2(1 - \varepsilon)} \frac{\boldsymbol{\mu}^\top \bar{\mathbf{Q}}(\gamma) \boldsymbol{\mu}}{(1 + \delta(\gamma))^2} \\ \bar{v} &\equiv \left(\frac{\delta(\gamma)}{1 + \delta(\gamma)}\right)^2 + \frac{\eta(\mathbf{C}_\varepsilon)}{(1 + \delta(\gamma))^4} + \frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)} \\ &\quad \times \boldsymbol{\mu}^\top \left(\frac{\delta(\gamma) \bar{\mathbf{Q}}(\gamma)}{(1 + \delta(\gamma))^3} + \frac{\Delta(\mathbf{C}_\varepsilon)}{(1 + \delta(\gamma))^4} - \frac{2\eta(\mathbf{C}_\varepsilon) \bar{\mathbf{Q}}}{(1 + \delta(\gamma))^5}\right) \boldsymbol{\mu}. \end{aligned}$$

Proof. See Subsection C.3.2.2. □

5.2.3 Experiments

Gaussian Approximation of the Decision Function. We complete this study by simulations to validate our theoretical findings. Figure 5.7 depicts histograms showing the distribution of $g(x)$ for both **(a)** training and **(b)** test data. As we can see, these distributions are well approximated by monivariate Gaussians as per Theorem 5.5 and Theorem 5.6. Since the α -Dropout removes features at random from the training data, the

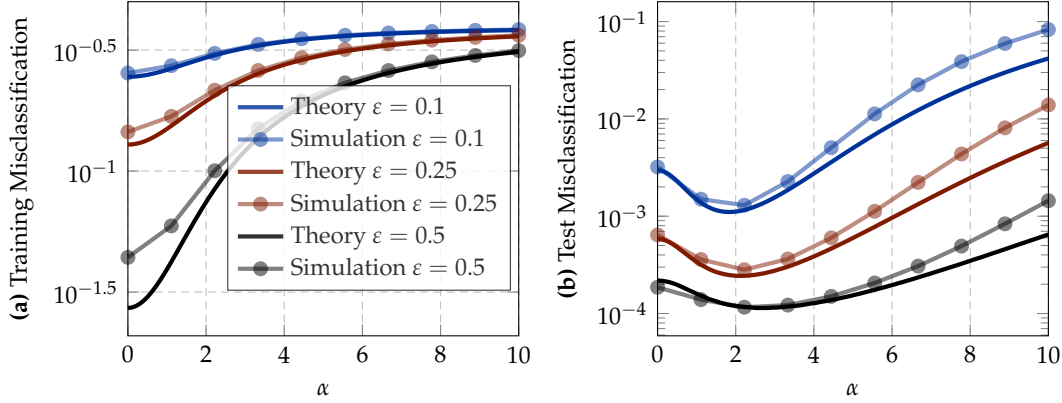


Figure 5.8: **(a)** Training and **(b)** Test misclassification errors as per Theorem 5.6 and Corollary 5.2 respectively. We used the parameters $\boldsymbol{\mu} = 4 \cdot [\frac{1}{\sqrt{4}}, \frac{1}{\sqrt{4}}, -\frac{1}{\sqrt{4}}, -\frac{1}{\sqrt{4}}, \mathbf{0}_{p-4}^\top]^\top$, $p = 125$, $n_1 = n_2 = 1000$ and $\gamma = 1 \cdot 10^{-2}$. Simulations are obtained through 100 Monte-Carlo runs of independent realizations of the matrix \mathbf{X} as in equation 5.15.

misclassification error happens to be larger on the training set compared to the test set. Notably, the difference between the training and test error arises theoretically from the term $\kappa \equiv \sqrt{\frac{\epsilon}{1+\alpha^2(1-\epsilon)}}$ as $\bar{m}_a \approx \kappa m_a$, therefore, for small values of ϵ the training error is larger than the test error, which shows the regularization effect of the Dropout.

Training and Test Performances. Figure 5.8 depicts the theoretical **(a)** training (through Theorem 5.6) and **(b)** test (through Corollary 5.2) misclassification errors, for different values of ϵ and in terms of α , and their simulated counterparts. We can notice from these plots that the training error increases with α and is minimal for $\alpha = 0$. In contrast, the test misclassification error is convex in terms of α and therefore the lowest generalization error corresponds to an optimal value $\alpha \neq 0$. We also remark that the optimal value of α increases in terms of ϵ which is counterintuitive since we expect and α near to 0 for large values of ϵ , but actually, the test misclassification error in terms of α gets more and more flatter as ϵ increases.

5.2.4 Central Contribution and Perspectives

Leveraging on random matrix theory, we have analyzed the effect of the α -Dropout layer on a one layer neural network, which allowed us to have a deeper understanding of the impact of this layer. We have notably exhibited an optimal Dropout operation (dropping our features with some $\alpha \neq 0$) in terms of the generalization error of the studied classifier. Although, our analysis was presented on a simple binary classification task, it can be straightforwardly generalized to a more realistic data model as the mixture of k -class model [LC18b, SLTC20]. Under a k -class model it may be beneficial to consider an α_ℓ per class \mathcal{C}_ℓ as the classes may be constructed with different statistics. Following the same approach one can derive the test misclassification error as per Corollary 5.2 in terms of scalar quantities involving the data statistics, and therefore exploit the formulas to find the optimal values of α_ℓ 's. Other perspectives of this work concern its application to real data and using multi-layer neural network architectures.

Chapter 6

Conclusions & Perspectives

6.1 Conclusions

The first notable result from this thesis, presented in Section 3.1, sheds light on the fact that artificially generated data through Generative Adversarial Nets (GANs) are random vectors which fall within the class of *concentrated vectors*. Consequently, real data can be modelled by this class of random vectors – and also since they generalize Gaussian vectors – if we can assimilate them to GAN data. Motivated by this first result, we have further investigated in this thesis three main ML methods, under the Mixture of Concentrated Vectors hypothesis for the input data model, where each data-class is described by its first and second order moments $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$. Indeed, relying on random matrix theory (RMT) and supposing the high-dimensional regime when the number of input data and their dimension are both large and comparable, we have first analyzed in Section 3.2 the spectral behavior of large Gram matrices which are at the core of various linear methods. As non-linearities appear in various ML methods, we have also investigated in Section 4.1 kernel methods through the analysis of large kernel matrices. Lastly, we have studied, through the example of the Softmax classifier in Section 5.1, ML methods which are implicitly defined by (convex) optimization problems.

As an aftermath of our analysis of these methods, we have emphasized that their effective performances – when applied to large dimensional data – *solely* depend on the class-wise *means* and *covariances* of the input data, namely the statistics $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ and $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$, a result which we have empirically verified through extensive experiments using CNN representations of GAN-generated images. This constitutes a major outcome of this thesis in the sense that it highlights the *universality aspect* of large ML classifiers regardless of their input data distribution. As a consequence, the universality aspect notably supports the validity of the Gaussian Mixture Model in the random matrix theory regime, thereby justifying the Gaussian mixture hypothesis on data as assumed in the works [Lia19, Mai19]. From a practical standpoint, RMT allows therefore for systematic analyzes and improvements of a wide range of ML classifiers which rely on the aforementioned methods – through the estimation of their asymptotic performances – on realistic data, such as deep learning representations of the surprising images generated artificially by GANs.

To provide an example of such analyzes and improvements through RMT, we have investigated in Section 4.2 the problem of sparse PCA where we have notably generalized recent ideas to tackle this problem by means of a broader class of random kernel matri-

ces. Specifically, our analysis has provided insights into how the principal components can be consistently recovered through a random kernel matrix $f(\hat{\Sigma})$, where $\hat{\Sigma}$ is a spiked covariance model and f a smooth function. Another example of such RMT analyzes and improvements concerns the analysis of the Dropout in neural networks. Indeed, as we presented in Section 5.2, we saw that the RMT analysis of a one-hidden-layer network with an α -Dropout layer has yielded insights about the effect of the Dropout operation, thereby improving it by suggesting dropping out features and replacing them with some $\alpha \neq 0$ in order to minimize the generalization error of the studied network.

6.2 Limitations and Perspectives

As we discussed in the previous section, the main outcomes of this thesis are twofold: The first outcome concerns the fact that GAN generated data are concentrated vectors by design, while the second outcome being the universality aspect of the studied ML methods, i.e., ML methods which rely on Gram matrices, kernel matrices or for which their solutions are explicit or implicitly defined through convex optimization problems. Yet, these findings have limitations, so we discuss in the following some potential future research directions in the same vein as our actual preliminary findings.

Validity of concentrated vectors for data modelling: So far in our experiments throughout this manuscript, we have considered as data CNN representations of GAN generated images. More fundamentally, we have shown that the framework of concentrated vectors is justified by the fact that GAN data fall within this class of random vectors. Hence, a first limitation concerns the fact that not all types of data can be generated by GANs, therefore modelling them as concentrated vectors is not straightforward. Indeed, texts (and thus their embeddings) are an example of such data types. However, in the natural language processing (NLP) paradigm, *text generation* is performed using Recurrent Neural Networks (RNNs) [SMH11]. Specifically, text data are generated sequentially starting from some hidden state $z_0 \sim \mathcal{N}(\mathbf{0}, I_d)$ and then computing at every time step t the new hidden state as $z_t = \text{RNN}(z_{t-1}, w_{t-1})$ where w_t stands for the embedding of t -th word in the processed sentence and RNN is the generator network. Consequently, studying this mechanism along with the RNN dynamics may yield to highlight that the learned word embeddings w_t are concentrated vectors, since the involved operations are Lipschitz transformations and the starting hidden state h_0 is a concentrated vector. Another and more revolutionary way of building text embeddings for NLP relies on more involved and huge neural networks architectures known as *transformers* which are based on the so-called attention modules [VSP⁺17]. The basic and fundamental operation of these architectures is the *self-attention* operation, which is a sequence-to-sequence operation and basically consists in transforming a sequence of input vectors w_1, \dots, w_t to another output sequence y_1, \dots, y_t as $y_i = \sum_j \alpha_{ij} w_j$ where the weights α_{ij} are function of the input sequences w_i 's (e.g., $\alpha_{ij} = w_i^\top w_j$ or as commonly used $\alpha_{ij} = \exp(w_i^\top w_j) / \sum_j \exp(w_i^\top w_j)$ involving the Softmax activation). Therefore, a first idea would be to study through the tools of RMT the mechanism of this operation which might explain the concentration property of the resulting learned word embeddings w_i 's.

Generalizations of the k -class Mixture of Concentrated Vectors Model: We assumed throughout this manuscript that data are made of distinct classes which are represented by their respective statistical means and covariances. This model does not take into account the between-class correlations or the fact that data have some hierarchical structure (e.g., the class animals is made of several sub-classes: cats, dogs, etc.). A possible generalization of our findings is to include these correlations/hierarchical structure and study their influence so as to understand the behavior of ML methods in some specific applications that are sensitive to between-class correlations. Concretely, a possible idea is to generalize the notion of hierarchical Gaussian mixture models [OP16, LM07] by relaxing the Gaussianity assumption to the class of concentrated vectors. Such generalization might also find applications in transfer learning in the vein of [YC08] as the target domain is very often different from the source domain, hence this difference shall be considered by modelling the underlying statistical dependencies.

On the analysis of more complex classifiers relying on non-convex optimization problems: As we discussed previously, the universality behavior concerns mostly ML methods which rely on Gram matrices, kernel matrices or for which their solutions are explicit or implicitly defined through convex optimization problems. However, the most successive ML methods nowadays rely on non-convex optimization problems which is the case of multi-layer neural networks. The main difficulty when studying these methods is the non-convex nature and the use of the backpropagation algorithm which make the analysis more challenging. However, using tools from *information theory* [SZT17], it has been shown that when building and training a neural network for a given task (e.g., classification), the *mutual information* between its input and successive hidden layers decreases across its internal layers. Roughly speaking, neural networks learn to produce hidden representations that are less *correlated* with the input by basically keeping the most relevant information across layers to solve the considered task, i.e., as long as there is a high correlation between the considered hidden layer and the desired target output. Relying on this intuition, efforts have been made to develop new alternatives to the end-to-end backpropagation algorithm [BEO19, PL20, MLK20, DYP20, LOV19]. Basically, these alternatives suggest to rather train the models layer-wise, i.e., training one layer at a time which might provide more flexibility and accessibility to theoretical analysis. A common approach for such alternatives relies on the *Information Bottleneck* (IB) principle which consists in minimizing the following objective

$$\min_{p(h|x)} \mathcal{I}(x; h) - \gamma \mathcal{I}(h; y)$$

where x stands for the input, y is the target label, h is the hidden feature vector, \mathcal{I} stands for the mutual information and $\gamma > 0$ is an hyperparameter. As such, optimizing IB is equivalent to minimizing the mutual information (i.e., dependence) between the input x and the hidden variable h while maximizing the mutual information between h and the desired target y . Specifically, h plays the role of the hidden representations of x in the neural network. Still, from a practical standpoint, a first issue appears from the fact that the mutual information is generally hard to compute. Authors in [MLK20] have very recently proposed an alternative to overcome this issue by considering the Hilbert-Schmidt Independent Criteria (HSIC) instead of the mutual information. Specifically, given some data matrix $\mathbf{X} = [x_1, \dots, x_n] \in \mathcal{M}_{p,n}$ and their corresponding targets $\mathbf{Y} = [y_1, \dots, y_n] \in \mathcal{M}_{k,n}$. Further supposing that the hidden features are given by some parametric transformation $\mathbf{H} = f(\mathbf{X}; \Theta) \in \mathcal{M}_{q,n}$ of the input matrix \mathbf{X} . The HSIC between \mathbf{X} and \mathbf{H} is particularly

given by

$$\text{HSIC}(\mathbf{X}, \mathbf{H}) = \frac{1}{(n-1)^2} \text{tr}(\mathbf{K}_X \mathbf{K}_H)$$

where $\mathbf{K}_X \in \mathcal{M}_n$ is the kernel matrix having entries $(\mathbf{K}_X)_{ij} = \kappa(x_i, x_j)$ for some kernel κ . Therefore, exploiting this formalism and RMT tools of the analysis of kernel matrices, one could gain new insights about the internal mechanism of this learning approach. In particular, two research directions might be interesting to explore:

1. A first direction concerns the analysis of the dynamics of a two hidden layers (or possibly networks with multiple layers) with the HSIC principle. Concretely, leveraging on RMT, an idea would be to study the involved functionals of the large kernel matrices \mathbf{K}_X and \mathbf{K}_H to get access to the dynamics of the studied network in the same vein as [LC18a]. Indeed, RMT analyses may provide statistical descriptions of the local or global extrema, thereby leading to a theoretical understanding of the model performances as well as improvements of the studied method through optimal hyperparameter tuning. Such an analysis could be first considered using a one-layer network with non-linear activation in order to analyze and understand the effect of non-linearity on the network dynamics, as investigated in [AS17] with the standard gradient descent approach.
2. Another direction of research, which might be more challenging, is to exploit the layer-wise learning approach to describe the *encoded information* or *encoded sufficient statistics* by neural networks layers weights in the same vein as our analysis of the Softmax layer in Section 5.1. A possible idea is to start by performing preliminary empirical studies to understand the correlations between the learned weights and the corresponding layers representations using the HSIC principle. Since we have identified, through the analysis of the Softmax layer in Section 5.1, that the last layer retrieves information from the network representations through their class-wise means and covariances, one could intuitively obtain the same conclusions layer-wise and characterize the *encoded statistics* in the internal layers weights. Such an analysis might be very useful for the understanding of the mechanism of neural networks and might explain their incredible ability to build hierarchical and discriminative representations.

As a concluding remark, we forcefully believe that the present manuscript scratches the surface of a much more ambitious endeavor. Most effective AI tools nowadays (e.g., deep neural networks) lack explainability even if they have demonstrated incredible super-human performances. However, there are several industrial domains (e.g., medical, transport or military) where explainability is foregrounded and is as crucial as performance, hence such inexplicable AI tools are left aside in these critical domains. To gain more confidence in these tools, RMT creates a new bridge, *from theory and progressively to practice*, which gradually explains and improves AI methods and makes them more accessible, understandable, reliable, hence opening a new way for research and industry along with the current explosion of AI by deep learning. In that sense, pushing forward the development of theoretical tools for the precise understanding of AI methods is of crucial interest to achieve more effective and trustworthy AI.

Appendix A

Synthèse de la thèse en Français

L'intelligence artificielle (IA) est connue comme l'ensemble des théories et des techniques utilisées pour créer des machines capables de simuler l'intelligence humaine. L'un des sous-domaines les plus intéressants de l'IA est l'apprentissage machine (ML) qui vise à fournir des algorithmes informatiques qui "apprennent" automatiquement par l'expérience afin de prendre des décisions futures sans être explicitement programmés. Fondamentalement, les algorithmes ML s'appuient sur la construction de modèles mathématiques - très souvent paramétriques - qui seront optimisés sur la base d'échantillons de données d'apprentissage et utilisés ensuite pour effectuer diverses tâches d'IA telles que la classification, la régression, le regroupement, etc.

Tout naturellement, l'IA trouve des applications dans divers domaines et, par conséquent, l'un des défis les plus importants de la ML est de fournir des algorithmes qui peuvent être appliqués à différents types de données (par exemple, des images, des textes, des graphes, etc.). Par construction, ces données peuvent être représentées sous différentes formes et, par conséquent, la performance des algorithmes de ML dépendra largement de la représentation choisie. Cette représentation devrait idéalement contenir des informations pertinentes sur les données afin de permettre l'apprentissage avec des modèles simples et une petite quantité de données. Historiquement, un grand nombre de travaux se sont concentrés sur la conception de représentations (ou de caractéristiques) artisanales, puis sur leur fourniture à des algorithmes de ML simples pour résoudre les tâches souhaitées. Mais pour la plupart des tâches et étant donné les différents types de données, ces approches ne sont pas facilement extensibles pour obtenir une IA efficace.

Depuis l'arrivée des réseaux neuronaux profonds (DNN), l'idée de développer des caractéristiques artisanales a été immédiatement écartée. En effet, les DNN ont surpassé la plupart des approches en démontrant leur incroyable capacité à apprendre automatiquement des représentations pertinentes à partir de données brutes dans un large éventail d'applications, y compris la vision par ordinateur, la reconnaissance de formes et le traitement du langage naturel. Malgré leur succès, de nombreuses questions restent sans réponse concernant les bases théoriques des DNN et qui sont très cruciales notamment pour leur explicabilité. Par exemple, la caractérisation complète de leurs représentations et/ou paramètres appris est encore un problème ouvert.

L'un des principaux aspects qui ont rendu les DNN efficaces dans la pratique est le fait qu'ils soient des modèles sur-paramétrés. En effet, il a été démontré que les architectures profondes de ces modèles surpassent les architectures peu profondes lorsqu'il s'agit de données multi-dimensionnelles (un échantillon de n données de dimension p) lorsque ces deux grandeurs n et p sont de grande taille, ce qui est souvent le cas dans les scénarios de

la vie réelle¹. En outre, les DNN les plus efficaces se trouvent à avoir un certain nombre de paramètres N qui sont au moins de l'ordre de p ou même beaucoup plus grands (par exemple, LeNet-5 [LeC98] contient des paramètres $N = 60000$).

En substance, ces grandes dimensions induisent de nombreux phénomènes contre-intuitifs qui font que les intuitions des petites dimensions s'effondrent complètement. Pour une meilleure compréhension de ces phénomènes, nous fournirons ultérieurement quelques exemples illustratifs qui révèlent ces aspects contre-intuitifs. Dans le cas particulier où les deux $p, n \rightarrow \infty$ avec $p/n \rightarrow 0 \in (0, \infty)$, la théorie des matrices aléatoires (RMT) fournit des outils puissants pour évaluer la performance de divers algorithmes de ML en tenant compte de l'effet de ces dimensions. En effet, la RMT donne accès au mécanisme interne d'un grand nombre de méthodes de ML, permettant ainsi une compréhension plus approfondie et des améliorations systématiques de ces méthodes. Nous renvoyons le lecteur à la thèse de Z. Liao [Lia19] pour les applications de RMT aux méthodes du noyau, aux réseaux neuronaux aléatoires peu profonds et à la dynamique des réseaux neuronaux ; la thèse de X. Mai [Mai19] qui traite des applications de RMT à l'apprentissage semi-supervisé et aux SVMs.

Les travaux susmentionnés s'appuient largement sur des hypothèses gaussiennes² concernant les données traitées. L'un des principaux résultats de cette thèse est d'aller au-delà de l'hypothèse gaussienne pour aborder l'applicabilité de RMT à des données réelles qui ne sont probablement pas proches des vecteurs gaussiens. En particulier, en travaillant sous le modèle statistique plus générique des vecteurs concentrés [LC20], nous fournissons des justifications – en nous appuyant sur les réseaux adversaires générateurs (GAN) – sur la pertinence d'un tel modèle pour la modélisation réaliste des données, et nous analysons en outre, sous l'hypothèse de concentration sur les données, le comportement des grandes matrices de noyaux (qui se trouvent être au cœur de divers algorithmes ML) ainsi que certains composants essentiels des réseaux neuronaux tels que la dernière couche Softmax. Un résultat majeur des travaux développés dans cette thèse est le résultat d'universalité énoncé comme :

"seulement les moments d'ordre un et deux importent pour décrire le comportement de ces méthodes"

justifiant ainsi l'hypothèse de la gaussianité des données selon les résultats de [Lia19, Mai19].

En effet, le premier résultat notable de cette thèse, présenté dans la section 3.1, met en lumière le fait que les données générées artificiellement par le biais des Generative Adversarial Nets (GAN) sont des vecteurs aléatoires qui entrent dans la classe des *vecteurs concentrés*. Par conséquent, les données réelles peuvent être modélisées par cette classe de vecteurs aléatoires – et aussi puisqu'ils généralisent les vecteurs gaussiens – si nous pouvons les assimiler aux données des GAN. Motivés par ce premier résultat, nous avons étudié plus en détail dans cette thèse trois méthodes principales de ML, sous l'hypothèse du mélange de vecteurs concentrés pour le modèle de données d'entrée, où chaque classe de données est décrite par ses moments de premier et de second ordre $\{\boldsymbol{\mu}_\ell\}_{\ell=1}^k$ et $\{\boldsymbol{\Sigma}_\ell\}_{\ell=1}^k$. En effet, en s'appuyant sur la théorie des matrices aléatoires (RMT) et en supposant le régime de grande dimension i.e., lorsque le nombre de données d'entrée et leur dimension sont à la fois grands et comparables, nous avons d'abord analysé dans la section 3.2 le comportement spectral des grandes matrices de Gram qui sont au cœur de diverses méthodes linéaires. Comme des non-linéarités apparaissent dans diverses méthodes de

¹À titre d'exemple, MNIST [LeC98] contient $n = 70000$ images de dimension $p = 28 \times 28 = 784$.

²Modélisant les données comme un mélange Gaussien de k classes.

ML, nous avons également étudié dans la section 4.1 les méthodes à noyaux par l’analyse de grandes matrices à noyaux. Enfin, nous avons étudié, à travers l’exemple du classificateur Softmax dans la section 5.1, les méthodes de ML qui sont définies par des problèmes d’optimisation (convexes) implicites.

Suite à notre analyse de ces méthodes, nous avons souligné que leurs performances effectives – lorsqu’elles sont appliquées à des données de grande dimension – *seulement* dépendent des classes *moyennes* et *covariances* des données d’entrée, à savoir les statistiques $\{\mu_\ell\}_{\ell=1}^k$ et $\{\Sigma_\ell\}_{\ell=1}^k$, un résultat que nous avons empiriquement vérifié par des expériences approfondies en utilisant les représentations CNN des images générées par un GAN. Ceci constitue un résultat majeur de cette thèse dans le sens où il met en évidence *l’aspect d’universalité* des grands classificateurs ML indépendamment de la distribution des données d’entrée. En conséquence, l’aspect d’universalité soutient notamment la validité du modèle de mélange gaussien dans le régime de la théorie des matrices aléatoires, justifiant ainsi l’hypothèse de mélange gaussien sur les données telle que supposée dans les travaux [Lia19, Mai19]. D’un point de vue pratique, la théorie des matrices aléatoires permet ainsi d’analyser et d’améliorer systématiquement un large éventail de classificateurs ML qui s’appuient sur les méthodes susmentionnées – par l’estimation de leurs performances asymptotiques – sur des données réalistes, telles que des représentations approfondi des images surprenantes générées artificiellement par les GANs.

Pour fournir un exemple de telles analyses et améliorations par RMT, nous avons étudié dans la section 4.2 le problème de l’ACP parcimonieuse où nous avons notamment généralisé des idées récentes pour s’attaquer à ce problème à travers l’analyse d’une classe plus large de matrices à noyaux aléatoires. Plus précisément, notre analyse a permis de comprendre comment les composantes principales peuvent être récupérés de manière cohérente grâce à une matrice à noyau aléatoire de la forme $f(\hat{\Sigma})$, où $\hat{\Sigma}$ est un modèle de covariance à pics et f une fonction lisse. Un autre exemple d’analyse et d’amélioration à travers la théorie des matrices aléatoires concerne l’analyse du Dropout dans les réseaux de neurones. En effet, comme nous l’avons présenté dans la section 5.2, nous avons vu que l’analyse RMT d’un réseau à une couche cachée avec une couche α -Dropout a permis de comprendre l’effet de l’opération Dropout, l’améliorant ainsi en suggérant d’abandonner des caractéristiques des données et de les remplacer par une valeur $\alpha \neq 0$ afin de minimiser l’erreur de généralisation du réseau étudié.

En conclusion générale, les outils d’IA les plus efficaces de nos jours (par exemple, les réseaux de neurones profonds) manquent d’explicabilité même s’ils ont démontré des performances surhumaines incroyables. Cependant, il existe plusieurs domaines industriels (par exemple, le médical, le transport ou le militaire) où l’explicabilité est au premier plan et est aussi cruciale que la performance, c’est pourquoi ces outils d’IA inexplicables sont laissés de côté dans ces domaines critiques. Pour gagner plus de confiance dans ces outils, la théorie des matrices aléatoires crée un nouveau pont, de la théorie et progressivement vers la pratique, qui explique et améliore progressivement les méthodes d’IA et les rend plus accessibles, compréhensibles, fiables, ouvrant ainsi une nouvelle voie pour la recherche et l’industrie en même temps que l’explosion actuelle de l’IA par les méthodes d’apprentissage profond. En ce sens, faire avancer le développement d’outils théoriques pour la compréhension précise des méthodes d’IA est d’un intérêt crucial pour parvenir à une IA plus efficace et plus fiable.

Appendix B

Practical Contributions

Contents

B.1 Generative Collaborative Networks for Super-resolution	133
B.1.1 Motivation	133
B.1.2 Proposed Methods	135
B.1.2.1 Proposed Framework	135
B.1.2.2 Existing Loss Functions	137
B.1.3 Experiments for Single Image Super-Resolution	137
B.1.3.1 Proposed Methods	137
B.1.3.2 Evaluation Metrics	138
B.1.3.3 Experiments	139
B.1.4 Central Contribution and Discussions	143
B.2 Neural Networks Compression	143
B.2.1 Motivation	143
B.2.2 Proposed Methods	144
B.2.2.1 Setting & Notations	144
B.2.2.2 Neural Nets PCA-based Distillation (Net-PCAD)	145
B.2.2.3 Neural Nets LDA-based Distillation (Net-LDAD)	145
B.2.3 Experiments	146
B.2.4 Central Contribution and Discussions	148

B.1 Generative Collaborative Networks for Super-resolution

This section is based on the following work:

(C6) MEA. Seddik, M. Tamaazousti, J. Lin, “*Generative Collaborative Networks for Single Image SuperResolution*”, Neurocomputing’2019.

B.1.1 Motivation

The super-resolution problem (\mathcal{P}_{sr}) consists in estimating a high resolution (HR) image from its corresponding low resolution (LR) counterpart. \mathcal{P}_{sr} finds a wide range of applications and has attracted much attention within the community of computer vision [NM14, YYDN07, ZY12]. Generally, the considered optimization objective of supervised

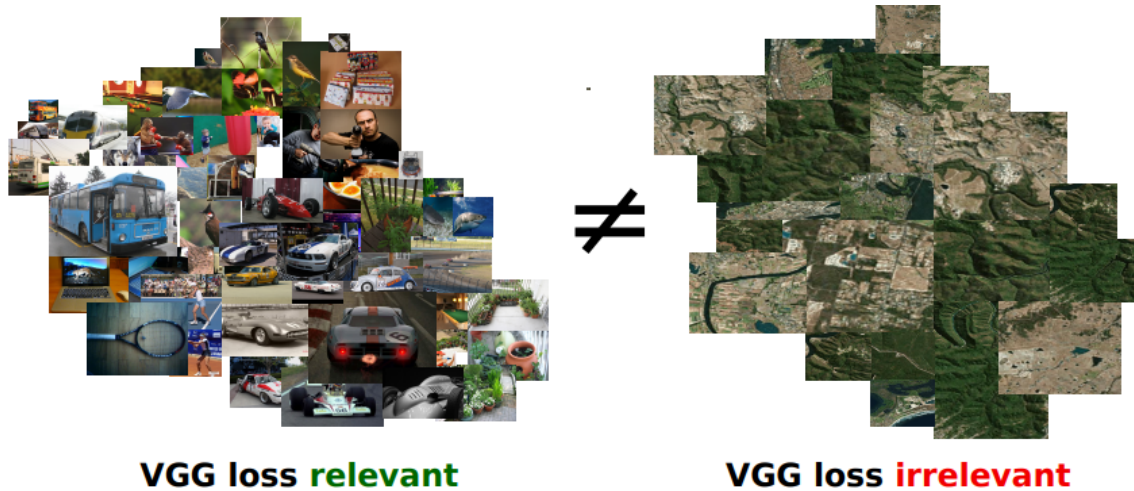


Figure B.1: When super-resolving images from a different domain (e.g., satellite images on the right) than the ImageNet domain (e.g., general objects on the left), the VGG loss introduced by [LTH⁺16] is no longer relevant. We propose a method that outperforms the SRGAN method [LTH⁺16] when super-resolving satellite images. Our method falls within a large class of methods which constitutes our proposed *Generative Collaborative Networks* framework.

methods to solve \mathcal{P}_{sr} is the minimization of the mean squared error (MSE) between the recovered HR image and ground-truth. This class of methods are known to be suboptimal to reconstruct texture details at large upscaling factors. In fact, since MSE consists in a pixel-wise images differences, its ability to recover high texture details is limited [LTH⁺16, GSBB11, WBS04, WSB03]. Furthermore, the minimization of MSE maximizes the Peak Signal-to-Noise-Ratio (PSNR) metric, which is commonly used for the evaluation of \mathcal{P}_{sr} methods [YMY14].

In order to correctly recover finer texture details when super-resolving at large upscaling factors, a recent (state-of-the-art) work [LTH⁺16] defined a perceptual loss which is a combination of an adversarial loss and a VGG loss. The former encourages solutions perceptually hard to distinguish from the HR ground-truth images, while the latter consists in using high-level feature maps of the VGG network [SZ14] pre-trained on ImageNet [DDS⁺09]. When considering the problem of super-resolving images from a target-domain *different* than ImageNet (e.g., satellite images), the features produced by the pre-trained VGG network on the source domain (ImageNet) are suboptimal and no longer relevant for the target domain. In fact, transfer-learning methods are known to be efficient only when the source and target domains are close enough [TLBH17a, TLBH⁺17b, KQD18]. In this work, we present a general framework which we call *Generative Collaborative Networks* (GCN), where the main idea consists in optimizing the generator (i.e., the mapping of interest) in the feature space of a network which we shall refer to as a *features extractor* network. The two networks are said to be *collaborative* in the sense that the features extractor network “helps” the generator by constructing (here, learning) relevant features. In particular, we applied our framework to the problem of single image super-resolution, and we demonstrated that it results in a method that is more adapted (compared to SRGAN [LTH⁺16]) when super-resolving images from a domain that is “far” from the ImageNet domain.

The problem of super-resolution has been tackled with a large range of approaches. In the following, we will consider the problem of *single* image super-resolution (\mathcal{P}_{sisr}) and thus the approaches that recover HR images from multiple images [BS98b, FREM04] are out of the scope of this paper. First approaches to solve \mathcal{P}_{sisr} were filtering-based methods (e.g., linear, bicubic or Lanczos [Duc79] filtering). Even if these methods are generally very fast, they usually yield overly smooth textures solutions [WBSS04]. Most promising and powerful approaches are learning-based methods which consist in establishing a mapping between LR images and their HR counterparts (supposed to be known). Initial work was proposed by Freeman *et al.* [FJP02]. This method has been improved in [DZSW11, ZEP10] by using compressed sensing approaches. Patch-based methods combined with machine learning algorithms were also proposed: in [TDVG13, TDSVG14] upsampling a LR image by finding similar LR training patches in a low dimensional space (using neighborhood embedding approaches) and a combination of the HR patches counterparts are used to reconstruct HR patches. A more general mapping of example pairs (using kernel ridge regression) was formulated by Kim and Kwon [KK10]. Similar approaches used Gaussian process regression [HS11], trees [SPP15] or Random Forests [SLB15] to solve the regression problem introduced in [KK10]. An ensemble method-based approach was adopted in [DTV15] by learning multiple patch regressors and selecting the most relevant ones during the test phase.

Convolutional neural networks (CNN)-based approaches outperformed other \mathcal{P}_{sisr} approaches, by showing excellent performance. Authors in [WLY⁺15] used an encoded sparse representation as a prior in a feed-forward CNN, based on the learned iterative shrinkage and thresholding algorithm of [GL10]. An end-to-end trained three layer deep fully convolutional network, based on bicubic interpolation to upscale the input images, was used in [DLHT14, DLHT16] and achieved good \mathcal{P}_{sisr} performances. Further works suggested that enabling the network to directly learn the upscaling filters, can remarkably increase performance in terms of both time complexity and accuracy [DLT16, SCH⁺16]. In order to recover visually more convincing HR images, Johnson *et al.* [JAFF16] and Bruna *et al.* bruna2015super used a closer loss function to perceptual similarity. More recently, authors in [LTH⁺16] defined a perceptual loss which is a combination of an adversarial loss and a VGG loss. The latter consists in minimizing the error between the recovered HR image and ground-truth in the high-level feature space of the pre-trained VGG network [SZ14] on ImageNet [DDS⁺09]. This method notably outperformed CNN-based methods for the problem \mathcal{P}_{sisr} .

B.1.2 Proposed Methods

B.1.2.1 Proposed Framework

Consider a problem \mathcal{P} of learning a mapping function \mathcal{F} , parameterized by $\theta_{\mathcal{F}}$, that transforms images from a domain \mathcal{X} to a domain \mathcal{Y} , given a training set of N pairs $\{(x_i, y_i)\}_{i=1}^N \in \mathcal{X} \times \mathcal{Y}$. Denote by $p_{\mathcal{X}}$ and $p_{\mathcal{Y}}$ the probability distributions respectively over \mathcal{X} and \mathcal{Y} . In addition, we introduce a given *features extractor* function denoted Φ , parameterized by θ_{Φ} , that maps an image $y \in \mathcal{Y}$ to a certain euclidean feature space \mathcal{S}_{Φ} of dimensionality d . The mappings \mathcal{F} and Φ are typically feed-forward Convolutional Neural Networks. The Generative Collaborative Networks (GCN) framework consists in learning the mapping function \mathcal{F} by minimizing a given loss function¹ in the space of

¹ ℓ_2 -loss is considered in the following.

features \mathcal{S}_Φ , between the generated images (through \mathcal{F}) and ground-truth. Formally,

$$\hat{\theta}_{\mathcal{F}} = \arg \min_{\theta_{\mathcal{F}}} \frac{\lambda_1}{N d} \sum_{i=1}^N \sum_{j=1}^d (\Phi_j(y_i) - \Phi_j(\mathcal{F}(x_i)))^2 + \lambda_2 \Omega(\theta_{\mathcal{F}}), \quad (\text{B.1})$$

where $\Omega(\theta_{\mathcal{F}})$ is a certain regularization term (detailed below) on the weights $\theta_{\mathcal{F}}$ and λ_1 and λ_2 are summation coefficients. The two networks \mathcal{F} and Φ are collaborative in the sense that, the latter learns specific features of the domain \mathcal{Y} and “helps” the former, as it is learned in the space \mathcal{S}_Φ . An important question arises about how to learn the mapping Φ . In following, we describe different classes of methods depending on the learning strategy of Φ . In fact, the features extractor function Φ can take different forms and be learned by different strategies. In particular, we distinguish two learning strategies (illustrated in Figure B.2), which we shall call *disjoint-learning* and *joint-learning*. The four following cases belong to the *disjoint-learning* strategy:

- (1.a) When Φ is the *identity operator* ($\Phi = \text{Id}$). In that case, the objective in equation B.1 becomes a simple pixel-wise MSE loss function. We refer to this class of methods by \mathcal{P}/mse .
- (1.b) When Φ corresponds to a *random feature map* neural network, that is to say, the weights θ_Φ are set randomly according to a given distribution μ . We refer to this class of methods by \mathcal{P}/ran .
- (1.c) When Φ is a part of a model that solves a *reconstruction problem* (jointly with an auxiliary mapping function $\Psi : \mathcal{S}_\Phi \rightarrow \mathcal{Y}$), by minimizing the pixel-wise ℓ_2 -loss function between the reconstructed images (through Ψ) and ground-truth:

$$(\hat{\theta}_\Phi, _) = \arg \min_{(\theta_\Phi, \theta_\Psi)} \frac{1}{N \dim(\mathcal{Y})} \sum_{i=1}^N \sum_{j=1}^{\dim(\mathcal{Y})} ((y_i)_j - (\Psi \circ \Phi(y_i)))_j^2. \quad (\text{B.2})$$

Notably, this strategy allows for the learning of reconstruction features which are different from classification-based features. We refer to this class of methods by \mathcal{P}/rec .

- (1.d) When Φ is trained to solve a *multi-label classification problem* [LTH⁺16], that is to say, when labels are available for the domain \mathcal{Y} . More precisely, it exists a dataset $\{(y_i, c_i)\}_{i=1}^n \in \mathcal{Y} \times \{1, \dots, m\}$ of n images labelled among m classes and Φ is learned to minimize the following objective:

$$(\hat{\theta}_\Phi, _) = \arg \max_{(\theta_\Phi, \theta_\Psi)} \mathbb{P} \{ \Psi \circ \Phi(y_i) = c_i \mid y_i; i \in \{1, \dots, m\} \}, \quad (\text{B.3})$$

where $\Psi : \mathcal{S}_\Phi \rightarrow \{1, \dots, m\}$. We refer to this class of methods by \mathcal{P}/cla .

The features extractor function Φ can also be trained *jointly* with the desired mapping function \mathcal{F} . Indeed, as in the GANs paradigm, one can use a discriminator to distinguish the generated images (through \mathcal{F}) and ground-truth, and thus learn more relevant and specific features for the problem of interest \mathcal{P} . In particular, the *joint-learning* strategy contains two cases:

- (2.a) When Φ is a part of a *discriminator*. $\mathcal{D} = \Psi \circ \Phi : \mathcal{Y} \rightarrow \{0, 1\}$ that classifies the generated images (through \mathcal{F}) and ground-truth. \mathcal{D} is optimized in an alternating manner along with \mathcal{F} to solve the adversarial min-max problem [SCT⁺16]:

$$\min_{\theta_{\mathcal{F}}} \max_{(\theta_\Phi, \theta_\Psi)} \mathbb{E}_{y \sim p_y} [\log \Psi \circ \Phi(y)] + \mathbb{E}_{x \sim p_x} [\log \{1 - \Psi \circ \Phi \circ \mathcal{F}(x)\}]. \quad (\text{B.4})$$

Standard methods	\mathcal{P}/mse	\mathcal{P}/cla	\mathcal{P}/rec	\mathcal{P}/dis	$\mathcal{P}/dis,rec$
Existence	✓[GSBB11]	✓[DB16]	✗	✗	✗
Adversarial methods	$\mathcal{P}/adv,mse$	$\mathcal{P}/adv,cla$	\mathcal{P}/adv	$\mathcal{P}/adv,rec$	
Existence	✓[YP16]	✓[LTH ⁺ 16]	✗	✗	

Table B.1: Existent loss functions of the proposed GCN framework.

The adversarial loss (second term of equation B.4) can thus be seen as a regularization of the parameters $\theta_{\mathcal{F}}$ by affecting this quantity to $\Omega(\theta_{\mathcal{F}})$ in equation B.1. This regularization “pushes” the solution of the problem in equation B.1 to the manifold of the images in the domain \mathcal{Y} . We refer to this class of methods by \mathcal{P}/adv . When $\lambda_2 = 0$, we refer to it by \mathcal{P}/dis .

- (2.b) When Φ is a part of a *discriminator* and an *auto-encoder*. Namely, by optimizing its weights θ_{Φ} to solve simultaneously, an *adversarial problem* as in equation B.4; through $\mathcal{D} = \Psi_1 \circ \Phi : \mathcal{Y} \rightarrow \{0, 1\}$, and a *reconstruction problem* as in equation B.2; through a mapping $\Psi_2 : \mathcal{S}_{\Phi} \rightarrow \mathcal{Y}$. We refer to this class of methods by $\mathcal{P}/adv,rec$ or $\mathcal{P}/dis,rec$ depending on the value of λ_2 in equation B.1.

B.1.2.2 Existing Loss Functions

The natural way to learn a mapping from a manifold to another is to use \mathcal{P}/mse methods. It is well known [GSBB11, LTH⁺16, WSB03, WBSS04] that this class of methods lead to overly-smooth and poor perceptual quality solutions. In order to handle the mentioned perceptual quality limitation, a variety of methods have been proposed in the literature. First methods used generative adversarial networks (GANs) for generating high perceptual quality images [DCF⁺15, MCL15], style transfer [LW16] and inpainting [YCL⁺16], namely the class of methods \mathcal{P}/adv with $\lambda_1 = 0$. Authors in [YP16] proposed to use \mathcal{P}/mse with an adversarial loss ($\lambda_1 > 0$ and $\lambda_2 > 0$) to train a network that super-resolves face images with large upscaling factors. Authors in [BSL15, JAFF16] and in [DB16] used \mathcal{P}/cla by considering respectively $\Phi = \text{VGG19}$ and $\Phi = \text{AlexNet}$ networks as fixed features extractors (learned *disjointly* from the mapping of interest), which result in a more perceptually convincing results for both super-resolution and artistic style-transfer [GEB15, GEB16]. More recently, authors in [LTH⁺16] used $\mathcal{P}/cla,adv$ by considering $\Phi = \text{VGG19}$ as a fixed features extractor combined with an adversarial loss ($\lambda_2 > 0$). To the best of our knowledge, as summarized in table B.1, the use of the other learning strategies of Φ ; namely (1.c), (2.a) and (2.b), have not been explored in the literature. We particularly apply these strategies in the context of Single Image Super-Resolution, which results in methods that are more suitable (comparing to the SRGAN method [LTH⁺16]) to super-resolution domains that differ from the ImageNet domain. The proposed methods as well as the corresponding experiments are presented in the following section.

B.1.3 Experiments for Single Image Super-Resolution

B.1.3.1 Proposed Methods

In this section, we consider the problem of Single Image Super-Resolution (\mathcal{P}_{sizr}). In particular, we suppose we are given N pairs $\{(I_i^{LR}, I_i^{HR})\}_{i=1}^N$ of low-resolution images and their high-resolution counterparts. Recalling our GCN framework (presented in the previous section) the proposed methods for the problem \mathcal{P}_{sizr} are: \mathcal{P}_{sizr}/rec , \mathcal{P}_{sizr}/dis , $\mathcal{P}_{sizr}/dis,rec$, \mathcal{P}_{sizr}/adv and $\mathcal{P}_{sizr}/adv,rec$. We show in the following that the most convincing results are given by $\mathcal{P}_{sizr}/adv,rec$. In particular, we show on a dataset of satellite

images (different from the ImageNet domain) that our method $\mathcal{P}_{sisr}/adv,rec$ outperforms the SRGAN method [LTH⁺16] by a large margin on the considered domain. Note that, as our goal is to show the irrelevance of the VGG loss for some visual domains (different from ImageNet), we do not consider the well-known SR benchmarks (*e.g.*, Set5, Set14, B100, Urban100) for the evaluation, as these benchmarks are relatively close to the ImageNet domain. The evaluation of the different methods is based on *perceptual metrics* [ZIE⁺18] which we recall in the following section.

B.1.3.2 Evaluation Metrics

The evaluation of super-resolution methods (more generally image regression-based methods) requires comparing visual patterns which remains an open problem in computer vision. In fact, classical metrics such as L2/PSNR, SSIM and FSIM often disagree with human judgments (*e.g.*, blurring causes large perceptual change but small L2 change). Thus, the definition of a *perceptual metric* which agrees with humans perception is an important aspect for the evaluation of \mathcal{P}_{sisr} methods. Zhang *et al.* [ZIE⁺18] recently evaluated deep features across different architectures (Squeeze [IHM⁺16], AlexNet [KSH12] and VGG [SZ14]) and tasks (supervised, self-supervised and unsupervised networks) and compared the resulting metrics with traditional ones. They found that deep features outperform all classical metrics (*e.g.*, L2/PSNR, SSIM and FSIM) by large margins on their introduced dataset. As a consequence, deep networks seem to provide an embedding of images which agrees surprisingly well with humans judgments.

Zhang *et al.* [ZIE⁺18] compute the distance between two images x, y with a network² Φ in the following way:

$$d_{\Phi}(x, y) = \sum_l \frac{1}{H_l W_l} \sum_{h,w} \|w_l \odot (\Phi^l(x)_{hw} - \Phi^l(y)_{hw})\|_2^2, \quad (\text{B.5})$$

where $\Phi^l(\cdot)$ are the extracted features from layer l and unit-normalized in the channel dimension. w_l is a re-scaling vector of the activations channel-wise at layer l . H_l and W_l are respectively the height and width of the l^{th} feature map.

Thus, we compute the *perceptual error* (PE) of a \mathcal{P}_{sisr} method (a mapping \mathcal{F}) on a given test-set of N low-resolution images and their high-resolution counterparts $\Pi = \{(I_i^{LR}, I_i^{HR})\}_{i=1}^N$ as the mean distances between the generated images (through \mathcal{F}) and ground-truth as follows:

$$\text{PE}_{\Phi}(\Pi) = \frac{1}{N} \sum_{i=1}^N d_{\Phi}(\mathcal{F}(I_i^{LR}), I_i^{HR}). \quad (\text{B.6})$$

Note that we use the implementation of [ZIE⁺18] to compute the perceptual distances $d_{\Phi}(\cdot, \cdot)$ using six variants which are based on the networks Squeeze [IHM⁺16], AlexNet [KSH12] and VGG [SZ14] and their “perceptual calibrated” versions. The best method is considered to be the one which minimizes the maximum amount of PEs across different networks $\Phi \in \{\text{Squ}, \text{Squ-l}, \text{Alex}, \text{Alex-l}, \text{VGG}, \text{VGG-l}\}$.

²The considered networks are Squeeze[IHM⁺16], AlexNet[KSH12] and VGG[SZ14] and their “perceptual calibrated” versions which we refer to respectively as Squeeze-l, AlexNet-l and VGG-l. See [ZIE⁺18] and the provided Github project within for further details.

B.1.3.3 Experiments

The overall goal of this section is to validate our statement about the relevance of the VGG loss when super-resolving images from a different domain than the ImageNet domain. To highlight this aspect, we first present the considered datasets, architectures and training details. Then we select the more appropriate method (across the GCN framework methods) for the \mathcal{P}_{sirs} problem based on perceptual metrics [ZIE⁺18]. Finally, we compare our proposed method to some baselines and the state-of-the-art SRGAN method [LTH⁺16], on three different datasets (detailed in the following section). We show in particular that our method outperforms SRGAN on the satellite images domain.

Datasets. The idea of replacing the MSE pixel-wise content loss on the image by a loss function that is closer to perceptual similarity is not new. Indeed, [LTH⁺16] defined a VGG loss on the feature map obtained by a specific layer of the pre-trained VGG19 network and shows that it fixes the inherent problem of overly smooth results which comes with the pixel-wise loss. Nevertheless, VGG19 being trained on ImageNet, their method would not perform particularly well on different images, the distribution of which is far away from that of ImageNet.

Therefore, we propose a similar method where the difference is that our features extractor is not pre-trained, but trained jointly with the generator. This removes the aforementioned limitation since the features extractor is trained on the same dataset as the generator and thus extract relevant features.

To show that, we trained our different networks (*i.e.*, with different features extractors) on three distinct datasets (examples of images of these datasets are shown in Figure B.3):

- A subset of *ImageNet* [DDS⁺09], for which we sampled 70,000 images. Since VGG19 was trained on ImageNet for many (more than 300K) iterations, we expect to have similar or worse results than the state-of-the-art method SRGAN from [LTH⁺16] on this database.
- *The Describable Textures Dataset (DTD)* [CMK⁺14], containing 5,600 images of textural patterns. These data are relatively close to ImageNet and we show that our method gives convincing results relatively close to SRGAN.
- A dataset containing satellite images³, which we generated by randomly cropping 256×256 images on a 7205×7205 satellite image which result in 235,183 images. We particularly show that our method significantly outperforms SRGAN on this dataset. We refer to this dataset by *Sat*.

All experiments are performed with a scale factor of $4\times$ between low- and high-resolutions images and the formers are obtained during the training by down-scaling the original images by a factor $1/4$.

Architectures. Our overall goal is to prove that the proposed GCN framework, is adapted to train a generative mapping model and that it surpasses the MSE loss in keeping perceptual similarity in the generated image (whereas the MSE loss tends to smooth things

³Can be found in http://www.terracolor.net/sample_imagery.html

out and lose high frequency details). As opposed to [LTH⁺16]’s work, our framework does not require to have a pre-trained network, like VGG, to extract helpful features for training. In this paper, we focus on the Super Resolution problem. Therefore, we chose our mapping function \mathcal{F} , or generator, to be that of Ledig *et al.* [LTH⁺16]: a feed-forward CNN parametrized by $\theta_{\mathcal{F}}$, composed of 10 residual blocks. These blocks are made of two convolutional layers with 3×3 kernels and 64 features maps, each followed by batch normalization and PReLU as activation. The image’s size is then increased of a factor 4 by two trained upsamplings. The architecture of all the used discriminators follows the guidelines of Radford *et al.* [RMC15] as it is composed of convolutional layers, followed by a batch normalization and a LeakyReLU ($\alpha = 0.2$) activation. This block is repeated eight times and each time the number of 3×3 kernels increases by a factor 2 (ranging from 64 to 512), a strided convolution is used to reduce the image resolution by 2. Two dense layers and a sigmoid activation then return the discrimination probability. In the case of an auto-encoder (every *Reconstruction* problem), we follow the same architecture for the encoder and a symmetric one for the decoder. Figure B.4 depicts an overview of the architectures for both the generator and the discriminator.

Training details and parameters. All networks were trained⁴ on a NVIDIA GeForce GTX 1070 GPU using the considered datasets, which do not contain the (1000) testing images shown as results. We scaled the range of both the LR input images and the HR images to $[-1, 1]$, which explains the tanh activation for the last layer of the generator. All variants of our networks, which differ in their features extractor, were trained from scratch (for the generator and the features extractor) with mini batches of 10 images. We used the Adam optimizer with a learning rate of $2 \cdot 10^{-4}$ and a decay of 0. The generator and the feature extractor are updated alternatively. As we realized training was stable and quite fast, we trained with only 5,000 update iterations to pinpoint the best method among the different GCNs. Finally, the regularization parameters in our global loss are set by default as $\lambda_1 = 1$ and $\lambda = 10^{-3}$. As a reminder, our goal here is, given a generator architecture (or mapping function \mathcal{F}), to find the best strategy to train it, following our GCNs paradigms. The best method is then further compared to baselines.

Features Extractor Selection. As we said above, we investigated the ability of different features extractor to construct relevant perceptual feature maps for training and improving the rendering quality of the generator. In order to select the best learning strategy given a certain dataset, we train the generator on each dataset using the different learning strategies: \mathcal{P}_{sirs}/rec , \mathcal{P}_{sirs}/dis , $\mathcal{P}_{sirs}/dis,rec$, \mathcal{P}_{sirs}/adv and $\mathcal{P}_{sirs}/adv,rec$. Note that, the features extractor for all the considered methods correspond to the first layer of the discriminators (or encoder-decoders). In fact, as the problem \mathcal{P}_{sirs} consists in recovering low-level perceptual cues, we limited our study to the first layer.

Table B.2 summarizes the results of the proposed \mathcal{P}_{sirs} methods in terms of low-level metrics (L2 and SSIM) and perceptual metrics [ZIE⁺18] which are given by Eq. equation B.6. We notice from this table that the method $\mathcal{P}_{sirs}/adv,rec$ performs relatively well on the datasets ImageNet and Sat in terms of perceptual metrics. While $\mathcal{P}_{sirs}/dis,rec$ gives better results on the DTD dataset. The main difference between these two methods is that the former considers an adversarial loss on the objective function while the latter does not consider the adversarial term. This explains the reason why $\mathcal{P}_{sirs}/adv,rec$ does

⁴A Keras implementation is provided in <https://github.com/melaseddik/GCN>

not perform well on DTD. In fact, texture images belong to a complex manifold and their distribution is relatively hard to fit by a generative model.

Figure B.5 shows qualitative results of the different proposed methods on the different presented datasets. Generally, the methods which were trained with an additional adversarial loss (\mathcal{P}_{sisr}/adv and $\mathcal{P}_{sisr}/adv,rec$) output images of higher quality (on the datasets ImageNet and Sat) as GANs were introduced to do just so: generate images that follow the distribution of the dataset. Among these two *adversarial* methods, it seems to us (as suggested by the quantitative results of table B.2) that $\mathcal{P}_{sisr}/adv,rec$ (column (c) of Figure B.5) is able to detect and render more details, due to its ability to generate more relevant features as the features extractor Φ is learned to solve a *multi-task* problem; namely a *discrimination* and a *reconstruction* problem, in particular, this method allows for the learning of both classification and reconstruction-based features. We will thus further investigate the $\mathcal{P}_{sisr}/adv,rec$ method for the comparison to the baseline and the state-of-the-art method SRGAN [LTH⁺16], on the satellite images domain.

$\mathcal{P}_{sisr}/adv,rec$ against baseline methods on the satellite images domain. Our main objective is to show that the VGG loss function (namely, the SRGAN method [LTH⁺16]) is no longer relevant when super-resolving images from a domain different than the ImageNet domain. In particular, by considering the satellite images domain, we show in this section that the selected method from the previous section ($\mathcal{P}_{sisr}/adv,rec$) outperforms some baselines, which are \mathcal{P}_{sisr}/mse (pixel-wise MSE loss) and $\mathcal{P}_{sisr}/adv,mse$ (pixel wise MSE loss combined with an adversarial loss), and the state-of-the-art super-resolution method, SRGAN [LTH⁺16]. Note that all the methods use the same architectures (depicted in figure B.4) for the generator and discriminator and are trained on the same domain (here, on satellite images). Our purpose being to show the relevance of the proposed method on a domain “far” from the ImageNet domain, we do not consider standard SR benchmarks, which are relatively “close” to the ImageNet domain.

Table B.3 presents quantitative results, in terms of classical metrics (L2 and SSIM) and perceptual metrics given by equation B.6, of the different methods on the Sat dataset. As we can notice, our method $\mathcal{P}_{sisr}/adv,rec$ outperforms the other methods in terms of perceptual metrics. Knowing that the perceptual metrics agree with human judgments [ZIE⁺18], these results validate the effectiveness of the $\mathcal{P}_{sisr}/adv,rec$ method. Note also that even if SRGAN [LTH⁺16] is optimized to minimize a VGG loss, it does not give the lowest perceptual errors in terms of the perceptual metrics VGG and VGG-l, this is due to the fact that the VGG features are not relevant for the satellite images domain. In addition, $\mathcal{P}_{sisr}/adv,rec$ gives the lowest perceptual errors in terms of the perceptual metrics Alex and Alex-l which agrees with a human perception. In fact, AlexNet network may more closely match the architecture of the human visual cortex [YD16].

Figure B.6 shows some qualitative results of different methods on a patch of an image from the Sat dataset. As we can notice, the $\mathcal{P}_{sisr}/adv,rec$ method gives the perceptually closest result to the ground-truth image, which agrees with the quantitative results of table B.3.

Further results. In this section, we provide further qualitative and quantitative comparisons to the considered baselines of the previous section. In particular, we consider

all the presented datasets for the comparisons. Qualitative results are provided in figure B.7. SRGAN performs better on ImageNet, which is not that surprising considering our features extractor was trained much less than VGG19 used in [LTH⁺16] and the VGG features being more relevant for images from the ImageNet domain. Nonetheless, we do have sharper images than the MSE based methods, although we show some artifact (especially on the boat) which we attribute to the competition between the content and adversarial losses. On DTD though, we can see the benefit of our method over a pre-trained VGG loss. Indeed, SRGAN is blurrier on both the house (first row) and the cliff (third row), in spite of having less artifacts than our method. On the “cracks” example (second row), SRGAN even totally obliterates the details in the center. Finally, results on the dataset Sat, which is the most different dataset compared to ImageNet, are the most compelling. Our method generates super resolved images that are really close to the real high resolution images, while we can clearly see imperfections on SRGAN’s results because of VGG19 which was not trained to detect perceptual features on satellite images.

Quantitative results are summarized in Table B.4. As shown in [LTH⁺16, ZIE⁺18], the standard quantitative measures such as L2 and SSIM fail to highlight image quality according to the human visual system. In fact, while the results of \mathcal{P}_{sirs}/mse are overly smooth perceptually, it has the lowest L2 and SSIM errors on Sat. However, perceptual metrics agree with what we assess qualitatively: SRGAN performs best on ImageNet but not on Sat, the distribution of which is the farthest from ImageNet. Actually, SRGAN ranks third of all four methods on Sat, just before $\mathcal{P}_{sirs}/adv,mse$, while still performing best on DTD which still is pretty close to ImageNet. This shows that the VGG features become less and less relevant as the dataset’s distribution part from ImageNet. On the other hand, our training framework allows to construct relevant features on any (never seen) dataset. Thus our method $\mathcal{P}_{sirs}/adv,rec$ performs best on Sat. Our method performing better than $\mathcal{P}_{sirs}/adv,mse$ also shows that our framework helps finding detail preserving features. Figure B.7 provides the results of the different baselines and our method on some examples of the considered datasets. We notice from these images that our method $\mathcal{P}_{sirs}/adv,rec$ recovers finer details on the different datasets while it outperforms the considered baselines on satellite images.

Table B.5 summarizes the results of the different methods on the considered datasets through the paper. From these results, we make the following conclusions:

- When the considered domain is far enough from the ImageNet domain, the VGG loss introduced by [LTH⁺16] is no longer relevant.
- The VGG network can not be fine-tuned when considering a domain for which there is no available labels for the images (*e.g.*, satellite images). Thus, the SRGAN method cannot be exploited efficiently in this case.
- Our framework results in a method ($\mathcal{P}_{sirs}/adv,rec$) that outperforms some baselines and the SRGAN method on the satellite images domain.
- Even on a domain close to the ImageNet domain (*e.g.*, texture images), one can find within our framework methods which give almost similar results to the SRGAN method, while the later is based on VGG features and thus need to train the VGG network on the whole ImageNet dataset.

B.1.4 Central Contribution and Discussions

In this work, we have proposed a general framework named Generative Collaborative Networks (GCN) which generalizes the existing methods for the problem of learning a mapping between two domains. The GCN framework highlights that there is a learning strategy of mappings that is not explored in the literature. In particular, the optimization of these mappings in the feature space of a features extractor network, which is mutually learned at the same time as the considered mapping (*joint-learning* strategy). The GCN framework was evaluated in the context of super-resolution on three datasets (ImageNet [DDS⁺09], DTD [CMK⁺14] and satellite images). We have shown that the proposed *joint-learning* strategy leads to a method that outperforms the state of the art [LTH⁺16] which uses a pre-trained features extractor network (VGG19 on ImageNet). Specifically, this holds when the domain of interest is “far” from the ImageNet domain (e.g., satellite images or images from the medical domain⁵). However, note that even for domains close to the ImageNet domain, the proposed method gives convincing (almost similar to [LTH⁺16]) results without using the whole ImageNet dataset to learn the features extractor network (as performed in [LTH⁺16]).

We systematically designed the proposed methods by using the first layer of the features extractor networks, while it could be interesting to evaluate in more detail the impact of this choice regarding the learning strategy. Moreover, the impact of the selected layer may also depend on the considered dataset. More generally, the GCN framework offers a large vision on the wide variety of existing loss functions used in the literature of learning mappings-based problems (e.g., super-resolution, image completion, artistic style transfer, etc.). In fact, we show that these loss functions can be simply reformulated, in the proposed framework, as a certain combination of a particular type of features extractor networks (\mathcal{P}/rec , \mathcal{P}/dis , $\mathcal{P}/dis,rec$, \mathcal{P}/adv and $\mathcal{P}/adv,rec$) and a particular learning strategies (*joint-learning* or *disjoint-learning*). Therefore it will be interesting to explore this promising framework in other learning mappings-based problems.

B.2 Neural Networks Compression

This section is based on the following work:

- (C7) MEA. Seddik, H.Essafi, A.Benzine, M.Tamaazousti, “Lightweight Neural Networks from PCA & LDA Based Distilled Dense Neural Networks”, International Conference on Image Processing (ICIP’20), Online, 2020.

B.2.1 Motivation

Neural networks are the most effective machine learning methods nowadays, and since they generally require millions of parameters, their implementation in an IoT environment is quite ineffective. Indeed, IoT requires machine learning models with limited amount of parameters since the edge devices have limited resources (computation capacity, storage, bandwidth, . . .). This requirement among others have triggered intensive research activities and led to the emergence of new computing paradigms, i.e. edge computing [Sat17] which has emerged as an answer to the need for shifting the computing

⁵This domain is particularly relevant for the proposed framework as it seems very far from the ImageNet domain. Unfortunately, we have not found a big amount of publicly available data (to the best of our knowledge) for medical images which prevented us from considering this domain through the paper.

from cloud to decentralized processing units close to the data sources [ASEA19]. The authors in [MMH⁺19] give a survey of works dealing with machine learning at the network edge. AI at the edge [LLP⁺19] is a concrete example of leveraging the computing and storage resources near the places where data is produced. The authors in [VTWA18] present a set of approaches proposed for embedding deep learning into the edge computing devices. They also present some applications that can fit with the edge computing paradigm and can take benefit from the edge network. However, the accuracy of deep learning models depends largely on the hyper-parameters of the network in particularly the number and the size of layers.

Nevertheless, big models are resources consuming which can impede the embedding of AI technologies in IoT devices with constrained resources. To tackle this issue some interesting solutions were proposed, model compression was among the first proposed methods. Model compression methods use well proved compression techniques for reducing the storage volume required by neural networks without impacting their performance. For instance the algorithm presented in [HMD15] is composed of three methods applied in pipeline: first pruning (selection of the important weights of the network) method is applied, followed by quantization and Huffman coding (compression without lost). Compression methods are useful for reducing the required storage space of network models but inefficient for reducing the computation power which is one of the main critical aspects of the IoT devices. Recent methods targeting the production of smaller models, with the accuracy near of lager ones, was investigated [MFLG19].

In this study, we present two methods for distilling a given large dense neural network into a smaller one. The proposed methods are based on knowledge distillation concept [YJBK17], where a large (teacher) pre-trained network is used to train a smaller (student) network. They have the advantage to produce models that consume less storage and computing resources; the knowledge distillation approach is a kind of transfer learning approach which is commonly used in various machine learning problems [TLBH⁺19, MTS⁺19]. In [MFLG19] the authors show that the accuracy of the student network model depends highly on the ratio size between the two network models (teacher and student), higher is this ratio (gap size between teacher and student is large), smaller is the accuracy. To alleviate this problem the authors suggest, instead of distilling the student model directly from the teacher model, to use succession of teacher assistants based knowledge distillation approach where the models are distilled step by step until obtaining the final model. The proposed methods in this paper are notably complementary to the methods in [HMD15, MFLG19].

B.2.2 Proposed Methods

B.2.2.1 Setting & Notations

Consider a dense neural network, which we refer to as the teacher network (TN), composed of L layers and constructed in the following way, for $\ell \in [L]$:

$$\text{(TN)} : \begin{cases} \mathbf{h}^{(0)} = \mathbf{x} \in \mathbb{R}^{p_0}, \\ \mathbf{h}^{(\ell)} = f_{\ell} \left(\mathbf{W}^{(\ell)} \mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)} \right) \in \mathbb{R}^{p_{\ell}}, \end{cases} \quad (\text{B.7})$$

where, \mathbf{x} corresponds to the input data features, f_{ℓ} denotes the ℓ -th layer activation, $\mathbf{h}^{(\ell)}$ stands for the features of \mathbf{x} extracted at layer ℓ , $\mathbf{W}^{(\ell)} \in \mathbb{R}^{p_{\ell} \times p_{\ell-1}}$ and $\mathbf{b}^{(\ell)} \in \mathbb{R}^{p_{\ell}}$ are respectively the weight matrix and bias at each layer ℓ . TN is typically of large size, meaning

that the hidden features dimensions p_ℓ are relatively large. In the following, we will present two methods that construct a small network size, which we refer to as the student network (SN), based on the TN learned features. The two methods target different learning problems, depending if TN solves a supervised problem or an unsupervised one.

B.2.2.2 Neural Nets PCA-based Distillation (Net-PCAD)

Given a set of n training samples $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p_0 \times n}$ on which TN was initially trained to perform some arbitrary learning task. The Net-PCAD method consists in performing a PCA [STC19b] at each hidden layer of TN, and then training a SN to perform the same learning task as TN along with the task of mapping its hidden features at each layer with the reduced features of TN. Formally, we denote by $\mathbf{H}_\ell = [\mathbf{h}_1^{(\ell)}, \dots, \mathbf{h}_n^{(\ell)}] \in \mathbb{R}^{p_\ell \times n}$ where $\mathbf{h}_i^{(\ell)}$ stands for the features of \mathbf{x}_i at layer ℓ . Therefore, a PCA is performed at each layer ℓ in order to reduce the dimension of the hidden features p_ℓ , relying on the top k_ℓ largest eigenvectors of the sample covariance matrix:

$$\mathbf{C}_\ell = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{h}}_i^{(\ell)} \bar{\mathbf{h}}_i^{(\ell)\top} \quad (\text{B.8})$$

where $\bar{\mathbf{h}}_i^{(\ell)} = \mathbf{h}_i^{(\ell)} - \frac{1}{n} \sum_{j=1}^n \mathbf{h}_j^{(\ell)}$ are the centred hidden features. We denote by $\mathbf{U}_\ell \in \mathbb{R}^{p_\ell \times k_\ell}$ the matrix containing the k_ℓ largest eigenvectors of \mathbf{C}_ℓ . Consequently, the student network (SN) is composed of L layers and has the following structure:

$$(\text{SN}) : \begin{cases} \tilde{\mathbf{h}}^{(0)} = \mathbf{x} \in \mathbb{R}^{p_0}, \\ \tilde{\mathbf{h}}^{(\ell)} = f_\ell \left(\tilde{\mathbf{W}}^{(\ell)} \tilde{\mathbf{h}}^{(\ell-1)} + \tilde{\mathbf{b}}^{(\ell)} \right) \in \mathbb{R}^{k_\ell}, \end{cases} \quad (\text{B.9})$$

with the convention $k_0 = p_0$ and where, \mathbf{x} corresponds to the input data features, $\tilde{\mathbf{h}}^{(\ell)}$ stands for the features of \mathbf{x} extracted at layer ℓ , $\tilde{\mathbf{W}}^{(\ell)} \in \mathbb{R}^{k_\ell \times k_{\ell-1}}$ and $\tilde{\mathbf{b}}^{(\ell)} \in \mathbb{R}^{k_\ell}$ are respectively the weight matrix and bias at each layer ℓ .

Given the initial learning problem of TN which corresponds to a loss function $\mathcal{L}_{\text{problem}}$, SN is therefore optimized with the following loss function, where the Homoscedastic loss [KGC18] is considered since the optimization problem for SN can be formulated as a multi-task problem.

$$\mathcal{L} = e^{-\sigma_{\text{problem}}} \mathcal{L}_{\text{problem}} + \sigma_{\text{problem}} + \sum_{\ell=1}^{L-1} e^{-\sigma_\ell} \mathcal{L}_{\text{mse}}(\tilde{\mathbf{h}}^{(\ell)}, \mathbf{U}_\ell^\top \mathbf{h}^{(\ell)}) + \sigma_\ell \quad (\text{B.10})$$

where \mathcal{L}_{mse} denotes the mean squared error loss function, σ_{problem} and σ_ℓ 's are the Homoscedastic loss parameters which are learned during the training of the student network. A full description of the Net-PCAD method is provided as a pseudo-code algorithm in Algorithm 1.

B.2.2.3 Neural Nets LDA-based Distillation (Net-LDAD)

If the initial learning problem of the TN is a supervised classification problem, one can take advantage of the fact that the data belong to K different classes $\{\mathcal{C}_j\}_{j=1}^K$ and therefore project the hidden features of the TN in structured low-dimensional spaces. Linear Discriminant Analysis (LDA) is a dimension reduction technique that specifically reduces

Algorithm 1: Net-PCAD description.

Input: A trained teacher network **TN**, a data matrix \mathbf{X} and the learning problem loss $\mathcal{L}_{\text{problem}}$.

Output: Trained student network **SN**.

for $\ell \leftarrow 1$ **to** $L - 1$ **do**

- 1. Extract the representations \mathbf{H}_ℓ of \mathbf{X} from **TN**;
- 2. Compute \mathbf{U}_ℓ through a PCA on \mathbf{H}_ℓ ;

end

Train the student network **SN** with \mathcal{L} as in equation B.10;

the dimension of data relying on their classes structure [TGIH17]. LDA is closely related to PCA but differs from the latter by the fact that it explicitly attempts to model the difference between the classes of the data, while PCA does not take into account any difference in class labels. Therefore, the idea behind Net-LDAD is to exploit the labels information layer-wise in the training of the student network from the teacher network. Specifically, we compute at each layer ℓ of the TN the within class scatter matrix as:

$$\mathbf{S}_w^{(\ell)} = \sum_{j=1}^K \sum_{x \in \mathcal{C}_j} (\mathbf{h}_x^{(\ell)} - \mathbf{m}_j^{(\ell)})(\mathbf{h}_x^{(\ell)} - \mathbf{m}_j^{(\ell)})^\top \quad (\text{B.11})$$

where $\mathbf{h}_x^{(\ell)}$ is the representation of \mathbf{x} at layer ℓ of the TN and $\mathbf{m}_j^{(\ell)} = \frac{1}{|\mathcal{C}_j|} \sum_{x \in \mathcal{C}_j} \mathbf{h}_x^{(\ell)}$. And the between class scatter matrix at each layer ℓ is given by:

$$\mathbf{S}_b^{(\ell)} = \sum_{j=1}^K |\mathcal{C}_j| (\mathbf{m}_j^{(\ell)} - \mathbf{m}^{(\ell)})(\mathbf{m}_j^{(\ell)} - \mathbf{m}^{(\ell)})^\top \quad (\text{B.12})$$

where $\mathbf{m}^{(\ell)} = \frac{1}{n} \sum_{x \in X} \mathbf{h}_x^{(\ell)}$. Therefore, the projection matrix of LDA at each layer ℓ is computed as the k_ℓ largest eigenvectors of $(\mathbf{S}_w^{(\ell)})^{-1} \mathbf{S}_b^{(\ell)}$. We denote by $\mathbf{V}_\ell \in \mathbb{R}^{p_\ell \times k_\ell}$ such a projection matrix. Similarly to the PCA case, the student network is therefore trained to minimize the following objective:

$$\begin{aligned} \mathcal{L} &= e^{-\sigma_{\text{classification}}} \mathcal{L}_{\text{classification}} + \sigma_{\text{classification}} \\ &+ \sum_{\ell=1}^{L-1} e^{-\sigma_\ell} \mathcal{L}_{\text{mse}}(\tilde{\mathbf{h}}^{(\ell)}, \mathbf{V}_\ell^\top \mathbf{h}^{(\ell)}) + \sigma_\ell \end{aligned} \quad (\text{B.13})$$

where $\mathcal{L}_{\text{classification}}$ is typically a categorical cross entropy loss function since the initial learning problem of the TN is supposed to be a classification problem. A full description of the Net-LDAD method is provided as a pseudo-code algorithm in Algorithm 2.

B.2.3 Experiments

In this section, we present experiments which highlight the effectiveness of the proposed methods to train student networks that are smaller in size w.r.t. a given large size dense teacher network. In particular, we consider three teacher networks composed of $L = 4$ dense layers which are trained to perform a classification problem respectively on the datasets MNIST [YCC98], Fashion-MNIST [XRV17] and CIFAR10 [KH10]. Therefore, we train the corresponding student networks by successively reducing their hidden dimensions using the presented methods Net-PCAD and Net-LDAD. Note that, for simplicity,

Algorithm 2: Net-LDAD description.

Input: A trained teacher network TN , a data matrix X and the learning problemloss $\mathcal{L}_{\text{problem}}$.**Output:** Trained student network SN .**for** $\ell \leftarrow 1$ **to** $L - 1$ **do**

1. Extract the representations H_ℓ of X from TN ;
2. Compute V_ℓ through a LDA on H_ℓ ;

endTrain the student network SN with \mathcal{L} as in equation B.13;

we reduce all the hidden dimensions to a constant value k and we vary k in all our experiments. Table B.6 presents the considered architectures for the teacher and student networks.

Figure B.8 depicts the training Homoscedastic loss of the student networks for different values of k and across the different considered datasets using our methods Net-PCAD and Net-LDAD. Note from this figure that both methods yield generally to a stable learning of the student networks, and the classification problem gets much easier as k increases. However, note that a careful choice of k must be made in order to get a smooth learning loss (e.g., see $k = 200$ on the Fashion-MNIST dataset).

Figure B.9 depicts the learned Homoscedastic parameters once the student networks have been trained for different values of k . We can observe from this figure (at least for the datasets MNIST and Fashion-MNIST), that the weight $e^{-\sigma_{\text{classification}}}$ is much larger than the weights on the hidden features for the Net-PCAD method, while all weights have the same order of magnitude for the Net-LDAD method. This can be interpreted by the fact that LDA finds layer wise a low dimensional space where data can be classified and therefore “helps” the classification learning problem. This is not the case regarding the curves for the CIFAR10 dataset, since it is a “hard” classification problem given the architecture of the teacher network (TN gets 45% accuracy).

In terms of the test accuracy, we note from Figure B.10 that learning the students networks with our methods improves largely their generalization capacities compared with learning them from scratch, and one can see that they even surpass the teacher network on the CIFAR10 dataset. As a summary, Net-PCAD and Net-LDAD yield to better generalization performances of the student networks as the learning problem gets harder in the sense of the teacher network test accuracy, knowing that classifying MNIST is an easy problem (TN gets 98%), Fashion-MNIST medium (TN gets 88%) while classifying CIFAR10 is a harder learning problem (TN gets 45%).

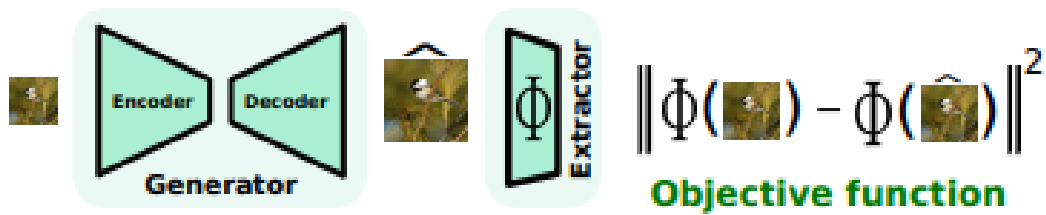
Table B.7 summarizes the performances of the learned student networks using the Net-PCAD method⁶, in terms of the forward execution time and test accuracy. As we can notice, Net-PCAD yields to an accurate speedup of inference time (depending on the choice of k) while not degrading the accuracy of the learned student networks w.r.t. the teacher network and even surpassing the teacher network’s accuracy in the case of hard classification problems (see CIFAR10).

⁶Net-LDAD gets similar results.

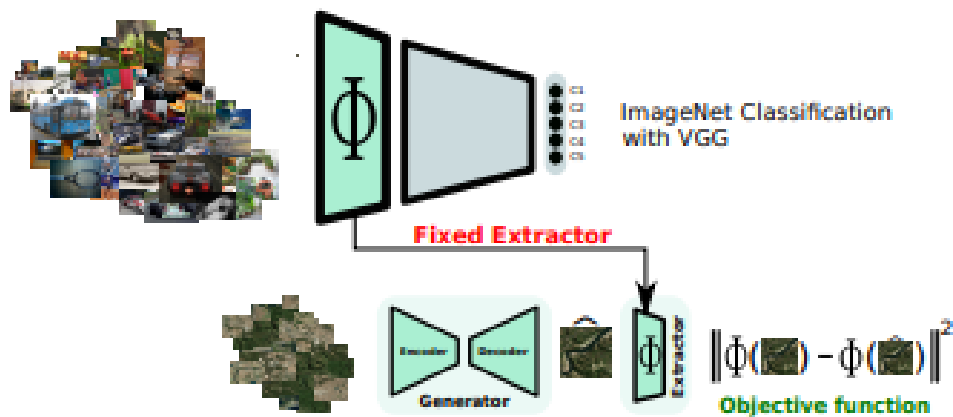
B.2.4 Central Contribution and Discussions

In this work we presented two methods to distillate a given teacher network (TN) into a student network (SN). Our methods improve the performance of SN compared to learning SN from scratch and even surpasses TN performances when the learning problem gets hard, therefore the resulting learned SN is suited to be implemented in an edge IoT device which requires limited resources. Note that the presented methods need to setup an hyper-parameter k_ℓ which will be addressed in an extended version of this work.

(a) GCN principle



(b) Disjoint-learning strategy example: SRGAN [1]



(c) Joint-learning strategy example: P/dis

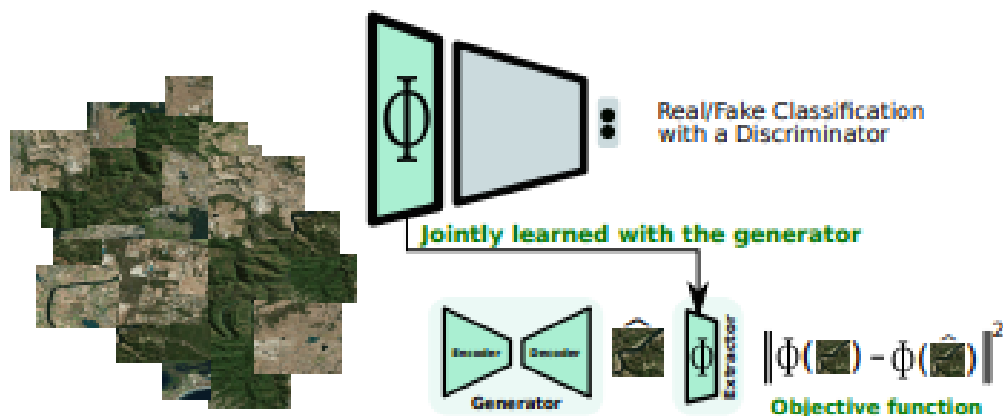


Figure B.2: Overview of the GCN framework with examples of the two learning strategies. The GCN framework consists in optimizing a *generator* in the feature space of an *extractor* as illustrated in (a). The extractor can be trained beforehand and used to optimize the generator, which we refer to as *disjoint-learning* strategy (b). The extractor can also be optimized jointly with the generator, i.e., using a *joint-learning* strategy (c).

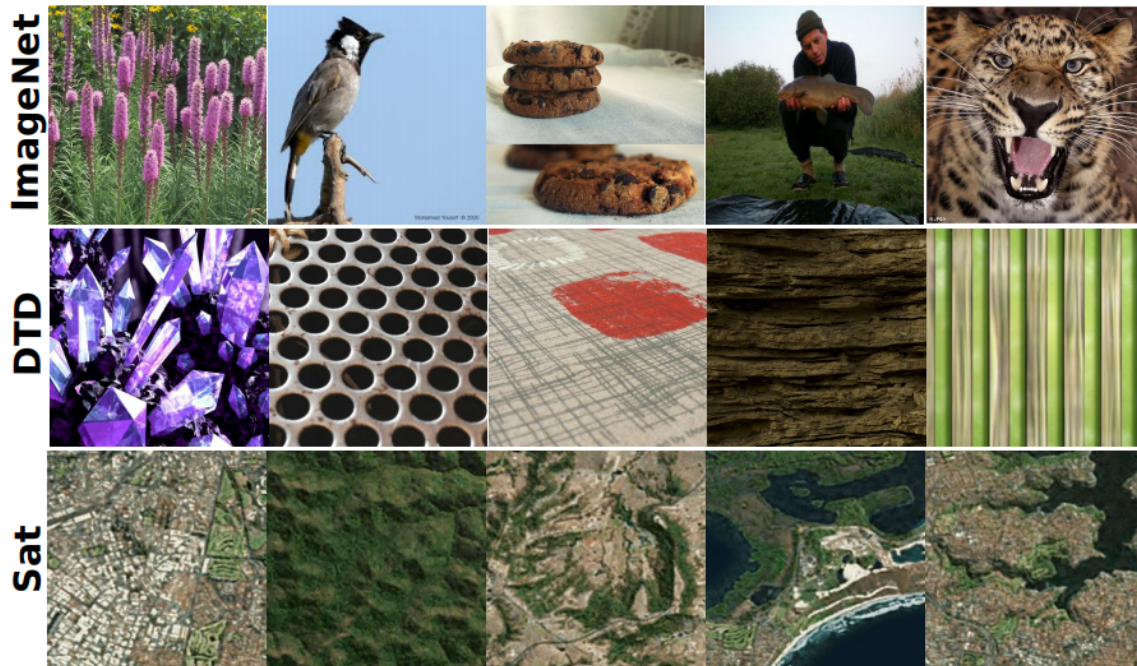
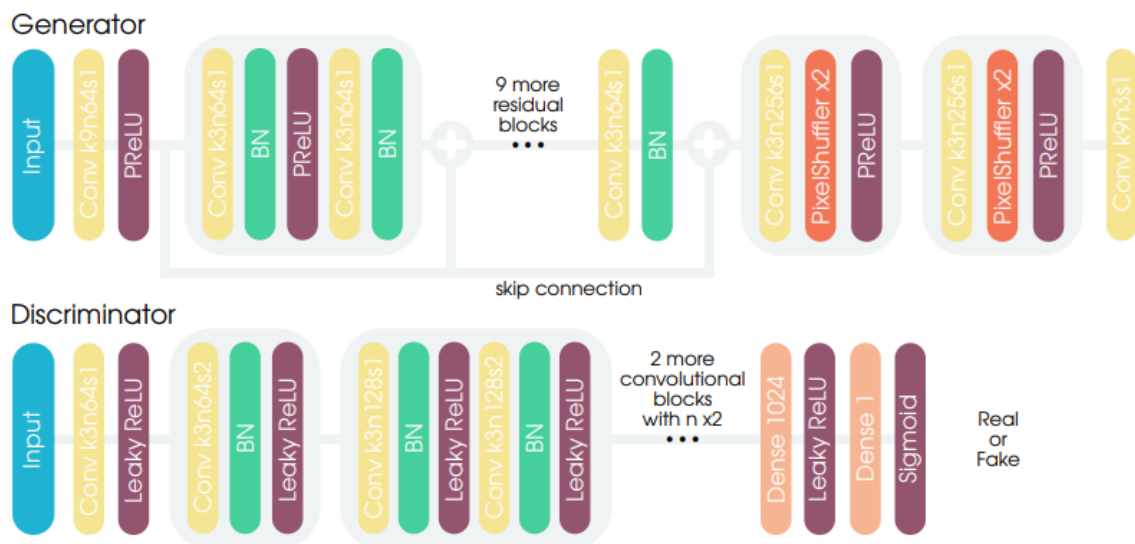


Figure B.3: Examples of images from the considered datasets.

Figure B.4: Overview of the used architectures for the generator and the discriminator. We have considered the same architectures as that of Ledig *et al.* [LTH⁺16].

		Low-level			Perceptual metrics				
	Methods	L2	SSIM	Squ	Squ-l	Alex	Alex-l	VGG	VGG-l
ImageNet	\mathcal{P}_{sirsr}/dis	0.018	<u>0.147</u>	1.606	0.279	1.470	0.398	2.088	0.358
	\mathcal{P}_{sirsr}/rec	0.020	0.162	1.723	0.301	1.595	0.425	2.243	0.388
	$\mathcal{P}_{sirsr}/dis,rec$	<u>0.017</u>	<u>0.147</u>	<u>1.587</u>	<u>0.279</u>	<u>1.420</u>	<u>0.382</u>	<u>2.052</u>	<u>0.353</u>
	\mathcal{P}_{sirsr}/adv	0.028	0.202	1.820	0.222	1.554	0.322	2.598	0.432
	$\mathcal{P}_{sirsr}/adv,rec$	0.016	0.141	1.533	<u>0.263</u>	1.362	<u>0.368</u>	1.994	0.340
DTD	\mathcal{P}_{sirsr}/dis	<u>0.027</u>	0.184	1.873	0.327	1.739	0.440	2.401	0.421
	\mathcal{P}_{sirsr}/rec	<u>0.027</u>	0.183	1.851	0.320	1.726	0.438	2.398	0.420
	$\mathcal{P}_{sirsr}/dis,rec$	0.023	0.167	1.703	0.292	1.576	0.404	2.260	0.392
	\mathcal{P}_{sirsr}/adv	0.036	0.227	2.077	<u>0.281</u>	1.812	<u>0.375</u>	2.770	0.473
	$\mathcal{P}_{sirsr}/adv,rec$	0.046	0.236	2.089	0.277	1.793	0.344	2.796	0.481
Sat	\mathcal{P}_{sirsr}/dis	0.011	0.129	<u>1.484</u>	0.210	1.508	0.356	2.121	<u>0.355</u>
	\mathcal{P}_{sirsr}/rec	0.060	0.168	1.705	0.245	1.762	0.423	2.260	0.395
	$\mathcal{P}_{sirsr}/dis,rec$	0.011	<u>0.138</u>	1.493	0.215	<u>1.435</u>	0.351	2.108	0.372
	\mathcal{P}_{sirsr}/adv	0.030	0.214	1.719	<u>0.181</u>	1.627	<u>0.306</u>	2.711	0.419
	$\mathcal{P}_{sirsr}/adv,rec$	<u>0.018</u>	0.183	1.359	0.140	1.310	0.220	<u>2.115</u>	0.344

Table B.2: Results of the proposed \mathcal{P}_{sirsr} methods in terms of traditional metrics (L2 and SSIM) and the *perceptual error* (PE) given by Eq. equation B.6 on different datasets. As we can notice, the method $\mathcal{P}_{sirsr}/adv,rec$ outperforms the other methods in the datasets ImageNet and Sat, while $\mathcal{P}_{sirsr}/dis,rec$ gives the best results on DTD.

		Low-level			Perceptual metrics				
	Methods	L2	SSIM	Squ	Squ-l	Alex	Alex-l	VGG	VGG-l
Sat	\mathcal{P}_{sirsr}/mse	0.011	0.134	1.873	0.245	1.855	0.411	2.536	0.419
	$\mathcal{P}_{sirsr}/adv,mse$	0.082	0.197	<u>1.458</u>	<u>0.205</u>	1.466	0.352	2.125	<u>0.347</u>
	SRGAN [LTH+16]	0.228	0.188	1.510	0.220	<u>1.361</u>	<u>0.282</u>	2.230	0.412
	$\mathcal{P}_{sirsr}/adv,rec$	<u>0.018</u>	<u>0.183</u>	1.359	0.140	1.310	0.220	2.115	0.344

Table B.3: Comparison of our method $\mathcal{P}_{sirsr}/adv,rec$ with baselines and the SRGAN method [LTH+16] on the satellite images domain, in terms of classical metrics (L2 and SSIM) and perceptual metrics [ZIE+18].

		Low-level			Perceptual metrics				
	Methods	L2	SSIM	Squ	Squ-l	Alex	Alex-l	VGG	VGG-l
ImageNet	\mathcal{P}_{sirsr}/mse	<u>0.017</u>	<u>0.146</u>	1.568	0.280	1.435	0.391	2.064	0.349
	$\mathcal{P}_{sirsr}/adv,mse$	0.020	0.156	1.634	<u>0.241</u>	1.397	<u>0.329</u>	2.223	0.384
	SRGAN	0.028	0.170	1.303	0.177	1.084	0.225	<u>2.045</u>	<u>0.342</u>
	$\mathcal{P}_{sirsr}/adv,rec$	0.016	0.141	<u>1.533</u>	0.263	<u>1.362</u>	0.368	1.994	0.340
DTD	\mathcal{P}_{sirsr}/mse	<u>0.029</u>	<u>0.185</u>	1.972	0.342	1.856	0.470	2.479	0.434
	$\mathcal{P}_{sirsr}/adv,mse$	0.025	0.188	1.880	<u>0.268</u>	1.586	0.349	2.512	0.430
	SRGAN	0.031	0.191	1.557	0.209	1.298	0.241	<u>2.308</u>	<u>0.393</u>
	$\mathcal{P}_{sirsr}/dis,rec$	0.023	0.167	<u>1.703</u>	0.292	<u>1.576</u>	0.404	2.260	0.392

Table B.4: Comparison of our methods $\mathcal{P}_{sirsr}/adv,rec$ and $\mathcal{P}_{sirsr}/dis,rec$ with baselines and the SRGAN method [LTH+16] on the datasets ImageNet (a subset of 200,000 randomly selected images) and DTD, in terms of classical metrics (L2 and SSIM) and perceptual metrics [ZIE+18].

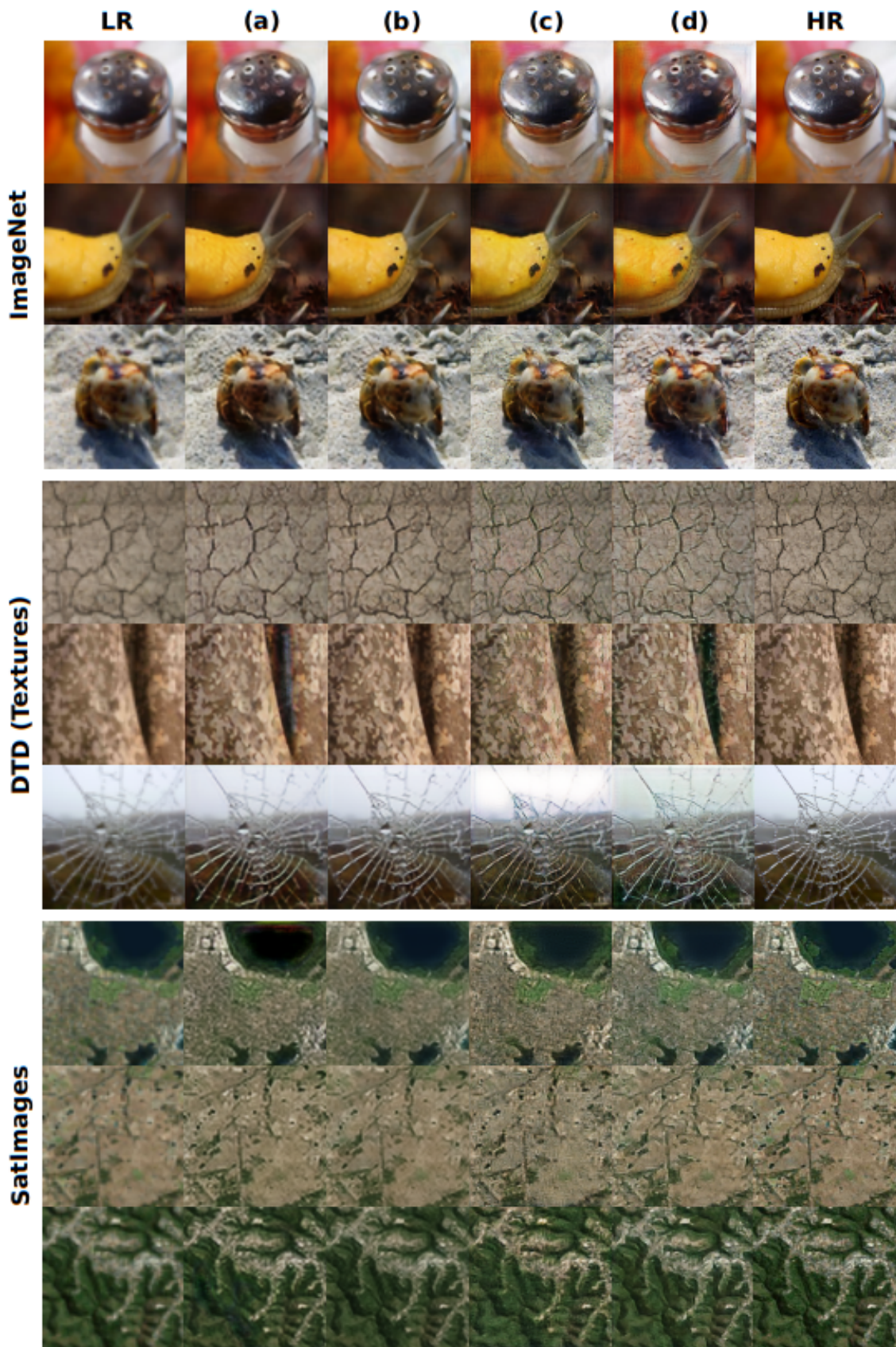
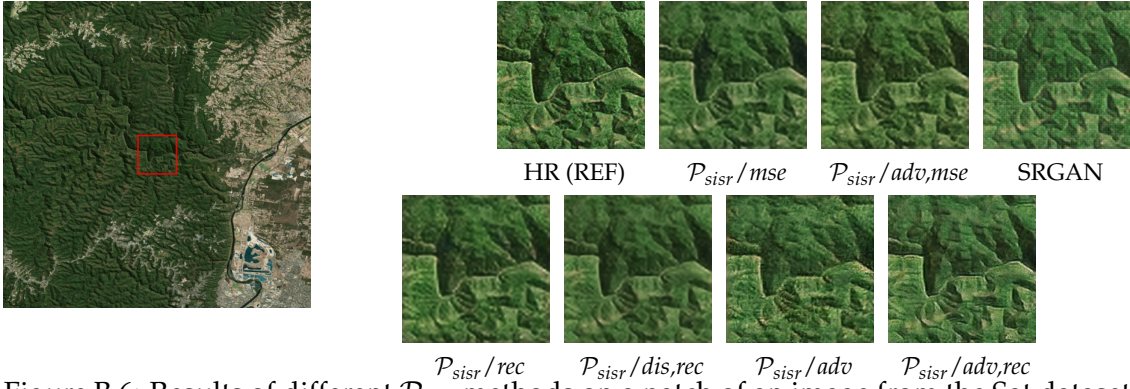


Figure B.5: Rows refer to the different considered Datasets. Columns refer to methods and ground-truth images: **LR** and **HR** refer to the low- and high-resolution pairs. The different used methods are: **(a)** \mathcal{P}_{sisr}/rec , **(b)** $\mathcal{P}_{sisr}/dis,rec$, **(c)** \mathcal{P}_{sisr}/adv and **(d)** $\mathcal{P}_{sisr}/adv,rec$. Best view in PDF.

Figure B.6: Results of different \mathcal{P}_{sisr} methods on a patch of an image from the Sat dataset.

		Low-level		Perceptual metrics					
	Methods	L2	SSIM	Squ	Squ-l	Alex	Alex-l	VGG	VGG-l
ImageNet	\mathcal{P}_{sisr}/mse	<u>0.017</u>	<u>0.146</u>	1.568	0.280	1.435	0.391	2.064	0.349
	$\mathcal{P}_{sisr}/adv,mse$	0.020	0.156	1.634	0.241	1.397	0.329	2.223	0.384
	SRGAN	0.028	0.170	1.303	0.177	1.084	0.225	<u>2.045</u>	<u>0.342</u>
	\mathcal{P}_{sisr}/dis	0.018	0.147	1.606	0.279	1.470	0.398	2.088	0.358
	\mathcal{P}_{sisr}/rec	0.020	0.162	1.723	0.301	1.595	0.425	2.243	0.388
	$\mathcal{P}_{sisr}/dis,rec$	<u>0.017</u>	0.147	1.587	0.279	1.420	0.382	2.052	0.353
	\mathcal{P}_{sisr}/adv	0.028	0.202	1.820	<u>0.222</u>	1.554	<u>0.322</u>	2.598	0.432
	$\mathcal{P}_{sisr}/adv,rec$	0.016	0.141	<u>1.533</u>	0.263	<u>1.362</u>	0.368	1.994	0.340
DTD	\mathcal{P}_{sisr}/mse	0.029	0.185	1.972	0.342	1.856	0.470	2.479	0.434
	$\mathcal{P}_{sisr}/adv,mse$	0.025	0.188	1.880	<u>0.268</u>	1.586	0.349	2.512	0.430
	SRGAN	0.031	0.191	1.557	0.209	1.298	0.241	<u>2.308</u>	<u>0.393</u>
	\mathcal{P}_{sisr}/dis	<u>0.027</u>	0.184	1.873	0.327	1.739	0.440	2.401	0.421
	\mathcal{P}_{sisr}/rec	<u>0.027</u>	<u>0.183</u>	1.851	0.320	1.726	0.438	2.398	0.420
	$\mathcal{P}_{sisr}/dis,rec$	0.023	0.167	<u>1.703</u>	0.292	<u>1.576</u>	0.404	2.260	0.392
	\mathcal{P}_{sisr}/adv	0.036	0.227	2.077	0.281	1.812	0.375	2.770	0.473
	$\mathcal{P}_{sisr}/adv,rec$	0.046	0.236	2.089	0.277	1.793	<u>0.344</u>	2.796	0.481
Sat	\mathcal{P}_{sisr}/mse	0.011	<u>0.134</u>	1.873	0.245	1.855	0.411	2.536	0.419
	$\mathcal{P}_{sisr}/adv,mse$	0.082	0.197	<u>1.458</u>	0.205	1.466	0.352	2.125	<u>0.347</u>
	SRGAN	0.228	0.188	1.510	0.220	1.361	0.282	2.230	0.412
	\mathcal{P}_{sisr}/dis	0.011	0.129	1.484	0.210	1.508	0.356	2.121	0.355
	\mathcal{P}_{sisr}/rec	0.060	0.168	1.705	0.245	1.762	0.423	2.260	0.395
	$\mathcal{P}_{sisr}/dis,rec$	0.011	0.138	1.493	0.215	<u>1.435</u>	0.351	2.108	0.372
	\mathcal{P}_{sisr}/adv	0.030	0.214	1.719	<u>0.181</u>	1.627	<u>0.306</u>	2.711	0.419
	$\mathcal{P}_{sisr}/adv,rec$	<u>0.018</u>	0.183	1.359	0.140	1.310	0.220	<u>2.115</u>	0.344

Table B.5: Comparison of the proposed \mathcal{P}_{sisr} methods in terms of traditional metrics (L2 and SSIM) and the *perceptual error* (PE) given by equation B.6 on all the considered datasets. In terms of perceptual metrics, the proposed \mathcal{P}_{sisr} methods rank in the second position after SRGAN [LTH⁺16] on the datasets ImageNet and DTD, while they outperform all the baselines on the satellite images domain which is far from the ImageNet domain.

Layer	Teacher	Student
Dense 1	$p_0 \times 1024$	$p_0 \times k$
Dense 2	1024×512	$k \times k$
Dense 3	512×256	$k \times k$
Dense 4	256×10	$k \times 10$

Table B.6: Architectures of the teacher and student networks. The dimensions of the weight matrix at each dense layer are shown for both networks.

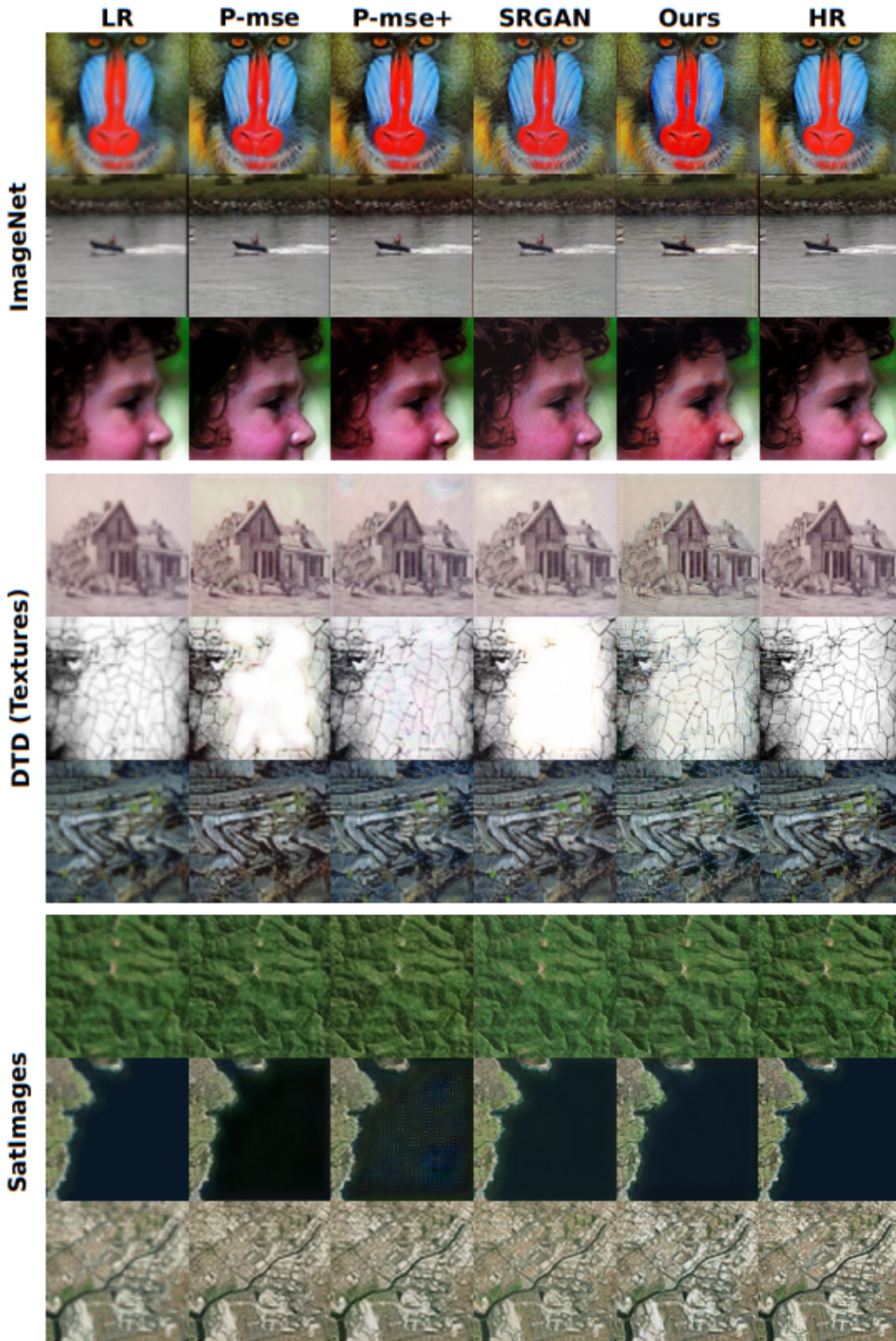


Figure B.7: Rows refer to the different Datasets. Columns refer to methods and ground-truth images: **LR** and **HR** refer to the low- and high-resolution pairs. **P-mse+** refers to the method $\mathcal{P}_{sirr/mse}$ with an adversarial loss ($\lambda_2 > 0$), **SRGAN** for the method in [LTH⁺16] and our method $\mathcal{P}_{sirr/adv,rec}$. Best view in PDF.

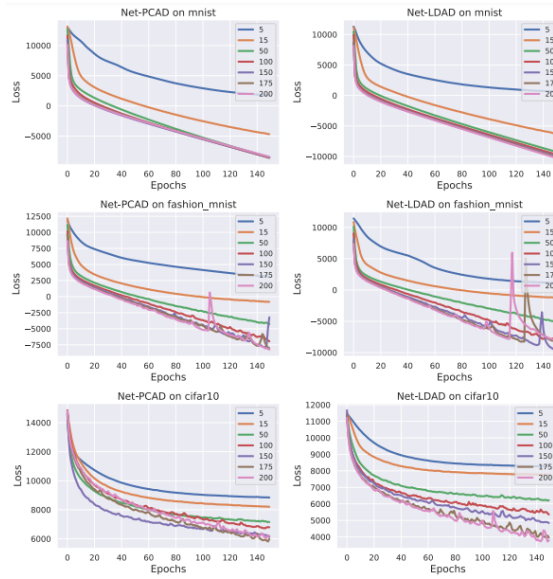


Figure B.8: Training loss of the student network when trained using Net-PCAD (**left**) and Net-LDAD (**right**) in terms of the training epochs, for different values of k , and across three different datasets.

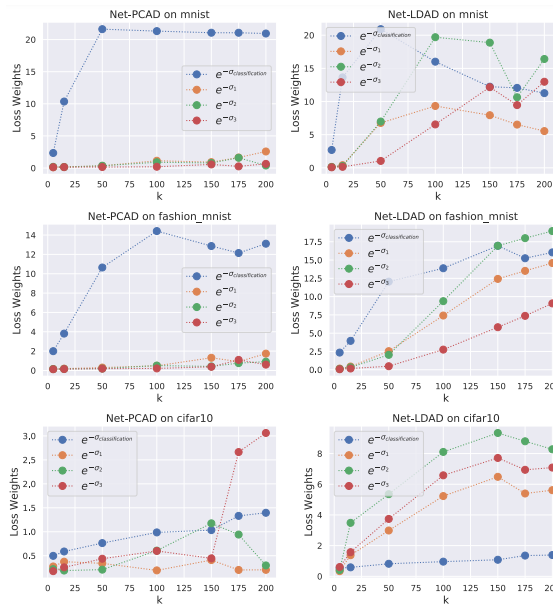


Figure B.9: The learned Homoscedastic loss parameters using Net-PCAD (**left**) and Net-LDAD (**right**) in terms of k and across three different datasets. The weights corresponding to the reduced features mapping loss are of the same order of magnitude for Net-LDAD as the classification loss, which is a consequence of the fact that LDA is classes dependent.

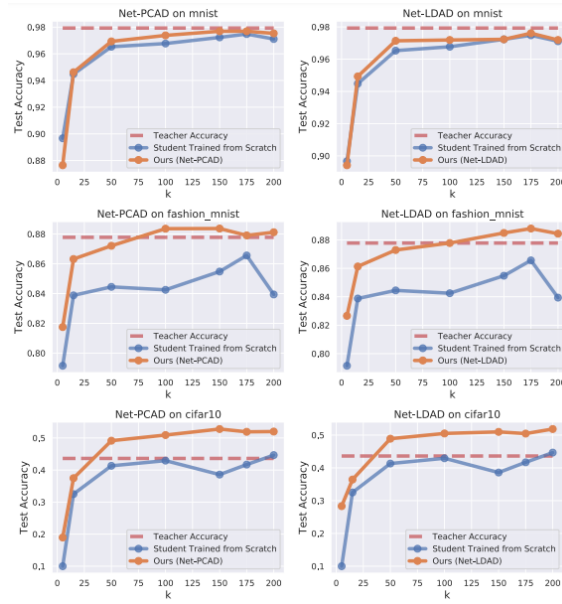


Figure B.10: Test accuracy of the student network in orange trained using Net-PCAD (left) and Net-LDAD (right), and test accuracy of the student network in blue trained from scratch, in terms of k and across three different datasets. The test accuracy of the teacher corresponds to the dashed red lines.

Dataset	Student			
	Teacher	$k = 50$	100	200
MNIST	2.23s	0.38s	0.45s	0.65s
	98%	97%	97.5%	97.8%
FASHION	2.23s	0.38s	0.45s	0.65s
	88%	87.5%	88.5%	88.5%
CIFAR10	4.63s	0.75s	0.92s	1.35s
	45%	50%	50.1%	50.3%

Table B.7: Forward execution time in seconds (and corresponding test accuracies in %) of the teacher network and the student network for different values of k , the forward pass is applied (on a i7-7700HQ CPU @ 2.80GHz) to the train set of the respective datasets using a batch size of 50000 images.

Appendix C

Proofs

Contents

C.1 Proofs of Chapter 3	157
C.1.1 Setting of the proof	157
C.1.2 Basic tools	158
C.1.3 Main body of the proof	159
C.2 Proofs of Chapter 4	162
C.2.1 Proofs of Section 4.1	162
C.2.2 Proofs of Section 4.2	165
C.2.2.1 Proof of Theorem 4.3	165
C.2.2.2 Proof of Theorem 4.4	167
C.3 Proofs of Chapter 5	168
C.3.1 Proofs of Section 5.1	168
C.3.2 Proofs of Section 5.2	172
C.3.2.1 Proof of Theorem 5.5	172
C.3.2.2 Proof of Theorem 5.6	173
C.3.2.3 Optimality	175

C.1 Proofs of Chapter 3

C.1.1 Setting of the proof

For simplicity, we will only suppose the case $k = 1$ and we consider the following notations that will be used subsequently.

$$\bar{x} = \mathbb{E}x_i, \quad \mathbf{C} = \mathbb{E}[x_i x_i^\top], \quad \mathbf{X}_0 = \mathbf{X} - \bar{x} \mathbf{1}_n^\top, \quad \mathbf{C}_0 = \mathbb{E}[\mathbf{X}_0 \mathbf{X}_0^\top / n].$$

Let

$$\mathbf{X}_{-i} = (x_1, \dots, x_{i-1}, 0, x_i, \dots, x_n)$$

the matrix \mathbf{X} with a vector of zeros at its i th column.

Denote the resolvents

$$\mathbf{R} = \left(\frac{\mathbf{X}^\top \mathbf{X}}{p} + z \mathbf{I}_n \right)^{-1}, \quad \mathbf{Q} = \left(\frac{\mathbf{X} \mathbf{X}^\top}{p} + z \mathbf{I}_p \right)^{-1}, \quad \mathbf{Q}_{-i} = \left(\frac{\mathbf{X} \mathbf{X}^\top}{p} - \frac{x_i x_i^\top}{p} + z \mathbf{I}_p \right)^{-1} \quad (\text{C.1})$$

And let

$$\tilde{\mathbf{Q}} = \left(\frac{1}{c} \frac{\mathbf{C}}{1 + \delta} + z \mathbf{I}_p \right)^{-1}, \quad (\text{C.2})$$

where δ is the solution to the fixed point equation

$$\delta = \frac{1}{p} \text{tr} \left(\mathbf{C} \left(\frac{1}{c} \frac{\mathbf{C}}{1 + \delta} + z \mathbf{I}_p \right)^{-1} \right).$$

C.1.2 Basic tools

Lemma C.1 ([Led05a]). *Let $z \in \mathcal{E}_q(1 | \mathbb{R}^p, \|\cdot\|)$ and $\mathbf{M} \in \mathcal{E}_q(1 | \mathbb{R}^{p \times n}, \|\cdot\|_F)$. Then, for some numerical constant $C > 0$*

- $\mathbb{E} \|z\| \leq \|\mathbb{E}z\| + C\sqrt{p}$, $\mathbb{E} \|z\|_\infty \leq \|\mathbb{E}z\|_\infty + C\sqrt{\log p}$.
- $\mathbb{E} \|\mathbf{M}\| \leq \|\mathbb{E}\mathbf{M}\| + C\sqrt{p+n}$, $\mathbb{E} \|\mathbf{M}\|_F \leq \|\mathbb{E}\mathbf{M}\|_F + C\sqrt{pn}$.

Lemma C.2. *Denote $\mathbf{Q}_{\bar{x}} = (\bar{x}\bar{x}^\top + z\mathbf{I}_p)^{-1}$, we have:*

$$\mathbf{Q}_{\bar{x}}\bar{x} = \frac{\bar{x}}{\|\bar{x}\|^2 + z} \quad \text{and} \quad \|\tilde{\mathbf{Q}}\bar{x}\|, \bar{x}\tilde{\mathbf{Q}}\bar{x} = \mathcal{O}(1).$$

Moreover, if $\|\bar{x}\| \geq \sqrt{p}$, $\|\tilde{\mathbf{Q}}\bar{x}\| = \mathcal{O}(p^{-1/2})$.

Proof. Since $z\mathbf{Q}_{\bar{x}} = \mathbf{I}_p - \mathbf{Q}_{\bar{x}}\bar{x}\bar{x}^\top$:

$$z\mathbf{Q}_{\bar{x}}\bar{x} = \bar{x} - \|\bar{x}\|^2\mathbf{Q}_{\bar{x}}\bar{x},$$

and we recover the first identity of the Lemma.

And since the matrix \mathbf{C}_0 is nonnegative symmetric, we have :

$$\tilde{\mathbf{Q}}\bar{x} = \left(\frac{1}{c} \frac{\mathbf{C}_0 + \bar{x}\bar{x}^\top}{1 + \delta} + z\mathbf{I}_p \right)^{-1} \bar{x} \leq \frac{c(1 + \delta)\bar{x}}{\|\bar{x}\|^2 + zc(1 + \delta)}.$$

Therefore, $\bar{x}\tilde{\mathbf{Q}}\bar{x} = \frac{c(1+\delta)\|\bar{x}\|^2}{\|\bar{x}\|^2 + zc(1+\delta)} = \mathcal{O}(1)$ and:

$$\|\tilde{\mathbf{Q}}\bar{x}\| = \frac{c(1 + \delta)\|\bar{x}\|}{\|\bar{x}\|^2 + zc(1 + \delta)} \leq \begin{cases} \frac{\|\bar{x}\|}{z} = \mathcal{O}(1) & \text{if } \|\bar{x}\| \leq 1, \\ \frac{c(1 + \delta)}{\|\bar{x}\|} = \mathcal{O}(1) & \text{if } \|\bar{x}\| \geq 1. \end{cases}$$

□

Proposition C.1. $\bar{x}^\top \mathbb{E}[\mathbf{Q}]\bar{x} = \bar{x}^\top \tilde{\mathbf{Q}}\bar{x} + \mathcal{O}\left(\sqrt{\frac{\log p}{p}}\right)$

Proof. Let us bound:

$$\left| \bar{x}^\top \mathbf{Q}\bar{x} - \bar{x}^\top \tilde{\mathbf{Q}}\bar{x} \right| \leq \frac{c^{-1}}{1 + \delta} \left| \mathbb{E} \left[\bar{x} \mathbf{Q} x_i x_i^\top \tilde{\mathbf{Q}}\bar{x} \left(\frac{1}{p} x_i^\top \mathbf{Q}_{-i} x_i - \delta \right) \right] + \frac{1}{p} \mathbb{E} \left[\bar{x}^\top \mathbf{Q}_{-i} x_i x_i^\top \mathbf{C} \tilde{\mathbf{Q}}\bar{x} \right] \right|$$

Now let us consider a supplementary random vector \mathbf{x}_{n+1} following the same law as the x_i 's and independent of \mathbf{X} . We divide the set $\mathbb{I} = [n + 1]$ into two sets $\mathbb{I}_{\frac{1}{2}}$ and $\mathbb{I}_{\frac{2}{2}}$

of same cardinality ($\lfloor \frac{n+1}{2} \rfloor \leq \#\mathbb{I}_{\frac{1}{2}}, \#\mathbb{I}_{\frac{2}{2}} \leq \lceil \frac{n+1}{2} \rceil$), we note $\mathbf{X}_{\frac{1}{2}} = (\mathbf{x}_i | i \in \mathbb{I}_{\frac{1}{2}})$, $\mathbf{X}_{\frac{2}{2}} = (\mathbf{x}_i | i \in \mathbb{I}_{\frac{2}{2}})$ and we introduce the diagonal matrices $\Delta = \text{diag} \left(\frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i - \delta | i \in \mathbb{I}_{\frac{1}{2}} \right)$, $D = \text{diag} \left(1 + \frac{1}{p+1} \mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i | i \in \mathbb{I}_{\frac{2}{2}} \right)$. We have the bound:

$$\begin{aligned}
& \left| \mathbb{E} \left[\bar{\mathbf{x}} \mathbf{Q} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \left(\frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i - \delta \right) \right] \right| \\
&= \left| \mathbb{E} \left[\left(1 + \frac{1}{p} \mathbf{x}_{n+1}^\top \mathbf{Q} \mathbf{x}_{n+1} \right) \mathbf{x}_{n+1} \mathbf{Q}_{+(n+1)} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \left(\frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i - \delta \right) \right] \right| \\
&= \frac{1}{p^2} \left| \mathbb{E} \left[\mathbf{1}^\top D \mathbf{X}_{\frac{2}{2}}^\top \mathbf{Q}_{+(n+1)} \mathbf{X}_{\frac{1}{2}} \Delta \mathbf{X}_{\frac{1}{2}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \right] \right| \\
&\leq \sqrt{\left| \mathbb{E} \left[\frac{1}{p^3} \mathbf{1}^\top D \mathbf{X}_{\frac{2}{2}}^\top \mathbf{Q}_{+(n+1)} \mathbf{X}_{\frac{1}{2}} \Delta^2 \mathbf{X}_{\frac{1}{2}}^\top \mathbf{Q}_{+(n+1)} \mathbf{X}_{\frac{2}{2}} D \mathbf{1} \right] \mathbb{E} \left[\frac{1}{p} \bar{\mathbf{x}}^\top \tilde{\mathbf{Q}} \mathbf{X}_{\frac{1}{2}} \mathbf{X}_{\frac{1}{2}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \right] \right|} \\
&\leq \sqrt{\left| \mathbb{E} \left[\left\| \frac{1}{p} \mathbf{X}_{\frac{2}{2}}^\top \mathbf{Q}_{+(n+1)} \mathbf{X}_{\frac{1}{2}} \right\|^2 \|D\|^2 \|\Delta\|^2 \right] \mathbb{E} [\bar{\mathbf{x}} \tilde{\mathbf{Q}} C \tilde{\mathbf{Q}} \bar{\mathbf{x}}] \right|} \leq \mathcal{O} \left(\sqrt{\frac{\log p}{p}} \right),
\end{aligned}$$

thanks to Lemma C.1 and Lemma C.2 (the spectral norm of Δ and D is just an infinity norm if we see them as random vectors of \mathbb{R}^n). We can bound $\frac{1}{p} \left| \mathbb{E} [\bar{\mathbf{x}}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q} C \tilde{\mathbf{Q}} \bar{\mathbf{x}}] \right|$ the same way to obtain the result of the proposition. \square

Proposition C.2. $\|\mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{X}_{-i}] - \frac{\bar{\mathbf{x}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \mathbf{1}^\top}{1+\delta}\| = \mathcal{O}(\sqrt{\log p})$

Proof. Considering $\mathbf{u} \in \mathbb{R}^n$ such that $\|\mathbf{u}\| = 1$:

$$\begin{aligned}
& \left| \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{X}_{-i} \mathbf{u}] - \frac{\bar{\mathbf{x}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \mathbf{1}^\top \mathbf{u}}{1+\delta} \right| \\
&= \left| \sum_{\substack{j=1 \\ j \neq i}}^n \mathbf{u}_j \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j}{1 + \frac{1}{p} \mathbf{x}_j^\top \mathbf{Q}_{-j} \mathbf{x}_j} - \frac{\mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_j}{1 + \delta} \right] \right| \\
&\leq \sqrt{n} \left| \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j}{1 + \frac{1}{p} \mathbf{x}_j^\top \mathbf{Q}_{-j} \mathbf{x}_j} - \frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j}{1 + \delta} \right] \right| + \left| \frac{1}{1 + \delta} \mathbb{E} [\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j - \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_j] \right| \quad (\text{where } i \neq j) \\
&\leq \sqrt{n} \left| \mathbb{E} \left[\bar{\mathbf{x}}^\top \mathbf{Q} \mathbf{x}_j \left(\frac{1}{p} \mathbf{x}_j^\top \mathbf{Q}_{-j} \mathbf{x}_j - \delta \right) \right] \right| + \sqrt{n} \left| \mathbb{E} [\bar{\mathbf{x}}^\top \mathbf{Q}_{-i} \bar{\mathbf{x}} - \bar{\mathbf{x}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}}] \right|,
\end{aligned}$$

where the first term is treated the same way as we did in the proof of Proposition C.1 and the second term is bounded thanks to Proposition C.1 \square

C.1.3 Main body of the proof

Proof of Theorem 3.1. Recall the definition of the resolvents \mathbf{R} and \mathbf{Q} in Equation (C.1). The first step of the proof is to show the concentration of \mathbf{R} . This comes from the fact that the application $\Phi : \mathbf{X} \mapsto (\mathbf{X}^\top \mathbf{X} + z \mathbf{I}_n)^{-1}$ is $2z^{-3/2}$ -Lipschitz w.r.t. the Frobenius norm. Indeed, by the matrix identity $\mathbf{A} - \mathbf{B} = \mathbf{A}(\mathbf{B}^{-1} - \mathbf{A}^{-1})\mathbf{B}$, we have

$$\Phi(\mathbf{X}) - \Phi(\mathbf{X} + \mathbf{H}) = \Phi(\mathbf{X})(\mathbf{H}^\top \mathbf{X} + (\mathbf{X} + \mathbf{H})^\top \mathbf{H})\Phi(\mathbf{X} + \mathbf{H})$$

And by the bounds $\|\mathbf{AB}\|_F \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|_F$, $\|\Phi(\mathbf{X})\mathbf{X}^\top\| \leq z^{-1/2}$ and $\|\Phi(\mathbf{X})\| \leq z^{-1}$, we have

$$\|\Phi(\mathbf{X} + \mathbf{H}) - \Phi(\mathbf{X})\|_F \leq \frac{2}{z^{3/2}} \|\mathbf{H}\|_F.$$

Therefore, given $\mathbf{X} \in \mathcal{E}_q(1 | \mathbb{R}^{p \times n}, \|\cdot\|_F)$ and since the application $\mathbf{X} \mapsto \mathbf{R} = \Phi(\mathbf{X}/\sqrt{p})$ is $2z^{-3/2}p^{-1/2}$ -Lipschitz, we have by Proposition 2.4 that $\mathbf{R} \in \mathcal{E}_q(p^{-1/2} | \mathbb{R}^{n \times n}, \|\cdot\|_F)$.

The second step consists in estimating $\mathbb{E}\mathbf{R}(z)$ through a deterministic matrix $\tilde{\mathbf{R}}$. Indeed, by the identity $(\mathbf{M}^\top\mathbf{M} + z\mathbf{I})^{-1}\mathbf{M}^\top = \mathbf{M}^\top(\mathbf{M}\mathbf{M}^\top + z\mathbf{I})^{-1}$, the resolvent \mathbf{R} can be expressed in function of \mathbf{Q} as follows

$$\mathbf{R} = \frac{1}{z} \left(\mathbf{I}_n - \frac{\mathbf{X}^\top\mathbf{Q}\mathbf{X}}{p} \right), \quad (\text{C.3})$$

thus a deterministic equivalent for \mathbf{R} can therefore be obtained through a deterministic equivalent of the matrix $\mathbf{X}^\top\mathbf{Q}\mathbf{X}$. However, as demonstrated in [LC19], the matrix \mathbf{Q} has as a deterministic equivalent the matrix $\tilde{\mathbf{Q}}$ defined in equation C.2. In the following, we aim at deriving a deterministic equivalent for $\frac{1}{p}\mathbf{X}^\top\mathbf{Q}\mathbf{X}$ in function of $\tilde{\mathbf{Q}}$. Let \mathbf{u} and \mathbf{v} be two unitary vectors in \mathbb{R}^n , and let us estimate

$$\Delta \equiv \mathbb{E} \left[\mathbf{u}^\top \left(\frac{\mathbf{X}^\top\mathbf{Q}\mathbf{X}}{p} - \frac{\mathbf{X}^\top\tilde{\mathbf{Q}}\mathbf{X}}{p} \right) \mathbf{v} \right] = \frac{1}{p} \mathbb{E} \left[\frac{\mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v}}{1+\delta} - \frac{1}{p} \mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}\mathbf{X}\mathbf{X}^\top\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v} \right]$$

With the following matrix identities (to explore the independence of the columns of \mathbf{X}):

$$\mathbf{Q} = \mathbf{Q}_{-i} - \frac{1}{p} \mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\mathbf{Q}, \quad \mathbf{Q}\mathbf{x}_i = \frac{\mathbf{Q}_{-i}\mathbf{x}_i}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i}, \quad \mathbf{A} - \mathbf{B} = \mathbf{A}(\mathbf{B}^{-1} - \mathbf{A}^{-1})\mathbf{B}$$

and the decomposition $\mathbf{Q}\mathbf{X}\mathbf{X}^\top = \sum_{i=1}^n \mathbf{Q}\mathbf{x}_i\mathbf{x}_i^\top$, we obtain:

$$\begin{aligned} \Delta &= \frac{1}{p^2} \mathbb{E} \left[\sum_{i=1}^n \frac{\mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}_{-i}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v}}{1+\delta} - \frac{\mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v}}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i} - \frac{1}{p} \frac{\mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\mathbf{Q}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v}}{1+\delta} \right] \\ &= \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{u}^\top\mathbf{X}_{-i}^\top\mathbf{Q}_{-i}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}_{-i}\mathbf{v}}{1+\delta} - \frac{\mathbf{u}^\top\mathbf{X}_{-i}^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\tilde{\mathbf{Q}}\mathbf{X}_{-i}\mathbf{v}}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i} \right. \\ &\quad + \frac{\mathbf{u}_i\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}_{-i}\mathbf{v}}{1+\delta} + \frac{\mathbf{v}_i\mathbf{u}^\top\mathbf{X}_{-i}^\top\mathbf{Q}_{-i}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{x}_i}{1+\delta} + \mathbf{u}_i\mathbf{v}_i \frac{\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{x}_i}{1+\delta} \\ &\quad - \frac{\mathbf{u}_i\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\tilde{\mathbf{Q}}\mathbf{X}_{-i}\mathbf{v}}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i} - \frac{\mathbf{v}_i\mathbf{u}^\top\mathbf{X}_{-i}^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\tilde{\mathbf{Q}}\mathbf{x}_i}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i} - \mathbf{u}_i\mathbf{v}_i \frac{\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\tilde{\mathbf{Q}}\mathbf{x}_i}{1 + \frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i} \\ &\quad \left. - \frac{1}{p} \frac{\mathbf{u}^\top\mathbf{X}^\top\mathbf{Q}_{-i}\mathbf{x}_i\mathbf{x}_i^\top\mathbf{Q}\mathbf{C}\tilde{\mathbf{Q}}\mathbf{X}\mathbf{v}}{1+\delta} \right] \end{aligned}$$

We can show with Holder's inequality and the concentration bounds (mainly the fact that $\frac{1}{p}\mathbf{x}_i^\top\mathbf{Q}_{-i}\mathbf{x}_i$ concentrates around δ) developed in [LC19], that most of the above quantities

vanish asymptotically. As a toy example, we consider the following term:

$$\begin{aligned}
& \left| \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{C} \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{1 + \delta} - \frac{\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] \right| \\
&= \left| \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{(1 + \delta) \left(1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \right)} \frac{\delta - \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] \right| \\
&\leq \left| \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\left(\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \right) \left(\mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v} \right) \left(\delta - \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \right) \right] \right| \\
&\leq \left| \frac{1}{p} \sum_{i=1}^n \left(\mathbb{E} \left[\left(\frac{1}{\sqrt{p}} \mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \right)^3 \right] \mathbb{E} \left[\left(\frac{1}{\sqrt{p}} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v} \right)^3 \right] \mathbb{E} \left[\left(\delta - \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \right)^3 \right] \right) \right|^{\frac{1}{3}} \\
&= \mathcal{O} \left(\frac{1}{\sqrt{p}} \right)
\end{aligned}$$

Similarly, we can show that:

$$\begin{aligned}
& \left| \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{u}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C} \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{1 + \delta} + \frac{\mathbf{v}_i \mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{C} \tilde{\mathbf{Q}} \mathbf{x}_i}{1 + \delta} \right. \right. \\
& \quad \left. \left. + \mathbf{u}_i \mathbf{v}_i \frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C} \tilde{\mathbf{Q}} \mathbf{x}_i}{1 + \delta} - \frac{1}{p} \frac{\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q} \mathbf{C} \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{1 + \delta} \right] \right| = \mathcal{O} \left(\frac{1}{\sqrt{p}} \right)
\end{aligned}$$

Finally, the remaining terms in Δ can be estimated as follows:

$$\begin{aligned}
\Delta &= \frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[- \frac{\mathbf{u}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{X}_{-i} \mathbf{v}}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right. \\
& \quad \left. - \frac{\mathbf{v}_i \mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i} - \mathbf{u}_i \mathbf{v}_i \frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] + \mathcal{O} \left(\frac{1}{\sqrt{p}} \right) \\
&= - \frac{2}{p} \frac{\delta \mathbf{u}^\top \mathbf{1} \bar{\mathbf{x}}^\top \tilde{\mathbf{Q}} \bar{\mathbf{x}} \mathbf{1}^\top \mathbf{v}}{1 + \delta} - \frac{\delta^2 \mathbf{u}^\top \mathbf{v}}{1 + \delta} + \mathcal{O} \left(\sqrt{\frac{\log p}{p}} \right)
\end{aligned}$$

Where the last equality is obtained through the following estimation:

$$\begin{aligned}
\frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{v}_i \mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i}{1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] &= \frac{1}{p} \sum_{i=1}^n \mathbb{E} \left[\frac{\mathbf{v}_i \mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i \left(\frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i (1 + \delta) - \delta \left(1 + \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i \right) \right)}{\left(1 + \frac{1}{p} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i \right) (1 + \delta)} \right] \\
& \quad + \frac{1}{p} \sum_{i=1}^n \frac{\mathbf{v}_i \delta \mathbb{E}[\mathbf{u}^\top \mathbf{X}_{-i}^\top \mathbf{Q}_{-i} \mathbf{x}_i]}{(1 + \delta)}
\end{aligned}$$

With the following bound:

$$\begin{aligned}
& \left| \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i (1 + \delta) - \delta \left(1 + \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i \right) \right| \\
&= \left| \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i (1 + \delta) - \delta (1 + \delta) + \delta (1 + \delta) - \delta \left(1 + \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i \right) \right| \\
&\leq \left| \frac{1}{p} \mathbf{x}_i^\top \tilde{\mathbf{Q}} \mathbf{x}_i - \delta \right| (1 + 2\delta),
\end{aligned}$$

we have again with Holder's inequality and Proposition C.2:

$$\frac{1}{p^2} \sum_{i=1}^n \mathbb{E} \left[\frac{v_i u^\top X_{-i}^\top Q_{-i} x_i x_i^\top \tilde{Q} x_i}{1 + \frac{1}{p} x_i^\top Q x_i} \right] = \frac{1}{p} \sum_{i=1}^n \frac{v_i \delta u^\top \mathbf{1} \bar{x}^\top \tilde{Q} \bar{x}}{1 + \delta} + \mathcal{O} \left(\sqrt{\frac{\log p}{p}} \right)$$

Now that we estimated Δ , it remains to estimate $\mathbb{E}[\frac{1}{p} X^\top \tilde{Q} X]$. Indeed, given two unit norm vectors $u, v \in \mathbb{R}^n$ we have:

$$\begin{aligned} \mathbb{E} \left[\frac{1}{p} u^\top X^\top \tilde{Q} X v \right] &= \frac{1}{p} \sum_{i,j=1}^n u_i v_j \mathbb{E}[x_i^\top \tilde{Q} x_j] = \frac{1}{p} \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n u_i v_j \bar{x}^\top \tilde{Q} \bar{x} + \sum_{i=1}^n u_i v_i \delta \\ &= \frac{1}{p} \bar{x}^\top \tilde{Q} \bar{x} u^\top \mathbf{1} \mathbf{1}^\top v + \left(\delta - \frac{1}{p} \bar{x}^\top \tilde{Q} \bar{x} \right) u^\top v = \frac{1}{p} \bar{x}^\top \tilde{Q} \bar{x} u^\top M_1 v^\top + \delta u^\top v + \mathcal{O} \left(\frac{1}{p} \right) \end{aligned}$$

since we have $\bar{x}^\top \tilde{Q} \bar{x} = \mathcal{O}(1)$ by Lemma C.2; we introduced the matrix $M_1 = \mathbf{1} \mathbf{1}^\top$. Therefore we have the following estimation:

$$\frac{1}{p} \mathbb{E} [X^\top Q X] = \frac{\delta}{1 + \delta} I_n + \frac{1}{p} \left(\frac{1 - \delta}{1 + \delta} \right) \bar{x}^\top \tilde{Q} \bar{x} M_1 + \mathcal{O}_{\|\cdot\|} \left(\sqrt{\frac{\log p}{p}} \right)$$

where $A = B + \mathcal{O}_{\|\cdot\|}(\alpha(p))$ means that $\|A - B\| = \mathcal{O}(\alpha(p))$. Finally, since R concentrates around its mean, we can then conclude:

$$R = \frac{1}{z} \left(I_n - \frac{1}{p} X^\top Q X \right) = \frac{1}{z} \frac{1}{1 + \delta} I_n + \frac{\delta - 1}{pz(\delta + 1)} \bar{x}^\top \tilde{Q} \bar{x} M_1 + \mathcal{O}_{\|\cdot\|} \left(\sqrt{\frac{\log p}{p}} \right).$$

□

C.2 Proofs of Chapter 4

C.2.1 Proofs of Section 4.1

Let us Taylor-expand K_{ij} in the vicinity of τ , i.e.,

$$\begin{aligned} K_{ij} &= f(\tau) + f'(\tau) \left(\frac{1}{p} \|x_i - x_j\|^2 - \tau \right) + \frac{1}{2} f''(\tau) \left(\frac{1}{p} \|x_i - x_j\|^2 - \tau \right)^2 \\ &\quad + \frac{1}{6} f^{(3)}(\zeta_{ij}^n) \left(\frac{1}{p} \|x_i - x_j\|^2 - \tau \right)^3, \end{aligned}$$

where $\zeta_{ij}^n \in [\frac{1}{p} \|x_i - x_j\|^2, \tau]$.

• *Control of the third order random matrix:* First, we show that the third order matrix term vanishes in operator norm. Indeed, exploiting the concentration

$$\frac{1}{p} \|x_i - x_j\|^2 - \tau \in \mathcal{O} \left(e^{-(n \cdot)^{q/2}} + e^{-(\sqrt{n} \cdot)^q} + e^{-(n \cdot)^q} \right),$$

we have for all $s \in 2\mathbb{N}^*$

$$\mathbb{E} \left| \frac{1}{p} \|x_i - x_j\|^2 - \tau \right|^s = \mathcal{O}(n^{-s/2}).$$

Let $K^{(3)}$ be the matrix with entries $\left(\frac{1}{p}\|x_i - x_j\|^2 - \tau\right)^3$, we have in particular

$$\mathbb{E}\|K^{(3)}\|_F^2 = \sum_{i,j=1}^n \mathbb{E}\left|\frac{1}{p}\|x_i - x_j\|^2 - \tau\right|^6 = \mathcal{O}(n^2 n^{-3}) = \mathcal{O}(n^{-1}).$$

On the other hand, if $f^{(3)}$ is bounded in a neighborhood of τ , we have with probability at least $1 - \delta$ that $\max_{i \neq j} f^{(3)}(\zeta_{ij}^n)$ is bounded. Indeed, since

$$\max_{1 \leq i \neq j \leq n} \left\{ \frac{1}{p}\|x_i - x_j\|^2 - \tau \right\} = \mathcal{O}\left(n^{-1/2} \log\left(\frac{n}{\sqrt{\delta}}\right)^{1/q}\right),$$

we also have

$$\max_{1 \leq i \neq j \leq n} |\zeta_{ij}^n - \tau| = \mathcal{O}\left(n^{-1/2} \log\left(\frac{n}{\sqrt{\delta}}\right)^{1/q}\right),$$

which shows that ζ_{ij}^n is bounded, so is $\max_{i \neq j} f^{(3)}(\zeta_{ij}^n)$ with probability at least $1 - \delta$. Thus

$$\mathbb{E}\|f^{(3)}(\zeta^n) \odot K^{(3)}\|_F^2 \leq \max_{i \neq j} f^{(3)}(\zeta_{ij}^n)^2 \mathbb{E}\|K^{(3)}\|_F^2 = \mathcal{O}(n^{-1}),$$

in particular, involving Markov's inequality and since the operator norm is bounded by the Frobenius norm, we conclude that $\|f^{(3)}(\zeta^n) \odot K^{(3)}\|_{op}$ is $\mathcal{O}(n^{-1/2})$ with probability at least $1 - \delta$.

• *Control of the second order random matrix:* First, let us show that the matrix with entries $E_{ij} = \mathbf{1}_{i \neq j}(\psi_i + \psi_j)z_i^\top z_j$ is vanishing in operator norm. Denote $\Sigma = ZZ - \text{diag}(ZZ)$, we can then express E as

$$E = \Sigma \text{diag}(\psi) + \text{diag}(\psi)\Sigma$$

Fix $\gamma \in (0, 1/2)$, we first show that $\|\Sigma\|_{op} = \mathcal{O}(\log(n)^\gamma)$. Indeed, since $\|\text{diag}(Z'Z)\|_{op}$ is bounded it remains to control $Z'Z$ in operator norm. Involving an ε -net argument, we have

$$\|Z'Z\|_{op} \leq \frac{1}{1 - 2\varepsilon} \max_{u \in \mathcal{E}} u'Z'Z u = \frac{1}{1 - 2\varepsilon} \max_{u \in \mathcal{E}} \|Zu\|_2^2$$

Recalling the q -exponential concentration, we have $Z \in \mathcal{O}(e^{-(\sqrt{n}\cdot)^q})$, and since $Z \mapsto \|Zu\|_2$ is 1-Lipschitz, $\|Zu\|_2 \in \mathcal{O}(e^{-(\sqrt{n}\cdot)^q})$ and thanks to Proposition 2.5, we have

$$\|Zu\|_2^2 \in \mathcal{O}(e^{-(n\cdot)^{q/2}} + e^{-(\sqrt{n}\cdot)^q})$$

Thus, there exists two absolute constants M and σ such that for all $t \geq 2\mathbb{E}\|Zu\|_2^2$, we have

$$\mathbb{P}\{\|Zu\|_2^2 \geq t\} \leq M(e^{-(nt/\sigma)^{q/2}} + e^{-(\sqrt{nt}/\sigma)^q})$$

To control the maximum over the ε -net \mathcal{E} , we need to bound its cardinality. In fact, we have the following Lemma

Lemma C.3 (Cardinality of an ε -net). *There exists an ε -net \mathcal{E} of the unit sphere in n dimensions, satisfying:*

$$|\mathcal{E}| \leq \left(1 + \frac{2}{\varepsilon}\right)^n$$

Thus, we have by the union bound

$$\mathbb{P} \left\{ \max_{u \in \mathcal{E}} \|Zu\|_2^2 \geq t \right\} \leq Me^{n \log(1+2/\varepsilon)} (e^{-(nt/\sigma)^{q/2}} + e^{-(\sqrt{n}t/\sigma)^q}) \equiv p(t)$$

In particular, for $t = \mathcal{O}(\log(n)^\gamma)$ and for all $q \geq 2$

$$\frac{n \log(1+2/\varepsilon)}{(n \frac{t}{\sigma})^{q/2}} = \mathcal{O} \left(\frac{n^{1-q/2}}{\log(n)^{q\gamma/2}} \right) \rightarrow 0 \text{ and } \frac{n \log(1+2/\varepsilon)}{(\sqrt{n} \frac{t}{\sigma})^q} = \mathcal{O} \left(\frac{n^{1-q}}{\log(n)^{q\gamma}} \right) \rightarrow 0$$

which shows that $\|Z^\top Z\|_{op} = \mathcal{O}(\log(n)^\gamma)$ with probability $1 - \mathcal{O}(e^{-(n \log(n)^\gamma)^{q/2}})$.

It remains to control $\text{diag}(\psi)$ in operator norm, we have for all $k \in 2\mathbb{N}$

$$\mathbb{E} \text{tr} \left(\text{diag}(\psi)^k \right) = \sum_{i=1}^n \mathbb{E} \psi_i^k = \mathcal{O}(n^{1-k/2}),$$

thus by Markov's inequality, we can show that $\|\text{diag}(\psi)\|_{op} = \mathcal{O}(n^{-1/2+1/k})$ for all $k \in 2\mathbb{N}$. Finally, since the operator norm is subadditive, we conclude that

$$\|\{\mathbf{1}_{i \neq j}(\psi_i + \psi_j)z_i^\top z_j\}_{i,j=1}^n\|_{op} = \mathcal{O}(n^{-1/2+\gamma} \log(n)^\gamma) \rightarrow 0$$

The kernel random matrix K can thus be expanded as

$$\begin{aligned} K &= f(\tau)' + f'(\tau) \left[\psi^\top + \psi^\top + \left\{ t_a \frac{\mathbf{1}_{n_a}}{\sqrt{p}} \right\}_{a=1}^k \quad \top + \left\{ t_b \frac{\mathbf{1}_{n_b}^\top}{\sqrt{p}} \right\}_{b=1}^k \right. \\ &\quad + \left\{ \|\bar{m}_a - \bar{m}_b\|^2 \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p} \right\}_{a,b=1}^k + 2 \left\{ \frac{1}{\sqrt{p}} Z_a^\top (\bar{m}_a - \bar{m}_b) \mathbf{1}_{n_b}^\top \right\}_{a,b=1}^k \\ &\quad \left. - 2 \left\{ \frac{1}{\sqrt{p}} \mathbf{1}_{n_a} (\bar{m}_a - \bar{m}_b)^\top Z_b \right\}_{a,b=1}^k - 2Z^\top Z \right] \\ &\quad + \frac{f''(\tau)}{2} \left[\psi^{2\top} + (\psi^2)^\top + \left\{ t_a^2 \frac{\mathbf{1}_{n_a}}{\sqrt{p}} \right\}_{a=1}^k \quad \top + \left\{ t_b^2 \frac{\mathbf{1}_{n_b}^\top}{\sqrt{p}} \right\}_{b=1}^k \right. \\ &\quad + 2 \left\{ t_a t_b \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p} \right\}_{a,b=1}^k + 2\psi \left\{ t_b \frac{\mathbf{1}_{n_b}^\top}{\sqrt{p}} \right\}_{b=1}^k + 2 \left\{ t_a \frac{\mathbf{1}_{n_a}}{\sqrt{p}} \right\}_{a=1}^k \psi^\top \\ &\quad + \frac{2}{\sqrt{p}} \text{diag}\{t_a \mathbf{1}_{n_a}\}_{a=1}^k \psi^\top + \frac{2}{\sqrt{p}} \psi^\top \text{diag}\{t_b \mathbf{1}_{n_b}\}_{b=1}^k \\ &\quad + 4 \left\{ \text{tr}(C_a C_b) \frac{\mathbf{1}_{n_a} \mathbf{1}_{n_b}^\top}{p^2} \right\}_{a,b=1}^k + 2\psi \psi^\top \left. \right] \\ &\quad + (f(0) - f(\tau) + \tau f'(\tau)) I_n + \mathcal{O}_{\|\cdot\|}(n^{-1/2+\gamma} \log(n)^\gamma) \end{aligned}$$

where $\psi \equiv [\psi_1, \dots, \psi_n]^\top \in \mathbb{R}^n$, $t_\ell \equiv p^{-\frac{1}{2}} \text{tr} \bar{C}_\ell$, $Z \equiv [z_1, \dots, z_n] \in \mathbb{R}^{p \times n}$ and $Z_\ell = [z_{n_1+\dots+n_{\ell-1}+1}, \dots, z_{n_1+\dots+n_\ell}]$ the restriction of Z to the mixture μ_ℓ .

In particular, the centered kernel matrix expresses as follows:

$$\begin{aligned}
\bar{K} &= PKP = f'(\tau)P \left[\left\{ \|\bar{m}_a - \bar{m}_b\|^2 \frac{1_{n_a} 1_{n_b}^\top}{p} \right\}_{a,b=1}^k + 2 \left\{ \frac{1}{\sqrt{p}} Z_a^\top (\bar{m}_a - \bar{m}_b) 1_{n_b}^\top \right\}_{a,b=1}^k \right. \\
&\quad \left. - 2 \left\{ \frac{1}{\sqrt{p}} 1_{n_a} (\bar{m}_a - \bar{m}_b)^\top Z_b \right\}_{a,b=1}^k - 2Z^\top Z \right] P \\
&\quad + \frac{f''(\tau)}{2} P \left[2 \left\{ t_a t_b \frac{1_{n_a} 1_{n_b}^\top}{p} \right\}_{a,b=1}^k + 2\psi \left\{ t_b \frac{1_{n_b}^\top}{\sqrt{p}} \right\}_{b=1}^k + 2 \left\{ t_a \frac{1_{n_a}}{\sqrt{p}} \right\}_{a=1}^k \psi^\top \right. \\
&\quad \left. + 4 \left\{ \text{tr}(C_a C_b) \frac{1_{n_a} 1_{n_b}^\top}{p^2} \right\}_{a,b=1}^k + 2\psi\psi^\top \right] P \\
&\quad + (f(0) - f(\tau) + \tau f'(\tau))P + \mathcal{O}_{\|\cdot\|}(n^{-1/2+\gamma} \log(n)^\gamma)
\end{aligned}$$

where $P = I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^\top$.

Introduce the following notations:

$$\begin{aligned}
M &= [\bar{m}_1, \dots, \bar{m}_k] \in \mathbb{R}^{p \times k}, \quad t = \left\{ \frac{1}{\sqrt{p}} \text{tr} \bar{C}_\ell \right\}_{\ell=1}^k \in \mathbb{R}^k \\
T &= \left\{ \frac{1}{p} \text{tr} \bar{C}_a \bar{C}_b \right\}_{a,b=1}^k \in \mathbb{R}^{k \times k}, \quad J = [j_1, \dots, j_k] \in \mathbb{R}^{n \times k} \\
\Phi &= Z^\top M - \{Z_\ell^\top \bar{m}_\ell\}_{\ell=1}^k \in \mathbb{R}^{n \times k}, \quad U = P \left[\frac{1}{\sqrt{p}} J, \Phi, \psi \right] \in \mathbb{R}^{n \times (2k+1)} \\
A_{11} &= M^\top M - \frac{f''(\tau)}{2f'(\tau)} t t^\top - \frac{f''(\tau)}{f'(\tau)} T - \left\{ \frac{\|\bar{m}_a\|^2 + \|\bar{m}_b\|^2}{2} \right\}_{a,b=1}^k \\
A &= \begin{bmatrix} A_{11} & I_k & -\frac{f''(\tau)}{2f'(\tau)} t \\ I_k & \mathbf{0}_{k \times k} & \mathbf{0}_{k \times 1} \\ -\frac{f''(\tau)}{2f'(\tau)} t^\top & \mathbf{0}_{1 \times k} & -\frac{f''(\tau)}{2f'(\tau)} \end{bmatrix}
\end{aligned}$$

The centred kernel matrix \bar{K} can be approximated by

$$\tilde{K} = -2f'(\tau) [PZ^\top ZP + UAU^\top] + (f(0) - f(\tau) + \tau f'(\tau))P \quad (\text{C.4})$$

such that

$$\|\bar{K} - \tilde{K}\|_2 = \mathcal{O}(n^{-1/2+\gamma} \log(n)^\gamma) \rightarrow 0 \quad (\text{C.5})$$

C.2.2 Proofs of Section 4.2

C.2.2.1 Proof of Theorem 4.3

Proof. We first need the following key Lemma.

Lemma C.4 (A Moment Result). *For $g_{ij}(X) \equiv [\Sigma_p^{1/2}]_{i \cdot} (\frac{1}{n} XX^\top - I_p) [\Sigma_p^{1/2}]_{\cdot j}$, we have, for all $k \in \mathbb{N}$ and for some absolute constant $C_k > 0$,*

$$\mathbb{E}|g_{ij}(X)|^{2k} \leq C_k \frac{\beta_p^{4k}}{n^k}. \quad (\text{C.6})$$

Proof. Given a random variable Z , we have

$$\forall m > 0, \mathbb{E}|Z|^m = \int_0^\infty m t^{m-1} \mathbb{P}\{|Z| \geq t\} dt,$$

whenever the right hand side is finite. Applying this identity to the random variable $g_{ij}(X)$ with $m = 2k$ and exploiting the concentration property in Lemma 4.2 yields the result. \square

The proof starts by a Taylor expansion of F_{ij} in the vicinity of $[\Sigma_p]_{ij}$, i.e.,

$$F_{ij} = \sum_{k=0}^2 \frac{f^{(k)}(\sigma_{ij})}{k!} F_{ij}^{(k)} + \frac{f^{(3)}(\zeta_{ij}^n)}{6} F_{ij}^{(3)}$$

where $\sigma_{ij} = [\Sigma_p]_{ij}$, $\zeta_{ij}^n \in [Y Y^\top / n]_{ij}, \sigma_{ij}]$,¹ and $F^{(k)}$ is the matrix with entries

$$F_{ij}^{(k)} \equiv [\Sigma_p^{1/2} (n^{-1} X X^\top - I_p) \Sigma_p^{1/2}]_{ij}^k = g_{ij}(X)^k.$$

We have by Lemma 4.2 that $[Y Y^\top / n]_{ij}$ concentrates around σ_{ij} , so that ζ_{ij}^n is bounded by $\sigma_{ij} + \varepsilon$, for all $\varepsilon > 0$, with high probability² (note that the condition $\max_{ij} |\sigma_{ij}| < B$ ensures that σ_{ij} is bounded and the condition on β_p ensures the quasi-exponential concentration of $[Y Y^\top / n]_{ij}$ around σ_{ij} ; see considered Assumptions above), formally

$$\begin{aligned} \mathbb{P}\left\{|\zeta_{ij}^n| \geq \sigma_{ij} + \varepsilon\right\} &\leq \mathbb{P}\left\{|g_{ij}(X)| \geq \varepsilon\right\} \leq K e^{-\frac{n}{\beta_p^2} \min(c_1 \varepsilon, \frac{c_2 \varepsilon^2}{\beta_p^2})} \\ &\leq K e^{-K' n^{\frac{1}{2} + 2\varepsilon} \min(c_1 \varepsilon, K' c_2 \varepsilon^2 n^{-\frac{1}{2} + 2\varepsilon})} \equiv p_n \rightarrow 0, \end{aligned}$$

where $K' > 0$. And since $f^{(3)}$ is continuous, we deduce that $f^{(3)}(\zeta_{ij}^n)$ is in particular bounded by

$$A \equiv \max_{x \in [\sigma_{ij} - \varepsilon, \sigma_{ij} + \varepsilon]} |f^{(3)}(x)|,$$

with probability $1 - p_n$. Knowing that the operator norm is bounded by the Frobenius norm, we look for a control of the Frobenius norm of the tailing term. We have

$$\|f^{(3)}(\zeta^n) \odot F^{(3)}\|_F^2 \leq A^2 \|F^{(3)}\|_F^2. \quad (\text{C.7})$$

By Lemma C.4, for all $k \in \mathbb{N}$

$$\mathbb{E} \|F^{(k)}\|_F^2 = \sum_{i,j=1}^p \mathbb{E} \left[|g_{ij}(X)|^{2k} \right] \leq C_k \frac{p^2 \beta_p^{4k}}{n^k},$$

for some absolute constant $C_k > 0$. Thus, by *Markov's inequality*, we have for all $\eta > 0$

$$\mathbb{P} \left\{ \|F^{(k)}\|_F \geq \frac{p \beta_p^{2k}}{n^{\frac{k}{2}}} \sqrt{\frac{C_k}{\eta}} \right\} \leq \eta.$$

Recalling Eq. equation C.7, we have with probability at least $1 - \eta$

$$\|f^{(3)}(\zeta^n) \odot F^{(3)}\|_F \leq C \frac{p \beta_p^6}{n^{\frac{3}{2}} \sqrt{\eta}}.$$

\square

¹The notation $[a, b]$ stands for the interval $[a, b]$ if $a < b$ or $[b, a]$ otherwise.

²For a given asymptotic variable n , we say that an event E_n occurs with high probability when it exist a function $\psi(n)$ quasi-exponentially decreasing in n such that $\mathbb{P}\{E_n\} \geq 1 - \psi(n)$.

C.2.2.2 Proof of Theorem 4.4

Proof. The proof needs the introduction of the following two lemmas, that can be found in [EK08] and which are a consequence of the ε -sparsity notion³

Lemma C.5. *Given an ε -sparse $p \times p$ real symmetric matrix M and calling $m = \max_{ij} |M_{ij}|$, we have, for all $k \in 2\mathbb{N}$*

$$\|M\|_{op} \leq (M^k)^{1/k} = \mathcal{O}(m p^{\varepsilon(1-1/k)+1/k}). \quad (\text{C.8})$$

Lemma C.6. *Given two real symmetric matrices M and N with $|M_{ij}| \leq N_{ij}$. Then, we have $\|M\|_{op} \leq \|N\|_{op}$.*

First, we show that when Σ_p is ε -sparse, the Hadamard product $f^{(k)}(\Sigma_p) \odot F^{(k)}$ is of vanishing operator norm for $k \geq 1$, precisely

Lemma C.7. *Let $\mu > 0$, suppose Σ_p is a $\frac{1}{2+\mu}$ -sparse matrix. For f a real and differentiable function, $k \in \{1, 2\}$ such that $f^{(k)}(0) = 0$ and for $\eta > 0$, we have for all $\varepsilon \in (0, \frac{k(2+\mu)-2}{2(3+2\mu)})$*

$$\|f^{(k)}(\Sigma_p) \odot F^{(k)}\|_{op} = \mathcal{O}_\eta^{[1/\varepsilon]} \left(n^{\frac{2-k(2+\mu)}{2(2+\mu)} + \varepsilon(2 - \frac{1}{2+\mu})} \right).$$

Proof. We start by proving that the matrix $F^{(k)}$ has entries of order $\mathcal{O}(n^{-k/2})$. In fact, we have by Lemma C.4, for all $m \in \mathbb{N}^*$

$$\mathbb{E}|F_{ij}^{(k)}|^{2m} = \mathbb{E}|g_{ij}(X)|^{2km} = \mathcal{O}(n^{-km}),$$

thus applying *Markov's inequality* to the random variable $|F_{ij}^{(k)}|^{2m}$ yields to the following tail control.

$$\mathbb{P}\{|F_{ij}^{(k)}| \geq t\} \leq \frac{\mathbb{E}|F_{ij}^{(k)}|^{2m}}{t^{2m}} \leq C n^{-km} t^{-2m},$$

where C is an absolute constant. Recalling Assumption A1 and by the union bound, we have

$$\mathbb{P}\{\max_{ij} |F_{ij}^{(k)}| \geq t\} \leq \sum_{i,j=1}^p \mathbb{P}\{|F_{ij}^{(k)}| \geq t\} \leq p^2 \mathbb{P}\{|F_{ij}^{(k)}| \geq t\} \leq C n^{2-km} t^{-2m},$$

which implies for $\eta > 0$ and for all $m > 0$

$$\max_{ij} |F_{ij}^{(k)}| = \mathcal{O}_\eta^m \left(n^{-\frac{k}{2} + \frac{1}{m}} \right) \quad (\text{C.9})$$

Besides, let M be the matrix defined as $M \equiv \max_{ij} |F_{ij}^{(k)}| \cdot f^{(k)}(\Sigma_p)$, we have

$$|[f^{(k)}(\Sigma_p) \odot F^{(k)}]_{ij}| \leq M_{ij},$$

³Through the identity $(M^k) \leq \max_{ij} |M_{ij}|^k \cdot |\mathcal{C}_p(k)|$.

thus, one has by Lemma C.6

$$\|f^{(k)}(\Sigma_p) \odot F^{(k)}\|_{op} \leq \|M\|_{op} = \max_{ij} |F_{ij}^{(k)}| \cdot \|f^{(k)}(\Sigma_p)\|_{op}.$$

In particular, since $f^{(k)}(\Sigma_p)$ is $\frac{1}{2+\mu}$ -sparse (by Remark 4.1), we have by Lemma C.5 and by equation C.9, for some $\eta > 0$

$$\|f^{(k)}(\Sigma_p) \odot F^{(k)}\|_{op} = \mathcal{O}_\eta^{2m} \left(n^{\frac{1}{2+\mu}(1-\frac{1}{2m}) + \frac{1}{2m} - \frac{k}{2} + \frac{1}{2m}} \right),$$

choosing $\epsilon = \frac{1}{2m} < \frac{k(2+\mu)-2}{2(3+2\mu)}$ yields the final result. \square

When considering f such that $f'(0) = f''(0) = 0$, the result holds by Corollary 4.1 and Lemma C.7. In fact, the dominant order corresponds to $k = 1$ in Lemma C.7. Which completes the proof. \square

C.3 Proofs of Chapter 5

C.3.1 Proofs of Section 5.1

This section provides the exact computation of the fixed point equation in Theorem 5.3. We start by introducing the main tools to perform the calculations.

Theorem C.1 giving the central limit theorem for concentrated vectors was originally proven for uniform distributions on convex subspaces of \mathbb{R}^p , but it was quickly understood that the result is true for a larger class of random vectors satisfying a so-called “thin shell property” (see [Fre19] for a simple and complete proof of this inference). The thin shell property expresses the fact that a random vector X lies principally on a thin shell around a sphere with the following inequality satisfied for some $\epsilon > 0$

$$\mathbb{P} \left(\left| \frac{\|X\|}{\sqrt{p}} - 1 \right| \geq \epsilon \right) \leq \epsilon.$$

For a concentrated random vector $X \propto \mathcal{E}_2$, such that $\mathbb{E}[XX^T] = I_p$ and $\sqrt{p} = O(\mathbb{E}[\|X\|])$ ⁴ the norm being a 1-Lipschitz observation, we know (as in equation 5.3) that there exist two constants $C, c > 0$ such that:

$$\mathbb{P} (|\|X\| - \mathbb{E}[\|X\||] \geq \epsilon) \leq C e^{-c\epsilon^2}.$$

Integrating this concentration inequality for $\epsilon \in [0, \infty)$ with Fubini Theorem, we have $\forall r \geq 2$:

$$\mathbb{E} \left[\left| \frac{\|X\|}{\mathbb{E}[\|X\|]} - 1 \right|^r \right] \leq C(r/2c)^{r/2} \mathbb{E}[\|X\|]^r.$$

Therefore, since $\mathbb{E}[\|X\|] \leq \sqrt{\mathbb{E}[\|X\|^2]} = \sqrt{p} = O(\mathbb{E}[\|X\|])$, we can deduce from Hölder’s inequality that:

$$\mathbb{E} \left[\left| \frac{\|X\|^2}{\mathbb{E}[\|X\|]^2} - 1 \right| \right] \leq \sqrt{\mathbb{E} \left[\left| \frac{\|X\|}{\mathbb{E}[\|X\|]} - 1 \right|^2 \right] \mathbb{E} \left[\left| \frac{\|X\|}{\mathbb{E}[\|X\|]} + 1 \right|^2 \right]} = O(1/\sqrt{p}),$$

⁴This is a very loose hypothesis needed to set that $\frac{\|X\|}{\mathbb{E}[\|X\|]}$ is sufficiently concentrated. Generally if $\sqrt{p} \ll \mathbb{E}[\|X\|]$, that means that one can obtain a better concentration than $X \propto \mathcal{E}_2$

from which we conclude that: $\mathbb{E} [\|X\|]^2 = \mathbb{E} [\|X\|^2] + O(\sqrt{p})$ and therefore $\mathbb{E} [\|X\|] = \sqrt{p} + O(p^{1/4})$. Choosing $\varepsilon = p^{-1/4}$ yields from the concentration of $\|X\|$ to the existence of some constant $K > 0$ such that:

$$\mathbb{P} \left(\left| \frac{\|X\|}{\sqrt{p}} - 1 \right| \geq Kp^{-1/4} \right) \leq Ce^{-cp^{1/2}} \leq Kp^{-1/4},$$

for p large enough. We can then infer (see [Fre19]), that the projections on small dimensional vector spaces of a concentrated vector are Gaussian vectors with high probability.

Theorem C.1 (CLT for concentrated vectors [Kla07, FGP07]). *Given a random vector $X \in \mathbb{R}^p$, and noting G , the cumulative distribution function of a Gaussian variable of zero mean and unit variance. If $X \propto \mathcal{E}_2$, $\mathbb{E}[X] = 0$ and $\mathbb{E}[XX^T] = I_p$, then for any integer $k, n \in \mathbb{N}$, small compared to p , for any $\eta \in (0, 1)$, there exists two constants $C, c > 0$ and a set $\Theta \subset \mathcal{S}^{p-1}$ such that $\sigma(\Theta) \geq 1 - \sqrt{p}Ce^{-c\sqrt{p}}$ and $\forall \theta = (\theta_1, \dots, \theta_k) \in \Theta^k$, there exists a Gaussian vector $Z \sim \mathcal{N}(0, \theta^T \theta)$ such that:*

$$\forall a \in \mathbb{R}^k : \sup_{t \in \mathbb{R}} |\mathbb{P}(a^T \theta^T X \geq t) - \mathbb{P}(a^T Z \geq t)| \leq Cp^{-1/4}.$$

We need a simple preliminary Lemma to state our Stein-like formula for concentrated vectors from Theorem C.1.

Lemma C.8. *Given two random variables $X, Y \in \mathbb{R}^k$, if:*

$$\sup_{t \in \mathbb{R}} |\mathbb{P}(a^T X \geq t) - \mathbb{P}(a^T Y \geq t)| \leq \varepsilon,$$

then for any differentiable mapping $f : \mathbb{R}^k \rightarrow \mathbb{R}$ integrable and bounded around ∞ by f_∞ :

$$|\mathbb{E} [f(X)] - \mathbb{E} [f(Y)]| \leq \frac{f_\infty + \int |f|}{\varepsilon}$$

Proof. Let us prove it in the case $k = 1$:

$$\begin{aligned} |\mathbb{E} [f(X)] - \mathbb{E} [f(Y)]| &\leq \left| \int_{t=-\infty}^{\infty} f(t)(d\mathbb{P}_X(t) - d\mathbb{P}_Y(t)) \right| \\ &\leq |f(t)(\mathbb{P}(X \geq t) - \mathbb{P}(Y \geq t))|_{-\infty}^{\infty} \\ &\quad + \left| \int_{t=-\infty}^{\infty} f'(t)((\mathbb{P}(X \geq t) - \mathbb{P}(Y \geq t))dt \right| \\ &\leq \frac{f_\infty + \int |f|}{\varepsilon} \end{aligned}$$

□

We have the following Stein-like [LN08] theorem for concentrated vectors which is the central tool to express the fixed point mappings in equation (8) of the Main Paper.

Proposition C.3. *Let $x \in \mathbb{R}^p$ be a random vector satisfying the concentration $x \propto \mathcal{E}_2$ (denote $m \equiv \mathbb{E}[x]$ and $C \equiv \mathbb{E}[xx^T]$) and let $f : \mathbb{R}^p \rightarrow \mathbb{R}$ be some three times differentiable function such that f, f' and f'' satisfy the hypotheses of Lemma C.8 with $f_\infty, \int |f|, f'_\infty, \int |f'|, f''_\infty, \int |f''| = O(1)$. Then, there exists a subset $\Theta \subset \mathcal{S}^{p-1}$ such that: $\sigma(\Theta) \geq 1 - Ce^{-cp/\log p}$ and $\forall w, v, u \in \Theta$:*

$$\begin{aligned} \mathbb{E}[f(w^T x)v^T x] &= \mathbb{E}[f(w^T x)]v^T m + \mathbb{E}[f'(w^T x)]v^T Cw + O\left(\frac{1}{n^{1/4}}\right) \\ \mathbb{E}[f(w^T x)v^T xx^T u] &= \mathbb{E}[f(w^T x)]v^T (mm^T + C)u + \mathbb{E}[f'(w^T x)]v^T (Cwm^T + mw^T C)u \\ &\quad + \mathbb{E}[f''(w^T x)]v^T Cww^T Cu + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

Proof. Let us first consider a random vector \mathbf{z} with zero mean and identity covariance satisfying $\mathbf{z} \propto \mathcal{E}_2$. Considering the subset $\Theta \subset \mathbb{S}^{p-1}$ mentioned in Theorem C.1 (for $k = 3$), we know that $\sigma(\Theta) \geq 1 - Ce^{-cp^{1-\varepsilon}}$ for two constants $C, c > 0$ and furthermore for any $\boldsymbol{\theta} = (\mathbf{w}, \mathbf{v}, \mathbf{u}) \in \Theta^3$, and some Gaussian random vector $\mathbf{Z} \sim \mathcal{N}(0, \boldsymbol{\theta}^\top \boldsymbol{\theta})$:

$$\forall a \in \mathbb{R}^2 : \sup_{t \in \mathbb{R}} |\mathbb{P}(a^\top \boldsymbol{\theta}_i^\top \mathbf{z} \geq t) - \mathbb{P}(a^\top \mathbf{Z} \geq t)| = O\left(\frac{1}{n^{1/4}}\right).$$

Given a mapping $g : \mathbb{R} \rightarrow \mathbb{R}$, we know thanks to Lemma C.8 and Stein's identity:

$$\mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] = \mathbb{E}[g(e_1^\top \mathbf{Z})e_1^\top \mathbf{Z}] + O\left(\frac{1}{n^{1/4}}\right) = \mathbb{E}[g'(\mathbf{w}^\top \mathbf{z})] \|\mathbf{w}\|^2 + O\left(\frac{1}{n^{1/4}}\right),$$

where $e_1 = (1, 0)$. Second, if we note $\tilde{\boldsymbol{\theta}} = (\mathbf{w}, \tilde{\mathbf{v}}, \tilde{\mathbf{u}}) = \mathbf{Q}\mathbf{R}$, the QR-decomposition of $\boldsymbol{\theta}$, and $\mathbf{Q} = (\mathbf{q}_1, \mathbf{q}_2, \mathbf{q}_3)$ (of course, $\mathbf{q}_1 = e_1$), we know that $e_1^\top \mathbf{Z}$ and $\mathbf{q}_2^\top \mathbf{Z}$ are independent (\mathbf{Z} is Gaussian and $\mathbb{E}[e_1^\top \mathbf{Z} \mathbf{q}_2^\top \mathbf{Z}] = e_1^\top \boldsymbol{\theta}^\top \boldsymbol{\theta} \mathbf{q}_2 = \mathbf{w}^\top \tilde{\mathbf{v}} = 0$). We can therefore estimate:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] = \mathbb{E}[g(e_1^\top \mathbf{Z})] \mathbb{E}[\mathbf{q}_2^\top \mathbf{Z}] + O\left(\frac{1}{n^{1/4}}\right) = O\left(\frac{1}{n^{1/4}}\right).$$

Combing those 2 estimations, we see that for any differentiable function g $\mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{v}^\top \mathbf{z}] = \mathbb{E}[g'(\mathbf{w}^\top \mathbf{z})]\mathbf{v}^\top \mathbf{w}$. Therefore if we take for g the mapping $t \mapsto f(\mathbf{w}^\top \mathbf{m} + t)$ (satisfying $f(\mathbf{w}^\top \mathbf{x}) = g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})$), we get the identity:

$$\begin{aligned} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})\mathbf{v}^\top \mathbf{x}] &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{v}^\top \mathbf{m}] + \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})]\mathbf{v}^\top \mathbf{C}^{1/2} \mathbf{w} + O\left(\frac{1}{n^{1/4}}\right) \\ &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})]\mathbf{v}^\top \mathbf{m} + \mathbb{E}[g'(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})]\mathbf{v}^\top \mathbf{C} \mathbf{w} + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

With the same method, let us first compute:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{w} = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{w}\mathbf{z}] \mathbb{E}[\mathbf{z}^\top \tilde{\mathbf{v}}] = O\left(\frac{1}{n^{1/4}}\right) = \mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \tilde{\mathbf{u}}.$$

Second:

$$\tilde{\mathbf{v}}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \tilde{\mathbf{u}} = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] \mathbb{E}[\tilde{\mathbf{v}}^\top \mathbf{z}\mathbf{z}^\top \tilde{\mathbf{u}}] + O\left(\frac{1}{n^{1/4}}\right) = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] \tilde{\mathbf{v}}^\top \tilde{\mathbf{u}} + O\left(\frac{1}{n^{1/4}}\right).$$

Third:

$$\begin{aligned} \mathbf{w}^\top \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{w} &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\mathbf{w}^\top \mathbf{z}\mathbf{w}^\top \mathbf{z}] = \mathbb{E}[g(\mathbf{w}^\top \mathbf{z}) + g'(\mathbf{w}^\top \mathbf{z})\mathbf{w}^\top \mathbf{z}] \|\mathbf{w}\|^2 + \\ &= \mathbb{E}[g(\mathbf{w}^\top \mathbf{z})] \|\mathbf{w}\|^2 + \mathbb{E}[g''(\mathbf{w}^\top \mathbf{z})] \|\mathbf{w}\|^4 + O\left(\frac{1}{n^{1/4}}\right). \end{aligned}$$

Therefore, $\mathbb{E}[g(\mathbf{w}^\top \mathbf{z})\tilde{\mathbf{v}}^\top \mathbf{z}\mathbf{z}^\top \tilde{\mathbf{v}}] = \mathbb{E}[(g(\mathbf{w}^\top \mathbf{z})]\tilde{\mathbf{v}}^\top \tilde{\mathbf{u}} + \mathbb{E}[g''(\mathbf{w}^\top \mathbf{z})]\tilde{\mathbf{v}}^\top \mathbf{w}\mathbf{w}^\top \tilde{\mathbf{u}}$, and we can conclude as before that:

$$\begin{aligned} \tilde{\mathbf{v}}^\top \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})\mathbf{x}\mathbf{x}^\top] \mathbf{u} &= \tilde{\mathbf{v}}^\top \mathbf{m} \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \mathbf{m}^\top \mathbf{u} + \tilde{\mathbf{v}}^\top \mathbf{m} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}^\top] \mathbf{C}^{1/2} \mathbf{u} \\ &\quad + \tilde{\mathbf{v}}^\top \mathbf{C}^{1/2} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}] \mathbf{m} \mathbf{u} + \tilde{\mathbf{v}}^\top \mathbf{C}^{1/2} \mathbb{E}[g(\mathbf{w}^\top \mathbf{C}^{1/2} \mathbf{z})\mathbf{z}\mathbf{z}^\top] \mathbf{C}^{1/2} \mathbf{u} + O\left(\frac{1}{n^{1/4}}\right) \\ &= \mathbb{E}[f(\mathbf{w}^\top \mathbf{x})] \tilde{\mathbf{v}}^\top (\mathbf{m}\mathbf{m}^\top + \mathbf{C}) \mathbf{u} \\ &\quad + \mathbb{E}[f'(\mathbf{w}^\top \mathbf{x})] \tilde{\mathbf{v}}^\top (\mathbf{C}\mathbf{w}\mathbf{m}^\top + \mathbf{m}\mathbf{C}\mathbf{w}^\top) \mathbf{u} \\ &\quad + \mathbb{E}[f''(\mathbf{w}^\top \mathbf{x})] \tilde{\mathbf{v}}^\top \mathbf{C}\mathbf{w}\mathbf{w}^\top \mathbf{C}\mathbf{u} + O\left(\frac{1}{n^{1/4}}\right) \end{aligned}$$

□

Let us now employ Proposition C.3 to express the mappings (defined here for a random vector x_i in the class \mathcal{C}_ℓ)

$$\begin{aligned} m_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \mathbf{\Lambda}^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i})] \\ c_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \frac{1}{n} \mathbf{\Lambda}^{-1} \mathbb{E}[\tilde{x}_i h_i(\tilde{x}_i^\top \mathbf{W}_{-i}) h_i(\tilde{x}_i^\top \mathbf{W}_{-i})^\top \tilde{x}_i^\top] \mathbf{\Lambda}^{-1}, \end{aligned}$$

from the class-wise means and covariances of x_1, \dots, x_n and \mathbf{m}_W and \mathbf{C}_W (that are respectively the mean and covariance of \mathbf{W} but also of \mathbf{W}_{-i}). We are going to fix successively x_i and \mathbf{W}_{-i} to be able to compute the expectations appearing in the formulations of m_ℓ and c_ℓ (it is made possible since x_i and \mathbf{W}_{-i} are independent). Although it is not fully rigorous, we employ Proposition C.3 as if the estimations of $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \mathbf{v}^\top \mathbf{x}]$ and $\mathbb{E}[f(\mathbf{w}^\top \mathbf{x}) \mathbf{v}^\top \mathbf{x} \mathbf{x}^\top \mathbf{u}]$ for \mathbf{w}, \mathbf{v} and \mathbf{u} belonging to a big subset of \mathbb{S}^{p-1} of measure bigger than $1 - Ce^{-cp-1/4}$ implied that the result would be true for all vectors $\mathbf{w}, \mathbf{v}, \mathbf{u} \in \mathbb{S}^{p-1}$. It is of course not rigorously correct, however, in practice, for the vectors \mathbf{w}, \mathbf{v} and \mathbf{u} we are considering, it appears to be valid.

To simplify the expression of the derivative in \tilde{x}_i of $h_i(\tilde{x}_i^\top \mathbf{W})$, let us replace our couple of variables $(\tilde{x}_i, \mathbf{W}) \in \mathbb{R}^{p \times k} \times \mathbb{R}^{pk}$ by the variables $(x_i, \tilde{\mathbf{W}}) \in \mathbb{R}^p \times \mathbb{R}^{p \times k}$ where $\tilde{\mathbf{W}} = (\mathbf{w}_1, \dots, \mathbf{w}_k)$. We have then $h_i(\tilde{x}_i^\top \mathbf{W}) = h_i(\tilde{\mathbf{W}}^\top x_i)$. Given a twice differentiable mapping $\phi : \mathbb{R}^k \rightarrow \mathbb{R}$ and $\psi : \mathbf{v} \mapsto \phi(\tilde{\mathbf{W}}^\top \mathbf{v})$ we have the identities:

$$\nabla \psi(\mathbf{v}) = \tilde{\mathbf{W}} \nabla \phi(\tilde{\mathbf{W}}^\top \mathbf{v}) \quad \text{and} \quad d^2 \mathbf{v} = \tilde{\mathbf{W}} d^2 \phi \tilde{\mathbf{W}}^\top \mathbf{v} \tilde{\mathbf{W}}^\top$$

Now, let us follow our new notations and try to compute $m_\ell(\mathbf{m}_W, \mathbf{C}_W) \equiv \mathbb{E}_{-i} [\sqrt{n} \mathbb{E}_i [x_i h^\ell(\tilde{\mathbf{W}}_{-i}^\top x_i)^\top]]$, where we noted $h^\ell : \mathbb{R}^k \rightarrow \mathbb{R}^k$, the mapping h_i for $k(i) = \ell$ (recall that $k(i)$ provides the class of x_i). Let us decompose the matrix $\mathbf{S}_W = \mathbb{E}[\mathbf{W} \mathbf{W}^\top]$ followingly:

$$\mathbf{S}_W = \left(\begin{array}{c|c|c} \mathbf{S}_W^{1,1} & \dots & \mathbf{S}_W^{1,k} \\ \vdots & & \vdots \\ \mathbf{S}_W^{k,1} & \dots & \mathbf{S}_W^{k,k} \end{array} \right) \in \mathcal{M}_{kp},$$

where for all $a, b \in [k]$, $\mathbf{S}_W^{a,b} \in \mathcal{M}_p$, so that we can introduce the low-dimensional random vector $z \sim \mathcal{N}(\ell, \mathbf{K}^\ell)$ with:

$$\ell \equiv \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \quad \text{and} \quad \mathbf{K}^\ell \equiv (\text{Tr}((\boldsymbol{\Sigma}_\ell + \boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top) \mathbf{S}_W^{a,b}))_{1 \leq a, b \leq k} - \boldsymbol{\mu}_\ell^\top \mathbf{m}_W \mathbf{m}_W^\top \boldsymbol{\mu}_\ell. \quad (\text{C.10})$$

With such a choice, z has the same distribution as $\tilde{x}_i^\top \mathbf{W}_{-i}$ for $k(i) = \ell$. Then we can compute thanks to Proposition C.3:

$$\begin{aligned} m_\ell(\mathbf{m}_W, \mathbf{C}_W) &= \boldsymbol{\mu}_\ell \mathbb{E} [h^\ell(\tilde{\mathbf{W}}^\top x_i)^\top] + \boldsymbol{\Sigma}_\ell \mathbb{E}_{-i} [\tilde{\mathbf{W}}_{-i} \mathbb{E}_i [dh_i^\top \tilde{\mathbf{W}}^\top x_i]] \\ &= \boldsymbol{\mu}_\ell \mathbb{E} [h^\ell(z)^\top] + \boldsymbol{\Sigma}_\ell \mathbf{m}_W \mathbb{E} [dh_i^\top z] + O_{\|\cdot\|} \left(\frac{1}{n^{1/4}} \right), \end{aligned}$$

since $\|\mathbf{C}_W\| \leq O(\sqrt{\log n/n})$ (because $\|\mathbf{W}\| \propto \mathcal{E}_2(\sqrt{\log n/n})$ on \mathcal{A}_X).

To estimate the mapping σ_ℓ , let us note for simplicity $H : \mathbf{v} \mapsto h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v}) h^\ell(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})^\top \in \mathcal{M}_{k,k}$, we know that: $\nabla H_{a,b}(\mathbf{v}) = \tilde{\mathbf{W}}_{-i} J(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})_{a,b}$ and $d^2 H_{a,b} \mathbf{v} = \tilde{\mathbf{W}}_{-i} K(\tilde{\mathbf{W}}_{-i}^\top \mathbf{v})_{a,b} \tilde{\mathbf{W}}_{-i}^\top$, where we introduce for any $\mathbf{u} \in \mathbb{R}^k$ the objects:

$$\begin{aligned} J(\mathbf{u})_{a,b} &= h_a^\ell(\mathbf{u}) \nabla h_b^\ell(\mathbf{u}) + h_b^\ell(\mathbf{u}) \nabla h_a^\ell(\mathbf{u}) \in \mathbb{R}^k \\ K(\mathbf{u})_{a,b} &= \nabla h_a^\ell(\mathbf{u}) \nabla h_b^\ell(\mathbf{u})^\top + h_b^\ell(\mathbf{u}) d^2 h_a^\ell \mathbf{u} + h_a^\ell(\mathbf{u}) d^2 h_b^\ell \mathbf{u} \in \mathcal{M}_k \end{aligned}$$

Following the same strategy as previously, we can show thanks to Proposition C.3 that with the decomposition:

$$c_\ell(\mathbf{m}_W, \mathbf{C}_W) = \left(\begin{array}{c|c|c} c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{1,1} & \dots & c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{1,k} \\ \vdots & & \vdots \\ \hline c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{k,1} & \dots & c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{k,k} \end{array} \right) \in \mathcal{M}_{kp},$$

for any $a, b \in [k]$ and $\ell \in [k]$ such that $k(i) = \ell$:

$$\begin{aligned} \Lambda c_\ell(\mathbf{m}_W, \mathbf{C}_W)_{a,b} \Lambda &= \frac{1}{n} \mathbb{E}[H_{a,b}(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top] \\ &= \frac{1}{n} \Sigma_\ell \tilde{\mathbf{m}}_W \mathbb{E}[J(z)_{a,b}] \boldsymbol{\mu}_\ell^\top + \frac{1}{n} \boldsymbol{\mu}_\ell \mathbb{E}[J(z)_{a,b}]^\top \tilde{\mathbf{m}}_W^\top \Sigma_\ell \\ &\quad + \frac{1}{n} \Sigma_\ell \left(\tilde{\mathbf{m}}_W \mathbb{E}[K(z)_{a,b}] \tilde{\mathbf{m}}_W^\top + \sum_{1 \leq c, d \leq k} S_W^{c,d} (\mathbb{E}[K(z)_{a,b}]_{c,d}) \right) \Sigma_\ell \\ &\quad + \frac{1}{n} \mathbb{E}[h_a^\ell(z) h_b^\ell(z)] (\boldsymbol{\mu}_\ell \boldsymbol{\mu}_\ell^\top + \Sigma_\ell) + O_{\|\cdot\|} \left(\frac{1}{n^{1/4}} \right) \end{aligned}$$

We are then left to estimating the derivatives of h to be able to compute m_ℓ and c_ℓ . From the implicit expression of g_i given by Proposition 4.3 one can deduce the formulas, for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^k$:

$$dg_i \mathbf{v} = -(\Delta_i df_i g(\mathbf{v}) - \mathbf{I}_k)^{-1} \quad \text{and} \quad d^2 g_i \mathbf{v} \cdot \mathbf{u} = dg_i \mathbf{v} \Delta_i (d^2 f_i g(\mathbf{v}) \cdot \mathbf{u}) dg_i \mathbf{v}$$

Then, recalling the identity $h_i = f_i \circ g_i$, one can derive from the upper formulas the expression of the differentiates of h_i thanks to the identities for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^k$:

$$dh_i \mathbf{v} = \Delta_i^{-1} (dg_i \mathbf{v} - \mathbf{I}_k) \quad \text{and} \quad d^2 h_i \mathbf{v} \cdot \mathbf{u} = \Delta_i^{-1} d^2 g_i \mathbf{v} \cdot \mathbf{u}$$

C.3.2 Proofs of Section 5.2

C.3.2.1 Proof of Theorem 5.5

Proof of Theorem 5.5. For clarity, we simply write $\mathbf{Q}(z) = \mathbf{Q}$ and $\delta(z) = \delta$ removing the dependence on z , and let the resolvent \mathbf{Q}_{-i} which is \mathbf{Q} without the i -th datum \mathbf{x}_i defined as

$$\mathbf{Q}_{-i} = \left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top - \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^\top + z \mathbf{I}_p \right)^{-1} \quad (\text{C.11})$$

And let

$$\mathbf{m} \equiv \sqrt{\frac{\varepsilon}{1 + a^2(1 - \varepsilon)}} \boldsymbol{\mu} \quad (\text{C.12})$$

Estimation of $\mathbb{E}[g(\mathbf{x})]$ Using the identity $\mathbf{Q} \mathbf{x}_i = \frac{\mathbf{Q}_{-i} \mathbf{x}_i}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}$, we have for $\mathbf{x} \in \mathcal{C}_a$

$$\begin{aligned} \mathbb{E}[\mathbf{x}^\top \mathbf{w}] &= \frac{1}{n} \mathbb{E}[\mathbf{x}^\top \mathbf{Q} \mathbf{X} \mathbf{y}] = \frac{1}{n} \sum_{i=1}^n y_i \mathbb{E}[\mathbf{x}^\top \mathbf{Q} \mathbf{x}_i] = \frac{1}{n} \sum_{i=1}^n y_i \mathbb{E} \left[\frac{\mathbf{x}^\top \mathbf{Q}_{-i} \mathbf{x}_i}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] \\ &= \frac{1}{n} \sum_{i=1}^n y_i \mathbb{E} \left[\frac{\mathbf{x}^\top \mathbf{Q}_{-i} \mathbf{x}_i}{1 + \delta} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) = (-1)^a \frac{\boldsymbol{\mu}^\top \mathbf{Q} \mathbf{m}}{1 + \delta} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \end{aligned}$$

Estimation of $\mathbb{E}[g(\mathbf{x})^2]$

$$\begin{aligned}
\mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] &= \frac{1}{n^2} \mathbb{E}[\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q} \mathbf{x} \mathbf{x}^\top \mathbf{Q} \mathbf{X} \mathbf{y}] = \frac{1}{n^2} \mathbb{E}[\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q} \mathbf{C}_1 \mathbf{Q} \mathbf{X} \mathbf{y}] \\
&= \frac{1}{n^2} \sum_{i,j=1}^n y_i y_j \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \mathbf{C}_1 \mathbf{Q} \mathbf{x}_j] = \frac{1}{n^2} \sum_{i=1}^n y_i^2 \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \mathbf{C}_1 \mathbf{Q} \mathbf{x}_i] + \frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \mathbf{C}_1 \mathbf{Q} \mathbf{x}_j] \\
&= \frac{1}{n^2} \sum_{i=1}^n y_i^2 \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C}_1 \mathbf{Q}_{-i} \mathbf{x}_i}{(1+\delta)^2}\right] + \frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C}_1 \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^2}\right] + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \\
&= \frac{1}{n} \frac{\text{Tr}(\mathbf{C} \mathbb{E}[\mathbf{Q}_{-i} \mathbf{C}_1 \mathbf{Q}_{-i}])}{(1+\delta)^2} + \frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C}_1 \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^2}\right] + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)
\end{aligned}$$

And using the identity $\mathbf{Q} = \mathbf{Q}_{-i} - \frac{\mathbf{Q}_{-i} \frac{1}{n} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i}}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}$, the second term develops as

$$\begin{aligned}
&\frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{C}_1 \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^2}\right] \\
&= \frac{1}{n^2} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{C}_1 \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^2}\right] - \frac{1}{n^3} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{C}_1 \mathbf{Q}_{-ji} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^3}\right] \\
&\quad - \frac{1}{n^3} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ji} \mathbf{x}_j \mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{C}_1 \mathbf{Q}_{-ji} \mathbf{x}_j}{(1+\delta)^3}\right] + \frac{1}{n^4} \sum_{i \neq j} y_i y_j \mathbb{E}\left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_j \mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{C}_1 \mathbf{Q}_{-ji} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-j} \mathbf{x}_j}{(1+\delta)^4}\right] + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \\
&= \frac{\mathbf{m}^\top \mathbb{E}[\mathbf{Q}_{-ij} \mathbf{C}_1 \mathbf{Q}_{-ji}] \mathbf{m}}{(1+\delta)^2} - \frac{2 \text{Tr}(\mathbb{E}[\mathbf{C} \mathbf{Q} \mathbf{C}_1 \mathbf{Q}])}{n(1+\delta)^3} \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m} + \frac{1}{n^2(1+\delta)^4} (\mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m})^2 \mathbf{m}^\top \mathbb{E}[\mathbf{Q} \mathbf{C}_1 \mathbf{Q}] \mathbf{m} + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)
\end{aligned}$$

where the term $\mathbb{E}[\mathbf{Q} \mathbf{A} \mathbf{Q}]$ is handled by

$$\begin{aligned}
\eta(\mathbf{A}) &\equiv \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \mathbb{E}[\mathbf{Q} \mathbf{A} \mathbf{Q}]) = \frac{(1+\delta) \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}})}{(1+\delta)^2 - \frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \bar{\mathbf{Q}} \mathbf{C}_\varepsilon \bar{\mathbf{Q}})} \\
\Delta(\mathbf{A}) &\equiv \mathbb{E}[\mathbf{Q} \mathbf{A} \mathbf{Q}] = \bar{\mathbf{Q}} \mathbf{A} \bar{\mathbf{Q}} + \frac{\eta(\mathbf{A})}{1+\delta} \bar{\mathbf{Q}} \mathbf{C}_\varepsilon \bar{\mathbf{Q}}
\end{aligned}$$

Putting all together yields to

$$\mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] = \frac{\eta(\mathbf{C}_1)}{(1+\delta)^2} + \frac{\mathbf{m}^\top \Delta(\mathbf{C}_1) \mathbf{m}}{(1+\delta)^2} - \frac{2\eta(\mathbf{C}_1) \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m}}{(1+\delta)^3} + \mathcal{O}\left(\frac{1}{n^2}\right)$$

□

C.3.2.2 Proof of Theorem 5.6

Proof of Theorem 5.6. Using the previous notations and matrix identities, for $\mathbf{x}_i \in \mathcal{C}_a$ a sample from the training set \mathbf{X} , we have:

Estimation of $\mathbb{E}[g(\mathbf{x}_i)]$

$$\begin{aligned}\mathbb{E}[\mathbf{x}_i^\top \mathbf{w}] &= \frac{1}{n} \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \mathbf{X} \mathbf{y}] = \frac{1}{n} \sum_{j=1}^n y_j \mathbb{E}[\mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_j] \\ &= \frac{1}{n} \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}{1 + \frac{1}{n} \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i} \right] + \frac{1}{n} \sum_{j \neq i} y_j \mathbb{E} \left[\frac{\mathbf{x}_i^\top \mathbf{Q}_{-j} \mathbf{x}_j}{(1 + \delta)^2} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \frac{\delta}{1 + \delta} + (-1)^a \frac{\mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m}}{(1 + \delta)^2} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right)\end{aligned}$$

Estimation of $\mathbb{E}[g(\mathbf{x}_i)^2]$

$$\begin{aligned}\mathbb{E}[(\mathbf{x}_i^\top \mathbf{w})^2] &= \frac{1}{n^2} \mathbb{E}[\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q} \mathbf{X} \mathbf{y}] = \frac{1}{n^2} \mathbb{E} \left[\frac{\mathbf{y}^\top \mathbf{X}^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{X} \mathbf{y}}{(1 + \delta)^2} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \frac{1}{n^2} \sum_{j,k=1}^n y_j y_k \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_k}{(1 + \delta)^2} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \frac{1}{n^2} \sum_{j=1}^n y_j^2 \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j}{(1 + \delta)^2} \right] + \frac{1}{n^2} \sum_{j \neq k} y_j y_k \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_k}{(1 + \delta)^2} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \frac{1}{n^2} \mathbb{E} \left[\frac{(\mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i)^2}{(1 + \delta)^2} \right] + \frac{1}{n^2} \sum_{j \neq i} y_j^2 \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_j}{(1 + \delta)^2} \right] + \frac{1}{n^2} \sum_{j \neq i} y_j y_i \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}{(1 + \delta)^2} \right] \\ &\quad + \frac{1}{n^2} \sum_{i \neq j \neq k} y_j y_k \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-i} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_k}{(1 + \delta)^2} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \left(\frac{\delta}{1 + \delta} \right)^2 + \frac{1}{n^2} \sum_{j \neq i} \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-ij} \mathbf{x}_j}{(1 + \delta)^4} \right] + \frac{1}{n^2} \sum_{j \neq i} y_j y_i \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-i} \mathbf{x}_i}{(1 + \delta)^3} \right] \\ &\quad + \frac{1}{n^2} \sum_{i \neq j \neq k} y_j y_k \mathbb{E} \left[\frac{\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{x}_i \mathbf{x}_i^\top \mathbf{Q}_{-ik} \mathbf{x}_k}{(1 + \delta)^4} \right] + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \left(\frac{\delta}{1 + \delta} \right)^2 + \frac{\frac{1}{n} \text{Tr}(\mathbf{C}_\varepsilon \mathbb{E}[\mathbf{Q}_{-ij} \mathbf{C}_\varepsilon \mathbf{Q}_{-ij}])}{(1 + \delta)^4} + \frac{\delta \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m}}{(1 + \delta)^3} \\ &\quad + \frac{1}{n^2} \sum_{i \neq j \neq k} y_j y_k \frac{\mathbb{E}[\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{C}_\varepsilon \mathbf{Q}_{-ik} \mathbf{x}_k]}{(1 + \delta)^4} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right)\end{aligned}$$

where we have previously estimated the term $\frac{1}{n^2} \sum_{i \neq j \neq k} y_j y_k \mathbb{E}[\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{C}_\varepsilon \mathbf{Q}_{-ik} \mathbf{x}_k]$ as

$$\frac{1}{n^2} \sum_{i \neq j \neq k} y_j y_k \mathbb{E}[\mathbf{x}_j^\top \mathbf{Q}_{-ij} \mathbf{C}_\varepsilon \mathbf{Q}_{-ik} \mathbf{x}_k] = \mathbf{m}^\top \mathbb{E}[\mathbf{Q}_{-ijk} \mathbf{C}_\varepsilon \mathbf{Q}_{-ijk}] \mathbf{m} - \frac{2 \text{Tr}(\mathbb{E}[\mathbf{C}_\varepsilon \mathbf{Q} \mathbf{C}_\varepsilon \mathbf{Q}])}{n(1 + \delta)} \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right)$$

Hence, putting all together we get

$$\begin{aligned}\mathbb{E}[g(\mathbf{x}_i)^2] &= \left(\frac{\delta}{1 + \delta} \right)^2 + \frac{\eta(\mathbf{C}_\varepsilon)}{(1 + \delta)^4} + \frac{\delta \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m}}{(1 + \delta)^3} + \frac{\mathbf{m}^\top \Delta(\mathbf{C}_\varepsilon) \mathbf{m}}{(1 + \delta)^4} - \frac{2\eta \mathbf{m}^\top \bar{\mathbf{Q}} \mathbf{m}}{(1 + \delta)^5} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right) \\ &= \left(\frac{\delta}{1 + \delta} \right)^2 + \frac{\eta(\mathbf{C}_\varepsilon)}{(1 + \delta)^4} + \mathbf{m}^\top \left(\frac{\delta \bar{\mathbf{Q}}}{(1 + \delta)^3} + \frac{\Delta(\mathbf{C}_\varepsilon)}{(1 + \delta)^4} - \frac{2\eta(\mathbf{C}_\varepsilon) \bar{\mathbf{Q}}}{(1 + \delta)^5} \right) \mathbf{m} + \mathcal{O} \left(\frac{1}{\sqrt{n}} \right)\end{aligned}$$

and the CLT is obtained with similar arguments than louart2018random. \square

C.3.2.3 Optimality

Denote

$$a = \frac{\varepsilon}{1 + \alpha^2(1 - \varepsilon)}, \quad b = 1 + cq, \quad d = \frac{1 - \varepsilon}{\varepsilon}$$

And suppose, for some $\beta > 0$ and some integer s

$$\boldsymbol{\mu} = \beta \sum_{i=1}^s \frac{(-1)^i}{\sqrt{s}} \mathbf{e}_i$$

where $\mathbf{e}_1, \dots, \mathbf{e}_s$ are the s -first canonical vectors of \mathbb{R}^p , with $[e_i]_i = 1$. The expression of the test missclassification error involves the following terms which depend on α :

$$\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \boldsymbol{\mu}, \quad \boldsymbol{\mu}^\top \Delta \boldsymbol{\mu}$$

Let us first expression $\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \boldsymbol{\mu}$, which develops as

$$q_1 \equiv \boldsymbol{\mu}^\top \bar{\mathbf{Q}} \boldsymbol{\mu} = \boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu} - \frac{a(\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu})^2}{b + a\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}}$$

where

$$d_1 \equiv \boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu} = qb \sum_{i=1}^p \frac{\mu_i^2}{b + aqd\mu_i^2} = \frac{qb\beta^2}{b + aqd\frac{\beta^2}{s}}$$

The term $\boldsymbol{\mu}^\top \Delta \boldsymbol{\mu}$ involves the quantity η which reduces to

$$\eta = \frac{(1 + \delta)cq^2}{(1 + \delta)^2 - cq^2}$$

And thus $\boldsymbol{\mu}^\top \Delta \boldsymbol{\mu}$ develops as

$$\boldsymbol{\mu}^\top \Delta \boldsymbol{\mu} = \boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_1 \bar{\mathbf{Q}} \boldsymbol{\mu} + \frac{\eta}{1 + \delta} \boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_\varepsilon \bar{\mathbf{Q}} \boldsymbol{\mu}$$

where $\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_1 \bar{\mathbf{Q}} \boldsymbol{\mu}$ is given by

$$\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_1 \bar{\mathbf{Q}} \boldsymbol{\mu} = \boldsymbol{\mu}^\top \bar{\mathbf{Q}}^2 \boldsymbol{\mu} + (\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \boldsymbol{\mu})^2$$

And

$$\begin{aligned} q_2 \equiv \boldsymbol{\mu}^\top \bar{\mathbf{Q}}^2 \boldsymbol{\mu} &= \boldsymbol{\mu}^\top \left(\mathcal{D}_z - \frac{a\mathcal{D}_z \boldsymbol{\mu} \boldsymbol{\mu}^\top \mathcal{D}_z}{b + a\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}} \right)^2 \boldsymbol{\mu} \\ &= \boldsymbol{\mu}^\top \mathcal{D}_z^2 \boldsymbol{\mu} - \frac{2a\boldsymbol{\mu}^\top \mathcal{D}_z^2 \boldsymbol{\mu} \boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}}{b + a\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu}} + \frac{a^2 \boldsymbol{\mu}^\top \mathcal{D}_z^2 \boldsymbol{\mu} (\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu})^2}{(b + a\boldsymbol{\mu}^\top \mathcal{D}_z \boldsymbol{\mu})^2} \end{aligned}$$

where

$$d_2 \equiv \boldsymbol{\mu}^\top \mathcal{D}_z^2 \boldsymbol{\mu} = q^2 b^2 \sum_{i=1}^p \frac{\mu_i^2}{(b + aqd\mu_i^2)^2} = \frac{q^2 b^2 \beta^2}{(b + aqd\frac{\beta^2}{s})^2}$$

And the remaining term $\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_\varepsilon \bar{\mathbf{Q}} \boldsymbol{\mu}$ develops as

$$\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \mathbf{C}_\varepsilon \bar{\mathbf{Q}} \boldsymbol{\mu} = \boldsymbol{\mu}^\top \bar{\mathbf{Q}}^2 \boldsymbol{\mu} + a(\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \boldsymbol{\mu})^2 + ad\boldsymbol{\mu}^\top \bar{\mathbf{Q}} \text{diag}(\boldsymbol{\mu}^{\odot 2} + 2\alpha\boldsymbol{\mu}) \bar{\mathbf{Q}} \boldsymbol{\mu}$$

where it remains to express $\mu^\top \bar{Q} \text{diag}(\mu^{\odot 2} + 2\alpha\mu) \bar{Q}\mu$ which develops as

$$\mu^\top \bar{Q} \text{diag}(\mu^{\odot 2} + 2\alpha\mu) \bar{Q}\mu = \mu^\top \bar{Q} \text{diag}(\mu^{\odot 2}) \bar{Q}\mu + 2\alpha\mu^\top \bar{Q} \text{diag}(\mu) \bar{Q}\mu$$

with

$$\begin{aligned} q_3 &\equiv \mu^\top \bar{Q} \text{diag}(\mu^{\odot 2}) \bar{Q}\mu = \mu^\top \left(\mathcal{D}_z - \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \right) \text{diag}(\mu^{\odot 2}) \left(\mathcal{D}_z - \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \right) \mu \\ &= \mu^\top \left(\mathcal{D}_z \text{diag}(\mu^{\odot 2}) - \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu^{\odot 2}) \right) \left(\mathcal{D}_z - \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \right) \mu \\ &= \mu^\top \left(\mathcal{D}_z \text{diag}(\mu^{\odot 2}) \mathcal{D}_z - \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu^{\odot 2}) \mathcal{D}_z - \mathcal{D}_z \text{diag}(\mu^{\odot 2}) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \right. \\ &\quad \left. + \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu^{\odot 2}) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \right) \mu \\ &= \mu^\top \mathcal{D}_z \text{diag}(\mu^{\odot 2}) \mathcal{D}_z\mu - \mu^\top \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu^{\odot 2}) \mathcal{D}_z\mu - \mu^\top \mathcal{D}_z \text{diag}(\mu^{\odot 2}) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \mu \\ &\quad + \mu^\top \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu^{\odot 2}) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \mu \end{aligned}$$

and

$$d_3 \equiv \mu^\top \mathcal{D}_z \text{diag}(\mu^{\odot 2}) \mathcal{D}_z\mu = q^2 b^2 \sum_{i=1}^p \frac{\mu_i^4}{(b + aqd\mu_i^2)^2} = \frac{q^2 b^2 \frac{\beta^2}{s}}{(b + aqd\frac{\beta^2}{s})^2}$$

Similarly $\mu^\top \bar{Q} \text{diag}(\mu) \bar{Q}\mu$ develops as

$$\begin{aligned} q_4 &\equiv \mu^\top \bar{Q} \text{diag}(\mu) \bar{Q}\mu = \mu^\top \mathcal{D}_z \text{diag}(\mu) \mathcal{D}_z\mu - \mu^\top \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu) \mathcal{D}_z\mu \\ &\quad - \mu^\top \mathcal{D}_z \text{diag}(\mu) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \mu + \mu^\top \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \text{diag}(\mu) \frac{a\mathcal{D}_z\mu\mu^\top\mathcal{D}_z}{b + a\mu^\top\mathcal{D}_z\mu} \mu \end{aligned}$$

with

$$d_4 \equiv \mu^\top \mathcal{D}_z \text{diag}(\mu) \mathcal{D}_z\mu = q^2 b^2 \sum_{i=1}^p \frac{\mu_i^3}{(b + aqd\mu_i^2)^2} = \frac{q^2 b^2 \frac{\beta^2}{\sqrt{s}}}{(b + aqd\frac{\beta^2}{s})^2}$$

Thus, putting all together, we have

$$\begin{aligned} k_\varepsilon &\equiv \mu^\top \bar{Q} C_\varepsilon \bar{Q}\mu = q_2 + aq_1^2 + ad \left(q_3 + 2q_4 \sqrt{\frac{\varepsilon - 1}{1 - \varepsilon}} \right) \\ k_1 &\equiv \mu^\top \bar{Q} C_1 \bar{Q}\mu = q_2 + q_1^2 \\ h &\equiv \mu^\top \Delta\mu = k_1 + \frac{\eta}{1 + \delta} k_\varepsilon \end{aligned}$$

where

$$\begin{aligned} q_1 &= d_1 - \frac{ad_1^2}{b + ad_1}, \quad q_2 = d_2 - \frac{2ad_2d_1}{b + ad_1} + \frac{a^2d_2d_1^2}{(b + ad_1)^2} \\ q_3 &= d_3 - \frac{2ad_1d_3}{b + ad_1} + \frac{a^2d_1^2d_3}{(b + ad_1)^2}, \quad q_4 = d_4 - \frac{2ad_1d_4}{b + ad_1} + \frac{a^2d_1^2d_4}{(b + ad_1)^2} \end{aligned}$$

And finally

$$d_1 = \frac{qb\beta^2}{b + aqd\frac{\beta^2}{s}}, \quad d_2 = \frac{q^2b^2\beta^2}{(b + aqd\frac{\beta^2}{s})^2}, \quad d_3 = \frac{q^2b^2\frac{\beta^2}{s}}{(b + aqd\frac{\beta^2}{s})^2}, \quad d_4 = \frac{q^2b^2\frac{\beta^2}{\sqrt{s}}}{(b + aqd\frac{\beta^2}{s})^2}$$

Therefore, the mean and variance of the decision are given by

$$m_\ell = (-1)^\ell \sqrt{a} \frac{q_1}{1 + \delta}$$

$$v = \frac{1}{(1 + \delta)^2} \left(\eta + a \left[h - q_1 - \frac{2\eta q_1}{1 + \delta} \right] \right)$$

with the change of variable $a \leftarrow \frac{\varepsilon}{1 + a^2(1 - \varepsilon)}$, the optimal α is provided by

$$\frac{\partial}{\partial a} Q \left(\frac{m_\ell}{\sqrt{v}} \right) = 0$$

Since $Q'(x) = -\frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$, the above equation reduces to

$$\sqrt{v} \frac{\partial m_\ell}{\partial a} - m_\ell \frac{\partial \sqrt{v}}{\partial a} = 0$$

Bibliography

- [A⁺02] Suykens Johan AK et al. *Least squares support vector machines*. World Scientific, 2002.
- [AKC18] Hafiz Tiomoko Ali, Abla Kammoun, and Romain Couillet. Random matrix-improved kernels for large dimensional spectral clustering. In *IEEE Statistical Signal Processing Workshop 2018*, 2018.
- [AMGS12] Konstantin Avrachenkov, Alexey Mishenin, Paulo Gonçalves, and Marina Sokol. Generalized optimization framework for graph-based semi-supervised learning. In *Proceedings of the 2012 SIAM International Conference on Data Mining*, pages 966–974. SIAM, 2012.
- [And63] Theodore Wilbur Anderson. Asymptotic theory for principal component analysis. *The Annals of Mathematical Statistics*, 34(1):122–148, 1963.
- [APD14] Megasthenis Asteris, Dimitris Papailiopoulos, and Alexandros Dimakis. Nonnegative sparse pca with provable guarantees. In *International Conference on Machine Learning*, pages 1728–1736, 2014.
- [AS17] Madhu S Advani and Andrew M Saxe. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- [ASD18] Joseph Antognini and Jascha Sohl-Dickstein. Pca of high dimensional random walks with comparison to neural network training. In *Advances in Neural Information Processing Systems*, pages 10307–10316, 2018.
- [ASE17] Antreas Antoniou, Amos Storkey, and Harrison Edwards. Data augmentation generative adversarial networks. *arXiv preprint arXiv:1711.04340*, 2017.
- [ASEA19] Ali Alnoman, Shree Krishna Sharma, Waleed Ejaz, and Alagan Anpalagan. Emerging edge computing technologies for distributed iot systems. *IEEE Network*, 33(6):140–147, 2019.
- [B⁺09] Yoshua Bengio et al. Learning deep architectures for ai. *Foundations and trends® in Machine Learning*, 2(1):1–127, 2009.
- [BAP⁺05] Jinho Baik, Gérard Ben Arous, Sandrine Péché, et al. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *The Annals of Probability*, 33(5):1643–1697, 2005.
- [BCV13] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, Aug 2013.

- [BDS18] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- [BEO19] Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can scale to imagenet. In *International Conference on Machine Learning*, pages 583–593, 2019.
- [BGC16] Florent Benaych-Georges and Romain Couillet. Spectral analysis of the gram matrix of mixture models. *ESAIM: Probability and Statistics*, 20:217–237, 2016.
- [BGN11] Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- [BGN12] Florent Benaych-Georges and Raj Rao Nadakuditi. The singular values and vectors of low rank perturbations of large rectangular random matrices. *Journal of Multivariate Analysis*, 111:120–135, 2012.
- [Bil08] Patrick Billingsley. *Probability and measure*. John Wiley & Sons, 2008.
- [BL08] Peter J Bickel and Elizaveta Levina. Covariance regularization by thresholding. *The Annals of Statistics*, pages 2577–2604, 2008.
- [BS⁺98a] Zhi-Dong Bai, Jack W Silverstein, et al. No eigenvalues outside the support of the limiting spectral distribution of large-dimensional sample covariance matrices. *The Annals of Probability*, 26(1):316–345, 1998.
- [BS98b] Sean Borman and Robert L Stevenson. Super-resolution from image sequences-a review. In *Circuits and Systems, 1998. Proceedings. 1998 Midwest Symposium on*, pages 374–378. IEEE, 1998.
- [BS08] Zhidong D Bai and Jack W Silverstein. Clt for linear spectral statistics of large-dimensional sample covariance matrices. In *Advances In Statistics*, pages 281–333. World Scientific, 2008.
- [BS10] Zhidong Bai and Jack W Silverstein. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.
- [BSL15] Joan Bruna, Pablo Sprechmann, and Yann LeCun. Super-resolution with deep convolutional sufficient statistics. *arXiv preprint arXiv:1511.05666*, 2015.
- [CBG⁺16] Romain Couillet, Florent Benaych-Georges, et al. Kernel spectral clustering of large dimensional data. *Electronic Journal of Statistics*, 10(1):1393–1454, 2016.
- [CDMF09] Mireille Capitaine, Catherine Donati-Martin, and Delphine Féral. The largest eigenvalues of finite rank deformation of large wigner matrices: convergence and nonuniversality of the fluctuations. *The Annals of Probability*, pages 1–47, 2009.

- [CHM⁺15] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surfaces of multilayer networks. In *Artificial intelligence and statistics*, pages 192–204, 2015.
- [CJ95] Jorge Cadima and Ian T Jolliffe. Loading and correlations in the interpretation of principle components. *Journal of Applied Statistics*, 22(2):203–214, 1995.
- [CMK⁺14] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [CS13] Xiuyuan Cheng and Amit Singer. The spectrum of random inner-product kernel matrices. *Random Matrices: Theory and Applications*, 2(04):1350010, 2013.
- [CSZ09] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- [DB16] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. In *Advances in Neural Information Processing Systems*, pages 658–666, 2016.
- [DCF⁺15] Emily L Denton, Soumith Chintala, Rob Fergus, et al. Deep generative image models using a laplacian pyramid of adversarial networks. In *Advances in neural information processing systems*, pages 1486–1494, 2015.
- [dCPS⁺18] Remi Tachet des Combes, Mohammad Pezeshki, Samira Shabanian, Aaron Courville, and Yoshua Bengio. On the learning dynamics of deep neural networks. *arXiv preprint arXiv:1809.06848*, 2018.
- [DDS⁺09] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [dGJL05] Alexandre d’Aspremont, Laurent E Ghaoui, Michael I Jordan, and Gert R Lanckriet. A direct formulation for sparse pca using semidefinite programming. In *Advances in neural information processing systems*, pages 41–48, 2005.
- [DK70] Chandler Davis and William Morton Kahan. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- [DLHT14] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European Conference on Computer Vision*, pages 184–199. Springer, 2014.
- [DLHT16] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2016.
- [DLT16] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *European Conference on Computer Vision*, pages 391–407. Springer, 2016.

- [DM14] Yash Deshpande and Andrea Montanari. Sparse pca via covariance thresholding. In *Advances in Neural Information Processing Systems*, pages 334–342, 2014.
- [DM16] Yash Deshpande and Andrea Montanari. Sparse pca via covariance thresholding. *J. Mach. Learn. Res.*, 17(1):4913–4953, January 2016.
- [DTVG15] Dengxin Dai, Radu Timofte, and Luc Van Gool. Jointly optimized regressors for image super-resolution. In *Computer Graphics Forum*, volume 34, pages 95–104. Wiley Online Library, 2015.
- [Duc79] Claude E Duchon. Lanczos filtering in one and two dimensions. *Journal of applied meteorology*, 18(8):1016–1022, 1979.
- [DYP20] Shiyu Duan, Shujian Yu, and Jose Principe. Modularizing deep learning via pairwise learning with kernels. *arXiv preprint arXiv:2005.05541*, 2020.
- [DZSW11] Weisheng Dong, Lei Zhang, Guangming Shi, and Xiaolin Wu. Image deblurring and super-resolution by adaptive sparse domain selection and adaptive regularization. *IEEE Transactions on Image Processing*, 20(7):1838–1857, 2011.
- [EI98] Stanley C Eisenstat and Ilse CF Ipsen. Three absolute perturbation bounds for matrix eigenvalues imply relative bounds. *SIAM Journal on Matrix Analysis and Applications*, 20(1):149–158, 1998.
- [EK08] Nouredine El Karoui. Operator norm consistent estimation of large-dimensional sparse covariance matrices. *The Annals of Statistics*, pages 2717–2756, 2008.
- [EK10a] Nouredine El Karoui. On information plus noise kernel random matrices. *The Annals of Statistics*, 38(5):3191–3216, 2010.
- [EK⁺10b] Nouredine El Karoui et al. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1):1–50, 2010.
- [FGP07] B. Fleury, O. Guédon, and G. Paouris. A stability result for mean width of l_p -centroid bodies. *Advances in Mathematics*, 214:865–877, 2007.
- [Fin92] Helmut Finner. A generalization of holder’s inequality and some probability inequalities. *The Annals of probability*, pages 1893–1901, 1992.
- [FJP02] William T Freeman, Thouis R Jones, and Egon C Pasztor. Example-based super-resolution. *IEEE Computer graphics and Applications*, 22(2):56–65, 2002.
- [FP07] Delphine Féral and Sandrine Péché. The largest eigenvalue of rank one deformation of large wigner matrices. *Communications in mathematical physics*, 272(1):185–228, 2007.
- [Fre19] Daniel J. Fresen. A simplified proof of clt for convex bodies. *arXiv preprint arXiv:1907.06785*, 2019.
- [FREM04] Sina Farsiu, M Dirk Robinson, Michael Elad, and Peyman Milanfar. Fast and robust multiframe super resolution. *IEEE transactions on image processing*, 13(10):1327–1344, 2004.

- [GAA⁺17] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems 30*, pages 5767–5777. 2017.
- [GCM18] Samantha Guerriero, Barbara Caputo, and Thomas Mensink. Deep nearest class mean classifiers. In *International Conference on Learning Representations, Worskhop Track*, 2018.
- [GEB15] Leon Gatys, Alexander S Ecker, and Matthias Bethge. Texture synthesis using convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 262–270, 2015.
- [GEB16] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Image style transfer using convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR).*, pages 2414–2423. IEEE, 2016.
- [GL10] Karol Gregor and Yann LeCun. Learning fast approximations of sparse coding. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, pages 399–406. Omnipress, 2010.
- [GP17] Bolin Gao and Lacra Pavel. On the properties of the softmax function with application in game theory and reinforcement learning. *arXiv preprint arXiv:1704.00805*, 2017.
- [GPAM⁺14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [GSBB11] Prateek Gupta, Priyanka Srivastava, Satyam Bhardwaj, and Vikrant Bhateja. A modified psnr metric based on hvs for quality assessment of color images. In *Communication and Industrial Application (ICCIA), 2011 International Conference on*, pages 1–4. IEEE, 2011.
- [HLM15] Elad Hazan, Roi Livni, and Yishay Mansour. Classification with low rank and missing data. In *ICML*, pages 257–266, 2015.
- [HLN⁺07] Walid Hachem, Philippe Loubaton, Jamal Najim, et al. Deterministic equivalents for certain functionals of large random matrices. *The Annals of Applied Probability*, 17(3):875–930, 2007.
- [HMD15] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015.
- [HRU⁺17] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.
- [HS11] He He and Wan-Chi Siu. Single image super-resolution using gaussian process regression. In *Computer Vision and Pattern Recognition (CVPR).*, pages 449–456. IEEE, 2011.
- [HZS06] Guang-Bin Huang, Qin-Yu Zhu, and Chee-Kheong Siew. Extreme learning machine: theory and applications. *Neurocomputing*, 70(1-3):489–501, 2006.

- [IHM⁺16] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- [IS15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [JAFF16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016.
- [JL09] Iain M Johnstone and Arthur Yu Lu. On consistency and sparsity for principal components analysis in high dimensions. *Journal of the American Statistical Association*, 104(486):682–693, 2009.
- [KC17] Abia J Kammoun and Romain Couillet. Subspace kernel clustering of large dimensional data. (*submitted to*) *Annals of Applied Probability*, 2017.
- [KGC18] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018.
- [KH10] Alex Krizhevsky and Geoff Hinton. Convolutional deep belief networks on cifar-10. *Unpublished manuscript*, 40(7), 2010.
- [KK10] Kwang In Kim and Younghee Kwon. Single-image super-resolution using sparse regression and natural image prior. *IEEE transactions on pattern analysis and machine intelligence*, 32(6):1127–1133, 2010.
- [Kla07] B. Klartag. A central limit theorem for convex sets. *Inventiones mathematicae volume*, 168:pages91–131, 2007.
- [KNV15] Robert Krauthgamer, Boaz Nadler, and Dan Vilenchik. Do semidefinite relaxations solve sparse pca up to the information limit? *Ann. Statist.*, 43(3):1300–1322, 06 2015.
- [KQD18] Alireza Karbalayghareh, Xiaoning Qian, and Edward R Dougherty. Optimal bayesian transfer learning. *IEEE Transactions on Signal Processing*, 2018.
- [KSH12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [KXR⁺19] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition, 2019.
- [KY13] Antti Knowles and Jun Yin. The isotropic semicircle law and deformation of wigner matrices. *Communications on Pure and Applied Mathematics*, 66(11):1663–1749, 2013.

- [LC17] Zhenyu Liao and Romain Couillet. Random matrices meet machine learning: A large dimensional analysis of ls-svm. In *ICASSP*, pages 2397–2401. IEEE, 2017.
- [LC18a] Zhenyu Liao and Romain Couillet. The dynamics of learning: A random matrix approach. *arXiv preprint arXiv:1805.11917*, 2018.
- [LC18b] Cosme Louart and Romain Couillet. Concentration of measure and large random matrices with an application to sample covariance matrices. *arXiv preprint arXiv:1805.08295*, 2018.
- [LC18c] Cosme Louart and Romain Couillet. A random matrix and concentration inequalities framework for neural networks analysis. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4214–4218. IEEE, 2018.
- [LC19] Cosme Louart and Romain Couillet. Concentration of measure and large random matrices with an application to sample covariance matrices. *submitted*, 2019.
- [LC20] Cosme Louart and Romain Couillet. Concentration of measure and large random matrices with an application to sample covariance matrices. *submitted to Random Matrices: Theory and Applications*, 2020.
- [LeC98] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [Led05a] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- [Led05b] Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical Soc., 2005.
- [Lia19] Zhenyu Liao. *A random matrix framework for large dimensional machine learning and neural networks*. Theses, Université Paris-Saclay, September 2019.
- [LLP⁺19] Lauri Lovén, Teemu Leppänen, Ella Peltonen, Juha Partala, Erkki Harjula, Pawani Porambage, Mika Ylianttila, and Jukka Riekk. Edge ai: A vision for distributed, edge-native artificial intelligence in future 6g networks. *The 1st 6G Wireless Summit*, pages 1–2, 2019.
- [LM07] Neil D Lawrence and Andrew J Moore. Hierarchical gaussian process latent variable models. In *Proceedings of the 24th international conference on Machine learning*, pages 481–488, 2007.
- [LN08] Zinoviy Landsman and Johanna Nešlehová. Stein’s lemma for elliptical random vectors. *Journal of Multivariate Analysis*, 99(5):912–927, 2008.
- [LOV19] Sindy Löwe, Peter O’Connor, and Bastiaan Veeling. Putting an end to end-to-end: Gradient-isolated learning of representations. In *Advances in Neural Information Processing Systems*, pages 3039–3051, 2019.
- [LTH⁺16] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. *arXiv preprint*, 2016.

- [LW16] Chuan Li and Michael Wand. Combining markov random fields and convolutional neural networks for image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2479–2486, 2016.
- [LWL⁺17] Xuezhi Liang, Xiaobo Wang, Zhen Lei, Shengcai Liao, and Stan Z Li. Soft-margin softmax for deep classification. In *International Conference on Neural Information Processing*, pages 413–421. Springer, 2017.
- [M⁺13] Zongming Ma et al. Sparse principal component analysis and iterative thresholding. *The Annals of Statistics*, 41(2):772–801, 2013.
- [Mai19] Xiaoyi Mai. *Methods of random matrices for large dimensional statistical learning*. PhD thesis, Université Paris-Saclay, 2019.
- [MBM11] Gregoire Montavon, Mikio L Braun, and Klaus-Robert Miller. Kernel analysis of deep networks. *Journal of Machine Learning Research*, 12(Sep):2563–2581, 2011.
- [MC17] Xiaoyi Mai and Romain Couillet. A random matrix analysis and improvement of semi-supervised learning for large dimensional data. *arXiv preprint arXiv:1711.03404*, 2017.
- [MCL15] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015.
- [MFLG19] Seyed-Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant: Bridging the gap between student and teacher. *arXiv preprint arXiv:1902.03393*, 2019.
- [MKKY18] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- [MLC19] Xiaoyi Mai, Zhenyu Liao, and Romain Couillet. A large scale analysis of logistic regression: Asymptotic performance and new insights. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3357–3361. IEEE, 2019.
- [MLK20] Kurt Wan-Duo Ma, JP Lewis, and W Bastiaan Kleijn. The hsic bottleneck: Deep learning without back-propagation. In *AAAI*, pages 5085–5092, 2020.
- [MMH⁺19] MG Murshed, Christopher Murphy, Daqing Hou, Nazar Khan, Ganesh Ananthanarayanan, and Faraz Hussain. Machine learning at the network edge: A survey. *arXiv preprint arXiv:1908.00080*, 2019.
- [MO14] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [MP67] Vladimir A Marčenko and Leonid Andreevich Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457, 1967.

- [MTS⁺19] Sara Meftah, Youssef Tamaazousti, Nasredine Semmar, Hassane Essafi, and Fatiha Sadat. Joint learning of pre-trained and random units for domain adaptation in part-of-speech tagging. *arXiv preprint arXiv:1904.03595*, 2019.
- [MVPC13] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2624–2637, 2013.
- [MWA06] Baback Moghaddam, Yair Weiss, and Shai Avidan. Spectral bounds for sparse pca: Exact and greedy algorithms. In *Advances in neural information processing systems*, pages 915–922, 2006.
- [NBA⁺18] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Sensitivity and generalization in neural networks: an empirical study. *arXiv preprint arXiv:1802.08760*, 2018.
- [NJW02] Andrew Y Ng, Michael I Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, pages 849–856, 2002.
- [NM14] Kamal Nasrollahi and Thomas B Moeslund. Super-resolution: a comprehensive survey. *Machine vision and applications*, 25(6):1423–1468, 2014.
- [OP16] Lukasz P Olech and Mariusz Paradowski. Hierarchical gaussian mixture model with objects attached to terminal and non-terminal dendrogram nodes. In *Proceedings of the 9th International Conference on Computer Recognition Systems CORES 2015*, pages 191–201. Springer, 2016.
- [Pau07] Debashis Paul. Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica*, 17:1617–1642, 2007.
- [PB17] Jeffrey Pennington and Yasaman Bahri. Geometry of neural network loss surfaces via random matrix theory. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 2798–2806. JMLR. org, 2017.
- [PDK13] Dimitris Papailiopoulos, Alexandros Dimakis, and Stavros Korokythakis. Sparse pca through low-rank approximations. In *International Conference on Machine Learning*, pages 747–755, 2013.
- [PGM⁺19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.
- [PL20] Roman Pogodin and Peter E Latham. Kernelized information bottleneck leads to biologically plausible 3-factor hebbian learning in deep networks. *arXiv preprint arXiv:2006.07123*, 2020.
- [PRU⁺18] Kristina Preuer, Philipp Renz, Thomas Unterthiner, Sepp Hochreiter, and Günter Klambauer. Fréchet chemblnet distance: A metric for generative models for molecules. *arXiv preprint arXiv:1803.09518*, 2018.

- [PW17] Jeffrey Pennington and Pratik Worah. Nonlinear random matrix theory for deep learning. In *Advances in Neural Information Processing Systems*, pages 2637–2646, 2017.
- [RHW⁺88] David E Rumelhart, Geoffrey E Hinton, Ronald J Williams, et al. Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1, 1988.
- [RLNH17] Kevin Roth, Aurelien Lucchi, Sebastian Nowozin, and Thomas Hofmann. Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems 30*, pages 2018–2028. 2017.
- [RMC15] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [Sat17] Mahadev Satyanarayanan. The emergence of edge computing. *Computer*, 50(1):30–39, 2017.
- [SB95] Jack W Silverstein and ZD Bai. On the empirical distribution of eigenvalues of a class of large dimensional random matrices. *Journal of Multivariate analysis*, 54(2):175–192, 1995.
- [SB16] Suraj Srinivas and R Venkatesh Babu. Generalized dropout. *arXiv preprint arXiv:1611.06791*, 2016.
- [SC95] Jack W Silverstein and Sang-Il Choi. Analysis of the limiting spectral distribution of large dimensional random matrices. *Journal of Multivariate Analysis*, 54(2):295–309, 1995.
- [SCH⁺16] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1874–1883, 2016.
- [SCT⁺16] Casper Kaae Sønderby, Jose Caballero, Lucas Theis, Wenzhe Shi, and Ferenc Huszár. Amortised map inference for image super-resolution. *arXiv preprint arXiv:1610.04490*, 2016.
- [SEBT20] Mohamed El Amine Seddik, Hassan Essafi, Abdallah Benzine, and Mohamed Tamaazousti. Lightweight neural networks from pca lda based distilled dense neural networks. In *International Conference on Image Processing*, 2020.
- [SH08] Haipeng Shen and Jianhua Z Huang. Sparse principal component analysis via regularized low rank matrix approximation. *Journal of multivariate analysis*, 99(6):1015–1034, 2008.
- [SHK⁺14] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.

- [SIVA17] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [SLB15] Samuel Schulter, Christian Leistner, and Horst Bischof. Fast and accurate image upscaling with super-resolution forests. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3791–3799, 2015.
- [SLCT20] Mohamed El Amine Seddik, Cosme Louart, Romain Couillet, and Mohamed Tamaazousti. The unexpected deterministic and universal behavior of large softmax classifiers. 2020.
- [SLTC20] Mohamed El Amine Seddik, Cosme Louart, Mohamed Tamaazousti, and Romain Couillet. Random matrix theory proves that deep learning representations of gan-data behave as gaussian mixtures. In *International Conference on Machine Learning*, 2020.
- [SMG13] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [SMH11] Ilya Sutskever, James Martens, and Geoffrey E Hinton. Generating text with recurrent neural networks. In *ICML*, 2011.
- [SPP15] Jordi Salvador and Eduardo Perez-Pellitero. Naive bayes super-resolution forest. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 325–333, 2015.
- [ŚST⁺18] Marek Śmieja, Lukasz Struski, Jacek Tabor, Bartosz Zieliński, and Przemysław Spurek. Processing of missing data by neural networks. In *Advances in Neural Information Processing Systems*, pages 2719–2729, 2018.
- [STC19a] Mohamed El Amine Seddik, Mohamed Tamaazousti, and Romain Couillet. Kernel random matrices of large concentrated data: The example of gan-generated images. In *International Conference on Acoustics, Speech, and Signal Processing*, 2019.
- [STC19b] Mohamed El Amine Seddik, Mohamed Tamaazousti, and Romain Couillet. A kernel random matrix-based approach for sparse PCA. In *International Conference on Learning Representations*, 2019.
- [STC19c] Mohamed El Amine Seddik, Mohamed Tamaazousti, and Romain Couillet. Pourquoi les matrices aléatoires expliquent l’apprentissage ? un argument d’universalité offert par les gans. In *Colloque francophone de traitement du signal et des images*, 2019.
- [STL19] Mohamed El Amine Seddik, Mohamed Tamaazousti, and John Lin. Generative collaborative networks for single image super-resolution. *arXiv:1902.10467*, 2019.
- [SV99] Johan AK Suykens and Joos Vandewalle. Least squares support vector machine classifiers. *Neural processing letters*, 9(3):293–300, 1999.

- [SZ14] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [SZS⁺13] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [SZT17] Ravid Shwartz-Ziv and Naftali Tishby. Opening the black box of deep neural networks via information. *arXiv preprint arXiv:1703.00810*, 2017.
- [TAKC18] Hafis Tiomoko Ali, Abla Kammoun, and Romain Couillet. Random matrix asymptotics of inner product spectral clustering. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018.
- [Tao12] Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.
- [TDSVG14] Radu Timofte, Vincent De Smet, and Luc Van Gool. A+: Adjusted anchored neighborhood regression for fast super-resolution. In *Asian Conference on Computer Vision*, pages 111–126. Springer, 2014.
- [TDVG13] Radu Timofte, Vincent De, and Luc Van Gool. Anchored neighborhood regression for fast example-based super-resolution. In *Computer Vision (ICCV)*, pages 1920–1927. IEEE, 2013.
- [TGIH17] Alaa Tharwat, Tarek Gaber, Abdelhameed Ibrahim, and Aboul Ella Hassanien. Linear discriminant analysis: A detailed tutorial. *AI Commun.*, 30:169–190, 2017.
- [TLBH17a] Youssef Tamaazousti, Hervé Le Borgne, and Céline Hudelot. Mucale-net: Multi categorical-level networks to generate more discriminating features. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6711–6720, 2017.
- [TLBH⁺17b] Youssef Tamaazousti, Hervé Le Borgne, Céline Hudelot, Mohamed El Amine Seddik, and Mohamed Tamaazousti. Learning more universal representations for transfer-learning. *arXiv:1712.09708*, 2017.
- [TLBH⁺19] Youssef Tamaazousti, Hervé Le Borgne, Céline Hudelot, Mohamed El Amine Seddik, and Mohamed Tamaazousti. Learning more universal representations for transfer-learning. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [VLBB08] Ulrike Von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering. *The Annals of Statistics*, pages 555–586, 2008.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [VTWA18] Sahar Voghoei, Navid Hashemi Tonekaboni, Jason G Wallace, and Hamid R Arabnia. Deep learning at the edge. In *2018 International Conference on*

- Computational Science and Computational Intelligence (CSCI)*, pages 895–901. IEEE, 2018.
- [WBSS04] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [WEG87] Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1-3):37–52, 1987.
- [WGR⁺09] John Wright, Arvind Ganesh, Shankar Rao, Yigang Peng, and Yi Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. In *Advances in neural information processing systems*, pages 2080–2088, 2009.
- [WLY⁺15] Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang. Deep networks for image super-resolution with sparse prior. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 370–378, 2015.
- [WSB03] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *Signals, Systems and Computers, 2004. Conference Record of the Thirty-Seventh Asilomar Conference on*, volume 2, pages 1398–1402. Ieee, 2003.
- [WWL13] Stefan Wager, Sida Wang, and Percy S Liang. Dropout training as adaptive regularization. In *Advances in neural information processing systems*, pages 351–359, 2013.
- [XRV17] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [YC08] Kai Yu and Wei Chu. Gaussian process models for link analysis and transfer learning. In *Advances in Neural Information Processing Systems*, pages 1657–1664, 2008.
- [YCC98] LeCun Yann, Cortes Corinna, and J Christopher. The mnist database of handwritten digits. URL <http://yhann. lecun. com/exdb/mnist>, 1998.
- [YCL⁺16] Raymond Yeh, Chen Chen, Teck Yian Lim, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with perceptual and contextual losses. *arXiv preprint arXiv:1607.07539*, 2016.
- [YD16] Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356, 2016.
- [YJBK17] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4133–4141, 2017.
- [YKYR18] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9291–9301, 2018.

- [YLK⁺19] Joonyoung Yi, Juhyuk Lee, Kwang Joon Kim, Sung Ju Hwang, and Eunho Yang. Why not to use zero imputation? correcting sparsity bias in training neural networks. In *International Conference on Learning Representations*, 2019.
- [YMY14] Chih-Yuan Yang, Chao Ma, and Ming-Hsuan Yang. Single-image super-resolution: A benchmark. In *European Conference on Computer Vision*, pages 372–386. Springer, 2014.
- [YP16] Xin Yu and Fatih Porikli. Ultra-resolving face images by discriminative generative networks. In *European Conference on Computer Vision*, pages 318–333. Springer, 2016.
- [YW19] Yaqiong Yao and HaiYing Wang. Optimal subsampling for softmax regression. *Statistical Papers*, 60(2):235–249, 2019.
- [YYDN07] Qingxiong Yang, Ruigang Yang, James Davis, and David Nistér. Spatial-depth super resolution for range images. In *Computer Vision and Pattern Recognition, CVPR.*, pages 1–8. IEEE, 2007.
- [YZ13] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research*, 14(Apr):899–925, 2013.
- [ZEP10] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *International conference on curves and surfaces*, pages 711–730. Springer, 2010.
- [ZHT06] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265–286, 2006.
- [ZIE⁺18] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *arXiv preprint*, 2018.
- [ZS07] Ron Zass and Amnon Shashua. Nonnegative sparse pca. In *Advances in Neural Information Processing Systems*, pages 1561–1568, 2007.
- [ZY12] Wilman WW Zou and Pong C Yuen. Very low resolution face recognition problem. *IEEE Transactions on Image Processing*, 21(1):327–340, 2012.