



Contributions to robust estimation: minimax optimality vs. computational tractability

Amir-Hossein Bateni

► To cite this version:

Amir-Hossein Bateni. Contributions to robust estimation: minimax optimality vs. computational tractability. Statistics [math.ST]. Institut Polytechnique de Paris, 2022. English. NNT: 2022IP-PAG004 . tel-03828199

HAL Id: tel-03828199

<https://theses.hal.science/tel-03828199>

Submitted on 25 Oct 2022

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.

Contributions to robust estimation: minimax optimality vs. computational tractability

Thèse de doctorat de l'Institut Polytechnique de Paris
préparée à l'École nationale de la statistique et de l'administration
économique

École doctorale n°574 École doctorale de mathématiques
Hadamard (EDMH)
Spécialité de doctorat: Mathématiques appliquées

Thèse présentée et soutenue à Palaiseau, le 6 juillet 2022, par

AMIR-HOSSEIN BATENI

Composition du Jury :

Alexandre Tsybakov CREST, ENSAE, IP Paris	Président
Pierre Alquier RIKEN AIP Tokyo	Rapporteur
Mahdi Soltanolkotabi University of Southern California	Rapporteur
Guillaume Lécué CREST, ENSAE, IP Paris	Examineur
Nicolas Verzelen INRAE, Université de Montpellier	Examineur
Arnak Dalalyan CREST, ENSAE, IP Paris	Directeur de thèse

Remerciements

Mes premières remerciements vont certainement à Arnak. Je ne saurais comment te remercier pour tout ce que tu as fait pour moi. Tu m'as fait confiance alors que j'avais le parcours d'un outlier dans le monde statistique. Tu as supporté mes bêtises et mes défauts, tu t'es engagé sérieusement pour que j'apprenne et que j'avance, et en même temps tu m'as laissé explorer mes propres idées en croyant en moi. Je te remercie vivement pour ta patience, ton engagement, ton esprit ouvert, et ta bienveillance. Tu m'as montré un éminent exemple de savoir-être et savoir-vivre.

Je tiens à exprimer ma reconnaissance aux membres du jury, dont la présence est une véritable honneur pour moi. Thank you Mahdi for taking your time to review thoughtfully this manuscript. Merci Pierre pour ton rapport détaillé. Tu a été un généreux et chaleureux collègue pour moi. Merci Sacha de m'avoir accueilli au sein du pôle statistique au CREST et d'avoir accepté d'être le président de mon jury. Tes travaux ont eu une influence majeure sur ma recherche. Merci Nicolas d'avoir eu la gentillesse d'accepter de faire partie de mon jury. Et merci Guillaume pour ton soutien constant et pour nos discussions que j'ai toujours appréciées.

Je tiens aussi à exprimer mes remerciements à mes collègues chercheurs à l'ENSAE. Merci Victor de tes efforts pour le bon déroulement de mon doctorat. Je regrette de ne pas avoir pu continué le projet intéressant qu'on avait ensemble. Merci Katia de tes grands soutiens pendant la bonne année que j'ai passée à Dauphine. Merci à Cristina, Nicolas, Mathieu, Vianney, Anna, et Jaouad. J'ai beaucoup appris de vous humainement et professionnellement lors de nos différents échanges. Un immense merci à Arshak. Tu es vraiment un bon ami et en plus un bon collaborateur. Je suis content d'avoir travaillé avec toi et j'espère qu'on pourra encore travailler ensemble.

Je remercie également les doctorants et les anciens doctorants de l'ENSAE et alentours. Merci Arya pour toutes nos conversations interminables. Tu es une source inépuisable de motivation, d'ambition, et d'inspiration. Notre amitié est une longue histoire et il n'y a pas de place ici pour tout évoquer. Merci Avetik ! C'était un sacré bonheur de t'avoir à côté de moi pendant plusieurs années. Tu te disposes d'une forte joie de vivre que tu transmets aussi aux autres. Merci Étienne, Clara, et Hugo de l'ambiance vivante que vous avez créée dans le bureau. J'ai passé de très bons moments avec vous et surtout vous m'avez appris plein d'expressions françaises. Un grand merci à Jules, Christophe et Lucie pour leur im-

portant travail de l'assistance d'enseignement. Merci aux différentes générations, à Lionel, Simo, Gautier, Alexis, Mamadou, Jérémy, Yannick, Evgenii, Léna, Geoffrey, Badr, Boris, Nicolas, Gabriel, Solenne, François-Pierre, Meyer, Julien, Suzanne, Dang, Flore, Nayel, Younes, Sasila, Yannis, Alexandre, Corentin, Maria, Théo, et tous ceux que j'ai éventuellement oubliés. Et pareillement merci aux post-doctorants, à Yann, Thomas, Martin, et Jiaying.

Je remercie grandement mes chers amis. Merci Ali et Sahar de votre gentillesse et générosité ! Vous êtes vraiment magnifiques ! Merci MohammadSaeid et Ebrahim de votre soutien et présence ! Sans vous, ces années auraient vraiment été difficiles.

Last but not least, I want to thank my family. Indeed, I can never thank you enough my lovely parents as you have made so many sacrifices for me. I also thank warmly my dear sisters, my brothers-in-law, and my sweet nephews. Your presence and your support have always been important for me.

1	Introduction	3
1.1	Problem definition	3
1.2	Robust approaches for one dimension	4
1.2.1	Filtering	4
1.2.2	Mean as center of symmetry	6
1.2.3	Minimum distance estimator	8
1.2.4	Lower bound	10
1.3	Robust approaches for high dimension	11
1.3.1	Filtering (multidimensional case)	12
1.3.2	Mean as center of symmetry (multidimensional case)	14
1.3.3	Minimum distance estimator (multidimensional case)	16
1.3.4	Leveraging covariance matrix	17
1.4	Prior work	19
1.5	Contributions	21
1.5.1	Contamination models and robust estimation on the probability simplex	22
1.5.2	Robust Estimation of Gaussian Mean	28
2	Contamination models and robust estimation on the probability simplex	37
2.1	Introduction	38
2.2	Various models of contamination	40
2.2.1	Huber's contamination	40
2.2.2	Huber's deterministic contamination	41
2.2.3	Oblivious contamination	42
2.2.4	Parameter contamination	42
2.2.5	Adversarial contamination	42
2.2.6	Minimax risk "in expectation" versus "in deviation"	43
2.3	Prior work	44
2.4	Minimax rates on the "sparse" simplex and confidence regions	45
2.4.1	Upper bounds: worst-case risk of the sample mean	45
2.4.2	Lower bounds on the minimax risk	46
2.4.3	Confidence regions	48
2.5	Instance based bounds	49
2.6	Illustration on a numerical example	50
2.7	Summary and conclusion	51
3	Robust Estimation of Gaussian Mean	53
3.1	Introduction	54
3.2	Adversarially corrupted sub-Gaussian model and spectral dimension reduction	57
3.2.1	Choice of the dimension reduction regime	59
3.2.2	Choice of the threshold	60

3.3	Assessing the error of the SDR estimator	60
3.4	The case of unknown covariance matrix	63
3.5	Numerical experiments	64
3.5.1	Implementation details	64
3.5.2	Experimental setup	65
3.5.3	Statistical accuracy	65
3.5.4	Computational efficiency	66
3.5.5	Breakdown point	66
3.6	Summary, related work and conclusion	68
4	Discussion	70
5	Introduction en français	72
5.1	Définition du problème	72
5.2	Variables discrètes	73
5.3	Variables gaussiennes	74
5.3.1	Problème en dimension 1	75
5.3.2	Problème en grande dimension	77
5.4	Organisation du manuscrit	81
A	Proofs for Chapter 2	82
A.1	Proofs of propositions	82
A.2	Minimax upper bounds over the sparse simplex	84
A.3	Minimax lower bounds over the sparse simplex	86
A.4	Proofs of bounds with high probability	88
A.5	Proofs of instance based bounds	90
B	Proofs for Chapter 3	92
B.1	Proof of Theorem 11	92
B.1.1	Bounding the projected error of the average of filtered observations	92
B.1.2	Bounding the error of the geometric median of projected observations	94
B.1.3	Bounding the number of filtered out observations	95
B.1.4	Estimating the mean from a low-dimensional adversarial projection	97
B.1.5	Bounding stochastic errors	98
B.1.6	Putting all the pieces together	102
B.2	Proof of Theorem 13	106
B.3	Extension to Sub-Gaussian distributions	111
B.3.1	Proof of Theorem 12	113

Chapter 1

Introduction

1.1	Problem definition	3
1.2	Robust approaches for one dimension	4
1.2.1	Filtering	4
1.2.2	Mean as center of symmetry	6
1.2.3	Minimum distance estimator	8
1.2.4	Lower bound	10
1.3	Robust approaches for high dimension	11
1.3.1	Filtering (multidimensional case)	12
1.3.2	Mean as center of symmetry (multidimensional case)	14
1.3.3	Minimum distance estimator (multidimensional case)	16
1.3.4	Leveraging covariance matrix	17
1.4	Prior work	19
1.5	Contributions	21
1.5.1	Contamination models and robust estimation on the probability simplex	22
1.5.2	Robust Estimation of Gaussian Mean	28

1.1 Problem definition

In statistics and learning theory, it is common to assume that samples are independently and identically distributed according to a reference probability distribution. A more realistic approach could be to relax this assumption by allowing a fraction of samples to not necessarily follow the reference distribution. These disobeying samples, called *outliers*, may drastically skew the classical estimators.

Consider the following setting. Let $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ be n i.i.d. samples drawn from $\mathcal{N}(\boldsymbol{\mu}, \sigma^2)$. Given $\varepsilon \in (0, 1/2)$, we observe ε -contaminated Gaussian random variables $\mathbf{X}_1, \dots, \mathbf{X}_n$ which

means that we have $X_i = Y_i$ for $i \in \mathcal{I}$ while X_i takes arbitrary values when $i \in \mathcal{O}$, where $|\mathcal{I}| = n(1 - \varepsilon)$ and $|\mathcal{O}| = n\varepsilon$. Let \mathcal{S} be $\mathcal{I} \cup \mathcal{O}$ which is indeed the set $\{1, \dots, n\}$. Our goal is to estimate μ and we measure the estimation error by absolute value. When there is no outlier ($\mathcal{O} = \emptyset$), the maximum likelihood estimator (MLE) which is the sample mean, estimates μ with optimal non-asymptotic error rate, while the presence of a single outlier with an extreme large value (far enough from μ) can cause an arbitrary large deviation of MLE from μ . So, we see that MLE is not robust to outliers.

In this work, we aim to construct estimators robust to outliers, and we focus on the fundamental task of mean estimation, the task to which many problems in statistics reduce. We are interested in the non-asymptotic behavior of estimators.

1.2 Robust approaches for one dimension

Let us start by exploring different ideas for estimating robustly μ in the setting of ε -corrupted Gaussian samples.

1.2.1 Filtering

As we have seen the problem with MLE is its high sensitivity to extreme far outliers from μ . To get around this problem, we can simply remove the largest and smallest points and then apply MLE. Indeed, we sort the observations $X_{(1)} \leq \dots \leq X_{(n)}$ and return $\hat{\mu}_F := \sum_{i \in F} X_i / |F|$ where $F = \{i \in \mathcal{I} \cup \mathcal{O} \mid X_{(2n\varepsilon)} < X_i < X_{(n-2n\varepsilon)}\}$ ¹. After filtering even if there are still outliers among the samples, we are sure that they are at most as far as $Y_{(n\varepsilon)}$ and $Y_{(n-n\varepsilon)}$ (where the initial samples are sorted: $Y_{(1)} \leq \dots \leq Y_{(n)}$). Hence, this estimator, called *Trimmed mean*, is not so much influenced by extreme outliers. See Figure 1.1.

We try to analyze informally the estimation error of $\hat{\mu}_F$ without entering into the details. Let ξ_i denote $Y_i - \mu$ for every $i \in \mathcal{S}$. By the triangle inequality, we have

$$\begin{aligned} |\hat{\mu}_F - \mu| &\leq \frac{1}{|F|} \left| \sum_{i \in \mathcal{I} \cap F} (X_i - \mu) \right| + \frac{1}{|F|} \left| \sum_{i \in \mathcal{O} \cap F} (X_i - \mu) \right| \\ &\leq \frac{1}{|F|} \left| \sum_{i \in \mathcal{I} \cap F} \xi_i \right| + \frac{|\mathcal{O} \cap F|}{|F|} \max(|\xi_{(n\varepsilon)}|, |\xi_{(n-n\varepsilon)}|). \end{aligned} \quad (1.1)$$

We aim to upper bound the two terms in (1.1).

Remark 1. While the random variables $(Y_i)_{i \in \mathcal{S}}$ are independent, the random variables $(Y_i \mathbb{1}(i \in \mathcal{I}))_{i \in \mathcal{S}}$ are generally not independent. Indeed, we do not have any particular assumption on outliers and this implies that worst scenarios are possible. For instance, we may suppose that there is an adversary who observes the initial samples $(Y_i)_{i \in \mathcal{S}}$, and knowing our estimator, replaces εn of them by εn arbitrary values in a malicious way. At the end, we receive

¹Without loss of generality we suppose that $n\varepsilon$ is an integer.

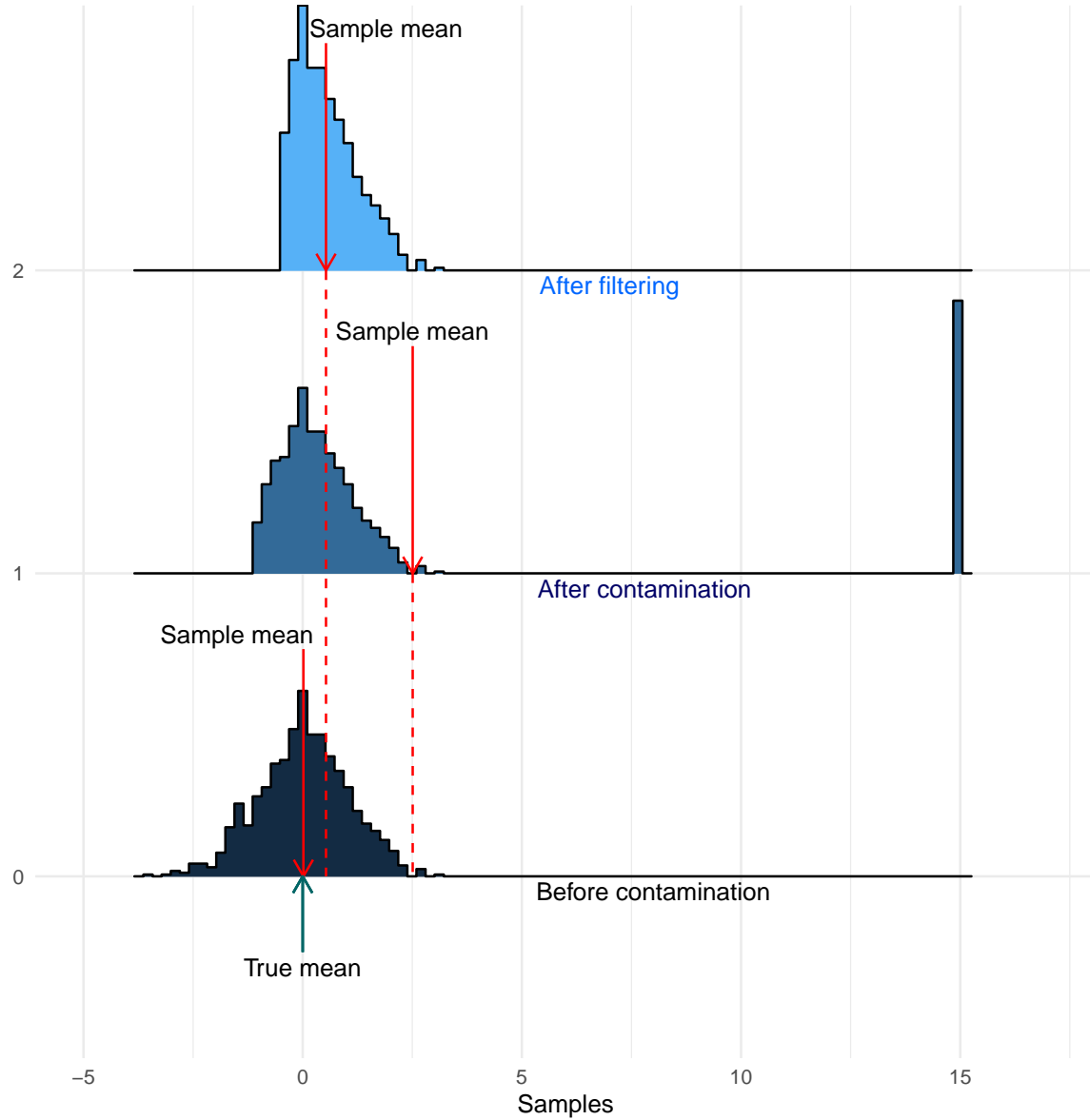


Figure 1.1: At the first floor the histogram of the initial samples $(Y_i)_{i \in \mathcal{S}}$ is depicted. At the second floor, we have the histogram of the contaminated samples $(X_i)_{i \in \mathcal{S}}$ where the smallest $n\varepsilon$ samples among $(Y_i)_{i \in \mathcal{S}}$ are replaced by $n\varepsilon$ samples equal to 15. At the third floor, we see the histogram of samples after the filtering procedure where the $2n\varepsilon$ greatest and $2n\varepsilon$ smallest samples among $(X_i)_{i \in \mathcal{S}}$ are removed. The sample means for each set are marked by red lines. The effect of the filtering procedure on the performance of the sample mean is apparent.

only X_1, \dots, X_n as described at the beginning of the chapter. This is called the adversarial contamination model. In this case, the nature of the subset \mathcal{I} might depend on many factors included all the initial samples $(Y_i)_{i \in S}$.

To bypass this problem of lack of independence in the first term of (1.1), we can state by the triangle inequality that

$$\begin{aligned} \left| \frac{1}{|F|} \sum_{i \in \mathcal{I} \cap F} \xi_i \right| &\leq \left| \frac{1}{|F|} \sum_{i \in S} \xi_i \right| + \left| \frac{1}{|F|} \sum_{i \in S \setminus (\mathcal{I} \cap F)} \xi_i \right| \\ &\leq \left| \frac{1}{|F|} \sum_{i=1}^n \xi_i \right| + \max_{|J| \leq 5n\varepsilon} \left| \frac{1}{|F|} \sum_{j \in J} \xi_j \right|, \end{aligned} \quad (1.2)$$

where in the second inequality we used $n - 5n\varepsilon \leq |\mathcal{I} \cap F|$. With high probability, the first term in (1.2) is of order $(1 - 4\varepsilon)^{-1}(\sigma/\sqrt{n})$ (using the Gaussian tail bounds and the fact that $|F| = n - 4n\varepsilon$), and the second term is bounded by a term of order $(1 - 4\varepsilon)^{-1}\sigma\varepsilon\sqrt{\log(e/\varepsilon)}$ (using the Gaussian properties and the union bound). For more details on these calculations, we refer the reader to the proof of Lemma 6.

To upper bound the second term in (1.1), we can use Hoeffding's inequality to prove that $\xi_{(n\varepsilon)}$ and $\xi_{(n-n\varepsilon)}$ respectively are concentrated around $\sigma\Phi^{-1}(\varepsilon)$ and $\sigma\Phi^{-1}(1 - \varepsilon)$ where Φ is the cumulative distribution function of $\mathcal{N}(0, 1)$. The absolute value of these two last terms is bounded by $\sigma\sqrt{\log(1/\varepsilon)}$ since we deal with Gaussian variables. This entails that the second term in (1.1) is bounded by a term of order $(1 - 4\varepsilon)^{-1}\sigma\varepsilon\sqrt{\log(1/\varepsilon)}$ with high probability. Therefore, we informally proved that $\hat{\mu}_F$ estimates μ with an error rate of order at most

$$\frac{(\sigma/\sqrt{n}) + \sigma\varepsilon\sqrt{\log(1/\varepsilon)}}{1 - 4\varepsilon}$$

with high probability.

The error rate of the trimmed mean is composed of the optimal error rate σ/\sqrt{n} in the case of non contamination (which is achieved by MLE), and $\sigma\varepsilon\sqrt{\log(1/\varepsilon)}$ which shows the impact of outliers in our estimation in worst case. If there is no outlier ($\varepsilon = 0$), the trimmed mean does not remove any point and returns the sample mean. Nevertheless, this estimator is tolerant to ε not larger than $1/4$.

By removing other quantities of points than $4n\varepsilon$ in the filtering step, one can change the performance of the trimmed mean. A special case is the sample median where $\lceil n/2 \rceil - 1$ points are filtered out from each side. We analyze the sample median as an estimator of the mean in the next section.

1.2.2 Mean as center of symmetry

The Gaussian distribution is a symmetric distribution, and the mean is not only the barycentre of the distribution, but also the center of symmetry. The idea is to estimate the mean, as a

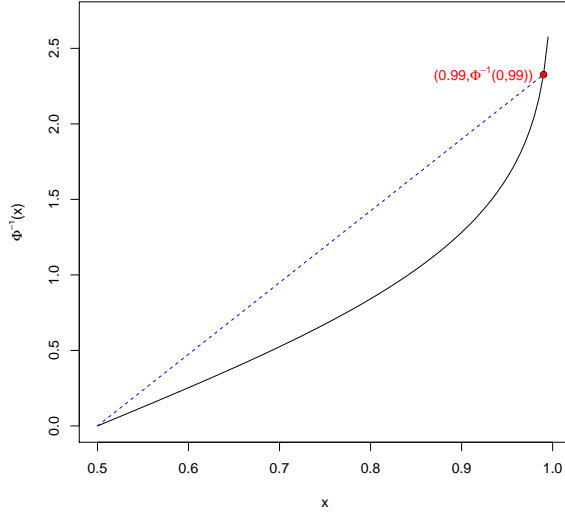


Figure 1.2: The behavior of the probit function: for $x \in (0, 0.49)$, $\Phi^{-1}(0.5 + x) < 5x$.

center of symmetry. There are various notions of symmetry. The one that we consider here is the symmetry induced by mass where there are same masses in the both sides of the center. More precisely, in this symmetry, the center is the median. Thus, we may estimate the center, *i.e.* the mean, by the sample median which is in fact the center of symmetry for the empirical mass.

To analyze the sample median, we exceptionally provide more details since it makes our presentation simpler. The sample median of the contaminated data², $X_{(\lceil n/2 \rceil)}$, is located between $Y_{(\lceil n/2 \rceil - n\varepsilon)}$ and $Y_{(\lceil n/2 \rceil + n\varepsilon)}$, and by Hoeffding's inequality we can prove that $|Y_{(\lceil n/2 \rceil - n\varepsilon)} - \mu|$ and $|Y_{(\lceil n/2 \rceil + n\varepsilon)} - \mu|$ are bounded by

$$\sigma \Phi^{-1}\left(\frac{1}{2} + \varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}}\right)$$

with probability at least $1 - \delta$ for any $\delta \in (0, 1)$ satisfying $\frac{1}{2} + \varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}} < 1$. Observing the behavior of the function Φ^{-1} (called *probit*) around 1/2 in Figure 1.2, we note that if $\varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}} \leq 0.49$, we have

$$\Phi^{-1}\left(\frac{1}{2} + \varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}}\right) \leq 5\left(\varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}}\right).$$

We can conclude that if $\varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}} \leq 0.49$, the inequality $|X_{(n/2)} - \mu| \leq 5\sigma\left(\varepsilon + \sqrt{\frac{\log(2/\delta)}{2n}}\right)$ holds with probability at least $1 - \delta$.

Remark 2. Let fix $\delta \in (0, 1)$. We have seen that for $\varepsilon < \frac{1}{2} - \sqrt{\frac{\log(2/\delta)}{2n}}$ the estimation error of

²For sake of simplicity, here we suppose that n is odd.

the sample median is bounded with probability at least $1 - \delta$. What if $\varepsilon \geq \frac{1}{2} - \sqrt{\frac{\log(2/\delta)}{2n}}$? For $\varepsilon \in \left[\frac{1}{2} - \sqrt{\frac{\log(2/\delta)}{2n}}, \frac{1}{2} - \frac{1}{n}\right]$, we can claim that the sample median $X_{(\lceil n/2 \rceil)}$ is located between $Y_{(1)}$ and $Y_{(n)}$ which means that $|X_{(\lceil n/2 \rceil)} - \mu|$ is bounded by a term of order $\sigma \sqrt{\log(n)}$ with high probability. For $\varepsilon > \frac{1}{2} - \frac{1}{n}$, the sample median could be an outlier (since the majority of data is outlier) and therefore, the estimation error is not bounded in worst case. So, we can state that for fixed n the breakdown point of the sample median is $\frac{1}{2} - \frac{1}{n}$. The breakdown point is the maximum proportion of outliers for which the estimation error is bounded in worst case. For instance, the breakdown point of the sample mean is $1/n$ and that of the trimmed mean is $\frac{1}{4} - \frac{1}{n}$. Given our non-asymptotic study framework, we are interested here in the notion of finite sample breakdown point. However, in Section 1.5.2 and Chapter 3, we will consider the asymptotic notion of the breakdown point.

As we emphasized, the sample median is characterized by the empirical mass, and thereby its robustness could be interpreted by the robustness of the empirical mass since the empirical mass is not sensitive to extreme far outliers (it can be changed at most by 2ε when data become contaminated). Another robust method could be to leverage the empirical mass to estimate the probability density function of the reference distribution. We will expose this method in the coming section.

1.2.3 Minimum distance estimator

We know the explicit expression of the probability density function (p.d.f.) for our reference distribution in terms of the mean parameter. Another approach for our problem might be to determine a value as mean in a way that the corresponding p.d.f. fits best the observations. To do so, we can select a set of possible candidates for the mean, and choose the one with the best fitting of associated p.d.f. to data.

First, we introduce, as it is presented in (Devroye and Lugosi, 2000), the *Scheffé's* estimate which chooses the best between two p.d.f. f_1 and f_2 . Let f be a p.d.f. and $\hat{\nu}$ be a probability distribution. The Scheffé's estimate \hat{f} is defined as

$$\hat{f} := \begin{cases} f_1 & \text{if } \left| \int_A f_1 - \hat{\nu}(A) \right| < \left| \int_A f_2 - \hat{\nu}(A) \right|, \\ f_2 & \text{otherwise,} \end{cases}$$

where $A = \{x : f_1(x) > f_2(x)\}$. (Devroye and Lugosi, 2000, Theorem 6.1) states that

$$\int |f - \hat{f}| \leq 3 \min \left(\int |f_1 - f|, \int |f_2 - f| \right) + 4 \max_{A \in \mathcal{A}} \left| \int_A f - \hat{\nu}(A) \right|, \quad (1.3)$$

where $\mathcal{A} = \{\{f_1 > f_2\}, \{f_2 > f_1\}\}$. Now, assume that f_μ is the p.d.f. of the reference distribution $\mathcal{N}(\mu, \sigma)$ and ν is the empirical measure of the non-contaminated random variables $(Y_i)_{i \in S}$, i.e., for every set A , $\nu(A) = \sum_{i=1}^n \mathbb{1}(Y_i \in A)/n$. However, for choosing the best p.d.f.

we have only access to the empirical measure $\hat{\nu}$ of the contaminated data which associates to every set A the mass $\hat{\nu}(A) = \sum_{i=1}^n \mathbb{1}(\mathbf{X}_i \in A)/n$. Using (1.3) and the triangle inequality, we obtain

$$\int |f_{\mu} - \hat{f}| \leq 3 \min \left(\int |f_1 - f_{\mu}|, \int |f_2 - f_{\mu}| \right) + 4 \max_{A \in \mathcal{A}} \left| \int_A f_{\mu} - \nu(A) \right| + 4 \max_{A \in \mathcal{A}} |\nu(A) - \hat{\nu}(A)|.$$

The term $\max_{A \in \mathcal{A}} \left| \int_A f_{\mu} - \nu(A) \right|$ is the difference between the empirical and theoretical measures. It is bounded by a term of order $\sqrt{1/n}$ with high probability via Hoeffding's inequality and the union bound. The term $\max_{A \in \mathcal{A}} |\nu(A) - \hat{\nu}(A)|$ is at most two times the maximum mass that the adversary is allowed to transport in order to contaminate data, namely 2ε . If f_1 and f_2 respectively are the p.d.f. of $\mathcal{N}(\mu_1, \sigma^2)$ and $\mathcal{N}(\mu_2, \sigma^2)$, for $i \in \{1, 2\}$ we have $\int |f_i - f_{\mu}| = 2\text{TV}(\mathcal{N}(\mu_i, \sigma^2), \mathcal{N}(\mu, \sigma^2))$ (Scheffé's theorem) where TV denotes the total variation distance. Furthermore, by Pinsker's inequality we know that $\text{TV}(\mathcal{N}(\mu_i, \sigma^2), \mathcal{N}(\mu, \sigma^2)) \leq |\mu_i - \mu|/(2\sigma)$. So, finally we deduce that there exists $C > 0$ such that the inequality

$$\int |f_{\mu} - \hat{f}| \leq \frac{3}{\sigma} \min(|\mu_1 - \mu|, |\mu_2 - \mu|) + C \frac{1}{\sqrt{n}} + 8\varepsilon,$$

is satisfied with high probability. As μ is located between $\mathbf{X}_{((\varepsilon+0.05)n)}$ and $\mathbf{X}_{((1-\varepsilon-0.05)n)}$ with high probability, given $\gamma > 0$ we can construct a covering set $M_{\gamma} = \{\mu_1, \dots, \mu_{N_{\gamma}}\}$ over $[\mathbf{X}_{((\varepsilon+0.05)n)}, \mathbf{X}_{((1-\varepsilon-0.05)n)}]$ where for all m in $[\mathbf{X}_{((\varepsilon+0.05)n)}, \mathbf{X}_{((1-\varepsilon-0.05)n)}]$ there exists $i \in \{1, \dots, N_{\gamma}\}$ such that $|\mu_i - m| \leq \gamma$. Now, we consider a finite family of p.d.f. candidates $(f_i)_{i \in \{1, \dots, N_{\gamma}\}}$ where f_i is the p.d.f. of $\mathcal{N}(\mu_i, \sigma^2)$ and extend the Scheffé's estimate in order to choose the best p.d.f. among $(f_i)_{i \in \{1, \dots, N_{\gamma}\}}$ by running a tournament between these candidates. We declare $f_{\hat{k}}$ as the winner of the tournament with

$$\hat{k} = \arg \min_{k \in \{1, \dots, N_{\gamma}\}} \max_{A \in \mathcal{A}_{\gamma}} \left| \int_A f_k - \hat{\nu}(A) \right|, \quad (1.4)$$

where $\mathcal{A}_{\gamma} = \{\{x : f_i(x) > f_j(x)\} : i, j \in \{1, \dots, N_{\gamma}\}\}$. When there are multiple minima for (1.4), \hat{k} is defined to be the smallest index. Using (Devroye and Lugosi, 2000, Theorem 6.3), Hoeffding's inequality and the union bound, we can show that there exists $C > 0$ such that the inequality

$$\int |f_{\mu} - f_{\hat{k}}| \leq \frac{3}{\sigma} \gamma + C \sqrt{\frac{\log(N_{\gamma}^2)}{n}} + 8\varepsilon,$$

holds with high probability. We set $\gamma = \frac{\sigma}{\sqrt{n}}$, and thereby N_{γ}^2 is of order $4n$ as the length of $[\mathbf{X}_{((\varepsilon+0.05)n)}, \mathbf{X}_{((1-\varepsilon-0.05)n)}]$ is less than 4σ with high probability. This yields the existence of $C > 0$ such that

$$\int |f_{\mu} - f_{\hat{k}}| \leq C \sqrt{\frac{\log(n)}{n}} + 8\varepsilon,$$

is true with high probability. Now, we argue that $\mu_{\hat{k}}$ is a robust estimator of μ . Assume without loss of generality $\mu < \mu_{\hat{k}}$. Then we have

$$\begin{aligned}
\int |f_{\mu} - f_{\hat{k}}| &= 2 \int_{-\infty}^{(\mu+\mu_{\hat{k}})/2} (f_{\mu} - f_{\hat{k}}) = \frac{2}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{(\mu+\mu_{\hat{k}})/2} e^{-(x-\mu)^2/(2\sigma^2)} dx \\
&\quad - \frac{2}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^{(\mu+\mu_{\hat{k}})/2} e^{-(x-\mu_{\hat{k}})^2/(2\sigma^2)} dx \\
&= 2\Phi\left(\frac{\mu_{\hat{k}} - \mu}{2\sigma}\right) - 2\Phi\left(\frac{\mu - \mu_{\hat{k}}}{2\sigma}\right) \\
&= 4\Phi\left(\frac{\mu_{\hat{k}} - \mu}{2\sigma}\right) - 2. \tag{1.5}
\end{aligned}$$

This leads to

$$|\mu - \mu_{\hat{k}}| \leq 2\sigma\Phi^{-1}\left(\frac{C}{4}\sqrt{\frac{\log(n)}{n}} + 2\varepsilon + \frac{1}{2}\right),$$

with high probability if $\frac{C}{4}\sqrt{\frac{\log(n)}{n}} + 2\varepsilon < \frac{1}{2}$. Again, given the behavior of the probit function around 1/2 (cf. Figure 1.2), we deduce that

$$|\mu - \mu_{\hat{k}}| \leq 10\sigma\left(\frac{C}{4}\sqrt{\frac{\log(n)}{n}} + 2\varepsilon\right),$$

is satisfied with high probability if $\frac{C}{4}\sqrt{\frac{\log(n)}{n}} + 2\varepsilon \leq 0.49$. This estimator $\mu_{\hat{k}}$, called *skeleton estimate*, has a breakdown point of at least $\frac{1}{4} - \frac{C}{8}\sqrt{\frac{\log(n)}{n}}$.

While the skeleton estimate appears to be based on the estimation of p.d.f., we should stress that its key tool is the empirical mass, and this provides robustness for this estimator.

1.2.4 Lower bound

We have obtained so far three different error rates for estimating robustly the Gaussian mean in one dimension, namely

- with the trimmed mean: $\frac{\sigma}{\sqrt{n}} + \sigma\varepsilon\log(1/\varepsilon)$,
- with the sample median: $\frac{\sigma}{\sqrt{n}} + \sigma\varepsilon$,
- and with the skeleton estimate: $\sigma\sqrt{\frac{\log(n)}{n}} + \sigma\varepsilon$.

The best rate is achieved by the sample median. At this stage, the question is whether it is possible to do better than the sample median. Here, we attempt to show a lower bound for *Huber's contamination model*, a contamination model weaker than the adversarial model. We assume that the observations $\mathbf{X}_1, \dots, \mathbf{X}_n$ are drawn from the mixture distribution defined by

$$(1 - \varepsilon)\mathcal{N}(\mu, \sigma^2) + \varepsilon Q,$$

where Q is an unknown distribution representing the corruption in our problem. In Chapter 2 (Proposition 3), we show that Huber's model is nearly contained in the adversarial model, and this implies that our lower bound for Huber's model holds for the adversarial model, too. This lower bound is originally proposed in (Chen et al., 2018).

Consider two distributions $\mathcal{N}(\mu_1, \sigma^2)$ and $\mathcal{N}(\mu_2, \sigma^2)$ with their respective p.d.f. f_1 and f_2 . We propose two distributions Q_1 and Q_2 defined respectively by the p.d.f. $\frac{1-\varepsilon}{\varepsilon}(f_1 - f_2)\mathbb{1}(f_1 \geq f_2)$ and $\frac{1-\varepsilon}{\varepsilon}(f_2 - f_1)\mathbb{1}(f_2 \geq f_1)$. It is now easy to verify that $(1 - \varepsilon)\mathcal{N}(\mu_1, \sigma^2) + \varepsilon Q_1$ and $(1 - \varepsilon)\mathcal{N}(\mu_2, \sigma^2) + \varepsilon Q_2$ are the same distribution while we have $\text{TV}(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) = \frac{\varepsilon}{1-\varepsilon}$. That is to say that the two different distributions $\mathcal{N}(\mu_1, \sigma^2)$ and $\mathcal{N}(\mu_2, \sigma^2)$ might be contaminated in a way that they are not distinguishable one from another. In this situation, μ_1 and μ_2 have the same chance to be the mean of the reference distribution and our best estimation of the mean would be the middle point $(\mu_2 + \mu_1)/2$. Thus, our estimation error is at least $|\mu_2 - \mu_1|/2$. Let δ denote $\mu_2 - \mu_1$, if $\mu_1 < \mu_2$ by similar calculations as in (1.5) we have

$$\text{TV}(\mathcal{N}(\mu_1, \sigma^2), \mathcal{N}(\mu_2, \sigma^2)) = \int_{-\infty}^{\mu_1 + \delta/2} (f_1 - f_2) = 2\Phi\left(\frac{\delta}{2\sigma}\right) - 1,$$

and this entails that $|\mu_1 - \mu_2| = 2\sigma\Phi^{-1}\left(\frac{1}{2} + \frac{1}{2}\frac{\varepsilon}{1-\varepsilon}\right)$. Adding to this, the lower bound for estimating the Gaussian mean with non-contaminated data, namely σ/\sqrt{n} , our error rate under contamination cannot be better than

$$\sigma\Phi^{-1}\left(\frac{1}{2} + \frac{\varepsilon}{2(1-\varepsilon)}\right) + \frac{\sigma}{\sqrt{n}}.$$

This lower bound matches our upper bound for the sample median $\sigma\Phi^{-1}\left(\frac{1}{2} + \varepsilon + \frac{1}{\sqrt{n}}\right)$ and we conclude that the sample median is a minimax estimator of the mean when data are corrupted. Note that when ε tends to 1/2, the minimax risk and the risk of the sample median explode to infinity.

Remark 3. *The three studied estimators in this section are computable in polynomial time in terms of sample size. The trimmed mean requires the knowledge of ε , the skeleton estimate requires the knowledge of σ whereas the sample median does not require any of them. Moreover, the sample median is tolerant to larger values for ε . Consequently, in addition to its optimal error rate, the sample median seems to be the most interesting estimator for our problem.*

1.3 Robust approaches for high dimension

Now, we assume that data are in \mathbb{R}^p and the initial random variables $(Y_i)_{i \in \{1, \dots, n\}}$ follow the multivariate Gaussian distribution $\mathcal{N}(\mu, \sigma^2 \mathbf{I}_p)$. As in the previous setting, we observe $\mathbf{X}_1, \dots, \mathbf{X}_n$, a ε -contaminated version of data under the adversarial model. We quantify the estimation error by the Euclidean distance $\|\cdot\|_2$.

A naive solution for the multidimensional setting would be to estimate each coordinate of μ separately via a robust method for one dimension. By the union bound, we may establish that applying for example the sample median to each coordinate gives an estimator $\hat{\mu}_{\text{CM}}$ satisfying

$$\|\mu - \hat{\mu}_{\text{CM}}\|_2 \leq 5\sigma \left(\sqrt{\frac{p \log(2p/\delta)}{2n}} + \varepsilon \sqrt{p} \right),$$

with probability at least $1 - \delta$ if $\varepsilon + \sqrt{\log(2/\delta)/(2n)} \leq 0.49$. This estimator is known as the *componentwise median*. The problem with such estimators is that the term $\varepsilon \sqrt{p}$ in the error could become problematic when p is large. [Chen et al. \(2018\)](#) prove that under Huber's contamination model, the error rate of the componentwise median cannot be better than $\sigma(\sqrt{p/n} + \varepsilon \sqrt{p})$ when the contamination distribution Q is defined as a Dirac on the point $\mu + \sigma(1, \dots, 1)^\top$. The same contamination distribution can be deployed to show that the other robust methods for one dimension will not have a better result in high dimension. Now, the question is whether we can have a better dependency on p .

Remark 4. *One might improve slightly the dependency on p for the componentwise median by applying it not in the canonical basis but in the basis formed by the eigenvectors of the sample covariance matrix. Against this method, the adversary may choose outliers $Z_1, \dots, Z_{n\varepsilon}$ where for Z_i the i first coordinates are σ and the rest is zero. In such situation, one can show that the error rate of the componentwise median is at least $\sigma(\sqrt{p/n} + \varepsilon \sqrt{\min(n\varepsilon, p)})$.*

1.3.1 Filtering (multidimensional case)

In the multidimensional setting, there is no notion of order. Therefore, we cannot apply the filtering method proposed previously for one dimension. However, we know that multivariate Gaussian samples are concentrated on a sphere of radius $\sigma \sqrt{p}$ with the mean as center (note the contrast with univariate Gaussian samples which are concentrated close to the mean). More precisely, there exists $c > 0$ such that for all $t > 0$

$$\mathbf{P}(|\|\xi_i\|_2 - \sigma \sqrt{p}| > t) \leq 2 \exp(-ct^2/\sigma^2), \quad (1.6)$$

(see e.g., [\(Vershynin, 2018, Theorem 3.1.1\)](#)). See Figure 1.3.

We can now design a filter based on the property described by (1.6). In fact, this property states that given $\tau \in (0, 1)$, at least $1 - \tau$ fraction of $(\xi_i)_{i \in \{1, \dots, n\}}$ satisfies with high probability

$$\|\xi_i\|_2 \leq C\sigma(\sqrt{p} + \sqrt{\log(1/\tau)}),$$

where C is a positive constant. For a formal statement of this property, see Lemma 3. On this event, given ε -contaminated data, we can claim that at least $1 - \tau - \varepsilon$ fraction of $(X_i)_{i \in \{1, \dots, n\}}$

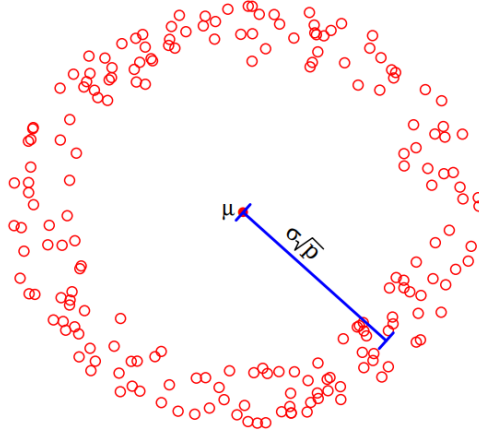


Figure 1.3: A schematic representation of the Gaussian point cloud in high dimension. Points are concentrated around a sphere with center μ and radius $\sigma\sqrt{p}$.

satisfies

$$\|X_i - \mu\|_2 \leq C\sigma(\sqrt{p} + \sqrt{\log(1/\tau)}). \quad (1.7)$$

For each point X_i , we consider the set $A_i := \{X_j : \|X_j - X_i\|_2 \leq 2C\sigma(\sqrt{p} + \sqrt{\log(1/\tau)})\}$. If the cardinality of this set is more than $n(\tau + \varepsilon)$, there is a point X_j in A_i satisfying (1.7). By the triangle inequality, this implies that X_i has a distance at most $3C\sigma(\sqrt{p} + \sqrt{\log(1/\tau)})$ from μ . Hence, all the points exhibited by the set $F := \{i : |A_i| \geq n(\varepsilon + \tau)\}$ are in a distance of order $\sigma(\sqrt{p} + \sqrt{\log(1/\tau)})$ from μ . Note that F contains all X_i satisfying (1.7) with $i \in \mathcal{I}$ if $n(1 - \tau - \varepsilon) > n(\tau + \varepsilon)$, i.e., if $\tau + \varepsilon < 1/2$. Otherwise, F could be of cardinality zero in worst case. So, F is of cardinality at least $n(1 - \tau - \varepsilon)$ if $\tau + \varepsilon < 1/2$. This procedure of filtering is called *Naive pruning* (Diakonikolas et al., 2016a).

The empirical mean $\hat{\mu}_F$ of the points in F satisfies

$$\begin{aligned} \|\hat{\mu}_F - \mu\|_2 &\leq \frac{1}{|F|} \left\| \sum_{i \in \mathcal{I} \cap F} (X_i - \mu) \right\|_2 + \frac{1}{|F|} \left\| \sum_{i \in \mathcal{O} \cap F} (X_i - \mu) \right\|_2 \\ &\leq \frac{1}{|F|} \left\| \sum_{i \in \mathcal{I} \cap F} \xi_i \right\|_2 + \frac{|\mathcal{O} \cap F|}{|F|} \max_{i \in \mathcal{O} \cap F} \|X_i - \mu\|_2. \end{aligned} \quad (1.8)$$

The first term of (1.8) is bounded by a term of order $(1 - \tau - \varepsilon)^{-1}(\sigma\sqrt{p/n} + \sigma(\varepsilon + \tau)\sqrt{\log(1/(\varepsilon + \tau))})$ via the same argument as in one dimension (cf. Lemma 6), and the second term is of order at most $(1 - \tau - \varepsilon)^{-1}\varepsilon\sigma(\sqrt{p} + \sqrt{\log(1/\tau)})$. Therefore, choosing $\tau = \varepsilon$, we conclude that if $\varepsilon < 1/4$ the error rate of the sample mean $\hat{\mu}_F$ after naive pruning is

$$\frac{\sigma}{1 - 2\varepsilon}(\sqrt{p/n} + \varepsilon\sqrt{\log(1/\varepsilon)} + \varepsilon\sqrt{p}).$$

At the end, \sqrt{p} unfortunately is still present behind ε in our estimation error.

1.3.2 Mean as center of symmetry (multidimensional case)

We may study three notions of symmetry for multivariate Gaussian distribution (borrowed from [Zuo and Serfling \(2000\)](#)). Each one proposes a different estimator for the center, a role played by the Gaussian mean under the three symmetries.

Central symmetry

A random vector Y is *centrally symmetric* if $Y - \theta$ and $\theta - Y$ have the same distribution where θ is the center of symmetry. This implies that $\theta = \mathbf{E}[Y]$ and suggests to estimate θ by the sample mean, however, the sample mean is not robust as we have already outlined.

Angular symmetry

A random vector Y is *angularly symmetric* about θ if $(Y - \theta)/\|Y - \theta\|_2$ is centrally symmetric. In our case, to estimate μ , the center of angular symmetry, we may compute a point $\hat{\mu}_{\text{GM}}$ at which the function $\phi : x \mapsto \sum_{i=1}^n (X_i - x)/\|X_i - x\|_2$ becomes null vector. It turns out that ϕ is the gradient of the function $x \mapsto \sum_{i=1}^n \|X_i - x\|_2$. The last function is convex and its minimum is attained somewhere in the convex hull of X_1, \dots, X_n . Consequently ϕ is equal to null vector at this point. That is to say

$$\hat{\mu}_{\text{GM}} = \arg \min_{x \in \mathbb{R}^p} \sum_{i=1}^n \|X_i - x\|_2, \quad (1.9)$$

and this formulation of $\hat{\mu}_{\text{GM}}$ is known as the *geometric median* of the samples. There are various optimization methods for approximating the minimum in (1.9). In particular, via these methods we may reach a point $\hat{\mu}'_{\text{GM}}$ satisfying

$$\frac{1}{n} \|\phi(\hat{\mu}'_{\text{GM}})\|_2 \leq \frac{1}{\sqrt{n}}. \quad (1.10)$$

This condition certifies the existence of a positive constant C such that with high probability we have

$$\|\mu - \hat{\mu}'_{\text{GM}}\|_2 \leq C(\sigma\sqrt{p/n} + \sigma\varepsilon\sqrt{p}). \quad (1.11)$$

Let us show this informally. Indeed, the random vector $(Y_i - \hat{\mu}'_{\text{GM}})/\|Y_i - \hat{\mu}'_{\text{GM}}\|_2$ enjoys the sub-Gaussian convergence rate, and since $\text{Tr}\left(\text{Var}\left[\frac{Y_i - \hat{\mu}'_{\text{GM}}}{\|Y_i - \hat{\mu}'_{\text{GM}}\|_2}\right]\right)$ is less than one, there exists $C' > 0$ such that

$$\left\| \mathbf{E} \left[\frac{Y_i - \hat{\mu}'_{\text{GM}}}{\|Y_i - \hat{\mu}'_{\text{GM}}\|_2} \right] \right\|_2 \leq \frac{1}{n} \|\phi(\hat{\mu}'_{\text{GM}})\|_2 + C' \frac{1}{\sqrt{n}} + 2\varepsilon, \quad (1.12)$$

is valid with high probability. For large values of p the random variable $\|Y_i - \hat{\mu}'_{\text{GM}}\|_2$ concentrates around its mean. This leads to

$$\begin{aligned} \left\| \mathbf{E} \left[\frac{Y_i - \hat{\mu}'_{\text{GM}}}{\|Y_i - \hat{\mu}'_{\text{GM}}\|_2} \right] \right\|_2 &\approx \left\| \mathbf{E} \left[\frac{Y_i - \hat{\mu}'_{\text{GM}}}{\mathbf{E}\|Y_i - \hat{\mu}'_{\text{GM}}\|_2} \right] \right\|_2 \\ &= \frac{\|\mu - \hat{\mu}'_{\text{GM}}\|_2}{\mathbf{E}\|Y_i - \hat{\mu}'_{\text{GM}}\|_2} \\ &\geq \frac{\|\mu - \hat{\mu}'_{\text{GM}}\|_2}{\|\mu - \hat{\mu}'_{\text{GM}}\|_2 + \sqrt{\mathbf{E}\|Y_i - \mu\|_2^2}} = \frac{\|\mu - \hat{\mu}'_{\text{GM}}\|_2}{\|\mu - \hat{\mu}'_{\text{GM}}\|_2 + \sigma\sqrt{p}}, \end{aligned} \quad (1.13)$$

where we used the triangle and Jensen's inequalities. Combining (1.10), (1.12) and (1.13), one obtains (1.11) with high probability if for some positive C'' we have $\frac{C''}{\sqrt{n}} + 2\varepsilon < 1$. Therefore, we obtain an estimator $\hat{\mu}'_{\text{GM}}$ of μ with the convergence rate $\sigma\sqrt{p/n} + \sigma\varepsilon\sqrt{p}$ in high dimension.

This analysis requires that $\varepsilon < \frac{1}{2} - \frac{C''}{2\sqrt{n}}$ which mean that the breakdown point of $\hat{\mu}'_{\text{GM}}$ is at least $\frac{1}{2} - \frac{C''}{2\sqrt{n}}$. [Lopuhaa and Rousseeuw \(1991\)](#) prove that the geometric median has the benefit of the maximum breakdown point value, namely $\frac{1}{2} - \frac{1}{n}$.

On an related note, under a more general setting, ([Lai et al., 2016](#), Proposition 2.1) construct an oblivious contamination (a contamination model weaker than the adversarial one, cf. [2.2.3](#)) for which the error rate of the geometric median cannot be better than $\sigma\varepsilon\sqrt{p}$.

Halfspace symmetry

A random vector Y is *halfspace symmetric* about θ if $\mathbf{P}(Y \in H) \geq 1/2$ for every closed halfspace H containing θ . In other words, θ is the center of the symmetry induced by mass in every direction. Indeed, this generalizes the notion of the median for higher dimensions. This means that for all unitary vector v , $v^\top \theta$ is the median of the random variable $v^\top Y$. To translate this notion of center for our sample distribution, we consider the projection of the samples in all the directions, and we look for the point minimizing the difference of empirical masses in both sides of the point for every direction. So, we may define the center as the point minimizing the maximum value of this difference over all the directions. This gives a point $\hat{\mu}_{\text{TM}}$ belonging to the set

$$\arg \min_{x \in \mathbb{R}^p} \max_{v \in \mathbb{S}^{p-1}} \left| \sum_{i=1}^n \mathbb{1}(v^\top X_i \leq v^\top x) - \sum_{i=1}^n \mathbb{1}(v^\top X_i > v^\top x) \right|,$$

where $\left| \sum_{i=1}^n \mathbb{1}(v^\top X_i \leq v^\top x) - \sum_{i=1}^n \mathbb{1}(v^\top X_i > v^\top x) \right|$ expresses the difference of empirical masses at the two sides of $v^\top x$ when the samples are projected on direction v . It can easily be shown that

$$\hat{\mu}_{\text{TM}} \in \arg \max_{x \in \mathbb{R}^p} \min_{v \in \mathbb{R}^{p-1}} \sum_{i=1}^n \mathbb{1}(v^\top X_i \leq v^\top x).$$

The last set is a convex polytope, and thus we can define $\hat{\mu}_{\text{TM}}$ in a unique way as the barycentre of this polytope. The point $\hat{\mu}_{\text{TM}}$ is known as *Tukey's median* ((Tukey, 1975)) of the samples, and represents a generalization of the sample median for higher dimensions. Given a point x , the value

$$\min_{v \in \mathbb{R}^{p-1}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(v^\top \mathbf{X}_i \leq v^\top x)$$

is called the *Tukey's halfspace depth* of x with respect to the samples, and so, Tukey's median is determined as the barycentre of the points of maximum depth with respect to the samples.

Chen et al. (2018) prove that under some conditions the convergence rate of Tukey's median is of order $\sigma \sqrt{p/n} + \sigma \varepsilon$. More precisely, under Huber's contamination, they establish that for some positive constants C_1 and C_2 , given $\delta \in (0, 1/2)$, $\hat{\mu}_{\text{TM}}$ satisfies

$$\|\mu - \hat{\mu}_{\text{TM}}\|_2 \leq \sigma \Phi^{-1} \left(\frac{1}{2} + \frac{\varepsilon}{1 - \varepsilon} + C_1 \sqrt{\frac{p}{n}} + C_2 \sqrt{\frac{\log(1/\delta)}{n}} \right),$$

with probability at least $1 - 2\delta$. The above formula is meaningful only if

$$\frac{1}{2} + \frac{\varepsilon}{1 - \varepsilon} + C_1 \sqrt{\frac{p}{n}} + C_2 \sqrt{\frac{\log(1/\delta)}{n}} < 1.$$

This implies that $\varepsilon < \frac{1}{3} - C \sqrt{\frac{p}{n}}$ for some positive constant C , and therefore the breakdown point of Tukey's median is at least $\frac{1}{3} - C \sqrt{\frac{p}{n}}$. Furthermore, results from (Liu et al., 2017) show that in general position the breakdown point of the Tukey's median is at most $\frac{1}{3} - \frac{2}{9} \frac{p-3}{n}$.

Remark 5. We finally succeeded in avoiding the factor \sqrt{p} behind ε in the error rate. However, computation of Tukey's median demands a number of operations of order n^p . Indeed, computation of the depth of any point in the space with respect to the samples costs already an exponential in p number of operations.

Remark 6. Central symmetry implies angular symmetry and angular symmetry implies halfspace symmetry. The foregoing material suggests that broader is the adopted notion of symmetry, more robust is the method estimating the mean as center of symmetry, and harder is the computation of the corresponding estimator.

Remark 7. The symmetry based estimators have two benefits: first, they do not require the knowledge of the covariance matrix and the contamination rate; second, they are translation, uniform scaling and orthogonal transformation equivariant.

1.3.3 Minimum distance estimator (multidimensional case)

The skeleton estimate can be built for higher dimensions in the same way as for one dimension. We need just to specify the construction of the candidates' set. Given a naive estimator

of μ , for instance the estimator $\hat{\mu}_F$ proposed by the filtering method, we know that μ belongs with high probability to the ball of center $\hat{\mu}_F$ and radius $C\sqrt{p}$ where C is a positive constant. A covering set $M_\gamma = \{\mu_1, \dots, \mu_{N_\gamma}\}$ can be constructed on this ball with $N_\gamma \leq (2C\sqrt{p}/\gamma)^p$. We consider the family $(f_i)_{i \in \{1, \dots, N_\gamma\}}$ where f_i is the p.d.f. of $\mathcal{N}(\mu_i, \sigma \mathbf{I}_p)$. Let $\hat{\nu}$ be the empirical measure defined by the samples X_1, \dots, X_n . We define

$$\hat{k} = \arg \min_{k \in \{1, \dots, N_\gamma\}} \max_{A \in \mathcal{A}_\gamma} \left| \int_A f_k - \hat{\nu}(A) \right|,$$

where $\mathcal{A}_\gamma = \{\{x : f_i(x) > f_j(x)\} : i, j \in \{1, \dots, N_\gamma\}\}$ with $\gamma = \sigma\sqrt{p/n}$. Using arguments similar to those of the one dimensional setting, the skeleton estimate $\mu_{\hat{k}}$ satisfies with high probability

$$\|\mu - \mu_{\hat{k}}\|_2 \leq 10\sigma \left(C\sqrt{\frac{p \log(n)}{n}} + 2\varepsilon \right),$$

where C is a positive constant and $C\sqrt{p \log(n)/n} + 2\varepsilon \leq 0.49$.

Remark 8. The set \mathcal{A}_γ is constructed in time exponential in p . For this reason, the skeleton estimate as Tukey's median is not computationally tractable. In fact, both estimators are determined by a minimax optimization problem and this suggests that exponential complexity is indispensable to achieve a error rate of order $\sigma\sqrt{p/n} + \sigma\varepsilon$ (without the factor \sqrt{p} multiplying ε). This raises the trade off challenge between the computational efficiency and statistical accuracy.

1.3.4 Leveraging covariance matrix

Applying the same technique introduced in Section 1.2.4, it is easily established that

$$\sigma\Phi^{-1}\left(\frac{1}{2} + \frac{\varepsilon}{2(1-\varepsilon)}\right) + \frac{\sigma\sqrt{p}}{\sqrt{n}}$$

is a lower bound for estimating the mean under Huber's contamination model when the non-contaminated samples follow $\mathcal{N}(\mu, \sigma \mathbf{I}_p)$. This shows that Tukey's median is a minimax estimator of the multivariate Gaussian mean when the samples are adversarially corrupted. As we have outlined, there is a significant gap between the convergence rate of the computational tractable methods presented here and that of the optimal method, *i.e.*, Tukey's median, which is not tractable. In recent years, one of the main challenges in robust estimation has been to fill this gap and various methods have been invented. The common point of almost all of these methods is to use the knowledge of the covariance matrix of the reference Gaussian distribution.

We try to give the high level scheme of these methods. Consider a more general setting where the initial samples $(Y_i)_{i \in \{1, \dots, n\}}$ are drawn from $\mathcal{N}(\mu, \Sigma)$. The key idea is the following.

If we have (contaminated) samples with a sample covariance matrix close to Σ , then we can certify that the empirical mean of these samples is close to μ . These methods use some or all of the following ingredients:

1. Majority of the inliers (X_i with $i \in \mathcal{I}$) lie in a distance of order \sqrt{p} from μ . Via a filtering method, we can obtain a set containing only samples with such property (see e.g., Lemma 4).
2. Let $\bar{\mu}_S$ be the sample mean of the set $\{Y_i : i \in S\}$. For all subsets $S \subset \{1, \dots, n\}$ of size at least $(1 - \varepsilon)n$, $\|\bar{\mu}_S - \mu\|_2$ is bounded by a term of order $\sqrt{p/n} + \varepsilon\sqrt{\log(1/\varepsilon)}$ (see e.g., Lemma 6).
3. Let $\bar{\Sigma}_S$ be the sample covariance matrix of the set $\{Y_i : i \in S\}$. For all subsets $S \subset \{1, \dots, n\}$ of size at least $(1 - \varepsilon)n$, $\|\bar{\Sigma}_S - \Sigma\|_{\text{op}}$ is bounded by a term of order $\max(\sqrt{p/n}, p/n) + \varepsilon \log(1/\varepsilon)$ where $\|\cdot\|_{\text{op}}$ denotes the operator norm for matrices (see e.g., Lemma 7).
4. Let $\hat{\mu}_S$ and $\hat{\Sigma}_S$ be the sample mean and sample covariance matrix of the set $\{X_i : i \in S\}$. We can show that $\|\mu - \hat{\mu}_S\|_2^2$ is controlled by a term depending on $\varepsilon\|\Sigma - \hat{\Sigma}_S\|_{\text{op}}$ (see e.g., Proposition 6).

Hence, now the goal would be to find a set S with a enough small $\|\Sigma - \hat{\Sigma}_S\|_{\text{op}}$. We can recapitulate various existing algorithms for realizing this goal in two categories:

- Filtering algorithms: they filter samples in order to control $\|\Sigma - \hat{\Sigma}_S\|_{\text{op}}$. For filtering they deploy spectral methods on the matrix $\Sigma - \hat{\Sigma}_S$ to detect in which directions there is more anomaly and then they remove the samples responsible for these anomalies. In this category, we can cite methods introduced in (Lai et al., 2016), (Diakonikolas et al., 2017), (Diakonikolas et al., 2018a) and our method presented in Chapter 3.
- Weighting algorithms: they attribute to each sample X_i a weight ω_i (such that $\sum_{i=1}^n \omega_i = 1$) and try to find weights in order to minimize $\|\Sigma - \hat{\Sigma}_\omega\|_{\text{op}}$ where $\hat{\Sigma}_\omega$ is the weighted sample covariance matrix. Finally, the weighted sample mean $\hat{\mu}_\omega$ determined by such weights appears to be a robust estimator of μ . For this problem, different optimization methods are proposed in (Diakonikolas et al., 2016a), (Cheng et al., 2019a), (Dalalyan and Minasyan, 2020), and (Cheng et al., 2020).

All of these algorithms succeed to achieve in polynomial time a error rate with dependency $\varepsilon\sqrt{\log(1/\varepsilon)}$ instead of $\varepsilon\sqrt{p}$. Diakonikolas et al. (2016b) provide a strong evidence that the dependency $\varepsilon\sqrt{\log(1/\varepsilon)}$ in error rate is necessary for any robust polynomial time algorithm under the adversarial model. Indeed, they demonstrate that any *Statistical Query* algorithm estimating μ under the adversarial model with dependency ε (instead of $\varepsilon\sqrt{\log(1/\varepsilon)}$) in error rate needs to have access to sample moments of degree at least $\log(1/\varepsilon)^{1/4}$, and this costs at least $p^{\log(1/\varepsilon)^{1/4}}$ operations, namely, an exponential number of operations in p . The *Statistical*

Query algorithms, first introduced in (Kearns, 1998), are a class of algorithms that are allowed to query expectations of functions of the reference distribution modulo an error rather than directly access the samples. This can model a wide range of algorithms in statistics and learning theory.

The four ingredients mentioned above can be generalized for samples drawn from other classes of distributions (such as the sub-Gaussian distributions or distributions with bounded moments), albeit with possibly different bounds corresponding to their respective concentration properties. For instance, Steinhardt et al. (2017) call a general form of the second ingredient *resilience*, Diakonikolas et al. (2020) call a general form of the second and third ingredients together *stability*, or the three first ingredients together with some other conditions are called *goodness* in (Diakonikolas et al., 2017). These generalizations make the algorithms with this approach extensible for robust mean estimation of a more general class of reference distributions.

1.4 Prior work

Robust estimation is an area of active research in Statistics since at least five decades (Donoho and Gasko, 1992; Donoho and Huber, 1983; Huber and Ronchetti, 2011; Huber, 1964; Maronna et al., 2006; Rousseeuw et al., 2011; Rousseeuw and Hubert, 1999; Tukey, 1975). Until very recently, theoretical guarantees were almost exclusively formulated in terms of the notions of breakdown point, sensitivity curve, influence function, etc. In 2015, a new line of research in this area started by the work of Chen et al. (2018) who considered the problem of estimating the mean and the covariance matrix of a Gaussian distribution in the high-dimensional setting under Huber’s contamination model by studying the rate of convergence of the minimax risk as a function of the sample size n , dimension p and the rate of contamination ε . The authors showed that the minimax optimal rate is attained by Tukey’s median which is not computationally tractable whereas classical tractable robust methods such as the componentwise median do not achieve this rate.

This phenomenon stimulated researches for computational tractable methods with optimal rate. First algorithms were proposed by computer science community. In this community, the error rates are expressed differently: it contains only the term depending on ε while the sample size necessary for obtaining this error rate is also considered. For example the error rates for the methods presented in the last section reduce to their term depending on ε , namely, $\varepsilon\sqrt{p}$, $\varepsilon\sqrt{\log(1/\varepsilon)}$ or ε and the other term which depends on p and n determines how many samples the methods need to achieve this error rate. Thus, there are two kinds of complexity in this approach: the computational complexity and the sample complexity.

For the Gaussian mean, Lai et al. (2016) proposed a recursive algorithm halving the dimension at each step based on a spectral analysis of the sample covariance matrix. Their algorithm achieves the error rate $\varepsilon\sqrt{\log(p)}$ under the oblivious contamination model (see the

definition in Section 2.2.3). Diakonikolas et al. (2016a) designed two algorithms with error rate $\varepsilon\sqrt{\log(1/\varepsilon)}$: one based on an iterative spectral filtering method and one based on a convex programming defined by the weighted samples. Cheng et al. (2019a) transformed this convex programming to a semi-definite programming (SDP) and proposed a faster iterative method with same error rate but with optimal sample complexity, *i.e.*, sample size of order d/ε^2 . Dalalyan and Minasyan (2020) introduce a different iterative method based on SDP by giving statistical guarantees for the error rate $\sqrt{\text{Tr}(\Sigma)/n} + \varepsilon\sqrt{\|\Sigma\|_{\text{op}}\log(1/\varepsilon)}$ in expectation and in probability. Finally, Cheng et al. (2021) prove that a suitable version of the gradient descent method on minimization of $\|\Sigma - \hat{\Sigma}_\omega\|$ converges to an estimator $\hat{\mu}_\omega$ of μ with same risk $\varepsilon\sqrt{\log(1/\varepsilon)}$. The only estimator giving a error rate of order ε in polynomial time (under oblivious contamination model) is proposed in (Diakonikolas et al., 2017) based on a spectral filtering method, however, it has a high sample complexity. Some other methods can be found in (Balakrishnan et al., 2017; Diakonikolas and Kane, 2019; Diakonikolas et al., 2020; Dong et al., 2019). All of these methods are extensible for sub-Gaussian distributions. In Chapter 3, we design a new estimator inspired by spectral filtering ideas in Lai et al. (2016) and Diakonikolas et al. (2017).

Robust mean estimation is studied also for the distributions with bounded moments. Many of the algorithms for Gaussian mean can be generalized for the distributions with bounded moments, see (Cheng et al., 2019a; Diakonikolas et al., 2017; Diakonikolas and Kane, 2019; Diakonikolas et al., 2020; Dong et al., 2019; Lai et al., 2016). For this problem, other algorithms with the idea of weighting the sample are proposed, see *e.g.*, (Depersin and Lecu , 2022; Liu et al., 2020a; Prasad et al., 2019; Steinhardt et al., 2017). In particular, Depersin and Lecu  (2022) give a computational tractable estimator under bounded second moment assumption with the optimal error rate $\sqrt{\text{Tr}(\Sigma)/n} + \sqrt{\|\Sigma\|_{\text{op}}\varepsilon}$, while the other mentioned methods are at best nearly optimal. All we have seen so far belong to the category of the approaches leveraging covariance matrix presented in Section 1.3.4. Minsker (2015) analyses the geometric median under bounded moments assumptions and proves that it has sub-optimal error rate. Lugosi and Mendelson (2021) extend the trimmed mean to high-dimensional setting and show that their extension enjoys optimal error rate, though, it cannot be computed in polynomial time. There are various other approaches, see *e.g.*, (Bahmani, 2021; Depersin and Lecu , 2021; Minsker, 2018b; Minsker and Ndaoud, 2021; Prasad et al., 2020).

A closely related problem is that of the mean estimation with sub-Gaussian rates for heavy tailed distributions. For one dimension, the so-called *median-of-means* estimator, introduced independently in (Nemirovsky and Yudin, 1983), (Jerrum et al., 1986) and (Alon et al., 1996), has the optimal performance. Other methods exist for one dimension, see (Catoni, 2012; Devroye et al., 2016; Minsker and Ndaoud, 2021) and for a more complete overview, we refer the reader to the survey (Lugosi and Mendelson, 2019a). For high-dimension, various estimators are proposed such as a median-of-means based estimator (Joly et al., 2017), geometric median (Hsu and Sabato, 2016; Minsker, 2015) or (Catoni and Giulini, 2017, 2018). The first estimator with optimal performance is introduced in (Lugosi and Mendelson, 2019b) which is

not computationally efficient. [Hopkins \(2020\)](#) formulates a SDP relaxation of the estimator of [Lugosi and Mendelson \(2019b\)](#) using *sum-of-squares* methods and designs a polynomial time algorithm with the optimal sub-Gaussian rates. [Cherapanamjeri et al. \(2019\)](#) combine this SDP formulation by a non-convex gradient descent procedure and improve the computational complexity of the last algorithm. [Depersin and Lecué \(2022\)](#) combine the idea of SDP and median-of-means and construct another tractable optimal estimator. [Lei et al. \(2019\)](#) propose a different optimal estimator not using SDP (SDP causes a high runtime). Furthermore, unified robust approaches against outliers and heavy tailed distributions are considered in ([Bahmani, 2021](#); [Depersin and Lecué, 2021](#); [Depersin and Lecué, 2022](#); [Minsker, 2015](#); [Minsker and Ndaoud, 2021](#); [Prasad et al., 2019](#)).

Concerning the robust mean estimation of discrete distributions, [Chen et al. \(2020\)](#); [Jain and Orlitsky \(2019\)](#); [Qiao and Valiant \(2018\)](#); [Steinhardt et al. \(2017\)](#) studied the case of group-contamination where the observations are batches of samples with ε -fraction of batches contaminated by an adversary. In Chapter 2, we study this problem under the normal adversarial model.

In addition, high-dimensional robust estimation is studied under a very different framework, Generative Adversarial Networks (GAN), and robust estimators based on GANs are constructed for various statistical tasks, see ([Gao et al., 2018, 2020](#); [Wang and Tan, 2021](#); [Zhu et al., 2022](#)). This line of work is close in spirit to the minimum distance estimator discussed in 1.3.3.

Beyond Huber and adversarial contamination models, other contamination models are considered in the literature such as parameter contamination models ([Bhatia et al., 2017](#); [Carpentier et al., 2018](#); [Collier and Dalalyan, 2019](#)) (cf. Section 2.2.4), oblivious contamination ([Diakonikolas et al., 2017](#); [Feng et al., 2014](#); [Lai et al., 2016](#)) (cf. Section 2.2.3), Wasserstein and total variation contaminations ([Diakonikolas et al., 2016b](#); [Zhu et al., 2019, 2020](#)).

Robust estimation is investigated for other statistical tasks as covariance matrix estimation (many paper studying mean estimation treat also covariance matrix estimation, see also ([Cheng et al., 2019b](#))), linear regression ([Alquier and Gerber, 2020](#); [Audibert and Catoni, 2011](#); [Bakshi and Prasad, 2021](#); [Cherapanamjeri et al., 2020](#); [Chinot, 2020](#); [Dalalyan and Thompson, 2019](#); [Depersin, 2020](#); [Diakonikolas et al., 2019b](#); [Gao, 2020](#); [Liu et al., 2018](#); [Pensia et al., 2020](#); [Sasai and Fujisawa, 2021](#)), risk minimization ([Chinot et al., 2018](#); [Chinot et al., 2020](#)), stochastic optimization ([Diakonikolas et al., 2019a](#)), classification ([Lecué et al., 2020](#)), density estimation ([Jain and Orlitsky, 2021](#); [Liu and Gao, 2019](#)), etc.

1.5 Contributions

We summarize our contributions in two sections. The first section is devoted to the problem of estimating the mean of a distribution supported by the k -dimensional probability simplex in the setting where an ε fraction of observations are subject to adversarial corruption. A

simple particular example is the problem of estimating the distribution of a discrete random variable. Assuming that the discrete variable takes k values, the unknown parameter θ is a k -dimensional vector belonging to the probability simplex. There we describe also the various models of contamination. Chapter 2 exposes a complete version of the contribution. This work is presented in (Bateni and Dalalyan, 2020).

The second section summarizes our contributions to the problem of robust estimation of the mean vector of a Gaussian distribution. We introduce an estimator based on spectral dimension reduction (SDR) and establish a finite sample upper bound on its error that is minimax-optimal up to a logarithmic factor. Furthermore, we prove that the (asymptotic) breakdown point of the SDR estimator is equal to $1/2$, the highest possible value of the breakdown point. In addition, the SDR estimator is equivariant by similarity transforms and has low computational complexity. This work is developed in Chapter 3 and is presented in (Bateni et al., 2022).

1.5.1 Contamination models and robust estimation on the probability simplex

Assume X_1, \dots, X_n are n independent and identically distributed random variables taking their values in the k -dimensional probability simplex $\Delta^{k-1} = \{v \in \mathbb{R}_+^k : v_1 + \dots + v_k = 1\}$. Our goal is to estimate the unknown vector $\theta = \mathbb{E}[X_i]$ in the case where the observations are contaminated by outliers. An important particular case is the estimation of the distribution of a discrete random variable X taking k distinct values. In this particular case, X_i 's take values in $\{e_1, \dots, e_k\}$, the set of the vectors of the canonical basis, which are also the extreme points of the simplex Δ^{k-1} .

In this introduction, to convey the main messages, we limit ourselves to the Huber contamination model, although our results apply to the more general adversarial contamination. Huber's contamination model assumes that there are two probability measures P, Q on Δ^{k-1} and a real $\varepsilon \in [0, 1/2)$ such that X_i is drawn from

$$P_i = (1 - \varepsilon)P + \varepsilon Q, \quad \forall i \in \{1, \dots, n\}.$$

This amounts to assuming that $(1 - \varepsilon)$ -fraction of observations, called inliers, are drawn from a reference distribution P , whereas ε -fraction of observations are outliers and are drawn from another distribution Q . In general, all the three parameters P, Q and ε are unknown. The parameter of interest is some functional (such as the mean, the standard deviation, etc.) of the reference distribution P , whereas Q and ε play the role of nuisance parameters.

When the unknown parameter lives on the probability simplex, there are many appealing ways of defining the risk. We focus on the following three metrics: total-variation, Hellinger

and \mathbb{L}^2 distances³

$$\begin{aligned} d_{\text{TV}}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) &:= 1/2 \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_1, \\ d_{\text{H}}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) &:= 1/\sqrt{2} \|\hat{\boldsymbol{\theta}}^{1/2} - \boldsymbol{\theta}^{1/2}\|_2, \\ d_{\mathbb{L}^2}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) &:= \|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2. \end{aligned}$$

The Hellinger distance above is well defined when the estimator $\hat{\boldsymbol{\theta}}$ is non-negative, which will be the case throughout this work. We will further assume that the dimension k may be large, but the vector $\boldsymbol{\theta}$ is s -sparse, for some $s \leq k$, i.e. $\#\{j : \theta_j \neq 0\} \leq s$. Our main interest is in constructing confidence regions and evaluating the minimax risk

$$\mathfrak{R}_{\square}(n, k, s, \varepsilon) := \inf_{\bar{\boldsymbol{\theta}}_n} \sup_{P, Q} \mathbf{E}[d_{\square}(\bar{\boldsymbol{\theta}}_n, \boldsymbol{\theta})],$$

where the *inf* is over all estimators $\bar{\boldsymbol{\theta}}_n$ built upon the observations $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} (1 - \varepsilon)\mathbf{P} + \varepsilon\mathbf{Q}$ and the *sup* is over all distributions P, Q on the probability simplex such that the mean $\boldsymbol{\theta}$ of P is s -sparse. The subscript \square of \mathfrak{R} above refers to the distance used in the risk, so that \square is TV, H, or \mathbb{L}^2 .

Various models of contamination

Different mathematical frameworks have been used in the literature to model the outliers. We present here five of them, from the most restrictive one to the most general, and describe their relationship. We present these frameworks in the general setting when the goal is to estimate the parameter $\boldsymbol{\theta}^*$ of a reference distribution $P_{\boldsymbol{\theta}^*}$ when ε proportion of the observations are outliers.

Huber's contamination The most popular framework for studying robust estimation methods is perhaps the one of Huber's contamination. In this framework, there is a distribution Q defined on the same space as the reference distribution $P_{\boldsymbol{\theta}^*}$ such that all the observations $\mathbf{X}_1, \dots, \mathbf{X}_n$ are independent and drawn from the mixture distribution $P_{\varepsilon, \boldsymbol{\theta}^*, Q} := (1 - \varepsilon)P_{\boldsymbol{\theta}^*} + \varepsilon Q$.

This corresponds to the following mechanism: one decides with probabilities $(1 - \varepsilon, \varepsilon)$ whether a given observation is an inlier or an outlier. If the decision is made in favor of being inlier, the observation is drawn from $P_{\boldsymbol{\theta}^*}$, otherwise it is drawn from Q . More formally, if we denote by \hat{O} the random set of outliers, then conditionally to $\hat{O} = O$,

$$\{\mathbf{X}_i : i \notin O\} \stackrel{\text{iid}}{\sim} P_{\boldsymbol{\theta}^*}, \{\mathbf{X}_i : i \in O\} \stackrel{\text{iid}}{\sim} Q, \{\mathbf{X}_i : i \in O\} \perp\!\!\!\perp \{\mathbf{X}_i : i \notin O\}, \quad (1.14)$$

for every $O \subset \{1, \dots, n\}$. Furthermore, for every subset O of the observations, we have

³We write $\|\mathbf{u}\|_q = (\sum_{j=1}^k |u_j|^q)^{1/q}$ and $\mathbf{u}^q = (u_1^q, \dots, u_k^q)$ for any $\mathbf{u} \in \mathbb{R}_+^k$ and $q > 0$.

$P(\widehat{O} = O) = (1-\varepsilon)^{n-|O|}\varepsilon^{|O|}$. We denote by⁴ $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ the set of joint probability distributions P_n of the random variables X_1, \dots, X_n satisfying the foregoing condition.

Huber's deterministic contamination The set of outliers as well as the number of outliers in Huber's model of contamination are random. This makes it difficult to compare this model to the others that will be described later in this section. To cope with this, we define here another model, termed Huber's deterministic contamination. As its name indicates, this new model has the advantage of containing a deterministic number of outliers, in the same time being equivalent to Huber's contamination in a sense that will be made precise below.

We say that the distribution P_n of X_1, \dots, X_n belongs to the Huber's deterministic contamination model denoted by $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \theta^*)$, if there are a set $O \subset \{1, \dots, n\}$ of cardinality at most $n\varepsilon$ and a distribution Q such that (1.14) is true. The apparent similarity of models $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ and $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \theta^*)$ can also be formalized mathematically in terms of the orders of magnitude of minimax risks. To ease notation, we let $R_d^\square(n, \varepsilon, \Theta, \widehat{\theta})$ to be the worst-case risk of an estimator $\widehat{\theta}$, where \square is either HC or HDC. More precisely, for $\mathcal{M}_n^\square(\varepsilon, \Theta) := \cup_{\theta \in \Theta} \mathcal{M}_n^\square(\varepsilon, \theta)$, we set⁵

$$R_d^\square(n, \varepsilon, \Theta, \widehat{\theta}) := \sup_{P_n \in \mathcal{M}_n^\square(\varepsilon, \Theta)} \mathbf{E}[d(\widehat{\theta}, \theta^*)].$$

This definition assumes that the parameter space Θ is endowed with a pseudo-metric $d : \Theta \times \Theta \rightarrow \mathbb{R}_+$. When $\Theta = \{\theta^*\}$ is a singleton, we write $R_{d,n}^\square(\varepsilon, \theta^*, \widehat{\theta})$ instead of $R_d^\square(n, \varepsilon, \{\theta^*\}, \widehat{\theta})$.

Proposition 1. *Let $\widehat{\theta}_n$ be an arbitrary estimator of θ^* . For any $\varepsilon \in (0, 1/2)$,*

$$\begin{aligned} R_d^{\text{HC}}(n, \varepsilon, \theta^*, \widehat{\theta}_n) &\leq R_{d,n}^{\text{HDC}}(2\varepsilon, \theta^*, \widehat{\theta}_n) + e^{-n\varepsilon/3} R_{d,n}^{\text{HDC}}(1, \theta^*, \widehat{\theta}_n), \\ \sup_{P_n \in \mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r) &\leq R_{d,n}^{\text{HDC}}(2\varepsilon, \theta^*, \widehat{\theta}_n) + r e^{-n\varepsilon/3}. \end{aligned}$$

Denote by \mathcal{D}_Θ the diameter of Θ , $\mathcal{D}_\Theta := \max_{\theta, \theta'} d(\theta, \theta')$. Proposition 1 implies that

$$\inf_{\widehat{\theta}_n} R_d^{\text{HC}}(n, \varepsilon, \Theta, \widehat{\theta}_n) \leq \inf_{\widehat{\theta}_n} R_{d,n}^{\text{HDC}}(n, 2\varepsilon, \Theta, \widehat{\theta}_n) + e^{-n\varepsilon/3} \mathcal{D}_\Theta.$$

When Θ is bounded, the last term is typically of smaller order than the minimax risk over $\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \Theta)$. Therefore, the minimax rate of estimation in Huber's model is not slower than the minimax rate of estimation in Huber's deterministic contamination model. This entails that a lower bound on the minimax risk established in HC-model furnishes a lower bound in HDC-model.

Oblivious contamination A third model of contamination that can be of interest is the oblivious contamination. In this model, it is assumed that the set O of cardinality o and the joint

⁴The superscript HC refers to the Huber's contamination

⁵The subscript d refers to the distance d used in the definition of the risk.

distribution Q_O of outliers are determined in advance, possibly based on the knowledge of the reference distribution P_{θ^*} . Then, the outliers $\{X_i : i \in O\}$ are drawn randomly from Q_O independently of the inliers $\{X_i : i \in O^c\}$. The set of all the joint distributions P_n of random variables X_1, \dots, X_n generated by such a mechanism will be denoted by $\mathcal{M}_n^{\text{OC}}(\varepsilon, \theta^*)$. The model of oblivious contamination is strictly more general than that of Huber's deterministic contamination, since it does not assume that the outliers are i.i.d. Therefore, the minimax risk over $\mathcal{M}_n^{\text{OC}}(\varepsilon, \Theta)$ is larger than the minimax risk over $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \Theta)$:

$$\inf_{\hat{\theta}_n} R_d^{\text{HDC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{OC}}(n, \varepsilon, \Theta, \hat{\theta}_n).$$

The last inequality holds true for any set Θ , any contamination level $\varepsilon \in (0, 1)$ and any sample size.

Parameter contamination In the three models considered above, the contamination acts on the observations. One can also consider the case where the parameters of the distributions of some observations are contaminated. More precisely, for some set $O \subset \{1, \dots, n\}$ selected in advance (but unobserved), the outliers $\{X_i : i \in O\}$ are independent and independent of the inliers $\{X_i : i \in O^c\}$. Furthermore, each outlier X_i is drawn from a distribution $Q_i = P_{\theta_i}$ belonging to the same family as the reference distribution, but corresponding to a contaminated parameter $\theta_i \neq \theta^*$. Thus, the joint distribution of the observations can be written as $(\bigotimes_{i \in O^c} P_{\theta^*}) \otimes (\bigotimes_{i \in O} P_{\theta_i})$. The set of all such distributions P_n will be denoted by $\mathcal{M}_n^{\text{PC}}(\varepsilon, \theta^*)$, where PC refers to "parameter contamination".

Adversarial contamination The last model of contamination we describe in this work, the adversarial contamination, is the most general one. It corresponds to the following two-stage data generation mechanism. In a first stage, iid random variables Y_1, \dots, Y_n are generated from a reference distribution P_{θ^*} . In a second stage, an adversary having access to Y_1, \dots, Y_n chooses a (random) set \hat{O} of (deterministic) cardinality s and arbitrarily modifies data points $\{Y_i : i \in \hat{O}\}$. The resulting sample, $\{X_i : i = 1, \dots, n\}$, is revealed to the Statistician. In this model, we have $X_i = Y_i$ for $i \notin \hat{O}$. However, since \hat{O} is random and potentially dependent of $Y_{1:n}$, it is not true that conditionally to $\hat{O} = O$, $\{X_i : i \in O^c\}$ are i.i.d. drawn from P_{θ^*} (for any deterministic set O of cardinality o).

We denote by $\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)$ the set of all the joint distributions P_n of all the sequences X_1, \dots, X_n generated by the aforementioned two-stage mechanism. This set $\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)$ is larger than all the four sets of contamination introduced in this section. Therefore, the following inequalities hold:

$$\inf_{\hat{\theta}_n} R_d^{\text{PC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{OC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{AC}}(n, \varepsilon, \Theta, \hat{\theta}_n),$$

for any n, ε, Θ and any distance d .

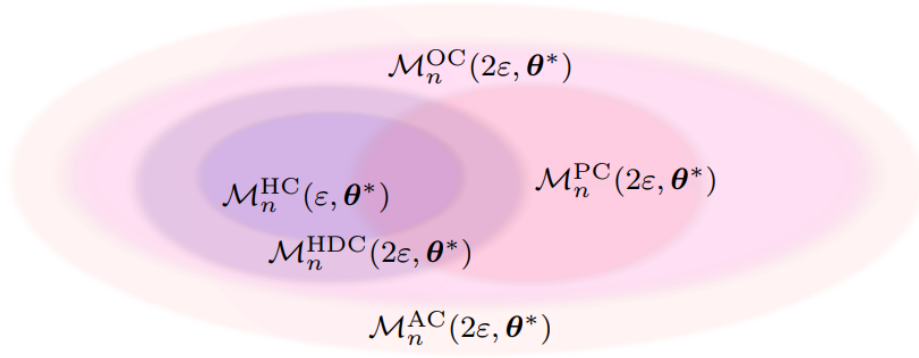


Figure 1.4: Visual representation of the hierarchy between various contamination models. Note that the inclusion of $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ in $\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)$ is somewhat heuristic, based on the relation on the worst-case risks reported in Proposition 1

Minimax rates on the “sparse” simplex and confidence regions

We now specialize the general setting of Section 1.5.1 to a reference distribution P , with expectation θ^* , defined on the simplex Δ^{k-1} . Along with this reference model describing the distribution of inliers, we will use different models of contamination. More precisely, we will establish upper bounds on worst-case risks of the sample mean in the most general, adversarial, contamination setting. Then, matching lower bounds will be provided for minimax risks under Huber’s contamination.

Upper bounds: worst-case risk of the sample mean We denote by Δ_s^{k-1} the set of all $v \in \Delta^{k-1}$ having at most s non-zero entries.

Theorem 1. *For every triple of positive integers (k, s, n) and for every $\varepsilon \in [0, 1]$, the sample mean $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$ satisfies*

$$\begin{aligned} R_{TV}^{\text{AC}}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (s/n)^{1/2} + 2\varepsilon, \\ R_H^{\text{AC}}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (s/n)^{1/2} + \sqrt{2} \varepsilon^{1/2}, \\ R_{\mathbb{L}^2}^{\text{AC}}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (1/n)^{1/2} + \sqrt{2} \varepsilon. \end{aligned}$$

An unexpected and curious phenomenon unveiled by this theorem is that all the three rates are different. As a consequence, the largest possible number of outliers, $o_d^*(n, s)$, that does not impact the minimax rate of estimation of θ^* crucially depends on the considered distance d . Taking into account the relation $\varepsilon = o/n$, we get

$$o_{TV}^*(n, s) \asymp (ns)^{1/2}, \quad o_H^*(n, s) \asymp s, \quad o_{\mathbb{L}^2}^*(n, s) \asymp n^{1/2}.$$

Lower bounds on the minimax risk A natural question, answered in the next theorem, is how tight are the upper bounds obtained in the last theorem. More importantly, one can wonder whether there is an estimator that has a worst-case risk of smaller order than that of the sample mean.

Theorem 2. *There are universal constants $c > 0$ and n_0 , such that for any integers $k \geq 3$, $s \leq k \wedge n$, $n \geq n_0$ and for any $\varepsilon \in [0, 1]$, we have*

$$\begin{aligned}\inf_{\bar{\theta}_n} R_{TV}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon\}, \\ \inf_{\bar{\theta}_n} R_H^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon^{1/2}\}, \\ \inf_{\bar{\theta}_n} R_{\mathbb{L}^2}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(1/n)^{1/2} + \varepsilon\},\end{aligned}$$

where $\inf_{\bar{\theta}_n}$ stands for the infimum over all measurable functions $\bar{\theta}_n$ from $(\Delta^{k-1})^n$ to Δ^{k-1} .

The main consequence of this theorem is that whatever the contamination model is (among those described in Section 1.5.1), the rates obtained for the MLE in Theorem 1 are minimax optimal. Indeed, Theorem 2 yields this claim for Huber's contamination, and in Chapter 2 we will show that the lower bounds obtained for HC remain valid for all the other contamination models and are minimax optimal.

Confidence regions In previous parts, we established bounds for the expected value of estimation error. The aim of this part is to present bounds on estimation error of the sample mean holding with high probability. This also leads to constructing confidence regions for the parameter vector θ^* . To this end, the contamination rate ε and the sparsity s are assumed to be known. It is an interesting open question whether one can construct optimally shrinking confidence regions for unknown ε and s .

Theorem 3. *Let $\delta \in (0, 1)$ be the tolerance level. If $\theta^* \in \Delta_s^{k-1}$, then under any contamination model, the regions of Δ^{k-1} defined by each of the following inequalities*

$$\begin{aligned}d_{\mathbb{L}^2}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}) &\leq (1/n)^{1/2} + \sqrt{2}\varepsilon + (\log(1/\delta)/n)^{1/2}, \\ d_{TV}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}) &\leq (s/n)^{1/2} + 2\varepsilon + (2\log(1/\delta)/n)^{1/2}, \\ d_H(\bar{\mathbf{X}}_n, \boldsymbol{\theta}) &\leq \sqrt{5}((s/n)\log(2s/\delta))^{1/2} + \varepsilon^{1/2} + ((1/2n)\log(2/\delta))^{1/2},\end{aligned}$$

contain θ^* with probability at least $1 - \delta$.

To illustrate the shapes of these confidence regions, we depicted them in Figure 1.5 for a three dimensional example, projected onto the plane containing the probability simplex. The sample mean in this example is equal to $(1/3, 1/2, 1/6)$.

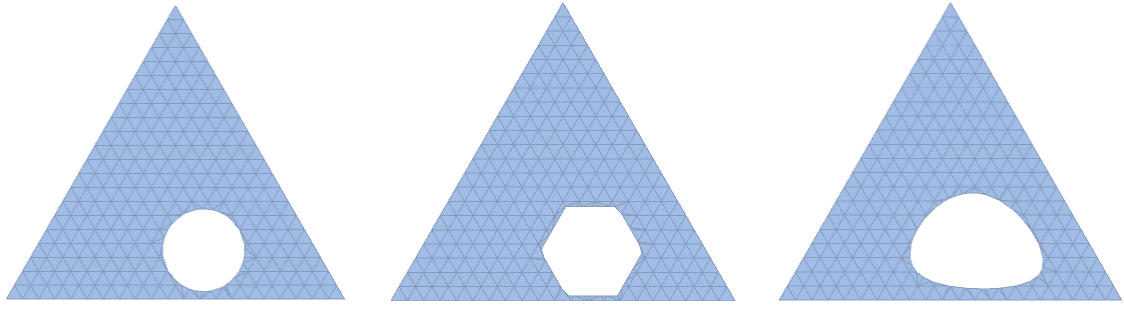


Figure 1.5: The shape of confidence sets (white regions) for the distances \mathbb{L}^2 (left), TV (center), and Hellinger (right) when the sample mean is $(1/3, 1/2, 1/6)$.

Instance based bounds When the dimension k is not finite, we can provide bounds which depend on the reference distribution θ^* or the sample mean \bar{X}_n . We restrict X_i 's to take value in $\{e_1, e_2, \dots\}$, and assume e_j occurs with probability θ_j^* . We define $\alpha_n(\theta) := 2 \sum_{\theta_j < 1/n} \theta_j$ and $\beta_n(\theta) := \frac{1}{\sqrt{n}} \sum_{\theta_j \geq 1/n} \sqrt{\theta_j}$. Using the results of [Berend and Kontorovich \(2013\)](#) the following upper and lower bounds are obtained for the error of the sample mean under the TV distance and adversarial model with ε -contamination.

Proposition 2. *Suppose X_i 's take value in $\{e_1, e_2, \dots\}$, and for $j \in \mathbb{N}$, e_j occurs with probability θ_j^* . For every n and for every $\varepsilon \in [0, 1]$, the sample mean \bar{X}_n satisfies*

$$\frac{\alpha_n(\theta_j^*) + \beta_n(\theta_j^*)}{4} - \frac{1}{4\sqrt{n}} - 2\varepsilon \leq \mathbf{E}d_{\text{TV}}(\bar{X}_n, \theta^*) \leq \alpha_n(\theta_j^*) + \beta_n(\theta_j^*) + 2\varepsilon.$$

These bounds need the knowledge of the reference distribution. The next theorem represent bounds based on the sample mean.

Theorem 4. *Suppose X_i 's take value in $\{e_1, e_2, \dots\}$, and for $j \in \mathbb{N}$, e_j occurs with probability θ_j^* . For every n and for every $\varepsilon \in [0, 1]$, the sample mean \bar{X}_n satisfies*

$$d_{\text{TV}}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \|\bar{X}_n^{1/2}\|_1 + 2\varepsilon + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

with probability at least $1 - \delta$, where $\delta \in (0, 1)$. We also have

$$\mathbf{E}d_{\text{TV}}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \mathbf{E} \|\bar{X}_n^{1/2}\|_1 + 2\varepsilon.$$

1.5.2 Robust Estimation of Gaussian Mean

The goal of this work is to make a step forward by designing an estimator which is not only nearly rate optimal and computationally tractable, but also has a (asymptotic) breakdown point equal to $1/2$, which is the highest possible value of the breakdown point. To construct the estimator, termed iterative spectral dimension reduction or SDR, we combine and suitably adapt

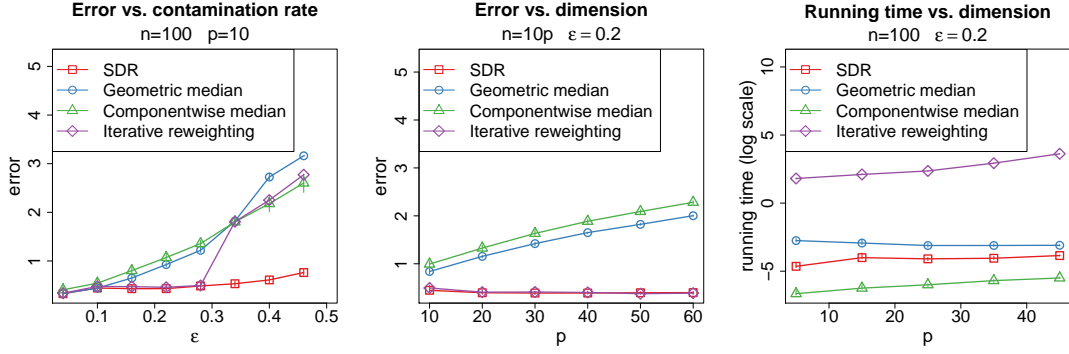


Figure 1.6: Plots that help to visually compare four robust estimators: SDR (our estimator), geometric median (GM) given by (3.1), componentwise median (CM), iteratively reweighted mean (IRM) of (Dalalyan and Minasyan, 2020). The first two plots show that SDR is as accurate as IRM for small ε , with SDR outperforming IRM for ε close to $1/2$. IRM and SDR are naturally much more accurate than GM and CM. The last plot shows that the running time of SDR is comparable to that of GM and is much smaller than that of IRM. More details on these experiments are provided in Section 3.5.

ideas from (Lai et al., 2016) and (Diakonikolas et al., 2017). The main underlying observation is that if we remove some clear outliers and restrict our attention to the subspace spanned by the eigenvectors of the sample covariance matrix corresponding to small eigenvalues, then the sample mean of the projected data points is a rate-optimal estimator. This allows us to iteratively reduce the dimension and eventually to estimate the remaining low-dimensional component of the mean by a standard robust estimator such as the componentwise median or the trimmed mean, see Algorithm 1.

The papers that are the closest to the present work are (Lai et al., 2016), (Diakonikolas et al., 2017) and (Dalalyan and Minasyan, 2020). The spectral dimension reduction scheme was proposed by (Lai et al., 2016) along with an initial sample splitting step ensuring the independence of the estimators over different subspaces. In the case of spherical Gaussian distribution contaminated by non-adversarial outliers, the paper states that the proposed estimator has a squared error at most of order $p \log^2 p \log(p/\varepsilon)/n + \varepsilon^2 \log p$. Compared to this, our results are valid in the more general setting of sub-Gaussian distribution, with arbitrary covariance matrix and adversarial contamination. In addition, our estimator does not rely on sample splitting and, therefore, has a risk with a better dependence on p . As compared to the filtering method of (Diakonikolas et al., 2017), our estimator has the advantage of being independent of ε and our error bound is valid for every covariance matrix and every confidence level. On the down side, our error bound has an extra factor $\log p$ in front of ε^2 . We believe that this factor is an artifact of the proof, but we were unable to remove it. Finally, compared to the iteratively reweighted mean (Dalalyan and Minasyan, 2020), the SDR estimator studied in the present work has a higher breakdown point, does not require the knowledge of ε and is much faster to compute. The advantages and shortcomings of these estimators are summarized in Table 1.1 and Figure 1.6.

Notation. For any pair of integers k and d such that $1 \leq k \leq d$, we denote by \mathcal{V}_k^d the set of all k -dimensional linear subspaces V of \mathbb{R}^d . For $V \in \mathcal{V}_k^d$, we write $k = \dim(V)$ and

	Comput. tractable	Breakdown point	known ε	Squared error rate	known Σ or $\Sigma \propto \mathbf{I}$
Gaussian distribution					
C./G. Median	yes	0.5	no	$\mathbf{r}_\Sigma/n + \varepsilon^2 p$	no
Tukey's Median	no	0.33	no	$\mathbf{r}_\Sigma/n + \varepsilon^2$	no
Agnostic Mean	yes	—	yes	$(p/n) \log^3 p + \varepsilon^2 \log p$	yes
Gaussian and sub-Gaussian distribution					
Iter. Rew. Mean	yes	0.28	yes	$(\mathbf{r}_\Sigma/n) + \varepsilon^2 \log(1/\varepsilon)$	yes
Iterative Filtering	yes	—	yes	$(p/n) \log^a p + \varepsilon^2 \log(1/\varepsilon)$	yes
SDR (this work)	yes	0.5	no	$(\mathbf{r}_\Sigma/n + \varepsilon^2 \log(1/\varepsilon)) \log p$	yes

Table 1.1: Properties of various robust estimators. Agnostic mean, iteratively reweighted mean and iterative filtering are the estimators studied in (Lai et al., 2016), (Dalalyan and Minasyan, 2020) and (Diakonikolas et al., 2017), respectively. The error rates reported for Tukey's median, componentwise median, geometric median and the agnostic mean have been proved for non-adversarial contamination. The squared error rate is provided in the case of a covariance matrix satisfying $\|\Sigma\|_{\text{op}} = 1$.

denote by P_V the orthogonal projection matrix onto V . \mathbb{S}^{d-1} stands for the unit sphere in \mathbb{R}^d . For a $d \times d$ symmetric matrix \mathbf{M} , we denote by $\lambda_1(\mathbf{M}), \dots, \lambda_d(\mathbf{M})$ its eigenvalues sorted in increasing order, and use the notation $\lambda_{\min}(\mathbf{M}) = \lambda_1(\mathbf{M})$, $\lambda_{\max}(\mathbf{M}) = \lambda_d(\mathbf{M})$, $\|\mathbf{M}\|_{\text{op}} = \max(|\lambda_{\min}(\mathbf{M})|, |\lambda_{\max}(\mathbf{M})|)$, $\text{Tr}(\mathbf{M}) = (\lambda_1 + \dots + \lambda_d)(\mathbf{M})$ and $\mathbf{r}_\mathbf{M} = \text{Tr}(\mathbf{M})/\|\mathbf{M}\|_{\text{op}}$. For any integer $n > 0$, we set $[n] = \{1, \dots, n\}$. We will denote by $\mathcal{O} \subset [n]$ the subscripts of the outliers and by $\mathcal{I} = [n] \setminus \mathcal{O}$ the subscripts of inliers. We also use notation $\log_+(x) = \max\{0, \log(x)\}$.

Algorithm 1 $\text{SDR}(\mathbf{X}_1, \dots, \mathbf{X}_n; \Sigma, t)$

```

1: let  $p$  be the dimension of  $\mathbf{X}_1$ 
2: let  $\hat{\mu}^{\text{GM}}$  be the geometric median of  $\mathbf{X}_1, \dots, \mathbf{X}_n$ 
3: let  $\mathcal{S} \leftarrow \{i : \|\mathbf{X}_i - \hat{\mu}^{\text{GM}}\| \leq t\sqrt{p}\}$ 
4: let  $\bar{\mathbf{X}}_\mathcal{S}$  be the sample mean of the filtered sample  $\{\mathbf{X}_i : i \in \mathcal{S}\}$ 
5: let  $\hat{\Sigma}_\mathcal{S}$  be the covariance matrix of the filtered sample  $\{\mathbf{X}_i : i \in \mathcal{S}\}$ 
6: if  $p > 1$  then
7:   let  $V$  be the span of the top  $\lceil p/e \rceil$  principal components of  $\hat{\Sigma}_\mathcal{S} - \Sigma$ 
8:   let  $P_V$  be the orth. projection onto  $V$ 
9:   let  $P_{V^\perp}$  be the orth. projection onto the orth. complement of  $V$ 
10:  let  $\hat{\mu} \leftarrow P_{V^\perp} \bar{\mathbf{X}}_\mathcal{S} + \text{SDR}(P_V \mathbf{X}_1, \dots, P_V \mathbf{X}_n; P_V \Sigma P_V, t)$ 
11: else
12:   let  $\hat{\mu} \leftarrow \hat{\mu}^{\text{GM}}$ 
13: end if
14: return  $\hat{\mu}$ 

```

Adversarially corrupted sub-Gaussian model and spectral dimension reduction

We assume that a set X_1, \dots, X_n of n data points drawn from a distribution P_n is given. This set is assumed to contain at least $n - [n\varepsilon]$ inliers, the remaining points being outliers. All the points lie in the p -dimensional Euclidean space and the inliers are independently drawn from a reference distribution, assumed to be sub-Gaussian with mean $\mu^* \in \mathbb{R}^p$ and covariance matrix Σ . To state the assumptions imposed on the observations in a more precise way, let us recall that the random vector ζ is said to be sub-Gaussian with zero mean and identity covariance matrix, if $\mathbb{E}[\zeta] = 0$, $\mathbb{E}[\zeta\zeta^\top] = \mathbf{I}_p$ and for some $\varsigma > 0$, we have

$$\mathbb{E}[e^{v^\top \zeta}] \leq \exp\{\varsigma \|v\|^2/2\}, \quad \forall v \in \mathbb{R}^p.$$

The parameter ς is commonly called the variance proxy and the writing $\zeta \sim \text{SG}_p(\varsigma)$ is used.

Definition 1. We say that the data generating distribution P_n is an adversarially corrupted sub-Gaussian distribution with mean μ^* , covariance matrix Σ , variance proxy ς and contamination rate ε , if there is a probability space on which we can define a sequence of random vectors $(X_1, Y_1), \dots, (X_n, Y_n)$ such that

1. Y_1, \dots, Y_n are independent and $\Sigma^{-1/2}(Y_i - \mu^*) \sim \text{SG}_p(\varsigma)$ for every $i \in [n]$.

2. the cardinality of $\mathcal{O} = \{i \in [n] : Y_i \neq X_i\}$ is at most equal to $n\varepsilon$.

3. the distribution of (X_1, \dots, X_n) is P_n .

We write then⁶ $P_n \in \text{SGAC}(\mu^*, \Sigma, \varsigma, \varepsilon)$. In the particular case where all Y_i are Gaussian, we will write $P_n \in \text{GAC}(\mu^*, \Sigma, \varepsilon)$.

The estimator we analyze in this work is termed iterative spectral dimension reduction and denoted by $\hat{\mu}^{\text{SDR}}$. It is closely related to the agnostic mean (Lai et al., 2016) and to iterative filtering (Diakonikolas et al., 2017) estimators. We will prove that SDR enjoys most of desired properties in the setting of robust estimation of the sub-Gaussian mean.

The parameters given as input to the iterative spectral dimension reduction algorithm are a strictly decreasing sequence of positive integers p_0, \dots, p_L such that $p_0 = p$ and a positive threshold $t > 0$. We recall that the geometric median is defined by

$$\hat{\mu}^{\text{GM}} \in \arg \min_{\mu \in \mathbb{R}^p} \sum_{i=1}^n \|X_i - \mu\|_2.$$

The algorithm for computing the SDR estimator reads as follows.

1. Start by setting $V_0 = \mathbf{I}_p$.

2. For $\ell = 0, \dots, L-1$ do

(a) Define $\bar{\mu}^{(\ell)} \in \mathbb{R}^{p_\ell}$ as the geometric median of $\{V_\ell^\top X_i : i \in [n]\}$.

⁶SGAC stands for sub-Gaussian with adversarial contamination.

- (b) Define the set $\mathcal{S}^{(\ell)} = \{i \in [n] : \|\mathbf{V}_\ell^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(\ell)}\|_2 \leq t\sqrt{p_\ell}\}$ of filtered data points.
- (c) Let $\bar{\mathbf{X}}^{(\ell)}$ and $\hat{\boldsymbol{\Sigma}}^{(\ell)}$ be the mean vector and the covariance matrix of the filtered sample $\{\mathbf{X}_i : i \in \mathcal{S}^{(\ell)}\}$, that is

$$\bar{\mathbf{X}}^{(\ell)} = \frac{1}{|\mathcal{S}^{(\ell)}|} \sum_{i \in \mathcal{S}^{(\ell)}} \mathbf{X}_i, \quad \hat{\boldsymbol{\Sigma}}^{(\ell)} = \frac{1}{|\mathcal{S}^{(\ell)}|} \sum_{i \in \mathcal{S}^{(\ell)}} (\mathbf{X}_i - \bar{\mathbf{X}})^{\otimes 2}.$$

- (d) Set $\hat{\boldsymbol{\mu}}^{(\ell)} = \mathbf{V}_\ell \mathbf{U}_\ell^\top \mathbf{U}_\ell \mathbf{V}_\ell^\top \bar{\mathbf{X}}^{(\ell)}$, where \mathbf{U}_ℓ is a $(p_\ell - p_{\ell+1}) \times p_\ell$ orthogonal matrix the rows of which are the eigenvectors of $\mathbf{V}_\ell^\top (\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}) \mathbf{V}_\ell$ corresponding to its $(p_\ell - p_{\ell+1})$ smallest eigenvalues.
- (e) Set $\mathbf{V}_{\ell+1} = \mathbf{V}_\ell (\mathbf{U}_\ell^\perp)^\top \in \mathbb{R}^{p \times p_{\ell+1}}$, where \mathbf{U}_ℓ^\perp is a $p_{\ell+1} \times p_\ell$ orthogonal matrix orthogonal to \mathbf{U}_ℓ , that is $\mathbf{U}_\ell^\perp \mathbf{U}_\ell^\top = \mathbf{0}$.
3. Define $\bar{\boldsymbol{\mu}}^{(L)}$ as the geometric median of $\mathbf{V}_L^\top \mathbf{X}_i$ for $i = 1, \dots, n$ and set $\mathcal{S}^{(L)} = \{i \in [n] : \|\mathbf{V}_L^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(L)}\|_2 \leq t\sqrt{p_L}\}$.
4. Define $\hat{\boldsymbol{\mu}}^{(L)} = \mathbf{V}_L \mathbf{V}_L^\top \bar{\mathbf{X}}^{(L)}$, the average of filtered and projected vectors.
5. Return $\hat{\boldsymbol{\mu}}^{\text{SDR}} = \hat{\boldsymbol{\mu}}^{(0)} + \hat{\boldsymbol{\mu}}^{(1)} + \dots + \hat{\boldsymbol{\mu}}^{(L)}$.

The steps described above can be summarized as follows. At each iteration $\ell < L$, we start by determining a filtered subsample $\mathcal{S}^{(\ell)}$ and a “nearly-outlier-orthogonal” subspace $\mathcal{U}_\ell = \text{Im}(\mathbf{V}_\ell \mathbf{U}_\ell^\top)$ of \mathbb{R}^p of dimension $p_\ell - p_{\ell+1}$. We define the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto \mathcal{U}_ℓ as the sample mean of the filtered and projected subsample, and we move to the next step for determining the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto the remaining part of the space. At the last iteration L , when the dimension is well reduced, the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto the subspace \mathcal{U}_L is defined as the average of the filtered subsample projected onto \mathcal{U}_L . The subspaces \mathcal{U}_ℓ are two-by-two orthogonal and span the whole space \mathbb{R}^p . Each subspace is determined from the spectral decomposition of the covariance matrix of the data points projected onto $(\mathcal{U}_0 \oplus \dots \oplus \mathcal{U}_{\ell-1})^\perp$, after removing the points lying at an abnormally large distance from the geometric median.

Choice of the dimension reduction regime The analysis of the error of the SDR estimator conducted in this work leads to an upper bound in which the sequence (p_0, \dots, p_L) is involved only through the expression

$$F(p_0, \dots, p_L) = \sum_{\ell=1}^L \frac{p_{\ell-1}}{p_\ell}.$$

Therefore, an appealing way of choosing this sequence is to minimize the function F under the constraint that the sequence is decreasing and $p_0 = p$ and $p_L = 1$. It follows from the inequality between the arithmetic and geometric means that $F(p_0, \dots, p_L) \geq L p^{1/L}$. Furthermore, the

equality is achieved⁷ in the case when all the terms in the definition of F are equal, *i.e.*, when for some $c > 0$ we have $p_{\ell-1} = cp_\ell$ for every $\ell \in [L]$. Since $p_0 = p$ and $p_L = 1$, this yields $c = p^{1/L}$ or, equivalently, $L = \log p / \log c$. Using these relations, we find that the function F is lower bounded by $Lc = (c/\log c) \log p$. The last step is to find the minimum of the function $c \mapsto c/\log c$ over the interval $(1, \infty)$. One easily checks that this function has a unique minimum at $c = e$. All these considerations advocate for using the dimension reduction regime defined by

$$p_0 = p, \quad p_\ell = \lfloor p_{\ell-1}/e \rfloor + 1, \quad \ell \in [L], \quad p_L = 1, \quad (1.15)$$

where $\lfloor x \rfloor$ is the largest integer strictly smaller than x . Such a definition of (p_ℓ) ensures that $p_{\ell-1}/p_\ell \leq e$ and that⁸ $L \leq 2 \log p$. In the rest of the paper, we assume that the sequence (p_ℓ) is chosen as in (1.15).

Choice of the threshold The SDR procedure has one important tuning parameter: the threshold t used to discard clearly outlying data points. Let us introduce the auxiliary notation

$$\bar{r}_n = \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n}}, \quad \text{and} \quad \tau = \frac{1}{4} \bigwedge \frac{\bar{r}_n}{\sqrt{\log_+(2/\bar{r}_n)}}. \quad (1.16)$$

Note that \bar{r}_n is essentially the quantile, up to a universal constant factor, of order $1 - \delta$ of the distribution of $\|\bar{\mathbf{Y}}_n - \boldsymbol{\mu}^*\|_2$ where \mathbf{Y}_i 's are independently drawn from $\mathcal{N}_p(\boldsymbol{\mu}^*, \Sigma)$ with $\|\Sigma\|_{\text{op}} = 1$. Our theoretical results advocate for using the value $t = t_1 + t_2$, where

$$t_1 = \frac{2(1 + \bar{r}_n)}{1 - 2\varepsilon^*}, \quad t_2 = 1 + \frac{\bar{r}_n}{\sqrt{\tau}} + \sqrt{2 + \log(2/\tau)},$$

where $\varepsilon^* < 1/2$ is the largest value of the contamination rate that the algorithm may handle.

Let $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_n$ be independent Gaussian with zero mean and covariance Σ . The expression of t_1 is obtained as an upper bound on the quantile of order $1 - \delta/2$ of the distribution of the random variable

$$T_1 = \sup_V \frac{2}{n(1 - 2\varepsilon) \dim(V)} \sum_{i=1}^n \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2,$$

see Lemma 2 and its proof for further details. Similarly, t_2 is defined so that the event

$$\sup_V \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i\|_2^2 > t_2^2 \dim(V)) \leq n\tau$$

has a probability at least $1 - \delta/2$. Although we tried to get sharp values for these thresholds t_1

⁷We relax here the assumption that all the entries p_ℓ are integers.

⁸To check this inequality, one can use the fact that $3 \leq p_{L-2} \leq pe^{2-L} + e/(e-1)$. This implies $L \leq 2 \log p$ for $p \geq 6$. For smaller values of p , the inequality can be checked by direct computations.

and t_2 , it is certainly possible to improve these values either by better mathematical arguments or by empirical considerations. Of course, smaller values of the thresholds t_1 and t_2 satisfying aforementioned conditions lead to an SDR estimator having smaller error.

Assessing the error of the SDR estimator

The iterative spectral dimension reduction estimator defined in previous sections has some desirable properties of a robust estimator that are easy to check. In particular, it is clearly equivariant by translation, orthogonal linear transform and global scaling. Furthermore, the breakdown point of the estimator is equal to that of the geometric median, that is to $1/2$. This means that even if almost the half of data points are chosen to be infinitely large, the estimator will not “break down” in the sense of becoming infinitely large. However, the fact that the estimated value does not become infinitely large, it might be not very close to the true mean. The next theorem shows that this is not the case and that the error of the SDR estimator has a nearly rate-optimal behavior even when the contamination rate is close to $1/2$. The adverb “nearly” is used here to reflect the presence of the $\sqrt{\log p}$ factor in the error bound, which is not present in the minimax rate.

Theorem 5. *Let $\varepsilon^* \in (0, 1/2)$, and $\delta \in (0, 1/2)$. Define \bar{r}_n and τ as in (1.16). For every $\varepsilon \leq \varepsilon^*$, let $\hat{\mu}^{\text{SDR}}$ be the estimator returned by Algorithm 1.5.2 with*

$$t = \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + \frac{\bar{r}_n}{\sqrt{\tau}} \right) + \sqrt{2 + 2 \log(1/\tau)}.$$

There exists a universal constant C such that for every $P_n \in \text{GAC}(\mu^, \Sigma, \varepsilon)$ with $\varepsilon \leq \varepsilon^*$ and $\|\Sigma\|_{\text{op}} = 1$, the probability of the event*

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_\Sigma}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\frac{\log(1/\delta)}{n}} \right)$$

is at least $1 - \delta$. Moreover, the constant C from the last display can be made explicit by replacing the effective rank \mathbf{r}_Σ by the dimension p in the definition of \bar{r}_n : That is, for every $\delta \in (0, 1/5)$ the inequality

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{156 \sqrt{2 \log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{2p}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\frac{3 \log(2/\delta)}{n}} \right)$$

holds with probability at least $1 - 5\delta$.

If we compare this result with its counterpart established in (Dalalyan and Minasyan, 2020) for the iteratively reweighted mean, besides the extra $\log p$ factor, we see that the above error bound does not reduce to the error of the empirical mean when the contamination rate goes

⁹Since in this theorem Σ is assumed to be known, we can always divide all the data points X_i by $\|\Sigma\|_{\text{op}}^{1/2}$ to get a data set with a covariance matrix satisfying $\|\Sigma\|_{\text{op}} = 1$.

to zero. We do not know whether this is just a drawback of our proof, or it is an intrinsic property of the estimator. Our numerical experiments reported later on suggest that it might be a property of the estimator.

There is another logarithmic factor, $\sqrt{\log(2/\varepsilon)}$, present in the second term of the error bounds provided by the last theorem, which does not appear in the minimax rate. There are computationally intractable robust estimators of the Gaussian mean, such as the Tukey median, that have an error bound free of this factor. However, all the known error bounds provably valid for polynomial time algorithms has this extra $\sqrt{\log(2/\varepsilon)}$ factor. Furthermore, this factor is known to be unavoidable in the case of sub-Gaussian model with adversarial contamination¹⁰, see (Lugosi and Mendelson, 2021, Section 2).

The case of unknown covariance matrix

The SDR estimator, as defined in Algorithm 1, requires the knowledge of covariance matrix Σ . In this section we consider the case where the matrix Σ is unknown, but an approximation of the latter is available. Namely, we assume that we have access to a matrix $\tilde{\Sigma}$ and to a real number $\gamma > 0$ such that $\|\tilde{\Sigma} - \Sigma\|_{\text{op}} \leq \gamma \|\Sigma\|_{\text{op}}$. In such a situation, we can replace in the SDR estimator the true covariance matrix by its approximation $\tilde{\Sigma}$. This will necessarily require to adjust the threshold t accordingly. The goal of the present section is to propose a suitable choice of t and to show the impact of the approximation error γ on the estimation accuracy.

As mentioned, the parameter t used in Algorithm 1 needs to be properly tuned in order to account for the approximation error in the covariance matrix. To this end, we introduce the following auxiliary notation similar to those presented in (1.16):

$$\tilde{r}_n = \frac{\sqrt{C_\gamma \mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n}} \quad \text{and} \quad \tilde{\tau} = \frac{1}{4} \bigwedge \frac{\tilde{r}_n}{\sqrt{\log_+(2/\tilde{r}_n)}}, \quad (1.17)$$

where $C_\gamma = (1 + \gamma)/(1 - \gamma)$. Compared to (1.16), the main difference here is the presence of the factor C_γ (which is equal to one if $\gamma = 0$) and the substitution of the effective rank of Σ by that of its approximation $\tilde{\Sigma}$. In the rest of this section, we assume that Σ is invertible.

Theorem 6. *Let $\varepsilon^* \in (0, 1/2)$, $\delta \in (0, 1/2)$ and define \tilde{r}_n and τ as in (1.17). Assume that $\tilde{\Sigma}$ satisfies $\|\Sigma^{-1/2} \tilde{\Sigma} \Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma$ for some $\gamma \in (0, 1/2]$. Let $\hat{\mu}^{\text{SDR}}$ be the output of $\text{SDR}(\mathbf{X}_1, \dots, \mathbf{X}_n; \tilde{\Sigma}, \tilde{t}_\gamma)$, see Algorithm 1, with*

$$\tilde{t}_\gamma = \frac{\|\tilde{\Sigma}\|_{\text{op}}}{1 - \gamma} \left\{ \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + \frac{\tilde{r}_n}{\sqrt{\tilde{\tau}}} \right) + \sqrt{2 + \log(2/\tau)} \right\}.$$

Then, there exists a universal constant C such that for every data generating distribution $P_n \in$

¹⁰Both sub-Gaussianity of the reference distribution and the adversarial nature of the contamination are important for getting the extra $\sqrt{\log(2/\varepsilon)}$ factor in the minimax rate.

$\text{GAC}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}, \varepsilon)$ with $\varepsilon \leq \varepsilon^*$, the probability of the event

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C \|\boldsymbol{\Sigma}\|_{\text{op}}^{1/2} \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_{\boldsymbol{\Sigma}}}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\varepsilon \gamma} + \sqrt{\frac{\log(1/\delta)}{n}} \right) \quad (1.18)$$

is at least $1 - \delta$.

On the one hand, if the value of γ is at most of order $\sqrt{(\mathbf{r}_{\boldsymbol{\Sigma}}/n) \log(1/\varepsilon)} + \varepsilon \log(1/\varepsilon)$ then Theorem 6 implies that the estimation error is of the same order as in the case of known covariance matrix $\boldsymbol{\Sigma}$ (Theorem 5). For instance, if the matrix $\boldsymbol{\Sigma}$ is assumed to be diagonal, one can defined $\tilde{\boldsymbol{\Sigma}}$ as the diagonal matrix composed of robust estimators of the variances of univariate contaminated Gaussian samples; see, for instance, Section 2 in (Comminges et al., 2021). For recent advances on robust estimation of (non-diagonal) covariance matrices by computationally tractable algorithms we refer the reader to (Cheng et al., 2019b).

On the other hand, if the value of γ for which the condition $\|\tilde{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}\|_{\text{op}} \leq \gamma \|\boldsymbol{\Sigma}\|_{\text{op}}$ is known to be true is of larger order than $\sqrt{(\mathbf{r}_{\boldsymbol{\Sigma}}/n) \log(1/\varepsilon)} + \varepsilon \log(1/\varepsilon)$, then $\sqrt{\varepsilon \gamma}$ dominates the other terms appearing in the error bound (1.18). Moreover, if γ is of constant order, then we get the error rate $\sqrt{\frac{\mathbf{r}_{\boldsymbol{\Sigma}}}{n}} + \sqrt{\varepsilon}$, which is in line with previously known bounds for computationally tractable estimators; see for example (Lai et al., 2016, Theorem 1.1), (Diakonikolas et al., 2017, Theorem 3.2), (Dalalyan and Minasyan, 2020, Theorem 4).

Chapter 2

Contamination models and robust estimation on the probability simplex

We consider the problem of estimating the mean of a distribution supported by the k -dimensional probability simplex in the setting where an ε fraction of observations are subject to adversarial corruption. A simple particular example is the problem of estimating the distribution of a discrete random variable. Assuming that the discrete variable takes k values, the unknown parameter θ is a k -dimensional vector belonging to the probability simplex. We first describe various settings of contamination and discuss the relation between these settings. We then establish minimax rates when the quality of estimation is measured by the total-variation distance, the Hellinger distance, or the \mathbb{L}^2 -distance between two probability measures. We also provide confidence regions for the unknown mean that shrink at the minimax rate. Our analysis reveals that the minimax rates associated to these three distances are all different, but they are all attained by the sample average. Furthermore, we show that the latter is adaptive to the possible sparsity of the unknown vector. Some numerical experiments illustrating our theoretical findings are reported.

2.1	Introduction	38
2.2	Various models of contamination	40
2.2.1	Huber's contamination	40
2.2.2	Huber's deterministic contamination	41
2.2.3	Oblivious contamination	42
2.2.4	Parameter contamination	42
2.2.5	Adversarial contamination	42
2.2.6	Minimax risk "in expectation" versus "in deviation"	43
2.3	Prior work	44
2.4	Minimax rates on the "sparse" simplex and confidence regions	45
2.4.1	Upper bounds: worst-case risk of the sample mean	45

2.4.2 Lower bounds on the minimax risk	46
2.4.3 Confidence regions	48
2.5 Instance based bounds	49
2.6 Illustration on a numerical example	50
2.7 Summary and conclusion	51

2.1 Introduction

Assume X_1, \dots, X_n are n independent and identically distributed random variables taking their values in the k -dimensional probability simplex $\Delta^{k-1} = \{v \in \mathbb{R}_+^k : v_1 + \dots + v_k = 1\}$. Our goal is to estimate the unknown vector $\theta = \mathbf{E}[X_i]$ in the case where the observations are contaminated by outliers. An important particular case is the estimation of the distribution of a discrete random variable X taking k distinct values. In this particular case, X_i 's take values in $\{e_1, \dots, e_k\}$, the set of the vectors of the canonical basis, which are also the extreme points of the simplex Δ^{k-1} .

In this introduction, to convey the main messages, we limit ourselves to the Huber contamination model, although our results apply to the more general adversarial contamination. Huber's contamination model assumes that there are two probability measures P, Q on Δ^{k-1} and a real $\varepsilon \in [0, 1/2)$ such that X_i is drawn from

$$P_i = (1 - \varepsilon)P + \varepsilon Q, \quad \forall i \in \{1, \dots, n\}.$$

This amounts to assuming that $(1 - \varepsilon)$ -fraction of observations, called inliers, are drawn from a reference distribution P , whereas ε -fraction of observations are outliers and are drawn from another distribution Q . In general, all the three parameters P, Q and ε are unknown. The parameter of interest is some functional (such as the mean, the standard deviation, etc.) of the reference distribution P , whereas Q and ε play the role of nuisance parameters.

When the unknown parameter lives on the probability simplex, there are many appealing ways of defining the risk. We focus on the following three metrics: total-variation, Hellinger and \mathbb{L}^2 distances¹

$$\begin{aligned} d_{\text{TV}}(\hat{\theta}, \theta) &:= 1/2 \|\hat{\theta} - \theta\|_1, \\ d_{\text{H}}(\hat{\theta}, \theta) &:= 1/\sqrt{2} \|\hat{\theta}^{1/2} - \theta^{1/2}\|_2, \\ d_{\mathbb{L}^2}(\hat{\theta}, \theta) &:= \|\hat{\theta} - \theta\|_2. \end{aligned}$$

The Hellinger distance above is well defined when the estimator $\hat{\theta}$ is non-negative, which will be the case throughout this work. We will further assume that the dimension k may be large,

¹We write $\|u\|_q = (\sum_{j=1}^k |u_j|^q)^{1/q}$ and $u^q = (u_1^q, \dots, u_k^q)$ for any $u \in \mathbb{R}_+^k$ and $q > 0$.

but the vector θ is s -sparse, for some $s \leq k$, *i.e.* $\#\{j : \theta_j \neq 0\} \leq s$. Our main interest is in constructing confidence regions and evaluating the minimax risk

$$\mathfrak{R}_{\square}(n, k, s, \varepsilon) := \inf_{\bar{\theta}_n} \sup_{P, Q} \mathbf{E}[d_{\square}(\bar{\theta}_n, \theta)], \quad (2.1)$$

where the *inf* is over all estimators $\bar{\theta}_n$ built upon the observations $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} (1-\varepsilon)\mathbf{P} + \varepsilon\mathbf{Q}$ and the *sup* is over all distributions \mathbf{P}, \mathbf{Q} on the probability simplex such that the mean θ of \mathbf{P} is s -sparse. The subscript \square of \mathfrak{R} above refers to the distance used in the risk, so that \square is TV, H, or \mathbb{L}^2 .

The problem described above arises in many practical situations. One example is an election poll: each participant expresses his intention to vote for one of k candidates. Thus, each θ_j is the true proportion of electors of candidate j . The results of the poll contain outliers, since some participants of the poll prefer to hide their true opinion. Another example related to elections, is the problem of counting votes across all constituencies. Each constituency communicates a vector of proportions to a central office, which is in charge of computing the overall proportions. However, in some constituencies (hopefully a small fraction only) the results are rigged. Hence, the set of observed vectors contains outliers.

We intend to provide non-asymptotic upper and lower bounds on the minimax risk that match up to numerical constants. In addition, we will provide confidence regions of the form $B_{\square}(\hat{\theta}_n, r_{n,\varepsilon,\delta}) = \{\theta : d_{\square}(\hat{\theta}_n, \theta) \leq r_{n,\varepsilon,\delta}\}$ containing the true parameter with probability at least $1-\delta$ and such that the radius $r_{n,\varepsilon,\delta}$ goes to zero at the same rate as the corresponding minimax risk.

When there is no outlier, *i.e.*, $\varepsilon = 0$, it is well known that the sample mean

$$\bar{\mathbf{X}}_n := \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

is minimax-rate-optimal and the rates corresponding to various distances are

$$\mathfrak{R}_{\mathbb{L}^2}(n, k, s, 0) \asymp (1/n)^{1/2} \quad \text{and} \quad \mathfrak{R}_{\square}(n, k, s, 0) \asymp (s/n)^{1/2} \quad \text{for } \square \in \{\text{TV}, \text{H}\}.$$

This raises several questions in the setting where data contains outliers. In particular, the following three questions will be answered in this work:

- Q1.** How do the risks \mathfrak{R}_{\square} depend on ε ? What is the largest proportion of outliers for which the minimax rate is the same as in the outlier-free case ?
- Q2.** Does the sample mean remain optimal in the contaminated setting?
- Q3.** What does happen if the unknown parameter θ is s -sparse ?

The most important step for answering these questions is to show that

$$\begin{aligned}\mathfrak{R}_{\text{TV}}(n, k, s, \varepsilon) &\asymp (s/n)^{1/2} + \varepsilon, \\ \mathfrak{R}_{\text{H}}(n, k, s, \varepsilon) &\asymp (s/n)^{1/2} + \varepsilon^{1/2}, \\ \mathfrak{R}_{\mathbb{L}^2}(n, k, s, \varepsilon) &\asymp (1/n)^{1/2} + \varepsilon.\end{aligned}$$

It is surprising to see that all the three rates are different leading to important discrepancies in the answers to the second part of question Q1 for different distances. Indeed, it turns out that the minimax rate is not deteriorated if the proportion of the outliers is smaller than $(s/n)^{1/2}$ for the TV-distance, s/n for the Hellinger distance and $(1/n)^{1/2}$ for the \mathbb{L}^2 distance. Furthermore, we prove that the sample mean is minimax rate optimal. Thus, even when the proportion of outliers ε and the sparsity s are known, it is not possible to improve upon the sample mean. In addition, we show that all these claims hold true for the adversarial contamination and we provide corresponding confidence regions.

The rest of the paper is organized as follows. Section 2.2 introduces different possible ways of modeling data sets contaminated by outliers. Pointers to relevant prior work are given in Section 2.3. Main theoretical results and their numerical illustration are reported in Section 2.4 and Section 2.6, respectively. Section 2.7 contains a brief summary of the obtained results and their consequences, whereas the proofs are postponed to the appendix.

2.2 Various models of contamination

Different mathematical frameworks have been used in the literature to model the outliers. We present here five of them, from the most restrictive one to the most general, and describe their relationship. We present these frameworks in the general setting when the goal is to estimate the parameter θ^* of a reference distribution P_{θ^*} when ε proportion of the observations are outliers.

2.2.1 Huber's contamination

The most popular framework for studying robust estimation methods is perhaps the one of Huber's contamination. In this framework, there is a distribution Q defined on the same space as the reference distribution P_{θ^*} such that all the observations X_1, \dots, X_n are independent and drawn from the mixture distribution $P_{\varepsilon, \theta^*, Q} := (1 - \varepsilon)P_{\theta^*} + \varepsilon Q$.

This corresponds to the following mechanism: one decides with probabilities $(1 - \varepsilon, \varepsilon)$ whether a given observation is an inlier or an outlier. If the decision is made in favor of being inlier, the observation is drawn from P_{θ^*} , otherwise it is drawn from Q . More formally, if we denote by \widehat{O} the random set of outliers, then conditionally to $\widehat{O} = O$,

$$\{X_i : i \notin O\} \stackrel{\text{iid}}{\sim} P_{\theta^*}, \{X_i : i \in O\} \stackrel{\text{iid}}{\sim} Q, \{X_i : i \in O\} \perp\!\!\!\perp \{X_i : i \notin O\}, \quad (2.2)$$

for every $O \subset \{1, \dots, n\}$. Furthermore, for every subset O of the observations, we have $P(\widehat{O} = O) = (1-\varepsilon)^{n-|O|}\varepsilon^{|O|}$. We denote by² $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ the set of joint probability distributions P_n of the random variables X_1, \dots, X_n satisfying the foregoing condition.

2.2.2 Huber's deterministic contamination

The set of outliers as well as the number of outliers in Huber's model of contamination are random. This makes it difficult to compare this model to the others that will be described later in this section. To cope with this, we define here another model, termed Huber's deterministic contamination. As its name indicates, this new model has the advantage of containing a deterministic number of outliers, in the same time being equivalent to Huber's contamination in a sense that will be made precise below.

We say that the distribution P_n of X_1, \dots, X_n belongs to the Huber's deterministic contamination model denoted by $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \theta^*)$, if there are a set $O \subset \{1, \dots, n\}$ of cardinality at most $n\varepsilon$ and a distribution Q such that (2.2) is true. The apparent similarity of models $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ and $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \theta^*)$ can also be formalized mathematically in terms of the orders of magnitude of minimax risks. To ease notation, we let $R_d^\square(n, \varepsilon, \Theta, \widehat{\theta})$ to be the worst-case risk of an estimator $\widehat{\theta}$, where \square is either HC or HDC. More precisely, for $\mathcal{M}_n^\square(\varepsilon, \Theta) := \cup_{\theta \in \Theta} \mathcal{M}_n^\square(\varepsilon, \theta)$, we set³

$$R_d^\square(n, \varepsilon, \Theta, \widehat{\theta}) := \sup_{P_n \in \mathcal{M}_n^\square(\varepsilon, \Theta)} \mathbf{E}[d(\widehat{\theta}, \theta^*)].$$

This definition assumes that the parameter space Θ is endowed with a pseudo-metric $d : \Theta \times \Theta \rightarrow \mathbb{R}_+$. When $\Theta = \{\theta^*\}$ is a singleton, we write $R_{d,n}^\square(\varepsilon, \theta^*, \widehat{\theta})$ instead of $R_d^\square(n, \varepsilon, \{\theta^*\}, \widehat{\theta})$.

Proposition 3. *Let $\widehat{\theta}_n$ be an arbitrary estimator of θ^* . For any $\varepsilon \in (0, 1/2)$,*

$$R_d^{\text{HC}}(n, \varepsilon, \theta^*, \widehat{\theta}_n) \leq R_{d,n}^{\text{HDC}}(2\varepsilon, \theta^*, \widehat{\theta}_n) + e^{-n\varepsilon/3} R_{d,n}^{\text{HDC}}(1, \theta^*, \widehat{\theta}_n), \quad (2.3)$$

$$\sup_{P_n \in \mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} rP(d(\widehat{\theta}_n, \theta^*) > r) \leq R_{d,n}^{\text{HDC}}(2\varepsilon, \theta^*, \widehat{\theta}_n) + re^{-n\varepsilon/3}. \quad (2.4)$$

Proof in the appendix, page 82

Denote by \mathcal{D}_Θ the diameter of Θ , $\mathcal{D}_\Theta := \max_{\theta, \theta'} d(\theta, \theta')$. Proposition 3 implies that

$$\inf_{\widehat{\theta}_n} R_d^{\text{HC}}(n, \varepsilon, \Theta, \widehat{\theta}_n) \leq \inf_{\widehat{\theta}_n} R_{d,n}^{\text{HDC}}(n, 2\varepsilon, \Theta, \widehat{\theta}_n) + e^{-n\varepsilon/3} \mathcal{D}_\Theta. \quad (2.5)$$

When Θ is bounded, the last term is typically of smaller order than the minimax risk over $\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \Theta)$. Therefore, the minimax rate of estimation in Huber's model is not slower than the minimax rate of estimation in Huber's deterministic contamination model. This entails that a lower bound on the minimax risk established in HC-model furnishes a lower bound in HDC-model.

²The superscript HC refers to the Huber's contamination

³The subscript d refers to the distance d used in the definition of the risk.

2.2.3 Oblivious contamination

A third model of contamination that can be of interest is the oblivious contamination. In this model, it is assumed that the set O of cardinality o and the joint distribution Q_O of outliers are determined in advance, possibly based on the knowledge of the reference distribution P_{θ^*} . Then, the outliers $\{X_i : i \in O\}$ are drawn randomly from Q_O independently of the inliers $\{X_i : i \in O^c\}$. The set of all the joint distributions P_n of random variables X_1, \dots, X_n generated by such a mechanism will be denoted by $\mathcal{M}_n^{\text{OC}}(\varepsilon, \theta^*)$. The model of oblivious contamination is strictly more general than that of Huber's deterministic contamination, since it does not assume that the outliers are iid. Therefore, the minimax risk over $\mathcal{M}_n^{\text{OC}}(\varepsilon, \Theta)$ is larger than the minimax risk over $\mathcal{M}_n^{\text{HDC}}(\varepsilon, \Theta)$:

$$\inf_{\hat{\theta}_n} R_d^{\text{HDC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{OC}}(n, \varepsilon, \Theta, \hat{\theta}_n).$$

The last inequality holds true for any set Θ , any contamination level $\varepsilon \in (0, 1)$ and any sample size.

2.2.4 Parameter contamination

In the three models considered above, the contamination acts on the observations. One can also consider the case where the parameters of the distributions of some observations are contaminated. More precisely, for some set $O \subset \{1, \dots, n\}$ selected in advance (but unobserved), the outliers $\{X_i : i \in O\}$ are independent and independent of the inliers $\{X_i : i \in O^c\}$. Furthermore, each outlier X_i is drawn from a distribution $Q_i = P_{\theta_i}$ belonging to the same family as the reference distribution, but corresponding to a contaminated parameter $\theta_i \neq \theta^*$. Thus, the joint distribution of the observations can be written as $(\bigotimes_{i \in O^c} P_{\theta^*}) \otimes (\bigotimes_{i \in O} P_{\theta_i})$. The set of all such distributions P_n will be denoted by $\mathcal{M}_n^{\text{PC}}(\varepsilon, \theta^*)$, where PC refers to "parameter contamination".

2.2.5 Adversarial contamination

The last model of contamination we describe in this work, the adversarial contamination, is the most general one. It corresponds to the following two-stage data generation mechanism. In a first stage, iid random variables Y_1, \dots, Y_n are generated from a reference distribution P_{θ^*} . In a second stage, an adversary having access to Y_1, \dots, Y_n chooses a (random) set \hat{O} of (deterministic) cardinality s and arbitrarily modifies data points $\{Y_i : i \in \hat{O}\}$. The resulting sample, $\{X_i : i = 1, \dots, n\}$, is revealed to the Statistician. In this model, we have $X_i = Y_i$ for $i \notin \hat{O}$. However, since \hat{O} is random and potentially dependent of $Y_{1:n}$, it is not true that conditionally to $\hat{O} = O$, $\{X_i : i \in O^c\}$ are iid drawn from P_{θ^*} (for any deterministic set O of cardinality o).

We denote by $\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)$ the set of all the joint distributions P_n of all the sequences

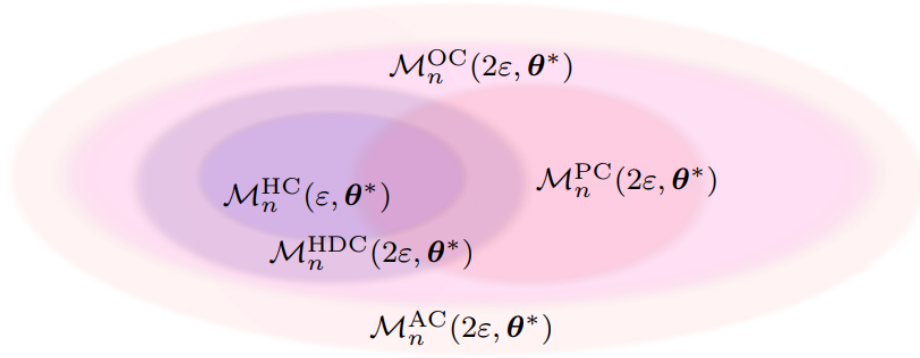


Figure 2.1: Visual representation of the hierarchy between various contamination models. Note that the inclusion of $\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$ in $\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)$ is somewhat heuristic, based on the relation on the worst-case risks reported in Proposition 3.

X_1, \dots, X_n generated by the aforementioned two-stage mechanism. This set $\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)$ is larger than all the four sets of contamination introduced in this section. Therefore, the following inequalities hold:

$$\inf_{\hat{\theta}_n} R_d^{\text{PC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{OC}}(n, \varepsilon, \Theta, \hat{\theta}_n) \leq \inf_{\hat{\theta}_n} R_d^{\text{AC}}(n, \varepsilon, \Theta, \hat{\theta}_n),$$

for any n, ε, Θ and any distance d .

2.2.6 Minimax risk “in expectation” versus “in deviation”

Most prior work on robust estimation focused on establishing upper bounds on the minimax risk in deviation⁴, as opposed to the minimax risk in expectation defined by (2.1). One of the reasons for dealing with the deviation is that it makes the minimax risk meaningful for models⁵ having random number of outliers and unbounded parameter space Θ . The formal justification of this claim is provided by the following result.

Proposition 4. *Let Θ be a parameter space such that $\mathcal{D}_\Theta = \sup_{\theta, \theta' \in \Theta} d(\theta, \theta') = +\infty$. Then, for every estimator $\hat{\theta}_n$, every $\varepsilon > 0$ and $n \in \mathbb{N}$, we have $R_d^{\text{HC}}(n, \varepsilon, \Theta, \hat{\theta}_n) = +\infty$.*

Proof in the appendix, page 83

This result shows, in particular, that the last term in (2.5), involving the diameter of Θ is unavoidable. Such an explosion of the minimax risk occurs because Huber’s model allows the number of outliers to be as large as $n/2$ with a strictly positive probability. One approach to overcome this shortcoming is to use the minimax risk in deviation. Another approach is to limit

⁴We call a risk bound in deviation any bound on the distance $d(\hat{\theta}, \theta^*)$ that holds true with a probability close to one, for any parameter value $\theta^* \in \Theta$.

⁵This is the case, for instance, of the Gaussian model with Huber’s contamination.

theoretical developments to the models HDC, PC, OC or AC, in which the number of outliers is deterministic.

2.3 Prior work

Robust estimation is an area of active research in Statistics since at least five decades (Donoho and Gasko, 1992; Donoho and Huber, 1983; Huber, 1964; Rousseeuw and Hubert, 1999; Tukey, 1975). Until very recently, theoretical guarantees were almost exclusively formulated in terms of the notions of breakdown point, sensitivity curve, influence function, etc. These notions are well suited for accounting for gross outliers, observations that deviate significantly from the data points representative of an important fraction of data set.

More recently, various authors investigated (Chen et al., 2013; Dalalyan and Chen, 2012; Nguyen and Tran, 2013) the behavior of the risk of robust estimators as a function of the rate of contamination ε . A general methodology for parametric models subject to Huber's contamination was developed in Chen et al. (2016, 2018). This methodology allowed for determining the rate of convergence of the minimax risk as a function of the sample size n , dimension k and the rate of contamination ε . An interesting phenomenon was discovered: in the problem of robust estimation of the Gaussian mean, classic robust estimators such as the coordinatewise median or the geometric median do not attain the optimal rate $(k/n)^{1/2} + \varepsilon$. This rate is attained by Tukey's median, the maximizer of Tukey's halfspace depth, the computation of which is costly in a high dimensional setting. Detailed analysis of Tukey's halfspace depth was conducted in (Brunel, 2019).

In the model analyzed in this paper, we find the same minimax rate, $(k/n)^{1/2} + \varepsilon$, only when the total-variation distance is considered. A striking difference is that this rate is attained by the sample mean which is efficiently computable in any dimension. This property is to some extent similar to the problem of robust density estimation (Liu and Gao, 2019), in which the standard kernel estimators are minimax optimal in contaminated setting. Note that the fact that in the sparse setting the improvement from $(k/n)^{1/2}$ to $(s/n)^{1/2}$ can be achieved without any penalization, and that the constraint of belonging to the probability simplex acts as a sparsity favoring penalty, was already known in the literature, see (Dalalyan and Sebban, 2018; Xia and Koltchinskii, 2016). Interestingly, similar phenomenon is observed in problems of estimation under shape constraints (Bellec, 2018; Guntuboyina and Sen, 2018). It is an interesting avenue of future research to analyze the robustness of the maximum likelihood estimator in this context.

Computational intractability of Tukey's median motivated a large number of studies that aimed at designing computationally tractable methods with nearly optimal statistical guarantees. Many of these works went beyond Huber's contamination by considering parameter contamination models (Bhatia et al., 2017; Carpentier et al., 2018; Collier and Dalalyan, 2019), oblivious contamination (Feng et al., 2014; Lai et al., 2016) or adversarial contami-

nation (Balakrishnan et al., 2017; Dalalyan and Minasyan, 2020; Diakonikolas et al., 2016a, 2017, 2018b). Interestingly, in the problem of estimating the Gaussian mean, it was proven that the minimax rates under adversarial contamination are within a factor at most logarithmic in n and k of the minimax rates under Huber’s contamination⁶. While each of the aforementioned papers introduced clearly the conditions on the contamination, to our knowledge, none of them described different possible models and the relationship between them.

Another line of growing literature on robust estimation aims at robustifying estimators and prediction methods to heavy tailed distributions, see (Audibert and Catoni, 2011; Chinot et al., 2018; Devroye et al., 2016; Donoho and Montanari, 2016; Joly et al., 2017; Lecué and Lerasle, 2017; Lugosi and Mendelson, 2019b; Minsker, 2015, 2018a). The results of those papers are of a different nature, as compared to the present work, not only in terms of the goals, but also in terms of mathematical and algorithmic tools.

In the case of the discrete distributions, Braess and Sauer (2004) established the minimax rates under the Kullback-Leibler divergence. More recently, Kamath et al. (2015) determined the minimax rates under other distances such as \mathbb{L}^2 , TV and the general family of f -divergence (including the χ^2 -divergence and the Hellinger distance). The estimator proposed in (Kamath et al., 2015), achieving the minimax rate for \mathbb{L}^2 and TV distances, is the sample mean while different estimators are proposed for the other distances. Concerning the robust estimation of discrete distributions, Chen et al. (2020); Jain and Orlitsky (2019); Qiao and Valiant (2018) studied the case of group-contamination. The distinctive feature of this setting is a better dependence of the minimax rate on the contamination rate ε . More precisely, if each group contains m samples, and ε fraction of groups are contaminated, the rates are obtained by replacing ε by ε/\sqrt{m} . The estimators achieving these rates are much more sophisticated than the sample mean.

2.4 Minimax rates on the “sparse” simplex and confidence regions

We now specialize the general setting of Section 2.2 to a reference distribution P , with expectation θ^* , defined on the simplex Δ^{k-1} . Along with this reference model describing the distribution of inliers, we will use different models of contamination. More precisely, we will establish upper bounds on worst-case risks of the sample mean in the most general, adversarial, contamination setting. Then, matching lower bounds will be provided for minimax risks under Huber’s contamination.

2.4.1 Upper bounds: worst-case risk of the sample mean

We denote by Δ_s^{k-1} the set of all $v \in \Delta^{k-1}$ having at most s non-zero entries.

⁶All these papers consider the risk in deviation, so that the minimax risk under Huber’s contamination is finite.

Theorem 7. *For every triple of positive integers (k, s, n) and for every $\varepsilon \in [0, 1]$, the sample mean $\bar{\mathbf{X}}_n := \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$ satisfies*

$$\begin{aligned} R_{TV}^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\mathbf{X}}_n) &\leq (s/n)^{1/2} + 2\varepsilon, \\ R_H^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\mathbf{X}}_n) &\leq (s/n)^{1/2} + \sqrt{2} \varepsilon^{1/2}, \\ R_{\mathbb{L}_2}^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\mathbf{X}}_n) &\leq (1/n)^{1/2} + \sqrt{2} \varepsilon. \end{aligned}$$

Proof in the appendix, page 84

An unexpected and curious phenomenon unveiled by this theorem is that all the three rates are different. As a consequence, the answer to the question “what is the largest possible number of outliers, $o_d^*(n, s)$, that does not impact the minimax rate of estimation of θ^* ?” crucially depends on the considered distance d . Taking into account the relation $\varepsilon = o/n$, we get

$$o_{TV}^*(n, s) \asymp (ns)^{1/2}, \quad o_H^*(n, s) \asymp s, \quad o_{\mathbb{L}_2}^*(n, s) \asymp n^{1/2}.$$

Furthermore, all the claims concerning the total variation distance, in the considered model, yield corresponding claims for the Wasserstein distances W_q , for every $q \geq 1$. Indeed, one can see an element $\theta \in \Delta^{k-1}$ as the probability distribution of a random vector \mathbf{X} taking values in the finite set $\mathcal{A} = \{e_1, \dots, e_k\}$ of vectors of the canonical basis of \mathbb{R}^k . Since these vectors satisfy $\|e_j - e_{j'}\|_2^2 = 2\mathbb{1}(j \neq j')$, we have

$$\begin{aligned} W_q^q(\theta, \theta') &= \inf_{\Gamma} \mathbf{E}_{(\mathbf{X}, \mathbf{X}') \sim \Gamma} [\|\mathbf{X} - \mathbf{X}'\|_2^q] \\ &= \inf_{\Gamma} 2^{q/2} \mathbf{P}(\mathbf{X} \neq \mathbf{X}') = 2^{q/2} \|\theta - \theta'\|_{TV}, \end{aligned} \tag{2.6}$$

where the *inf* is over all joint distributions Γ on $\mathcal{A} \times \mathcal{A}$ having marginal distributions θ and θ' . This implies that

$$R_{W_q}^{AC}(n, \varepsilon, \Delta_s^{k-1}) \leq \sqrt{2} \{(s/n)^{1/2} + 2\varepsilon\}^{1/q}, \quad \forall q \geq 1. \tag{2.7}$$

In addition, since the \mathbb{L}_2 norm is an upper bound on the \mathbb{L}_∞ -norm, we have $R_{\mathbb{L}_\infty}^{AC}(n, \varepsilon, \Delta^{k-1}) \leq (1/n)^{1/2} + \sqrt{2} \varepsilon$. Thus, we have obtained upper bounds on the risk of the sample mean for all commonly used distances on the space of probability measures.

2.4.2 Lower bounds on the minimax risk

A natural question, answered in the next theorem, is how tight are the upper bounds obtained in the last theorem. More importantly, one can wonder whether there is an estimator that has a worst-case risk of smaller order than that of the sample mean.

Theorem 8. *There are universal constants $c > 0$ and n_0 , such that for any integers $k \geq 3$,*

$s \leq k \wedge n, n \geq n_0$ and for any $\varepsilon \in [0, 1]$, we have

$$\begin{aligned}\inf_{\bar{\theta}_n} R_{TV}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon\}, \\ \inf_{\bar{\theta}_n} R_H^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon^{1/2}\}, \\ \inf_{\bar{\theta}_n} R_{\mathbb{L}^2}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(1/n)^{1/2} + \varepsilon\},\end{aligned}$$

where $\inf_{\bar{\theta}_n}$ stands for the infimum over all measurable functions $\bar{\theta}_n$ from $(\Delta^{k-1})^n$ to Δ^{k-1} .

Proof in the appendix, page 86

The main consequence of this theorem is that whatever the contamination model is (among those described in Section 2.2), the rates obtained for the MLE in Theorem 7 are minimax optimal. Indeed, Theorem 8 yields this claim for Huber's contamination. For Huber's deterministic contamination and the TV-distance, on the one hand, we have

$$\begin{aligned}R_{TV}^{HDC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq R_{TV}^{HDC}(n, 0, \Delta_s^{k-1}, \bar{\theta}_n) \\ &\stackrel{(1)}{=} R_{TV}^{HC}(n, 0, \Delta_s^{k-1}, \bar{\theta}_n) \stackrel{(2)}{\geq} c(s/n)^{1/2},\end{aligned}$$

where (1) uses the fact that for $\varepsilon = 0$ all the sets $\mathcal{M}_n^\square(\varepsilon, \theta^*)$ are equal, while (2) follows from the last theorem. On the other hand, in view of Proposition 3, for $\varepsilon \geq (6/n) \log(8n/c)$ (implying that $2e^{-n\varepsilon/6} \leq (c/4)\varepsilon$),

$$\begin{aligned}R_{TV}^{HDC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq R_{TV}^{HC}(n, \varepsilon/2, \Delta_s^{k-1}, \bar{\theta}_n) - 2e^{-n\varepsilon/6} \\ &\geq (c/4)\{(s/n)^{1/2} + \varepsilon\}.\end{aligned}$$

Combining these two inequalities, for $n \geq (10 + 2 \log(1/c))^2$, we get

$$R_{TV}^{HDC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) \geq (c/4)\{(s/n)^{1/2} + \varepsilon\}$$

for every $k \geq 1$ and every $\varepsilon \in [0, 1]$. The same argument can be used to show that all the inequalities in Theorem 8 are valid for Huber's deterministic contamination model as well. Since the inclusions $\mathcal{M}_n^{HDC}(\varepsilon, \theta^*) \subset \mathcal{M}_n^{OC}(\varepsilon, \theta^*) \cap \mathcal{M}_n^{PC}(\varepsilon, \theta^*) \subset \mathcal{M}_n^{AC}(\varepsilon, \theta^*)$ hold true, we conclude that the lower bounds obtained for HC remain valid for all the other contamination models and are minimax optimal.

The main tool in the proof of Theorem 8 is the following result (Chen et al., 2018, Theorem 5.1). There is a universal constant $c_1 > 0$ such that

$$\inf_{\bar{\theta}_n} \sup_{\mathcal{M}_n^{HC}(\varepsilon, \Delta)} P(d(\bar{\theta}_n, \theta^*) \geq w_d(\varepsilon, \Delta)) \geq c_1, \quad \forall \varepsilon \in [0, 1],$$

where $w_d(\varepsilon, \Delta)$ is the modulus of continuity defined by $w_d(\varepsilon, \Delta) = \sup\{d(\theta, \theta') : d_{TV}(\theta, \theta') \leq \varepsilon\}$.

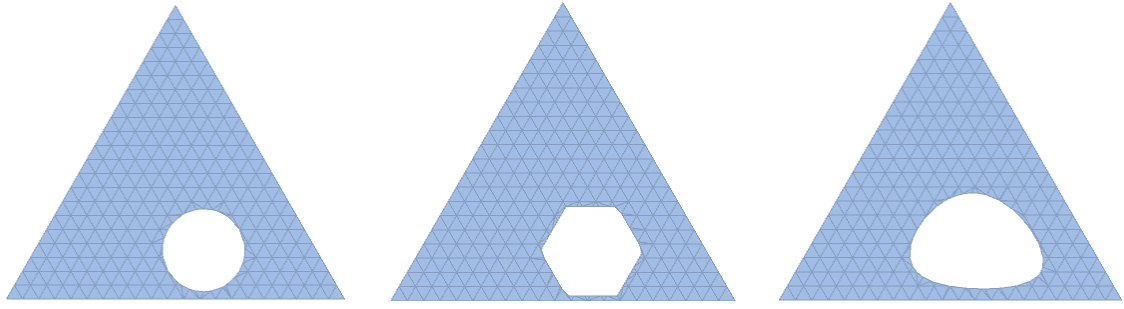


Figure 2.2: The shape of confidence sets (white regions) for the distances \mathbb{L}^2 (left), TV (center), and Hellinger (right) when the sample mean is $(1/3, 1/2, 1/6)$.

$\varepsilon/(1 - \varepsilon)\}$. Choosing θ and θ' to differ on only two coordinates, one can check that, for any $\varepsilon \leq 1/2$, $w_{\text{TV}}(\varepsilon, \Delta_s^{k-1}) \geq \varepsilon$, $w_{\text{H}}(\varepsilon, \Delta_s^{k-1}) \geq \varepsilon^{1/2}$ and $w_{\mathbb{L}^2}(\varepsilon, \Delta_s^{k-1}) \geq \sqrt{2}\varepsilon$. Combining with the lower bounds in the non-contaminated setting, this result yields the claims of Theorem 8. In addition, (2.6) combined with the results of this section implies that the rate in (2.7) is minimax optimal.

2.4.3 Confidence regions

In previous sections, we established bounds for the expected value of estimation error. The aim of this section is to present bounds on estimation error of the sample mean holding with high probability. This also leads to constructing confidence regions for the parameter vector θ^* . To this end, the contamination rate ε and the sparsity s are assumed to be known. It is an interesting open question whether one can construct optimally shrinking confidence regions for unknown ε and s .

Theorem 9. *Let $\delta \in (0, 1)$ be the tolerance level. If $\theta^* \in \Delta_s^{k-1}$, then under any contamination model, the regions of Δ_s^{k-1} defined by each of the following inequalities*

$$\begin{aligned} d_{\mathbb{L}^2}(\bar{\mathbf{X}}_n, \theta) &\leq (1/n)^{1/2} + \sqrt{2}\varepsilon + (\log(1/\delta)/n)^{1/2}, \\ d_{\text{TV}}(\bar{\mathbf{X}}_n, \theta) &\leq (s/n)^{1/2} + 2\varepsilon + (2\log(1/\delta)/n)^{1/2}, \\ d_{\text{H}}(\bar{\mathbf{X}}_n, \theta) &\leq \sqrt{5}((s/n)\log(2s/\delta))^{1/2} + \varepsilon^{1/2} + ((1/2n)\log(2/\delta))^{1/2}, \end{aligned}$$

contain θ^ with probability at least $1 - \delta$.*

Proof in the appendix, page 88

To illustrate the shapes of these confidence regions, we depicted them in Figure 2.2 for a three dimensional example, projected onto the plane containing the probability simplex. The sample mean in this example is equal to $(1/3, 1/2, 1/6)$.

While the estimator, $\bar{\mathbf{X}}_n$, used in the construction of confidence regions is adaptive both to the sparsity s and to the contamination rate ε , the confidence region itself is not adaptive.

Indeed, the radius of the confidence region depends on s and/or on ε . Constructing adaptive confidence regions and adaptive tests is an important question in statistics, we refer to (Cai and Low, 2006; Hoffmann and Nickl, 2011) for some precise results.

In the setting of shape constrained regression, Bellec (2016) proposed a method of constructing adaptive confidence regions that could be of interest for our setting as well. In particular, if there is no contamination, *i.e.*, $\varepsilon = 0$, we can define $\hat{s} = \text{Card}(j : (\bar{X}_n)_j \neq 0)$, the observed sparsity. It is clear that $\hat{s} \leq s$. The analog of Bellec's method in our setting would consist in replacing s by \hat{s} in the radius of the ball given by Theorem 9. Of course, this does not inflate the ball. However, we did not manage to adapt the argument in (Bellec, 2016, Theorem 4.1) for proving that the modified region (with \hat{s} instead of s) has still the required coverage property and, therefore, is adaptive to s . The main difficulty, at a very high-level, comes from the fact that the distance we consider, d_{TV} , is not induced by an inner product. Thus, constructing adaptive confidence regions even when there is no contamination is an open question.

2.5 Instance based bounds

When the dimension k is not finite, we can provide bounds which depend on the reference distribution θ^* or the sample mean \bar{X}_n . We restrict X_i 's to take value in $\{e_1, e_2, \dots\}$, and assume e_j occurs with probability θ_j^* . We define $\alpha_n(\theta) := 2 \sum_{\theta_j < 1/n} \theta_j$ and $\beta_n(\theta) := \frac{1}{\sqrt{n}} \sum_{\theta_j \geq 1/n} \sqrt{\theta_j}$. Using the results of Berend and Kontorovich (2013) the following upper and lower bounds are obtained for the error of the sample mean under the TV distance and adversarial model with ε -contamination.

Proposition 5. *Suppose X_i 's take value in $\{e_1, e_2, \dots\}$, and for $j \in \mathbb{N}$, e_j occurs with probability θ_j^* . For every n and for every $\varepsilon \in [0, 1]$, the sample mean \bar{X}_n satisfies*

$$\frac{\alpha_n(\theta_j^*) + \beta_n(\theta_j^*)}{4} - \frac{1}{4\sqrt{n}} - 2\varepsilon \leq \mathbf{E}d_{TV}(\bar{X}_n, \theta^*) \leq \alpha_n(\theta_j^*) + \beta_n(\theta_j^*) + 2\varepsilon.$$

Proof in the appendix, page 90

These bounds need the knowledge of the reference distribution. The next theorem represent bounds based on the sample mean.

Theorem 10. *Suppose X_i 's take value in $\{e_1, e_2, \dots\}$, and for $j \in \mathbb{N}$, e_j occurs with probability θ_j^* . For every n and for every $\varepsilon \in [0, 1]$, the sample mean \bar{X}_n satisfies*

$$d_{TV}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \|\bar{X}_n^{1/2}\|_1 + 2\varepsilon + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

with probability at least $1 - \delta$, where $\delta \in (0, 1)$. We also have

$$\mathbf{E}d_{\text{TV}}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \mathbf{E} \|\bar{X}_n^{1/2}\|_1 + 2\varepsilon.$$

Proof in the appendix, page 91

2.6 Illustration on a numerical example

We provide some numerical experiments which illustrate theoretical results of Section 2.4. The data set is the collection of 38 books written by Alexandre Dumas (1802-1870) and 38 books written by Emile Zola (1840-1902)⁷. To each author, we assign a parameter vector corresponding to the distribution of the number of words contained in the sentences used in the author's books. To be more clear, a sentence containing l words is represented by vector e_l , and if the parameter vector of an author is $(\theta_1, \dots, \theta_k)$, it means that a sentence used by the author is of size $l \in \{1, \dots, k\}$ with probability θ_l . We carried out synthetic experiments in which the reference parameter to estimate is the probability vector of Dumas, while the distribution of outliers is determined by the probability vector of Zola. Ground truths for these parameters are computed from the aforementioned large corpus of their works. Only the dense case $s = k$ were considered. For various values of ε and n , a contaminated sample was generated by randomly choosing n sentences either from Dumas' works (with probability $1 - \varepsilon$) or from Zola's works (with probability ε). The sample mean was computed for this corrupted sample, and the error with respect to Dumas' parameter vector was measured by the three distances TV, \mathbb{L}^2 and Hellinger. This experiment was repeated 10^4 times for each special setting to obtain information on error's distribution. Furthermore, by grouping nearby outcomes we created samples of different dimensions for illustrating the behavior of the error as a function of k .

The error of \bar{X}_n as a function of the sample size n , dimension k , and contamination rate ε is plotted in Figures 2.3 and 2.4. These plots are conform to the theoretical results. Indeed, the first plot in Figure 2.3 shows that the errors for the three distances is decreasing w.r.t. n . Furthermore, we see that up to some level of n this decay is of order $n^{-1/2}$. The second plot in Figure 2.3 confirms that the risk grows linearly in k for the TV and Hellinger distances, while it is constant for the \mathbb{L}^2 error.

Left panel of Figure 2.4 suggests that the error grows linearly in terms of contamination rate. This is conform to our results for the TV and \mathbb{L}^2 errors. But it might seem that there is a disagreement with the result for the Hellinger distance, for which the risk is shown to increase at the rate $\varepsilon^{1/2}$ and not ε . This is explained by the fact that the rate $\varepsilon^{1/2}$ corresponds to the worst-case risk, whereas here, the setting under experiment does not necessarily represent the worst case. When the parameter vectors of the reference and contamination distributions,

⁷The works of both authors are available from <https://www.gutenberg.org/>

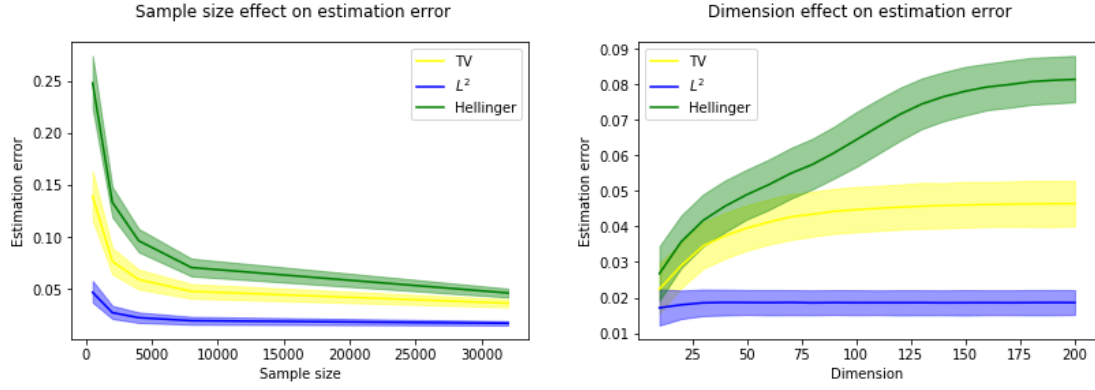


Figure 2.3: Estimation error of \bar{X}_n measured by total variation, Hellinger, and \mathbb{L}^2 distances as a function of (left panel) number of observations with contamination rate 0.2 and dimension 10^2 and (right panel) dimension with contamination rate 0.2 and 10^4 samples. The interval between 5th and 95th quantiles of the error, obtained from 10^4 repetitions, is also depicted for every graph.

respectively, are e_i and e_j with $i \neq j$ (i.e., when these two distributions are at the largest possible distance, which we call an extreme case), the graph of the error as a function of ε (right panel of Figure 2.4) is similar to that of square-root function.

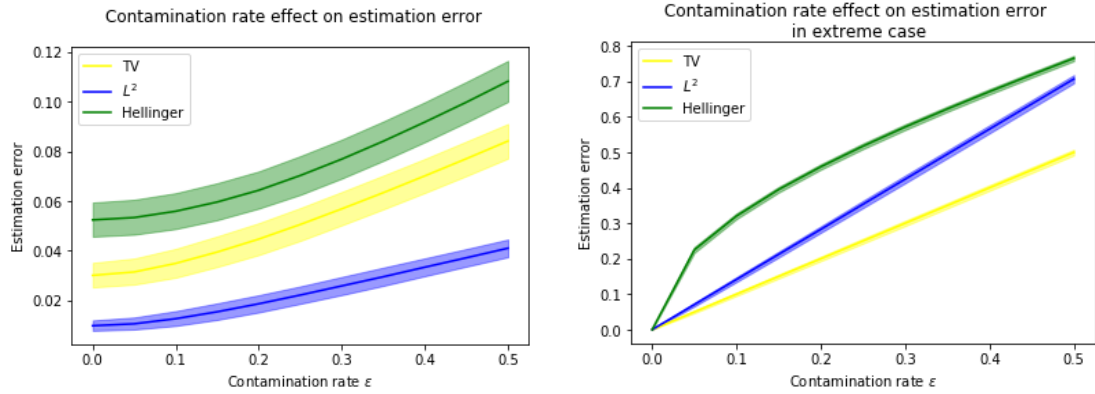


Figure 2.4: The estimation error of \bar{X}_n measured by total variation, Hellinger, and \mathbb{L}^2 distances in terms of the contamination rate (with dimension 10^2 and 10^4 samples). At right, we plotted the error with respect to the contamination rate for an extreme case, where the reference and contamination distributions have the largest distance. The interval between 5th and 95th quantiles of the error, obtained from 10^4 trials, is also depicted.

2.7 Summary and conclusion

We have analyzed the problem of robust estimation of the mean of a random vector belonging to the probability simplex. Four measures of accuracy have been considered: total variation, Hellinger, Euclidean and Wasserstein distances. In each case, we have established the min-imax rates of the expected error of estimation under the sparsity assumption. In addition,

confidence regions shrinking at the minimax rate have been proposed.

An intriguing observation is that the choice of the distance has much stronger impact on the rate than the nature of contamination. Indeed, while the rates for the aforementioned distances are all different, the rate corresponding to one particular distance is not sensitive to the nature of outliers (ranging from Huber's contamination to the adversarial one). While the rate obtained for the TV-distance coincides with the previously known rate of robustly estimating a Gaussian mean, the rates we have established for the Hellinger and for the Euclidean distances appear to be new. Interestingly, when the error is measured by the Euclidean distance, the quality of estimation does not get deteriorated with increasing dimension.

Chapter 3

Robust Estimation of Gaussian Mean

We study the problem of robust estimation of the mean vector of a sub-Gaussian distribution. We introduce an estimator based on spectral dimension reduction (SDR) and establish a finite sample upper bound on its error that is minimax-optimal up to a logarithmic factor. Furthermore, we prove that the breakdown point of the SDR estimator is equal to $1/2$, the highest possible value of the breakdown point. In addition, the SDR estimator is equivariant by similarity transforms and has low computational complexity. More precisely, in the case of n vectors of dimension p —at most εn out of which are adversarially corrupted—the SDR estimator has a squared error of order $(r_{\Sigma}/n + \varepsilon^2 \log(1/\varepsilon)) \log p$ and a running time of order $p^3 + np^2$. Here, $r_{\Sigma} \leq p$ is the effective rank of the covariance matrix of the reference distribution. Another advantage of the SDR estimator is that it does not require knowledge of the contamination rate and does not involve sample splitting. We also investigate extensions of the proposed algorithm and of the obtained results in the case of (partially) unknown covariance matrix.

3.1	Introduction	54
3.2	Adversarially corrupted sub-Gaussian model and spectral dimension reduction	57
3.2.1	Choice of the dimension reduction regime	59
3.2.2	Choice of the threshold	60
3.3	Assessing the error of the SDR estimator	60
3.4	The case of unknown covariance matrix	63
3.5	Numerical experiments	64
3.5.1	Implementation details	64
3.5.2	Experimental setup	65
3.5.3	Statistical accuracy	65
3.5.4	Computational efficiency	66
3.5.5	Breakdown point	66
3.6	Summary, related work and conclusion	68

3.1 Introduction

Robust estimation of a finite-dimensional parameter is a classical problem in statistics. The broad goal of robust estimation is to design statistical procedures that are not very sensitive to small changes in data or to small departures from the modeling assumptions. A typical example, extensively studied in the literature, and considered in the present work, is when the data set contains outliers.

The literature on robustness to outliers in parametric estimation is very rich; it would be impossible to review here all the important contributions. For an in-depth exposition of by now classical results and approaches, such as the influence function, the breakdown point and the efficiency, we refer to the books (Huber and Ronchetti, 2011; Maronna et al., 2006; Rousseeuw et al., 2011). In their vast majority, these well-established approaches treated the dimension of the parameter as a fixed and small constant. This simple setting was convenient for mathematical analysis and for computational purposes, but somewhat disconnected from many practical situations. Furthermore, it was hiding some fascinating phenomena that emerge only when the dimension is considered as a parameter that might be large, in the same way as the sample size.

More recently, (Chen et al., 2018) considered the problem of estimating the mean and the covariance matrix of a Gaussian distribution in the high-dimensional setting. The authors namely uncovered a new phenomenon: under the Huber contamination, the componentwise median is not minimax-rate optimal whereas the Tukey median is. More precisely, if a p -dimensional mean vector is to be estimated from n independent vectors drawn from the mixture distribution $(1 - \varepsilon)\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \varepsilon\mathbf{Q}$ (where $\varepsilon \in (0, 1/2)$ is the rate of contamination and \mathbf{Q} is the unknown distribution of outliers), then the mean squared error of the componentwise median is of order $p/n + p\varepsilon^2$ while that of Tukey's median is of order $p/n + \varepsilon^2$. This extra factor p in front of ε^2 has been proven in (Lai et al., 2016) to be present in the error of another widely used robust estimator of the mean, the geometric median (Minsker, 2015). Thus, as long as only statistical properties of the estimators are considered, Tukey's median is thus superior to its competitors, the componentwise and the geometric medians. However, the componentwise and the geometric medians are better than the Tukey's median in terms of the breakdown point: their breakdown point is equal to $1/2$ (Lopuhaa and Rousseeuw, 1991) whereas that of Tukey's median is $1/3$ (Donoho and Gasko, 1992). This is one of the appealing phenomena taking place in the high dimensional setting.

Another specificity of the high dimensional setting uncovered by (Chen et al., 2018) was the lack of computational tractability of the estimators that are statistically optimal. Indeed, Tukey's median is computationally intractable, but minimax-rate optimal, whereas the componentwise and the geometric medians are computationally tractable but statistically suboptimal. This observation led to the development of a number of computationally tractable estimators having an error with a better dependence on dimension than that of Tukey's median (Dalalyan and Minasyan, 2020; Diakonikolas et al., 2016a, 2017, 2018a; Dong et al., 2019; Lai et al., 2016).

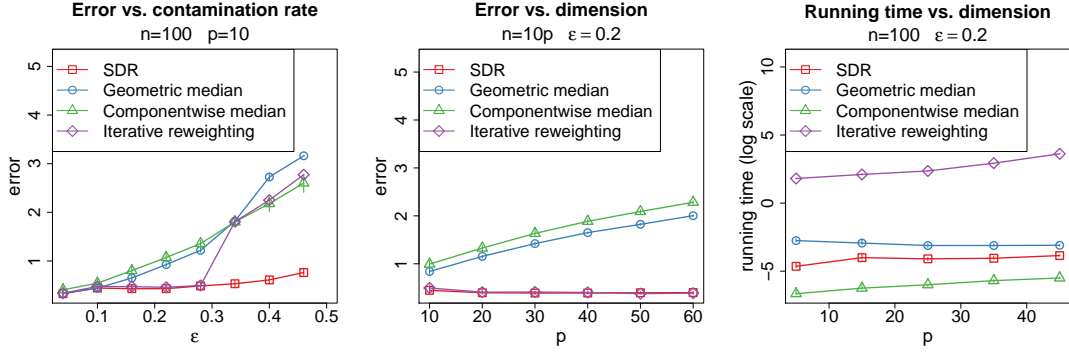


Figure 3.1: Plots that help to visually compare four robust estimators: SDR (our estimator), geometric median (GM) given by (3.1), componentwise median (CM), iteratively reweighted mean (IRM) of (Dalalyan and Minasyan, 2020). The first two plots show that SDR is as accurate as IRM for small ϵ , with SDR outperforming IRM for ϵ close to $1/2$. IRM and SDR are naturally much more accurate than GM and CM. The last plot shows that the running time of SDR is comparable to that of GM and is much smaller than that of IRM. More details on these experiments are provided in Section 3.5.

In particular, most estimators introduced in these papers allow to conciliate computational tractability (*i.e.*, are computable in time polynomial in $n, p, 1/\epsilon$) and statistical optimality up to logarithmic factors.

The goal of this paper is to make a step forward by designing an estimator which is not only nearly rate optimal and computationally tractable, but also has a breakdown point equal to $1/2$, which is the highest possible value of the breakdown point. To construct the estimator, termed iterative spectral dimension reduction or SDR, we combine and suitably adapt ideas from (Lai et al., 2016) and (Diakonikolas et al., 2017). The main underlying observation is that if we remove some clear outliers and restrict our attention to the subspace spanned by the eigenvectors of the sample covariance matrix corresponding to small eigenvalues, then the sample mean of the projected data points is a rate-optimal estimator. This allows us to iteratively reduce the dimension and eventually to estimate the remaining low-dimensional component of the mean by a standard robust estimator such as the componentwise median or the trimmed mean, see Algorithm 2.

The main contributions of this paper are methodological and theoretical. The SDR estimator, thoroughly defined in Section 3.2, is a fast and accurate method for robustly estimating the mean of a set of points. It depends on one tuning parameter, the threshold used for identifying and removing clear outliers, and on the dimension reduction regime. Our theoretical considerations provide some recommendations for their choices and our numerical experiments reported in Section 3.5 confirm the relevance of these choices. Importantly, the SDR estimator does not require as input the rate of contamination ϵ but only an upper bound on ϵ . As for theoretical contributions of this paper, we state in Section 3.3 an upper bound on the error of the SDR estimator, showing that it is nearly minimax-rate optimal and has a breakdown point equal to $1/2$. This is done in the general case of a sub-Gaussian distribution with heterogeneous covariance matrix contaminated by adversarial noise. In Section 3.4, we further investigate the error of the SDR estimator in the case where only an approximation to the

	Comput. tractable	Breakdown point	known ε	Squared error rate	known Σ or $\Sigma \propto \mathbf{I}$
Gaussian distribution					
C./G. Median	yes	0.5	no	$\mathbf{r}_\Sigma/n + \varepsilon^2 p$	no
Tukey's Median	no	0.33	no	$\mathbf{r}_\Sigma/n + \varepsilon^2$	no
Agnostic Mean	yes	—	yes	$(p/n) \log^3 p + \varepsilon^2 \log p$	yes
Gaussian and sub-Gaussian distribution					
Iter. Rew. Mean	yes	0.28	yes	$(\mathbf{r}_\Sigma/n) + \varepsilon^2 \log(1/\varepsilon)$	yes
Iterative Filtering	yes	—	yes	$(p/n) \log^a p + \varepsilon^2 \log(1/\varepsilon)$	yes
SDR (this work)	yes	0.5	no	$(\mathbf{r}_\Sigma/n + \varepsilon^2 \log(1/\varepsilon)) \log p$	yes

Table 3.1: Properties of various robust estimators. Agnostic mean, iteratively reweighted mean and iterative filtering are the estimators studied in (Lai et al., 2016), (Dalalyan and Minasyan, 2020) and (Diakonikolas et al., 2017), respectively. The error rates reported for Tukey’s median, componentwise median, geometric median and the agnostic mean have been proved for non-adversarial contamination. The squared error rate is provided in the case of a covariance matrix satisfying $\|\Sigma\|_{\text{op}} = 1$.

covariance matrix is available.

The papers that are the closest to the present one are (Lai et al., 2016), (Diakonikolas et al., 2017) and (Dalalyan and Minasyan, 2020). The spectral dimension reduction scheme was proposed by (Lai et al., 2016) along with an initial sample splitting step ensuring the independence of the estimators over different subspaces. In the case of spherical Gaussian distribution contaminated by non-adversarial outliers, the paper states that the proposed estimator has a squared error at most of order $p \log^2 p \log(p/\varepsilon)/n + \varepsilon^2 \log p$. Compared to this, our results are valid in the more general setting of sub-Gaussian distribution, with arbitrary covariance matrix and adversarial contamination. In addition, our estimator does not rely on sample splitting and, therefore, has a risk with a better dependence on p . As compared to the filtering method of (Diakonikolas et al., 2017), our estimator has the advantage of being independent of ε and our error bound is valid for every covariance matrix and every confidence level. On the down side, our error bound has an extra factor $\log p$ in front of ε^2 . We believe that this factor is an artifact of the proof, but we were unable to remove it. Finally, compared to the iteratively reweighted mean (Dalalyan and Minasyan, 2020), the SDR estimator studied in the present paper has a higher breakdown point, does not require the knowledge of ε and is much faster to compute. The advantages and shortcomings of these estimators are summarized in Table 3.1 and Figure 3.1.

Notation. For any pair of integers k and d such that $1 \leq k \leq d$, we denote by \mathcal{V}_k^d the set of all k -dimensional linear subspaces V of \mathbb{R}^d . For $V \in \mathcal{V}_k^d$, we write $k = \dim(V)$ and denote by P_V the orthogonal projection matrix onto V . \mathbb{S}^{d-1} stands for the unit sphere in \mathbb{R}^d . For a $d \times d$ symmetric matrix \mathbf{M} , we denote by $\lambda_1(\mathbf{M}), \dots, \lambda_d(\mathbf{M})$ its eigenvalues sorted in increasing order, and use the notation $\lambda_{\min}(\mathbf{M}) = \lambda_1(\mathbf{M})$, $\lambda_{\max}(\mathbf{M}) = \lambda_d(\mathbf{M})$, $\|\mathbf{M}\|_{\text{op}} =$

$\max(|\lambda_{\min}(\mathbf{M})|, |\lambda_{\max}(\mathbf{M})|)$, $\text{Tr}(\mathbf{M}) = (\lambda_1 + \dots + \lambda_d)(\mathbf{M})$ and $\mathbf{r}_M = \text{Tr}(\mathbf{M})/\|\mathbf{M}\|_{\text{op}}$. For any integer $n > 0$, we set $[n] = \{1, \dots, n\}$. We will denote by $\mathcal{O} \subset [n]$ the subscripts of the outliers and by $\mathcal{I} = [n] \setminus \mathcal{O}$ the subscripts of inliers. We also use notation $\log_+(x) = \max\{0, \log(x)\}$. For a matrix \mathbf{M} , we denote by $\sigma_{\min}(\mathbf{M})$ and $\sigma_{\max}(\mathbf{M})$ respectively its smallest and largest singular values.

Algorithm 2 $\text{SDR}(\mathbf{X}_1, \dots, \mathbf{X}_n; \Sigma, t)$

```

1: let  $p$  be the dimension of  $\mathbf{X}_1$ 
2: let  $\hat{\mu}^{\text{GM}}$  be the geometric median of  $\mathbf{X}_1, \dots, \mathbf{X}_n$ 
3: let  $\mathcal{S} \leftarrow \{i : \|\mathbf{X}_i - \hat{\mu}^{\text{GM}}\| \leq t\sqrt{p}\}$ 
4: let  $\bar{\mathbf{X}}_{\mathcal{S}}$  be the sample mean of the filtered sample  $\{\mathbf{X}_i : i \in \mathcal{S}\}$ 
5: let  $\hat{\Sigma}_{\mathcal{S}}$  be the covariance matrix of the filtered sample  $\{\mathbf{X}_i : i \in \mathcal{S}\}$ 
6: if  $p > 1$  then
7:   let  $V$  be the span of the top  $\lceil p/e \rceil$  principal components of  $\hat{\Sigma}_{\mathcal{S}} - \Sigma$ 
8:   let  $P_V$  be the orth. projection onto  $V$ 
9:   let  $P_{V^\perp}$  be the orth. projection onto the orth. complement of  $V$ 
10:  let  $\hat{\mu} \leftarrow P_{V^\perp} \bar{\mathbf{X}}_{\mathcal{S}} + \text{SDR}(P_V \mathbf{X}_1, \dots, P_V \mathbf{X}_n; P_V \Sigma P_V, t)$ 
11: else
12:   let  $\hat{\mu} \leftarrow \hat{\mu}^{\text{GM}}$ 
13: end if
14: return  $\hat{\mu}$ 

```

3.2 Adversarially corrupted sub-Gaussian model and spectral dimension reduction

We assume that a set $\mathbf{X}_1, \dots, \mathbf{X}_n$ of n data points drawn from a distribution P_n is given. This set is assumed to contain at least $n - \lfloor n\varepsilon \rfloor$ inliers, the remaining points being outliers. All the points lie in the p -dimensional Euclidean space and the inliers are independently drawn from a reference distribution, assumed to be sub-Gaussian with mean $\mu^* \in \mathbb{R}^p$ and covariance matrix Σ . To state the assumptions imposed on the observations in a more precise way, let us recall that the random vector ζ is said to be sub-Gaussian with zero mean and identity covariance matrix, if $\mathbb{E}[\zeta] = 0$, $\mathbb{E}[\zeta \zeta^\top] = \mathbf{I}_p$ and for some $\varsigma > 0$, we have

$$\mathbb{E}[e^{v^\top \zeta}] \leq \exp\{\varsigma \|v\|^2/2\}, \quad \forall v \in \mathbb{R}^p.$$

The parameter ς is commonly called the variance proxy and the writing $\zeta \sim \text{SG}_p(\varsigma)$ is used.

Definition 2. We say that the data generating distribution P_n is an adversarially corrupted sub-Gaussian distribution with mean μ^* , covariance matrix Σ , variance proxy ς and contamination rate ε , if there is a probability space on which we can define a sequence of random vectors $(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)$ such that

1. $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ are independent and $\Sigma^{-1/2}(\mathbf{Y}_i - \mu^*) \sim \text{SG}_p(\varsigma)$ for every $i \in [n]$.

2. the cardinality of $\mathcal{O} = \{i \in [n] : \mathbf{Y}_i \neq \mathbf{X}_i\}$ is at most equal to $n\varepsilon$.

3. the distribution of $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ is P_n .

We write then¹ $P_n \in \text{SGAC}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}, \varepsilon, \varepsilon)$. In the particular case where all \mathbf{Y}_i are Gaussian, we will write $P_n \in \text{GAC}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}, \varepsilon)$.

For an overview of various kind of contamination models we refer the interested reader to (Bateni and Dalalyan, 2020). The adversarial contamination considered throughout this work is perhaps the most general one considered in the literature as the elements of the set \mathcal{O} —called outliers—may be chosen using $\boldsymbol{\mu}^*, \boldsymbol{\Sigma}, \varepsilon$ but also $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ by an omniscient adversary. Note that in this setting, even the set \mathcal{O} is random and depends on $\mathbf{Y}_1, \dots, \mathbf{Y}_n$. Therefore, the inliers $\{\mathbf{X}_i : i \in \mathcal{I}\}$ cannot be considered as independent random variables. The problem studied in this work consists in estimating the mean $\boldsymbol{\mu}^*$ of the reference distribution from the adversarially corrupted observations $\mathbf{X}_1, \dots, \mathbf{X}_n$.

The estimator we analyze in this work is termed iterative spectral dimension reduction and denoted by $\hat{\boldsymbol{\mu}}^{\text{SDR}}$. It is closely related to the agnostic mean (Lai et al., 2016) and to iterative filtering (Diakonikolas et al., 2017) estimators. We will prove that SDR enjoys most of desired properties in the setting of robust estimation of the sub-Gaussian mean.

The parameters given as input to the iterative spectral dimension reduction algorithm are a strictly decreasing sequence of positive integers p_0, \dots, p_L such that $p_0 = p$ and a positive threshold $t > 0$. We recall that the geometric median is defined by

$$\hat{\boldsymbol{\mu}}^{\text{GM}} \in \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^p} \sum_{i=1}^n \|\mathbf{X}_i - \boldsymbol{\mu}\|_2. \quad (3.1)$$

The algorithm for computing the SDR estimator reads as follows.

1. Start by setting $\mathbf{V}_0 = \mathbf{I}_p$.
2. For $\ell = 0, \dots, L - 1$ do
 - (a) Define $\bar{\boldsymbol{\mu}}^{(\ell)} \in \mathbb{R}^{p_\ell}$ as the geometric median of $\{\mathbf{V}_\ell^\top \mathbf{X}_i : i \in [n]\}$.
 - (b) Define the set $\mathcal{S}^{(\ell)} = \{i \in [n] : \|\mathbf{V}_\ell^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(\ell)}\|_2 \leq t\sqrt{p_\ell}\}$ of filtered data points.
 - (c) Let $\bar{\mathbf{X}}^{(\ell)}$ and $\hat{\boldsymbol{\Sigma}}^{(\ell)}$ be the mean vector and the covariance matrix of the filtered sample $\{\mathbf{X}_i : i \in \mathcal{S}^{(\ell)}\}$, that is

$$\bar{\mathbf{X}}^{(\ell)} = \frac{1}{|\mathcal{S}^{(\ell)}|} \sum_{i \in \mathcal{S}^{(\ell)}} \mathbf{X}_i, \quad \hat{\boldsymbol{\Sigma}}^{(\ell)} = \frac{1}{|\mathcal{S}^{(\ell)}|} \sum_{i \in \mathcal{S}^{(\ell)}} (\mathbf{X}_i - \bar{\mathbf{X}})^{\otimes 2}.$$

- (d) Set $\hat{\boldsymbol{\mu}}^{(\ell)} = \mathbf{V}_\ell \mathbf{U}_\ell^\top \mathbf{U}_\ell \mathbf{V}_\ell^\top \bar{\mathbf{X}}^{(\ell)}$, where \mathbf{U}_ℓ is a $(p_\ell - p_{\ell+1}) \times p_\ell$ orthogonal matrix the rows of which are the eigenvectors of $\mathbf{V}_\ell^\top (\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}) \mathbf{V}_\ell$ corresponding to its $(p_\ell - p_{\ell+1})$ smallest eigenvalues.

¹SGAC stands for sub-Gaussian with adversarial contamination.

(e) Set $\mathbf{V}_{\ell+1} = \mathbf{V}_\ell(\mathbf{U}_\ell^\perp)^\top \in \mathbb{R}^{p \times p_{\ell+1}}$, where \mathbf{U}_ℓ^\perp is a $p_{\ell+1} \times p_\ell$ orthogonal matrix orthogonal to \mathbf{U}_ℓ , that is $\mathbf{U}_\ell^\perp \mathbf{U}_\ell^\top = \mathbf{0}$.

3. Define $\bar{\boldsymbol{\mu}}^{(L)}$ as the geometric median of $\mathbf{V}_L^\top \mathbf{X}_i$ for $i = 1, \dots, n$ and set $\mathcal{S}^{(L)} = \{i \in [n] : \|\mathbf{V}_L^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(L)}\|_2 \leq t\sqrt{p_L}\}$.
4. Define $\hat{\boldsymbol{\mu}}^{(L)} = \mathbf{V}_L \mathbf{V}_L^\top \bar{\mathbf{X}}^{(L)}$, the average of filtered and projected vectors.
5. Return $\hat{\boldsymbol{\mu}}^{\text{SDR}} = \hat{\boldsymbol{\mu}}^{(0)} + \hat{\boldsymbol{\mu}}^{(1)} + \dots + \hat{\boldsymbol{\mu}}^{(L)}$.

The steps described above can be summarised as follows. At each iteration $\ell < L$, we start by determining a filtered subsample $\mathcal{S}^{(\ell)}$ and a “nearly-outlier-orthogonal” subspace $\mathcal{U}_\ell = \text{Im}(\mathbf{V}_\ell \mathbf{U}_\ell^\top)$ of \mathbb{R}^p of dimension $p_\ell - p_{\ell+1}$. We define the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto \mathcal{U}_ℓ as the sample mean of the filtered and projected subsample, and we move to the next step for determining the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto the remaining part of the space. At the last iteration L , when the dimension is well reduced, the projection of $\hat{\boldsymbol{\mu}}^{\text{SDR}}$ onto the subspace \mathcal{U}_L is defined as the average of the filtered subsample projected onto \mathcal{U}_L . The subspaces \mathcal{U}_ℓ are two-by-two orthogonal and span the whole space \mathbb{R}^p . Each subspace is determined from the spectral decomposition of the covariance matrix of the data points projected onto $(\mathcal{U}_0 \oplus \dots \oplus \mathcal{U}_{\ell-1})^\perp$, after removing the points lying at an abnormally large distance from the geometric median.

3.2.1 Choice of the dimension reduction regime

The analysis of the error of the SDR estimator conducted in this work leads to an upper bound in which the sequence (p_0, \dots, p_L) is involved only through the expression

$$F(p_0, \dots, p_L) = \sum_{\ell=1}^L \frac{p_{\ell-1}}{p_\ell}.$$

Therefore, an appealing way of choosing this sequence is to minimize the function F under the constraint that the sequence is decreasing and $p_0 = p$ and $p_L = 1$. It follows from the inequality between the arithmetic and geometric means that $F(p_0, \dots, p_L) \geq Lp^{1/L}$. Furthermore, the equality is achieved² in the case when all the terms in the definition of F are equal, *i.e.*, when for some $c > 0$ we have $p_{\ell-1} = cp_\ell$ for every $\ell \in [L]$. Since $p_0 = p$ and $p_L = 1$, this yields $c = p^{1/L}$ or, equivalently, $L = \log p / \log c$. Using these relations, we find that the function F is lower bounded by $Lc = (c / \log c) \log p$. The last step is to find the minimum of the function $c \mapsto c / \log c$ over the interval $(1, \infty)$. One easily checks that this function has a unique minimum at $c = e$. All these considerations advocate for using the dimension reduction regime defined by

$$p_0 = p, \quad p_\ell = \lfloor p_{\ell-1}/e \rfloor + 1, \quad \ell \in [L], \quad p_L = 1, \quad (3.2)$$

²We relax here the assumption that all the entries p_ℓ are integers.

where $\lfloor x \rfloor$ is the largest integer strictly smaller than x . Such a definition of (p_ℓ) ensures that $p_{\ell-1}/p_\ell \leq e$ and that³ $L \leq 2 \log p$. In the rest of the paper, we assume that the sequence (p_ℓ) is chosen as in (3.2).

3.2.2 Choice of the threshold

The SDR procedure has one important tuning parameter: the threshold t used to discard clearly outlying data points. Let us introduce the auxiliary notation

$$\bar{r}_n = \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n}}, \quad \text{and} \quad \tau = \frac{1}{4} \bigwedge \frac{\bar{r}_n}{\sqrt{\log_+(2/\bar{r}_n)}}. \quad (3.3)$$

Note that \bar{r}_n is essentially the quantile, up to a universal constant factor, of order $1 - \delta$ of the distribution of $\|\bar{\mathbf{Y}}_n - \boldsymbol{\mu}^*\|_2$ where \mathbf{Y}_i 's are independently drawn from $\mathcal{N}_p(\boldsymbol{\mu}^*, \Sigma)$ with $\|\Sigma\|_{\text{op}} = 1$. Our theoretical results advocate for using the value $t = t_1 + t_2$, where

$$t_1 = \frac{2(1 + \bar{r}_n)}{1 - 2\varepsilon^*}, \quad t_2 = 1 + \frac{\bar{r}_n}{\sqrt{\tau}} + \sqrt{2 + \log(2/\tau)},$$

where $\varepsilon^* < 1/2$ is the largest value of the contamination rate that the algorithm may handle.

Let ξ_1, \dots, ξ_n be independent Gaussian with zero mean and covariance Σ . The expression of t_1 is obtained as an upper bound on the quantile of order $1 - \delta/2$ of the distribution of the random variable

$$T_1 = \sup_V \frac{2}{n(1 - 2\varepsilon) \dim(V)} \sum_{i=1}^n \|\mathbf{P}_V \xi_i\|_2,$$

see Lemma 2 and its proof for further details. Similarly, t_2 is defined so that the event

$$\sup_V \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \xi_i\|_2^2 > t_2^2 \dim(V)) \leq n\tau$$

has a probability at least $1 - \delta/2$. The related computations are deferred to Appendix B.1.3. Although we tried to get sharp values for these thresholds t_1 and t_2 , it is certainly possible to improve these values either by better mathematical arguments or by empirical considerations. Of course, smaller values of the thresholds t_1 and t_2 satisfying aforementioned conditions lead to an SDR estimator having smaller error.

3.3 Assessing the error of the SDR estimator

The iterative spectral dimension reduction estimator defined in previous sections has some desirable properties of a robust estimator that are easy to check. In particular, it is clearly

³To check this inequality, one can use the fact that $3 \leq p_{L-2} \leq pe^{2-L} + e/(e-1)$. This implies $L \leq 2 \log p$ for $p \geq 6$. For smaller values of p , the inequality can be checked by direct computations.

equivariant by translation, orthogonal linear transform and global scaling. Furthermore, the breakdown point of the estimator is equal to that of the geometric median, that is to $1/2$. This means that even if almost the half of data points are chosen to be infinitely large, the estimator will not “break down” in the sense of becoming infinitely large. However, the fact that the estimated value does not become infinitely large, it might be not very close to the true mean. The next theorem shows that this is not the case and that the error of the SDR estimator has a nearly rate-optimal behavior even when the contamination rate is close to $1/2$. The adverb “nearly” is used here to reflect the presence of the $\sqrt{\log p}$ factor in the error bound, which is not present in the minimax rate.

Theorem 11. *Let $\varepsilon^* \in (0, 1/2)$, and $\delta \in (0, 1/2)$. Define \bar{r}_n and τ as in (3.3). For every $\varepsilon \leq \varepsilon^*$, let $\hat{\mu}^{\text{SDR}}$ be the estimator returned by Algorithm 3.2 with*

$$t = \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + \frac{\bar{r}_n}{\sqrt{\tau}} \right) + \sqrt{2 + 2 \log(1/\tau)}.$$

There exists a universal constant C such that for every $P_n \in \text{GAC}(\mu^, \Sigma, \varepsilon)$ with $\varepsilon \leq \varepsilon^*$ and⁴ $\|\Sigma\|_{\text{op}} = 1$, the probability of the event*

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_\Sigma}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\frac{\log(1/\delta)}{n}} \right)$$

is at least $1 - \delta$. Moreover, the constant C from the last display can be made explicit by replacing the effective rank \mathbf{r}_Σ by the dimension p in the definition of \bar{r}_n : That is, for every $\delta \in (0, 1/5)$ the inequality

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{156 \sqrt{2 \log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{2p}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\frac{3 \log(2/\delta)}{n}} \right)$$

holds with probability at least $1 - 5\delta$.

Proof in the appendix, page 92

If we compare this result with its counterpart established in (Dalalyan and Minasyan, 2020) for the iteratively reweighted mean, besides the extra $\log p$ factor, we see that the above error bound does not reduce to the error of the empirical mean when the contamination rate goes to zero. We do not know whether this is just a drawback of our proof, or it is an intrinsic property of the estimator. Our numerical experiments reported later on suggest that it might be a property of the estimator.

There is another logarithmic factor, $\sqrt{\log(2/\varepsilon)}$, present in the second term of the error bounds provided by the last theorem, which does not appear in the minimax rate. There are computationally intractable robust estimators of the Gaussian mean, such as the Tukey median, that have an error bound free of this factor. However, all the known error bounds

⁴Since in this theorem Σ is assumed to be known, we can always divide all the data points X_i by $\|\Sigma\|_{\text{op}}^{1/2}$ to get a data set with a covariance matrix satisfying $\|\Sigma\|_{\text{op}} = 1$.

provably valid for polynomial time algorithms has this extra $\sqrt{\log(2/\varepsilon)}$ factor. Furthermore, this factor is known to be unavoidable in the case of sub-Gaussian model with adversarial contamination⁵, see (Lugosi and Mendelson, 2021, Section 2).

As shows the next theorem, the claims of Theorem 11 carry over the sub-Gaussian reference distributions with some slight modifications. These modifications mainly stem from the following lemma assessing the tail behavior of the singular values of a matrix having independent and sub-Gaussian columns.

Lemma 1 (Vershynin (2012), Theorem 5.39). *Let $\xi_{1:n}$ be a matrix consisting of sub-Gaussian vectors with variance proxy s . There is a universal constant C_0 such that for every $t > 0$ and for every pair of positive integers n and p , we have*

$$\begin{aligned} \mathbf{P}(\sigma_{\min}(\xi_{1:n}) \leq \sqrt{n} - C_0 s(\sqrt{p} + t)) &\leq e^{-t^2}, \\ \mathbf{P}(\sigma_{\max}(\xi_{1:n}) \geq \sqrt{n} + C_0 s(\sqrt{p} + t)) &\leq e^{-t^2}. \end{aligned}$$

Note that in the Gaussian case $s = 1$ and the constant C_0 can be chosen equal to $\sqrt{2}$. The last lemma leads to the following adaptations in the values of the thresholds used in the SDR estimator. First, we introduce auxiliary definitions

$$\tau = \frac{1}{4} \bigwedge \frac{\bar{r}_{n,s}}{\sqrt{\log_+(2/\bar{r}_{n,s})}}, \quad \text{with} \quad \bar{r}_{n,s} = \frac{3\sqrt{s}(\sqrt{p} + 2\sqrt{\log(2/\delta)})}{\sqrt{n}}. \quad (3.4)$$

Then, we set $t = t_1 + t_2$ with

$$t_1 = \frac{2(1 + C_0 \bar{r}_{n,s} \sqrt{s})}{1 - 2\varepsilon^*}, \quad t_2 = 1 + C_0 \sqrt{s} \left(\frac{\bar{r}_{n,s}}{\sqrt{\tau}} + \sqrt{2 + \log(1/\tau)} \right),$$

where C_0 is the same as in Lemma 1.

Now we are ready to state the theorem for the sub-Gaussian distributions showing that the SDR estimator with the threshold depending on the variance proxy s yields the same upper bound on ℓ_2 distance between our estimator $\hat{\mu}^{\text{SDR}}$ and the true value μ^* replacing the effective rank \mathbf{r}_Σ with the dimension p .

Theorem 12 (Sub-Gaussian version). *Let $\varepsilon^* \in (0, 1/2)$, and $\delta \in (0, 1/2)$. Define $\bar{r}_{n,s}$ and τ as in (3.4). For every $\varepsilon \leq \varepsilon^*$, let $\hat{\mu}^{\text{SDR}}$ be the estimator returned by Algorithm 3.2 with*

$$t = \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + C_0 \bar{r}_{n,s} \sqrt{\frac{s}{\tau}} \right) + C_0 s \sqrt{2 + 2 \log(1/\tau)},$$

where C is a universal constant. Then, there exists a constant C_s depending only on the variance proxy s such that for every $P_n \in \text{SGAC}(\mu^*, \Sigma, s, \varepsilon)$ with $\varepsilon \leq \varepsilon^*$ and $\|\Sigma\|_{\text{op}} = 1$, the

⁵Both sub-Gaussianity of the reference distribution and the adversarial nature of the contamination are important for getting the extra $\sqrt{\log(2/\varepsilon)}$ factor in the minimax rate.

probability of the event

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C_s \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{p}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\frac{\log(1/\delta)}{n}} \right)$$

is at least $1 - \delta$.

Proof in the appendix, page 111

3.4 The case of unknown covariance matrix

The SDR estimator, as defined in Algorithm 2, requires the knowledge of covariance matrix Σ . In this section we consider the case where the matrix Σ is unknown, but an approximation of the latter is available. Namely, we assume that we have access to a matrix $\tilde{\Sigma}$ and to a real number $\gamma > 0$ such that $\|\tilde{\Sigma} - \Sigma\|_{\text{op}} \leq \gamma \|\Sigma\|_{\text{op}}$. In such a situation, we can replace in the SDR estimator the true covariance matrix by its approximation $\tilde{\Sigma}$. This will necessarily require to adjust the threshold t accordingly. The goal of the present section is to propose a suitable choice of t and to show the impact of the approximation error γ on the estimation accuracy.

As mentioned, the parameter t used in Algorithm 2 needs to be properly tuned in order to account for the approximation error in the covariance matrix. To this end, we introduce the following auxiliary notation similar to those presented in (3.3):

$$\tilde{r}_n = \frac{\sqrt{C_\gamma \mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n}} \quad \text{and} \quad \tilde{\tau} = \frac{1}{4} \bigwedge \frac{\tilde{r}_n}{\sqrt{\log_+(2/\tilde{r}_n)}}, \quad (3.5)$$

where $C_\gamma = (1 + \gamma)/(1 - \gamma)$. Compared to (3.3), the main difference here is the presence of the factor C_γ (which is equal to one if $\gamma = 0$) and the substitution of the effective rank of Σ by that of its approximation $\tilde{\Sigma}$. In the rest of this section, we assume that Σ is invertible.

Theorem 13. *Let $\varepsilon^* \in (0, 1/2)$, $\delta \in (0, 1/2)$ and define \tilde{r}_n and τ as in (3.5). Assume that $\tilde{\Sigma}$ satisfies $\|\Sigma^{-1/2} \tilde{\Sigma} \Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma$ for some $\gamma \in (0, 1/2]$. Let $\hat{\mu}^{\text{SDR}}$ be the output of $\text{SDR}(\mathbf{X}_1, \dots, \mathbf{X}_n; \tilde{\Sigma}, \tilde{t}_\gamma)$, see Algorithm 2, with*

$$\tilde{t}_\gamma = \frac{\|\tilde{\Sigma}\|_{\text{op}}}{1 - \gamma} \left\{ \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + \frac{\tilde{r}_n}{\sqrt{\tilde{\tau}}} \right) + \sqrt{2 + \log(2/\tau)} \right\}.$$

Then, there exists a universal constant C such that for every data generating distribution $P_n \in \text{GAC}(\mu^, \Sigma, \varepsilon)$ with $\varepsilon \leq \varepsilon^*$, the probability of the event*

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C \|\Sigma\|_{\text{op}}^{1/2} \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_\Sigma}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\varepsilon \gamma} + \sqrt{\frac{\log(1/\delta)}{n}} \right) \quad (3.6)$$

is at least $1 - \delta$.

Proof in the appendix, page 106

On the one hand, if the value of γ is at most of order $\sqrt{(\mathbf{r}_\Sigma/n) \log(1/\varepsilon)} + \varepsilon \log(1/\varepsilon)$ then Theorem 13 implies that the estimation error is of the same order as in the case of known covariance matrix Σ (Theorem 11). For instance, if the matrix Σ is assumed to be diagonal, one can defined $\tilde{\Sigma}$ as the diagonal matrix composed of robust estimators of the variances of univariate contaminated Gaussian samples; see, for instance, Section 2 in (Comminges et al., 2021). For recent advances on robust estimation of (non-diagonal) covariance matrices by computationally tractable algorithms we refer the reader to (Cheng et al., 2019b).

On the other hand, if the value of γ for which the condition $\|\tilde{\Sigma} - \Sigma\|_{\text{op}} \leq \gamma \|\Sigma\|_{\text{op}}$ is known to be true is of larger order than $\sqrt{(\mathbf{r}_\Sigma/n) \log(1/\varepsilon)} + \varepsilon \log(1/\varepsilon)$, then $\sqrt{\varepsilon\gamma}$ dominates the other terms appearing in the error bound (3.6). Moreover, if γ is of constant order, then we get the error rate $\sqrt{\frac{\mathbf{r}_\Sigma}{n}} + \sqrt{\varepsilon}$, which is in line with previously known bounds for computationally tractable estimators; see for example (Lai et al., 2016, Theorem 1.1), (Diakonikolas et al., 2017, Theorem 3.2), (Dalalyan and Minasyan, 2020, Theorem 4).

3.5 Numerical experiments

We conducted numerical experiments on synthetic contaminated data to corroborate our theoretical results. The main goal of these experiments is to display statistical and computational features of the SDR and their dependence on various parameters. Moreover, we compared SDR to some other estimators proposed in the literature as well as to the oracle (empirical mean of the inliers). To do so, we selected componentwise median (CM), geometric median (GM) and Tukey’s median (TM) as the three classic estimators of the context, and the iteratively reweighted mean (IRM), introduced in (Dalalyan and Minasyan, 2020), as an example of optimization based method.

3.5.1 Implementation details

The experiments were run on a laptop with a 1.8 GHz Intel Core i7 and 8 GB of RAM. R codes of the experiments are freely available on the last author’s website. For GM and TM the R packages Gmedian⁶ (Cardot et al. (2013)) and TukeyRegion⁷ (Liu et al. (2019)) were used. IRM had been already implemented in R using Mosek⁸.

To optimize SDR, several choices were made. First, since geometric median is used in SDR as a rough estimator of the location, we limited it to at most 15 iterations and to stop at an accuracy of order 1. See the reference manual of Gmedian to have more details on these parameters. Second, since at the last step of the SDR one can use any estimator which is robust in low-dimensional setting, we chose to use the median of the projected data points.

⁶<https://cran.r-project.org/package=Gmedian>

⁷<https://cran.r-project.org/web/packages/TukeyRegion/index.html>

⁸www.mosek.com

Finally, we adjusted the numerical constant in the threshold t . In all the experiments, we assumed that the true value of ε is known and used $\varepsilon^* = \varepsilon$.

3.5.2 Experimental setup

Experiments were conducted on synthetic data sets obtained by applying a contamination scheme to n i.i.d. samples drawn from $\mathcal{N}_p(\mathbf{0}, \mathbf{I}_p)$. The following contamination schemes were considered.

- *Contamination by uniform outliers* (CUO): the locations of $n\varepsilon$ outliers are chosen at random independently of the inliers. The outliers are independent Gaussian with identity covariance matrix and with means having coordinates independently drawn from the uniform in $[0, 3]$ distribution.
- *Gaussian mixture contamination* (GMC): the locations of $n\varepsilon$ outliers are chosen at random independently of the inliers. The outliers are independent $\mathcal{N}_p(\boldsymbol{\mu}, \mathbf{I}_p)$. In our experiments, we chose $\boldsymbol{\mu}$ such that $\|\boldsymbol{\mu}\| = 15$.
- *Contamination by "smallest" eigenvector* (CSE): We replace the $n\varepsilon$ samples most correlated with the smallest principal eigenvector \mathbf{v}_p of the sample covariance matrix, by $n\varepsilon$ vectors all equal to $\sqrt{p}\mathbf{v}_p$ (\mathbf{v}_p is assumed to be a unit vector). In contrast with the two previous schemes, this one is adversarial.

Each experiment was repeated 50 times for SDR, CM, GM, the oracle and 10 times for IRM and TM. The tolerance probability δ was set to 0.1 in all the experiments. In the figures, points on the curves are median values of the error or of the running time for these trials whereas vertical bars overlaid on the points show the spread between the first and third quartiles. Since the computation of TM is prohibitively costly and is possible only for small sample sizes and dimensions, it is excluded from most of the experiments.

3.5.3 Statistical accuracy

At the upper left panel of Figure 3.2, we illustrate the behavior of the risk when the sample size increases for four different contamination levels: $\varepsilon \in \{0.1, 0.2, 0.3, 0.4\}$. The data are of dimension 60 and generated by the GMC scheme. The median estimation error converges respectively to the values 0.18, 0.36, 0.62 and 1.06. According to our theoretical result, the limit of the error should be proportional to $\frac{\varepsilon \log(1/\varepsilon)}{1-2\varepsilon}$. This is confirmed by the experimental results, since the ratio between the empirical limit of the median error and $\frac{\varepsilon \log(1/\varepsilon)}{1-2\varepsilon}$ for each ε is between 0.58 and 0.69.

At the upper right panel of Figure 3.2, the dependence of the error on the dimension is displayed. To better illustrate the effect of the dimension on the estimation error, we carried out our experiment on data sets of small sample size $n = 100$ with CUO contamination.

We compared the error of GM, CM, IRM and the oracle. In this plot, we clearly observe the supremacy of SDR and IRM as compared to GM and CM, which is in line with theoretical results. An important observation is that the error of the SDR estimator is very close to those of the IRM estimator and the oracle. This suggests that the factor $\sqrt{\log p}$ present in our theoretical results might be an artifact of the proof rather than an intrinsic property of the estimator, at least for nonadversarial contamination.

The last experiment aiming to display the behavior of the estimation error is depicted in the lower left panel of Figure 3.2. The examined synthetic datasets were generated by the CSE scheme with $\varepsilon = 0.2$. We measured the error for different values of the dimension and for sample size $n = 10p$ proportional to the dimension. In this case, the term $\sqrt{p/n}$ in the risk bound remains unchanged and we may perceive if the dimension virtually effects the term dependent on ε in the bound. The plot clearly confirms that the error is stable for SDR as it is for IRM and the oracle, in sharp contrast with GM and CM. The last point, of course, is not surprising since the risks of GM and CM scale as $\varepsilon\sqrt{p}$. Once again, this plot suggests that the factor $\sqrt{\log p}$ present in the SDR's risk bound might be unnecessary.

3.5.4 Computational efficiency

We conducted another experiment in order to better understand the computational complexity of SDR. Note that the computational cost of SDR comes from two operations done at each iterations: SVD of sample covariance matrix and computation of geometric median. We see that SDR can be computed in a reasonable time even in high dimensions. For instance, for $n = 10000$ and $p = 1000$ it takes nearly 26 seconds (tested over 20 trials).

At the lower right panel of Figure 3.2, we plotted the running times (in seconds) of GM, CM, IRM and SDR for various dimensions. Sample size in this experiment was set to 100, contamination rate was $\varepsilon = 0.2$ and CSE contamination scheme was used. As expected, IRM has substantially larger running time compared to SDR, GM and CM; this is due to the semidefinite programming solver running at each iteration of IRM. The fact that SDR is faster than GM (even though GM is deployed at each iteration of SDR) is explained by our choice of computing only a rough approximation of GM within SDR (limiting to 15 iterations and a tolerance parameter set to 1).

3.5.5 Breakdown point

A natural measure of robustness of an estimator is its resistance to a large fraction of outliers. The goal here is to demonstrate empirically our theoretical result showing that the breakdown point of the SDR estimator is $1/2$.

In Figure 3.3, at the left panel, we evaluated the error of the estimators on samples of size 100 and dimension 10 generated by the CSE scheme, for various values of ε . We can observe that SDR preserves its robustness with large contamination rates and outperforms

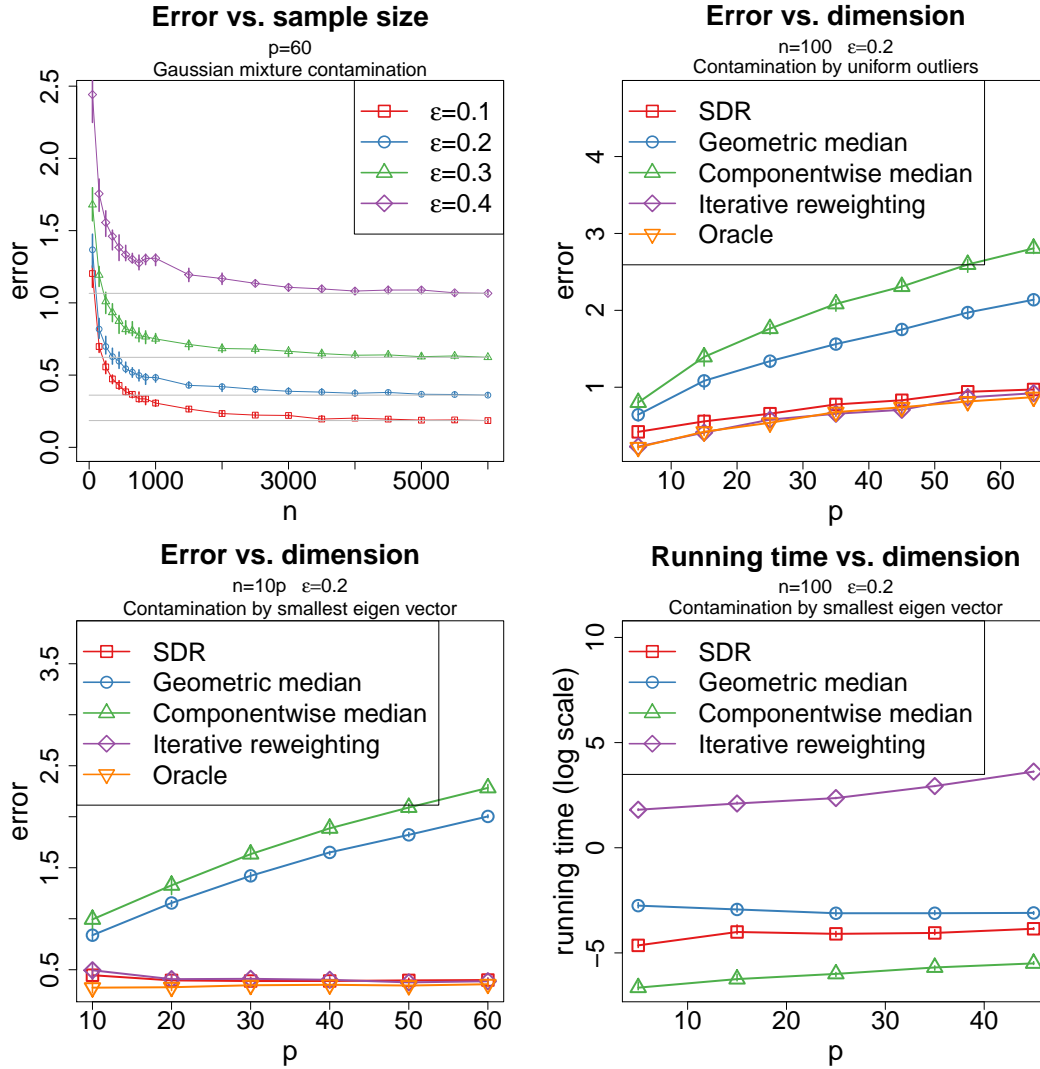


Figure 3.2: The upper left panel illustrates the convergence of SDR's median error when the sample size tends to infinity, for various contamination rates. The limiting values are shown by gray lines. The upper right panel shows the effect of dimension on the error. We see that in the case of SDR this effect is almost the same as for IRM and the oracle. The lower left panel plots the quantities when the sample-size increases proportionally to the dimension. Once again, we see that SDR is almost as accurate as IRM and the oracle. The lower right panel plots the running times of different estimators for various dimensions. It shows the huge computational gain of the SDR estimator as compared to IRM.

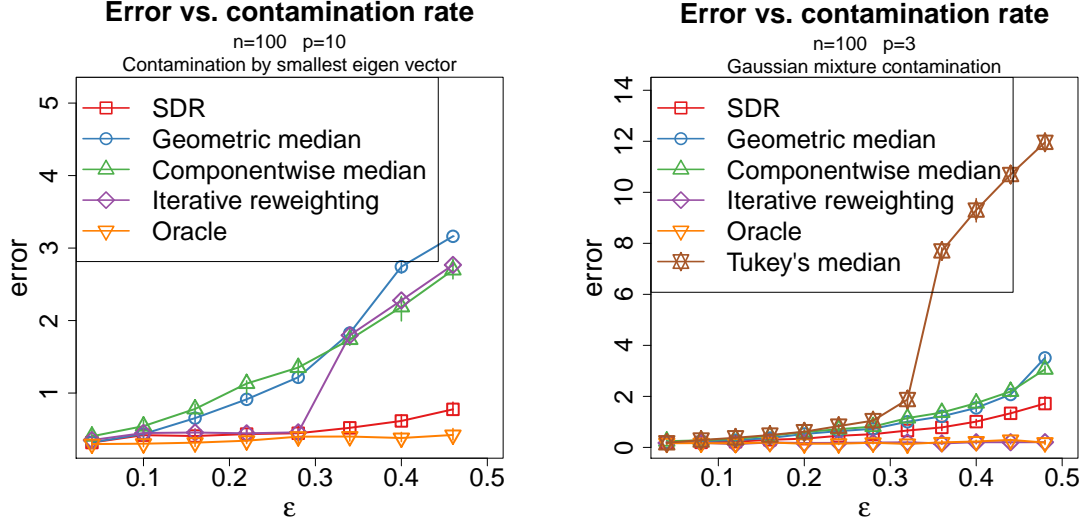


Figure 3.3: The left panel compares the robustness of various estimators by displaying the estimation error for different contamination rates under SCE scheme. SDR outperforms other estimators (even the oracle). Results displayed in the right panel are obtained by a similar experiment conducted for the GMC scheme. SDR is remarkably stable for different contamination schemes, while we see that SDR and Tukey’s median may behave poorly for $\varepsilon > 1/3$.

other estimators, excepted the oracle. More precisely, SDR and IRM have roughly the same error up to $\varepsilon = 0.28$. Starting from this value, the error of IRM starts a steep deterioration joining CM and GM.

At the right panel of Figure 3.3, we plotted the error as a function of the contamination rate for TM, CM, GM, IRM and SDR. Data used in this experiment were of size 100 and dimension 3, corrupted by GMC scheme. For this type of contamination, we observe that the IRM estimator remains robust even for ε close to $1/2$, whereas the error of TM deteriorates significantly for $\varepsilon > 1/3$. As a conclusion, for two contamination schemes which are challenging for iteratively reweighted mean and Tukey’s median, SDR shows very stable behavior.

3.6 Summary, related work and conclusion

We have proved that the multivariate mean estimator obtained by the iterative spectral dimension reduction method enjoys several appealing properties in the setting of sub-Gaussian observations subject to adversarial contamination. More precisely, in addition to being rigid transform equivariant and having breakdown point equal to $1/2$, the estimator has been shown to achieve the nearly minimax rate. Furthermore, the SDR estimator has low computational complexity, confirmed by reported numerical experiments. Indeed, its computational complexity is of the same order as that of computing the sample covariance matrix and performing a SVD on it. Presumably, at the cost of a moderate drop in accuracy, further speed-ups can be obtained by randomization (Halko et al., 2011) in the spirit of the prior work (Cheng et al., 2019a; Depersin and Lecué, 2022).

Notably, we have proved that the SDR estimator achieves the nearly optimal error rate without requiring the precise knowledge of the contamination rate. It however requires the knowledge of the covariance matrix. To alleviate this constraint, we have also established estimation guarantees in the case where an approximation of the covariance matrix is used instead of the true one. We have conducted numerical tests that show that the SDR is both fast and accurate.

Many recent works studied the problem of robust estimation in more complex high dimensional settings such linear regression or sparse mean and covariance estimation; see (Balakrishnan et al., 2017; Cheng et al., 2021; Chinot, 2020; Chinot et al., 2020; Collier and Dalalyan, 2019; Dalalyan and Thompson, 2019; Goes et al., 2020; Liu et al., 2020b; Pensia et al., 2020) and the references therein. It is under current investigation whether the results of the present paper can be extended to these settings.

Another interesting avenue for future research is to find an estimator that is rate-optimal, computationally tractable, with breakdown point equal to $1/2$ and, in the same time, asymptotically optimal in the sense that its risk is of order $\sqrt{p/n}$ when ε tends to zero. On a related note, it would be interesting to push further the exploration of second-order properties of the risk started in (Minasyan, 2020). Finally, an open question is how the minimax risk blows-up when the contamination rate tends to $1/2$. For the SDR estimator studied in this work, we established an upper bound of order $1/(1 - 2\varepsilon)$. However, we have no clue whether this is optimal. Our intuition is that it is not as the lower bound provided in Chapter 1 is of order $\log^{1/2}(1/(1 - 2\varepsilon))$ when ε is close to $1/2$.

Chapter 4

Discussion

In this manuscript, we investigated the problem of the robust estimation of the mean against outliers for two types of distributions: distributions supported by probability simplex and multivariate sub-Gaussian distributions. At the beginning, we presented existing approaches for dealing with outliers. We reviewed some folklore results and in particular we outlined the main challenge of the problem in high-dimensional settings which is the reconciliation of the statistical accuracy and computational efficiency.

We started our study by the simpler case of the distributions supported by the probability simplex, a special case of which is the category of the discrete distributions. We established the minimax rates under three different distances: total variation, Hellinger and Euclidean distances. We discovered that the convergence rate is different under each distance. Then, we proposed confidence regions shrinking at the minimax rate and instance based bounds. Finally, we provided some experimental results corroborating the established theoretical rates. In addition, throughout this study, we presented various contamination models and relations among them.

In the next stage, we considered the problem of the robust estimation of the multivariate sub-Gaussian mean under the adversarial contamination model. Unlike to the case of the distributions supported by probability simplex, the sample mean is not a robust estimator for the Gaussian mean. We designed a new robust estimator (called SDR), based on the idea of spectral dimension reduction proposed in (Lai et al., 2016). SDR enjoys a near optimal risk rate, a high breakdown point ($1/2$), a low computational complexity and is equivariant by similarity transforms. While SDR has the advantage of not requiring the knowledge of the contamination rate, it requires the knowledge of the covariance matrix, as it is the case for all the proposed computational tractable estimators attaining a near optimal error rate (we outlined this fact in Section 1.3.4). However, to alleviate this constraint, we established estimation guarantees in the case where an approximation of the covariance matrix is used instead of the true one. At the end, we presented some practical evidences for the theoretical properties of SDR by carrying out different experiments on synthetic data and comparing the performance of SDR to that of some other existing estimators in the literature.

An avenue for future research is to extend the SDR for robust realization of other statistical tasks such as linear regression, covariance matrix estimation, etc. It is also interesting to investigate the properties of SDR with respect to the pseudo-norms of the form $x \in \mathbb{R}^p \mapsto \|x\|_S = \sup_{v \in S} \langle v, x \rangle$ where S is a symmetric subset of \mathbb{R}^p .

Despite the existence several tractable robust estimators achieving a near optimal risk rate (with dependency $\varepsilon \sqrt{\log(1/\varepsilon)}$), and the evidence provided by [Diakonikolas et al. \(2016b\)](#) on the necessity of the factor $\varepsilon \sqrt{\log(1/\varepsilon)}$ for tractable robust estimators under the adversarial contamination, there are still some open questions. Is there an estimator for the adversarial contamination model that is rate-optimal, computationally tractable, with breakdown point equal to $1/2$, not requiring the knowledge of the contamination rate, and in the same time, asymptotically optimal in the sense that its risk is equivalent to $\sqrt{p/n}$ when ε tends to zero? Is there a tractable estimator for Huber's contamination which attains the optimal error rate $\sqrt{p/n} + \varepsilon$?

On another note, Tukey's median which is considered as an optimal estimator in the robust estimation context has still unknown properties to be explored. For example, in the case of the non-spherical Gaussian distributions, where the covariance matrix of the reference distribution is the matrix Σ , does Tukey's median enjoy the optimal risk rate $\sqrt{\text{Tr}(\Sigma)/n} + \|\Sigma\|_{\text{op}}^{1/2} \varepsilon$? How does it perform in robust estimation of the mean for other distributions such as distributions with bounded moments? Let ε^* denote the finite-sample breakdown point of Tukey's median. In Chapter 1, we showed that for multivariate spherical Gaussian data, ε^* satisfies

$$C \frac{p}{n} \leq \frac{1}{3} - \varepsilon^* \leq C' \sqrt{\frac{p}{n}},$$

where C and C' are positive constants. Can we improve the lower bound of $\frac{1}{3} - \varepsilon^*$ to $\sqrt{\frac{p}{n}}$? And there are questions on the computational aspects of Tukey's median. Tukey's median is a NP-hard problem, but, is there an approximation algorithm returning a point with an acceptable high Tukey's halfspace depth? We hope that these questions yield new directions and creative approaches in statistics and computer science.

Chapter 5

Introduction en français

5.1	Définition du problème	72
5.2	Variables discrètes	73
5.3	Variables gaussiennes	74
5.3.1	Problème en dimension 1	75
5.3.2	Problème en grande dimension	77
5.4	Organisation du manuscrit	81

5.1 Définition du problème

On observe des données indépendamment identiquement distribuées (iid) selon une loi de référence. On suppose que cette loi admet une moyenne et on souhaite estimer cette moyenne à partir des observations. Dans le cas des distributions sous-gaussiennes, on sait que la moyenne empirique est un estimateur optimal dans le sens minimax, c'est-à-dire que l'erreur d'estimation de la moyenne empirique dans son pire cas est au moins aussi bonne que celle de tout autre estimateur dans son pire cas respectif. Mais que se passe-t-il si parmi nos données, il y a quelques-unes qui ne suivent pas la loi de référence ? Ces données désobéissantes que l'on appellera *données aberrantes* ou *outlier*, peuvent avoir des effets significatifs et indésirables sur notre estimations. Par exemple, la présence d'une donnée aberrante extrêmement éloignée de la moyenne peut considérablement dévier la moyenne empirique et entraîner une grande erreur d'estimation. Ainsi, la moyenne empirique n'est pas robuste contre les données aberrantes.

Dans ce travail, on étudie le problème de l'estimation robuste de la moyenne. On suppose que Y_1, \dots, Y_n sont n échantillons iid selon une loi P_θ où θ est la moyenne de cette loi. Étant donné $\varepsilon \in [0, 1/2)$, on observe des échantillons ε -contaminés X_1, \dots, X_n ce qu'il veut dire

que $X_i = Y_i$ si $i \in \mathcal{I}$ et que X_i prend une valeur arbitraire si $i \in \mathcal{O}$ où $|\mathcal{I}| = n(1 - \varepsilon)$ et $|\mathcal{O}| = n\varepsilon$ ¹. On pose $\mathcal{S} := \mathcal{I} \cup \mathcal{O} = \{1, \dots, n\}$. Le but est d'estimer θ à partir de X_1, \dots, X_n .

5.2 Variables discrètes

Pour commencer, on considère le cas où la loi de référence P_θ est une distribution discrète à k valeurs distinctes. On représente les variables aléatoires issue de cette loi par des vecteurs de dimension k en utilisant l'encodage one-hot. Ainsi, Y_1, \dots, Y_n et X_1, \dots, X_n appartiennent à l'ensemble $E_k := \{e_1, \dots, e_k\}$, l'ensemble des vecteurs de la base canonique. Pour mesurer l'erreur de notre estimation $\hat{\theta}$ de θ , on utilise trois différentes distances : celle² de variation totale, de Hellinger, et de \mathbb{L}^2

$$\begin{aligned} d_{\text{TV}}(\hat{\theta}, \theta) &:= 1/2 \|\hat{\theta} - \theta\|_1, \\ d_{\text{H}}(\hat{\theta}, \theta) &:= 1/\sqrt{2} \|\hat{\theta}^{1/2} - \theta^{1/2}\|_2, \\ d_{\mathbb{L}^2}(\hat{\theta}, \theta) &:= \|\hat{\theta} - \theta\|_2. \end{aligned}$$

On s'intéresse à l'erreur minimax définie par

$$\mathfrak{R}_\square(n, k, \varepsilon) := \inf_{\bar{\theta}_n} \sup_{P_n \in \mathcal{M}_n(\varepsilon, \theta)} \mathbf{E}[d_\square(\bar{\theta}_n, \theta)],$$

où *inf* porte sur tous les estimateurs $\bar{\theta}_n$ construits à partir des observations X_1, \dots, X_n et *sup* porte sur toutes les distributions jointes P_n de X_1, \dots, X_n déterminées par notre modèle de contamination $\mathcal{M}_n(\varepsilon, \theta)$. L'indice \square de \mathfrak{R} ci-dessus fait référence à la distance employée dans le risque, et donc \square est TV, H, ou \mathbb{L}^2 .

$\mathcal{M}_n(\varepsilon, \theta)$ détermine la nature des outliers dans un modèle où la loi de référence est P_θ et les outliers constituent une ε fraction des données. Le modèle de contamination le plus général est le *modèle adversarial* dans lequel un adversaire omniscient observe les variables initiales Y_1, \dots, Y_n et remplace $n\varepsilon$ d'entre eux par des valeurs arbitraires. Le modèle de contamination le plus restrictif dans la littérature est celui de *Huber* dans lequel les observations sont des variables iid issues de la loi de mélange suivante

$$(1 - \varepsilon)P_\theta + \varepsilon Q,$$

où Q désigne la distribution des outliers. On détaille les différents modèles de contamination dans la section 2.2.

¹ Sans perte de généralité, on suppose que $n\varepsilon$ est un entier.

² On écrit $\|u\|_q = (\sum_{j=1}^k |u_j|^q)^{1/q}$ et $u^q = (u_1^q, \dots, u_k^q)$ pour tout $u \in \mathbb{R}_+^k$ et $q > 0$.

Concentrons maintenant sur la distance euclidienne. On sait que la moyenne empirique

$$\bar{\mathbf{X}}_n := \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

est l'estimateur minimax quand il n'y a pas de donnée aberrante. On va vérifier si cet estimateur reste robuste lors que des outliers contaminent les données. Nous avons

$$\begin{aligned} d_{\mathbb{L}^2}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}^*) &= \|\bar{\mathbf{X}}_n - \boldsymbol{\theta}^*\|_2 = \|\bar{\mathbf{Y}}_n - \boldsymbol{\theta}^* + \frac{1}{n} \sum_{i \in O} (\mathbf{X}_i - \mathbf{Y}_i)\|_2 \\ &\leq \|\bar{\mathbf{Y}}_n - \boldsymbol{\theta}^*\|_2 + \frac{|O|}{n} \sup_{\mathbf{x}, \mathbf{y} \in E_k} \|\mathbf{x} - \mathbf{y}\|_2 \\ &= \|\bar{\mathbf{Y}}_n - \boldsymbol{\theta}^*\|_2 + \sqrt{2}\varepsilon. \end{aligned}$$

On voit que l'erreur dans le cas contaminé est majorée par l'erreur dans le cas non-contaminé plus un terme de l'ordre de ε . Ceci montre que la contamination a au plus un impact de l'ordre de ε sur notre erreur d'estimation. Remarquons que cette relation est valable sous n'importe quel modèle de contamination. Par ailleurs, on peut montrer qu'il y a une borne inférieure de même ordre de grandeur pour l'erreur en question dans le modèle de Huber, *i.e.* dans le modèle le plus restrictif de la littérature. Cela nous aide à établir les vitesses de convergence minimax pour les différents modèles de contaminations. En effet, dans le chapitre 2, en généralisant cet argument, on prouve que

$$\begin{aligned} \mathfrak{R}_{\text{TV}}(n, k, \varepsilon) &\asymp (k/n)^{1/2} + \varepsilon, \\ \mathfrak{R}_{\text{H}}(n, k, \varepsilon) &\asymp (k/n)^{1/2} + \varepsilon^{1/2}, \\ \mathfrak{R}_{\mathbb{L}^2}(n, k, \varepsilon) &\asymp (1/n)^{1/2} + \varepsilon. \end{aligned}$$

C'est curieux de voir que toutes ces vitesses de convergence sont distinctes. En fait, il s'avère que l'erreur minimax n'est pas détériorée si la proportion des outliers est plus petite que $(k/n)^{1/2}$ pour la distance de variation totale, k/n pour la distance de Hellinger, et $(1/n)^{1/2}$ pour la distance euclidienne. De plus, on en déduit que la moyenne empirique est l'estimateur minimax pour les trois distances.

Dans le même chapitre, on étudie la relation entre les différents modèles de contamination, et on généralise ces résultats pour toutes les distributions à support dans le simplexe de dimension k , soit $\Delta^{k-1} = \{\mathbf{v} \in \mathbb{R}_+^k : v_1 + \dots + v_k = 1\}$. En outre, on y établit des régions de confiance qui rétrécissent à vitesses minimax.

5.3 Variables gaussiennes

Le cas des variables à support dans un simplexe est simple tandis que si la loi de référence est à support infini, il est un peu délicat de trouver un estimateur robuste contre les outliers.

La raison est que la présence d'un seul outlier avec une valeur arbitrairement grande peut dévier largement un estimateur tel que la moyenne empirique. Dans cette partie, on étudie le problème de l'estimation de la moyenne d'une loi normale à partir des observations contaminées par des outliers sous le modèle adversarial. On commence par le cas uni-dimensionnel et puis on aborde la version multidimensionnelle du problème.

5.3.1 Problème en dimension 1

On garde le cadre introduit dans la section 5.1, et on suppose désormais que la loi de référence est $\mathcal{N}(\mu, \sigma^2)$. On regarde brièvement deux approches possibles pour estimer μ de manière robuste. L'erreur d'estimation est mesurée par la valeur absolue.

Filtrage

Comme nous avons évoqué le problème avec la moyenne empirique est sa sensibilité aux outliers extrêmement éloignés de μ . Pour régler ce problème, on peut calculer la moyenne empirique sans considérer les plus petits et les plus grands points. En fait, on trie les observations $\mathbf{X}_{(1)} \leq \dots \leq \mathbf{X}_{(n)}$ et on retourne $\hat{\mu}_F := \sum_{i \in F} \mathbf{X}_i / |F|$ où $F = \{i \in \mathcal{I} \cup \mathcal{O} \mid \mathbf{X}_{(2n\varepsilon)} < \mathbf{X}_i < \mathbf{X}_{(n-2n\varepsilon)}\}$. Même s'il y a des outliers parmi les échantillons après ce filtrage, on peut être sûr qu'ils sont au plus aussi éloignés que $\mathbf{Y}_{(n\varepsilon)}$ et $\mathbf{Y}_{(n-n\varepsilon)}$ (où les échantillons initiaux sont triés : $\mathbf{Y}_{(1)} \leq \dots \leq \mathbf{Y}_{(n)}$). Ainsi, cet estimateur, appelé la *moyenne tronquée*, n'est pas tant influencé par les outliers avec des valeurs extrêmes. Regarder la figure 5.1.

Dans la section 1.2.1, on donne des clés de preuve pour montrer que $\hat{\mu}_F$ estime μ avec une erreur de l'ordre de

$$\frac{\sigma}{\sqrt{n}} + \sigma\varepsilon\sqrt{\log(1/\varepsilon)}$$

avec grande probabilité, si $\varepsilon < 1/4$.

Moyenne en tant que le centre de symétrie

La loi normale est une loi symétrique, et la moyenne n'est pas seulement le barycentre de la distribution, mais elle est également le centre de symétrie. Ainsi, une autre idée pour estimer la moyenne serait d'estimer le centre de symétrie. Il y a de différentes notions de symétrie. La notion que l'on adopte ici c'est celle induite par la masse : il existe la même quantité de masse dans les deux côtés du centre. Plus précisément, dans cette symétrie, le centre est la médiane. Donc, on peut estimer le centre, *i.e.* la moyenne, par la médiane de l'échantillon, qui en fait le centre de la symétrie induite par la masse empirique.

Dans la section 1.2.2, on prouve que la médiane de l'échantillon, sous une condition non

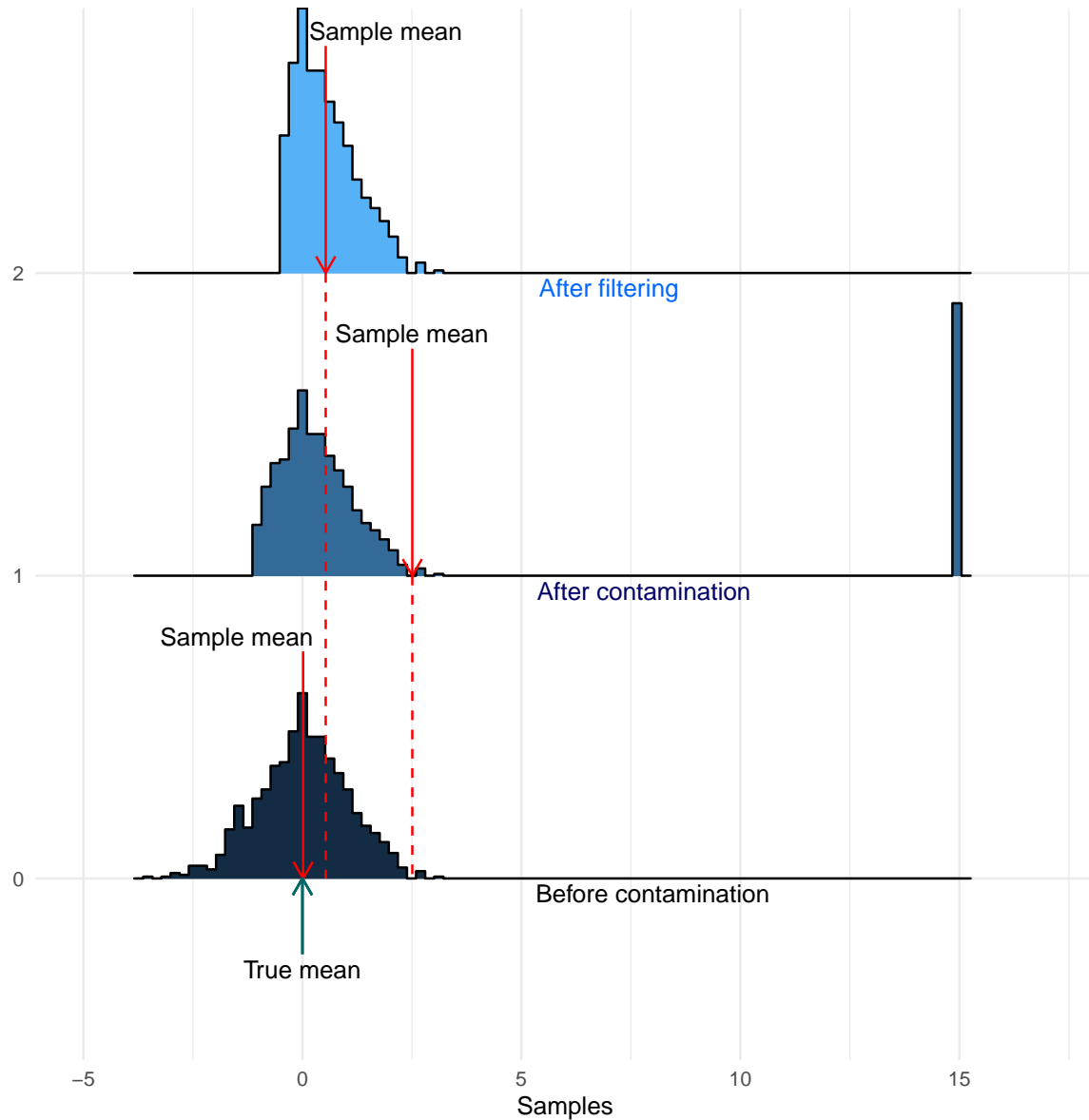


Figure 5.1: Au premier étage, l'histogramme des échantillons initiaux $(Y_i)_{i \in \mathcal{S}}$ est affiché. Au deuxième étage, nous avons l'histogramme des échantillons contaminés $(X_i)_{i \in \mathcal{S}}$ où les $n\varepsilon$ plus petits échantillons sont remplacées par $n\varepsilon$ échantillons de valeur 15. Au troisième étage, on voit l'histogramme des échantillons après le filtrage où les $2n\varepsilon$ plus grands et les $2n\varepsilon$ plus petits échantillons parmi $(X_i)_{i \in \mathcal{S}}$ sont mis de côté. La moyenne empirique pour chaque jeu de données est marquée par des lignes rouges. L'impact du filtrage sur la performance de la moyenne empirique est visible.

très restrictive, estime μ avec une erreur de l'ordre de

$$\frac{\sigma}{\sqrt{n}} + \sigma\varepsilon$$

avec grande probabilité. Par ailleurs, [Chen et al. \(2018\)](#) montre une borne inférieure du même ordre pour l'erreur minimax, ce qui implique que la médiane de l'échantillon est un estimateur minimax pour ce problème. Pour plus de détails sur cette borne inférieure regarder la section [1.2.4](#).

À part les deux méthodes mentionnées ci-dessus, il y a aussi la méthode de tournoi qui consiste à sélectionner parmi un ensemble fini de fonctions de densité, celle qui s'adapte mieux aux observations. On aborde cette méthode dans la section [1.2.3](#).

5.3.2 Problème en grande dimension

On suppose que nos données vivent dans \mathbb{R}^p et que les échantillons initiaux $(Y_i)_{i \in \{1, \dots, n\}}$ suivent la loi normale multidimensionnelle $\mathcal{N}(\mu, \sigma^2 \mathbf{I}_p)$. Tout comme le cas précédant, nous observons X_1, \dots, X_n , une version ε -contaminée des données sous le modèle adversarial. On quantifie l'erreur d'estimation par la distance euclidienne $\|\cdot\|_2$.

Une solution naïve pour le cas multidimensionnel serait d'estimer séparément chaque coordonnée de μ par une méthode robuste unidimensionnelle. Par la borne d'union, on pourrait prouver qu'en appliquant par exemple la médiane de l'échantillon à chaque coordonnée, on obtient un estimateur avec une erreur de l'ordre de

$$\sigma \sqrt{\frac{p \log(p)}{n}} + \sigma \varepsilon \sqrt{p}$$

avec grande probabilité. Le problème lié à ce genre d'estimateurs est que le terme $\varepsilon \sqrt{p}$ pourrait être problématique lorsque la valeur de p est grande. [Chen et al. \(2018\)](#) montre que sous le modèle de contamination de Huber, l'erreur de la médiane coordonnée par coordonnée ne peut pas être meilleure que $\sigma(\sqrt{p/n} + \varepsilon \sqrt{p})$ quand la distribution de contamination Q est définie comme un Dirac sur le point $\mu + \sigma(1, \dots, 1)^\top$. La même contamination peut être utilisée pour montrer que les autres méthodes robustes pour dimension 1 n'auront pas de meilleur comportement en grande dimension si on les applique coordonnée par coordonnée. Maintenant, la question est si on peut avoir une meilleure dépendance en p dans l'expression de l'erreur de notre estimation.

Filtrage (cas multidimensionnel)

Dans le cas multidimensionnel, il n'y a pas de notion d'ordre. Par conséquent, ce n'est pas possible d'employer la méthode de filtrage uni-dimensionnelle proposée ci-dessus. Pourtant, on sait que les échantillons gaussiens multidimensionnels se concentrent sur une sphère autour de la moyenne de rayon $\sigma \sqrt{p}$ (remarquer le contraste avec le cas unidimensionnel où

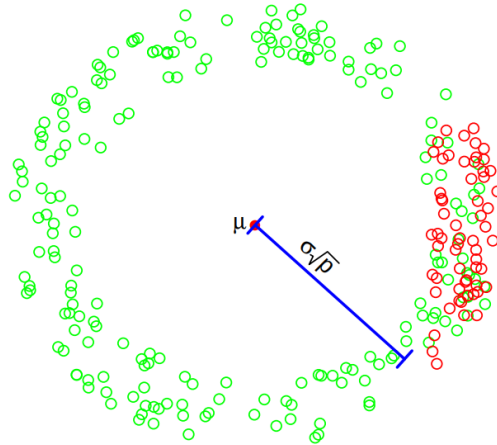


Figure 5.2: En grande dimension, les échantillons se concentrent sur une sphère de rayon $\sigma\sqrt{p}$ autour de la moyenne. Ici, les points rouges sont les outliers et les points verts sont les bons échantillons. Les simples outils de filtrage ne peuvent pas détecter ces outliers, et cette contamination dévie la moyenne empirique par un terme de l'ordre de $\varepsilon\sqrt{p}$.

les échantillons gaussiens se concentrent près de la moyenne). Plus précisément, il existe $c > 0$ tel que pour tout $t > 0$

$$\mathbf{P}(\|Y_i - \mu\|_2 - \sigma\sqrt{p} > t) \leq 2 \exp(-ct^2/\sigma^2),$$

(regarder e.g., ([Vershynin, 2018](#), Theorem 3.1.1)).

Cette propriété nous aide à construire des procédures de filtrage capables de détecter les échantillons extrêmement éloignés de la moyenne. En voir un exemple dans la section 1.3.1. Le problème est que ce genre de filtres ne peuvent enlever que les outliers extrêmes. Dans ce cas, les outliers peuvent se concentrer sur une même direction dans la sphère indiquée plus haut sans être filtrés, et après dévier la moyenne empirique par un terme de l'ordre de $\varepsilon\sqrt{p}$. Regarder la figure 5.2. En conséquence, on voit qu'un filtrage simple ne nous permet pas non plus de nous débarrasser du terme $\varepsilon\sqrt{p}$ dans l'expression de l'erreur de notre estimation. Pourtant, comme on verra c'est un outil utile.

La moyenne en tant que le centre de symétrie (cas multidimensionnel)

Pour généraliser la notion de la médiane de l'échantillon au cas multivarié, il faut adopter une notion de symétrie. Les différentes notions de symétrie donnent lieu à de différentes définitions de la médiane. Comme on analyse dans la section 1.3.2, la notion de la *symétrie centrale* définit la moyenne empirique comme la médiane de l'échantillon, ce qui n'est pas un estimateur robuste de la moyenne ; la notion de la *symétrie angulaire* définit la médiane géométrique, et ce dernier a le même défaut que les estimateurs précédents, à savoir le terme

$\varepsilon\sqrt{p}$ dans l'expression de son erreur dans le pire scénario.

En dimension 1, on a défini la médiane en tant que le centre de la symétrie induite par la masse. Cette notion de symétrie se généralise au cas multidimensionnel par la notion de la *symétrie de demi-espace*. Un vecteur aléatoire Y respecte la symétrie de demi-espace autour le centre θ si l'inégalité $\mathbf{P}(Y \in H) \geq 1/2$ est valide pour tout demi-espace fermé H contenant θ . Autrement dit, θ est le centre de la symétrie induite par la masse dans toute direction. Cette notion de symétrie définit la *médiane de Tukey* qui est un point appartenant à l'ensemble

$$\arg \max_{\mathbf{x} \in \mathbb{R}^p} \min_{\mathbf{v} \in \mathbb{R}^{p-1}} \sum_{i=1}^n \mathbb{1}(\mathbf{v}^\top \mathbf{X}_i \leq \mathbf{v}^\top \mathbf{x}).$$

Pour plus de précisions sur cette définition, regarder la section 1.3.2.

Chen et al. (2018) prouvent que sous certaines conditions, le risque de la médiane de Tukey est de l'ordre de

$$\sigma \sqrt{p/n} + \sigma \varepsilon$$

avec grande probabilité, et qu'en plus il s'agit de l'erreur optimale. C'est à dire que la médiane de Tukey est un estimateur optimal dans le sens minimax. Cependant, la médiane de Tukey n'est pas calculable en temps polynomial car son calcul nécessite le traitement de toutes les directions, soit un nombre d'opérations qui est exponentiel en p .

Utiliser la matrice de covariance

Comparant $\varepsilon\sqrt{p}$ et ε , on remarque qu'en grande dimension, il y a un écart significatif entre l'erreur des méthodes calculables déjà citées (telles que la médiane coordonnée par coordonnée, la médiane géométrique, le filtrage, etc.) et celle de la médiane de Tukey ou de la méthode de tournoi expliquée dans la section 1.3.3, qui ne sont pas calculables efficacement. Ce phénomène a stimulé beaucoup de recherches afin de trouver un compromis entre la précision statistique et l'efficacité de calcul. Le point commun entre presque toutes les méthodes proposant un tel compromis, c'est qu'elles se servent de la connaissance de la matrice de covariance de la loi de référence.

On suppose désormais que la loi de référence est $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Soit $\hat{\boldsymbol{\Sigma}}$ la matrice de covariance empirique. En nous appuyant sur notre connaissance de $\boldsymbol{\Sigma}$, on essaie de filtrer ou pondérer les observations de sorte que $\hat{\boldsymbol{\Sigma}}$ s'approche de $\boldsymbol{\Sigma}$. En effet, on se sert ici d'une propriété essentielle : $\|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}\|_2$ est contrôlé par $\varepsilon \|\boldsymbol{\Sigma} - \hat{\boldsymbol{\Sigma}}\|_{\text{op}}$ où $\hat{\boldsymbol{\mu}}$ et $\hat{\boldsymbol{\Sigma}}$ sont respectivement la moyenne empirique et la matrice de covariance empirique des observations filtrées ou pondérées. Ainsi, notre but consiste à minimiser $\|\boldsymbol{\Sigma} - \hat{\boldsymbol{\Sigma}}\|_{\text{op}}$ en filtrant ou pondérant les observations.

La méthode que nous construisons dans le chapitre 3, appelée *Spectral Dimension Reduction* (SDR), appartient à ce paradigme des méthodes employant la matrice de covariance

de la loi de référence. Il s'agit d'une méthode qui filtre les observations de manière itérative où dans chaque itération elle réduit la dimension en regardant le spectre de la matrice $\hat{\Sigma} - \Sigma$. Elle est basée sur l'observation qu'il y a des directions dans lesquelles l'impact de la contamination n'est pas significatif. Autrement dit, au lieu de regarder toutes les directions comme c'est nécessaire pour le calcul de la médiane de Tukey, on se focalise sur les directions les plus impactées par la contamination.

En fait, nous pouvons partitionner les vecteurs propres de la matrice $\hat{\Sigma} - \Sigma$ en deux groupes : les vecteurs propres supérieurs (correspondant aux grandes valeurs propres) et les vecteurs propres inférieurs (correspondant aux petites valeurs propres). Puis, nous projetons orthogonalement les échantillons sur les sous-espaces engendrés par chaque groupe. Il s'avère que la moyenne empirique est un estimateur robuste dans le sous-espace engendré par les vecteurs propres inférieurs. Donc, nous pouvons nous concentrer sur le sous-espace engendré par les vecteurs propres supérieurs. Sur ce sous-espace, nous répétons récursivement la même démarche. L'idée de la réduction itérative de dimension pour l'estimation robuste est originellement introduite dans (Lai et al., 2016). Cette procédure peut se schématiser dans l'algorithme "SDR" suivant :

- Si $p = 1$ on retourne la médiane de l'échantillon
- Sinon :
 1. Étape de filtrage : on supprime des échantillons afin d'avoir tous les échantillons dans une distance de l'ordre de \sqrt{p} de μ
 2. Soit V le sous-espace engendré par les $p/2$ vecteurs propres supérieurs de $\hat{\Sigma} - \Sigma$
 3. On projette orthogonalement les échantillons sur V et V^\top
 4. Sur V^\top on applique la moyenne empirique
 5. Sur V on applique récursivement SDR

La version détaillée de cet algorithme se trouve dans le chapitre 3. On y prouve que le risque de SDR est de l'ordre de

$$\sqrt{\log p} \left(\sqrt{\text{Tr}(\Sigma)/n} + \|\Sigma\|_{\text{op}}^{1/2} \varepsilon \sqrt{\log(1/\varepsilon)} \right)$$

avec grande probabilité. Donc, SDR est une méthode calculable en temps polynomial avec une erreur presque optimale (optimale modulo des facteurs logarithmiques). Cette méthode peut être également généralisée pour toutes les distributions sous-gaussiennes.

L'avantage de la méthode SDR, ce qu'elle est rapide, qu'elle ne requière pas la connaissance de la proportion de contamination (ε), et que son breakdown point vaut $1/2$ (breakdown point est la proportion maximale des outliers pour laquelle le risque de notre estimateur reste borné). De plus, on montre à la fin que SDR est adaptable au cas où notre connaissance de Σ est partielle.

5.4 Organisation du manuscrit

Ce manuscrit est constitué de quatre chapitres et deux appendices. Le chapitre 1 est une introduction sur l'estimation robuste de la moyenne gaussienne. On y présente les approches générales et on prouve informellement quelques résultats folklores du domaine. À la fin du chapitre, on expose les contributions principales de ce travail. Le chapitre 2 aborde le problème de l'estimation robuste des distributions à support dans le simplexe de probabilité. Dans ce chapitre, on étudie aussi les différents modèles de contamination. Cette partie du manuscrit est publiée dans (Bateni and Dalalyan, 2020). Dans le chapitre 3, on propose une méthode calculatoirement efficace avec une erreur presque optimale pour estimer la moyenne d'une loi gaussienne sous le modèle adversarial. Cette partie du travail se trouve dans (Bateni et al., 2022). On finit avec le chapitre 4 où on présente une conclusion et une perspective sur les directions possibles pour les prochains travaux. Les deux appendices contiennent les preuves formelles des résultats exposés dans les chapitres 2 et 3.

Appendix A

Proofs for Chapter 2

A.1 Proofs of propositions	82
A.2 Minimax upper bounds over the sparse simplex	84
A.3 Minimax lower bounds over the sparse simplex	86
A.4 Proofs of bounds with high probability	88
A.5 Proofs of instance based bounds	90

A.1 Proofs of propositions

Proof of Proposition 3 on page 41. Recall that \widehat{O} is the set of outliers in the Huber model. Let O be any subset of $\{1, \dots, n\}$. It follows from the definition of Huber's model that if P_n^O stands for the conditional distribution of (X_1, \dots, X_n) given $\widehat{O} = O$, when (X_1, \dots, X_n) is drawn from $P^n \in \mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)$, then $P_n^O \in \mathcal{M}_n^{\text{HDC}}(|O|/n, \theta^*)$. Therefore, for every O of cardinality $o \geq 2\varepsilon n$, we have

$$\begin{aligned}
 \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*) \mathbf{1}(\widehat{O} = O)] &= \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*) | \widehat{O} = O] \mathbf{P}(\widehat{O} = O) \\
 &\leq \sup_{\mathcal{M}_n^{\text{HDC}}(o/n, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \mathbf{P}(\widehat{O} = O) \\
 &\stackrel{(1)}{\leq} \sup_{\mathcal{M}_n^{\text{HDC}}(1, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \mathbf{P}(\widehat{O} = O).
 \end{aligned} \tag{A.1}$$

Inequality (1) above is a direct consequence of the inclusion $\mathcal{M}_n^{\text{HDC}}(o/n, \theta^*) \subset \mathcal{M}_n^{\text{HDC}}(1, \theta^*)$. Summing the obtained inequality over all sets O of cardinality $\geq 2\varepsilon n$, we get

$$\sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*) \mathbf{1}(|\widehat{O}| \geq 2\varepsilon n)] \leq \sup_{\mathcal{M}_n^{\text{HDC}}(1, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \mathbf{P}(|\widehat{O}| \geq 2\varepsilon n).$$

It follows from the multiplicative form of Chernoff's inequality that $\mathbf{P}(|\widehat{O}| \geq 2\varepsilon n) \leq e^{-n\varepsilon/3}$. This leads to the last term in inequality (2.3).

Using the same argument as for (A.1), for any O of cardinality $o < 2n\varepsilon$, we get

$$\begin{aligned} \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*) \mathbb{1}(|\widehat{O}| < 2n\varepsilon)] &\leq \sup_{\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \sum_{|O| \leq 2n\varepsilon} \mathbf{P}(\widehat{O} = O) \\ &= \sup_{\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)]. \end{aligned}$$

This completes the proof of (2.3).

One can use the same arguments along with the Tchebychev inequality to establish (2.4). Indeed, for every S of cardinality $o \leq 2\varepsilon n$, we have

$$\begin{aligned} \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r \text{ and } \widehat{O} = O) \\ &= \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r \mid \widehat{O} = O) \mathbf{P}(\widehat{O} = O) \\ &\leq \sup_{\mathcal{M}_n^{\text{HDC}}(o/n, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r) \mathbf{P}(\widehat{O} = O) \\ &\leq \sup_{\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \mathbf{P}(\widehat{O} = O). \end{aligned}$$

Summing the obtained inequality over all sets O of cardinality $o \leq 2\varepsilon n$, we get

$$\begin{aligned} \sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r \text{ and } |\widehat{O}| \leq 2\varepsilon n) &\leq \sup_{\mathcal{M}_n^{\text{HDC}}(2\varepsilon, \theta^*)} \mathbf{E}[d(\widehat{\theta}_n, \theta^*)] \\ &= R_d^{\text{HDC}}(n, 2\varepsilon, \theta^*, \widehat{\theta}). \end{aligned}$$

On the other hand, it holds that

$$\sup_{\mathcal{M}_n^{\text{HC}}(\varepsilon, \theta^*)} r \mathbf{P}(d(\widehat{\theta}_n, \theta^*) > r \text{ and } |\widehat{O}| > 2\varepsilon n) \leq r \mathbf{P}(|\widehat{O}| > 2\varepsilon n) \leq r e^{-n\varepsilon/3},$$

and the claim of the proposition follows. \square

Proof of Proposition 4 on page 43. Let θ_1 and θ_2 be two points in Θ . We have

$$\begin{aligned} 2R_d^{\text{HC}}(n, \varepsilon, \Theta, \widehat{\theta}_n) &\geq R_d^{\text{HC}}(n, \varepsilon, \theta_1, \widehat{\theta}_n) + R_d^{\text{HC}}(n, \varepsilon, \theta_2, \widehat{\theta}_n) \\ &\geq \mathbf{E}_{(1-\varepsilon)\mathbf{P}_{\theta_1} + \varepsilon\mathbf{P}_{\theta_2}}[d(\widehat{\theta}_n, \theta_1)] + \mathbf{E}_{(1-\varepsilon)\mathbf{P}_{\theta_1} + \varepsilon\mathbf{P}_{\theta_2}}[d(\widehat{\theta}_n, \theta_2)]. \end{aligned}$$

To ease writing, assume that n is an even number. Let O be any fixed set of cardinality $n/2$. It is clear that the set of outliers \widehat{O} satisfies

$$p_O = \mathbf{P}(\widehat{O} = O) = \mathbf{P}(\widehat{O} = O^c) > 0.$$

Furthermore, if $\mathbf{X}_{1:n}$ is drawn from $((1 - \varepsilon)\mathbf{P}_{\theta_1} + \varepsilon\mathbf{P}_{\theta_2})^{\otimes n} := \mathbf{P}_\varepsilon^{\otimes n}$, then its conditional distribution given $\widehat{O} = O$ is exactly the same as the conditional distribution of $\mathbf{X}_{1:n} \sim \mathbf{P}_\varepsilon^{\otimes n}$ given $\widehat{O} = O^c$. This implies that

$$\begin{aligned} 2R_d^{\text{HC}}(n, \varepsilon, \Theta, \widehat{\theta}_n) &\geq p_O(\mathbf{E}_{\mathbf{P}_\varepsilon}[d(\widehat{\theta}_n, \theta_1)|\widehat{O} = O] + \mathbf{E}_{\mathbf{P}_\varepsilon}[d(\widehat{\theta}_n, \theta_2)|\widehat{O} = O^c]) \\ &= p_O \mathbf{E}_{\mathbf{P}_\varepsilon}[d(\widehat{\theta}_n, \theta_1) + d(\widehat{\theta}_n, \theta_2)|\widehat{O} = S] \geq p_O d(\theta_1, \theta_2), \end{aligned}$$

where in the last step we have used the triangle inequality. The obtained inequality being true for every $\theta_1, \theta_2 \in \Theta$, we can take the supremum to get

$$R_d^{\text{HC}}(n, \varepsilon, \Theta, \widehat{\theta}_n) \geq (p_O/2) \sup_{\theta_1, \theta_2 \in \Theta} d(\theta_1, \theta_2) = +\infty.$$

This completes the proof. □

A.2 Minimax upper bounds over the sparse simplex

This section is devoted to the proof of the upper bounds on minimax risks in the discrete model with respect to various distances.

Proof of Theorem 7 on page 46. To ease notation, we set

$$\bar{\mathbf{X}}_n = \frac{1}{n} \sum_i \mathbf{X}_i, \quad \bar{\mathbf{Y}}_n = \frac{1}{n} \sum_i \mathbf{Y}_i, \quad \bar{\mathbf{X}}_O = \frac{1}{o} \sum_{i \in O} \mathbf{X}_i, \quad \bar{\mathbf{Y}}_O = \frac{1}{o} \sum_{i \in O} \mathbf{Y}_i$$

In the adversarial model, we have $\mathbf{X}_i = \mathbf{Y}_i$ if $i \notin O$ where $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ are generated from the reference distribution θ^* .

$$\begin{aligned} d_{\mathbb{L}^2}(\bar{\mathbf{X}}_n, \theta^*) &= \|\bar{\mathbf{X}}_n - \theta^*\|_2 = \|\bar{\mathbf{Y}}_n - \theta^* + \frac{1}{n} \sum_{i \in O} (\mathbf{X}_i - \mathbf{Y}_i)\|_2 \\ &\leq \|\bar{\mathbf{Y}}_n - \theta^*\|_2 + \frac{|O|}{n} \sup_{\mathbf{x}, \mathbf{y} \in \Delta^{k-1}} \|\mathbf{x} - \mathbf{y}\|_2 \\ &= \|\bar{\mathbf{Y}}_n - \theta^*\|_2 + \sqrt{2\varepsilon}, \end{aligned}$$

which gives us

$$\sup_{\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)} \mathbf{E}[d_{\mathbb{L}^2}(\bar{\mathbf{X}}_n, \theta^*)] \leq \sup_{\theta^*} \mathbf{E}[d_{\mathbb{L}^2}(\bar{\mathbf{Y}}_n, \theta^*)] + \sqrt{2\varepsilon}.$$

And for a fixed θ^* it is well known that

$$\begin{aligned}\mathbf{E}[d_{\mathbb{L}^2}^2(\bar{\mathbf{Y}}_n, \theta^*)] &= \sum_{j=1}^k \mathbf{Var}[\bar{\mathbf{Y}}_{n,j}] = \sum_{j=1}^k \mathbf{Var}\left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}(Y_i = \mathbf{e}_j)\right] \\ &= \frac{1}{n} \sum_{j=1}^k \mathbf{Var}[\mathbb{1}(Y_1 = \mathbf{e}_j)] \leq \frac{1}{n} \sum_{j=1}^k \mathbf{E}[\mathbb{1}(Y_1 = \mathbf{e}_j)] = \frac{1}{n}.\end{aligned}$$

Hence, we obtain $R_{\mathbb{L}^2}^{\text{AC}}(n, \varepsilon, \Delta^{k-1}) \leq (1/n)^{1/2} + \varepsilon$. Similarly,

$$\begin{aligned}d_{\text{TV}}(\bar{\mathbf{X}}_n, \theta^*) &\leq \|\bar{\mathbf{Y}}_n - \theta^*\|_1 + \frac{|O|}{n} \sup_{\mathbf{x}, \mathbf{y} \in \Delta^{k-1}} \|\mathbf{x} - \mathbf{y}\|_1 \\ &= \|\bar{\mathbf{Y}}_n - \theta^*\|_1 + 2\varepsilon.\end{aligned}$$

This gives

$$\sup_{\mathcal{M}_n^{\text{AC}}(\varepsilon, \theta^*)} \mathbf{E}[d_{\text{TV}}(\bar{\mathbf{X}}_n, \theta^*)] \leq \sup_{\theta^*} \mathbf{E}[d_{\text{TV}}(\bar{\mathbf{Y}}_n, \theta^*)] + 2\varepsilon.$$

In addition, for every θ^* ,

$$\begin{aligned}\mathbf{E}[d_{\text{TV}}(\bar{\mathbf{Y}}_n, \theta^*)] &= \frac{1}{2} \sum_{j=1}^k \mathbf{E}[|\bar{\mathbf{Y}}_{n,j} - \theta_j^*|] \\ &\leq \frac{1}{2} \sum_{j=1}^k \left(\mathbf{E}[|\bar{\mathbf{Y}}_{n,j} - \theta_j^*|^2] \right)^{1/2} \\ &= \frac{1}{2} \sum_{j=1}^k \left(\mathbf{Var}[\bar{\mathbf{Y}}_{n,j}] \right)^{1/2} \\ &\stackrel{(1)}{=} \frac{1}{2} \sum_{j \in J} \left(\frac{1}{n} \theta_j^* (1 - \theta_j^*) \right)^{1/2} \\ &\stackrel{(2)}{\leq} \frac{1}{2} s^{1/2} \left(\sum_{j=1}^k \frac{1}{n} \theta_j^* (1 - \theta_j^*) \right)^{1/2} \leq \frac{1}{2} (s/n)^{1/2},\end{aligned}$$

where in (1) we have used the notation $J = \{j : \theta_j^* \neq 0\}$ and in (2) we have used the Cauchy-Schwarz inequality. This leads to

$$R_{\text{TV}}^{\text{AC}}(n, \varepsilon, \Delta^{k-1}) \leq (k/n)^{1/2} + \varepsilon.$$

Finally, for the Hellinger distance

$$d_{\text{H}}(\bar{\mathbf{X}}_n, \theta^*) \leq d_{\text{H}}(\bar{\mathbf{X}}_n, \bar{\mathbf{Y}}_n) + d_{\text{H}}(\bar{\mathbf{Y}}_n, \theta^*) \leq d_{\text{TV}}(\bar{\mathbf{X}}_n, \bar{\mathbf{Y}}_n)^{1/2} + d_{\text{H}}(\bar{\mathbf{Y}}_n, \theta^*),$$

where we have already seen that

$$d_{\text{TV}}(\bar{\mathbf{X}}_n, \bar{\mathbf{Y}}_n) = \frac{o}{n} \|\bar{\mathbf{X}}_O - \bar{\mathbf{Y}}_O\|_1 \leq 2\varepsilon.$$

This yields

$$\sup_{\mathcal{M}_n^{\text{AC}}(\varepsilon, \boldsymbol{\theta}^*)} \mathbf{E}[d_{\text{H}}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}^*)] \leq \sup_{\boldsymbol{\theta}^*} \mathbf{E}[d_{\text{H}}(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*)] + \sqrt{2\varepsilon}.$$

Note that $\mathbf{E}[\bar{\mathbf{Y}}_{n,j}] = \theta_j^*$ implies that $\bar{\mathbf{Y}}_{n,j} = \theta_j^* = 0$ for every j that does not belong to the sparsity pattern J . Furthermore, for every $j \in J$, we have

$$|\sqrt{\bar{\mathbf{Y}}_{n,j}} - \sqrt{\theta_j^*}| = \frac{|\bar{\mathbf{Y}}_{n,j} - \theta_j^*|}{\sqrt{\bar{\mathbf{Y}}_{n,j}} + \sqrt{\theta_j^*}} \leq \frac{|\bar{\mathbf{Y}}_{n,j} - \theta_j^*|}{\sqrt{\theta_j^*}}.$$

This implies that for every $\boldsymbol{\theta}^* \in \Delta_s^{k-1}$,

$$\begin{aligned} \mathbf{E}[d_{\text{H}}^2(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*)] &= \mathbf{E}\left[\frac{1}{2} \sum_{j=1}^k \left(\sqrt{\bar{\mathbf{Y}}_{n,j}} - \sqrt{\theta_j^*}\right)^2\right] \leq \frac{1}{2} \mathbf{E}\left[\sum_{j \in J} \frac{(\bar{\mathbf{Y}}_{n,j} - \theta_j^*)^2}{\theta_j^*}\right] \\ &= \frac{1}{2} \sum_{j \in J} \frac{1}{\theta_j^*} \text{Var}[\bar{\mathbf{Y}}_{n,j}] = \frac{1}{2} \sum_{j \in J} \frac{1}{\theta_j^*} \times \frac{\theta_j^*(1 - \theta_j^*)}{n} = \frac{s-1}{2n}. \end{aligned}$$

Hence, by Jensen's inequality $\mathbf{E}[d_{\text{H}}(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*)] < \sqrt{s/n}$. Therefore, we infer that

$$R_{\text{H}}^{\text{AC}}(n, \varepsilon, \Delta_s^{k-1}) \leq (s/n)^{1/2} + \sqrt{2\varepsilon}^{1/2}$$

and the last claim of the theorem follows. \square

A.3 Minimax lower bounds over the sparse simplex

This section is devoted to the proof of the lower bounds on minimax risks in the discrete model with respect to various distances. Note that the rates over the high-dimensional “sparse” simplex Δ_s^{k-1} coincide with those for the dense simplex Δ^{s-1} . For this reason, all the lower bounds will be proved for Δ^{s-1} only (for $s \geq 2$ an even integer). In addition, we will restrict our attention to the distributions \mathbf{P}, \mathbf{Q} over Δ^{s-1} that are supported by the set \mathcal{A} of the elements of the canonical basis that is $\mathbf{P}(\mathcal{A}) = \mathbf{Q}(\mathcal{A}) = 1$.

Proof of Theorem 8 on page 46. We denote by \mathbf{e}_j the vector in \mathbb{R}^s having all the coordinates equal to zero except the j th coordinate which is equal to one. Setting

$$\boldsymbol{\theta} = \mathbf{e}_1, \quad \text{and} \quad \boldsymbol{\theta}' = \left(1 - \frac{\varepsilon}{1 - \varepsilon}\right) \mathbf{e}_1 + \frac{\varepsilon}{1 - \varepsilon} \mathbf{e}_2$$

we have

$$d_{\text{TV}}(\boldsymbol{\theta}, \boldsymbol{\theta}') = \frac{\varepsilon}{1 - \varepsilon}, \quad d_{\mathbb{L}^2}(\boldsymbol{\theta}, \boldsymbol{\theta}') \geq \sqrt{2}\varepsilon, \quad \text{and} \quad d_{\text{H}}(\boldsymbol{\theta}, \boldsymbol{\theta}') \geq \varepsilon^{1/2}.$$

Therefore, modulus of continuity defined by

$$w_d(\varepsilon, \Delta) = \sup \{d(\boldsymbol{\theta}, \boldsymbol{\theta}') : \boldsymbol{\theta}, \boldsymbol{\theta}' \in \Delta, d_{\text{TV}}(\boldsymbol{\theta}, \boldsymbol{\theta}') \leq \varepsilon/(1 - \varepsilon)\}$$

for a distance d and a set Δ , satisfies for any $\varepsilon \leq 1/2$,

$$w_{\text{TV}}(\varepsilon, \Delta^{k-1}) \geq \varepsilon, \quad w_{\mathbb{L}^2}(\varepsilon, \Delta^{k-1}) \geq \sqrt{2}\varepsilon, \quad \text{and} \quad w_{\text{H}}(\varepsilon, \Delta^{k-1}) \geq \varepsilon^{1/2}.$$

These bounds on the modulus of continuity are the first ingredient we need for lower bounding the minimax risk using (Chen et al., 2018, Theorem 5.1).

The second ingredient is the minimax rate in the non-contaminated case. are well known. For each of the considered distances, this rate is well-known. However, for the sake of completeness, we provide below a proof those lower bounds using Fano's method. For this, we use the Varshamov-Gilbert lemma (see e.g. Lemma 2.9 in (Tsybakov, 2009)) and Theorem 2.5 in (Tsybakov, 2009). The Varshamov-Gilbert lemma guarantees the existence of a set $\boldsymbol{\omega}^{(1)}, \dots, \boldsymbol{\omega}^{(M)} \in \{0, 1\}^{\lfloor s/2 \rfloor}$ of cardinality $M \geq 2^{s/16}$ such that

$$\rho(\boldsymbol{\omega}^{(i)}, \boldsymbol{\omega}^{(j)}) \geq \frac{s}{16}, \quad \text{for all } i \neq j,$$

where $\rho(\cdot, \cdot)$ stands for the Hamming distance. Using these binary vectors $\{\boldsymbol{\omega}_j\}$, a parameter $\beta \in [0, 1/s]$ to be specified later and the “baseline” vector $\boldsymbol{\theta}^{(0)} = (1/s, \dots, 1/s)$, we define hypotheses $\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(M)}$ by the relations

$$\boldsymbol{\theta}_{2j-1}^{(i)} = \boldsymbol{\theta}_{2j-1}^{(0)} + \boldsymbol{\omega}_j^{(i)}\beta \quad \text{and} \quad \boldsymbol{\theta}_{2j}^{(i)} = \boldsymbol{\theta}_{2j}^{(0)} - \boldsymbol{\omega}_j^{(i)}\beta \quad \forall j \in \{0, \dots, \lfloor s/2 \rfloor\}.$$

Remark that $\boldsymbol{\theta}^{(0)}, \dots, \boldsymbol{\theta}^{(M)}$ are all probability vectors of dimension s . Denoting the Kullback-Leibler divergence by $d_{\text{KL}}(\cdot, \cdot)$, one can check the conditions of Theorem 2.5 in Tsybakov (2009):

$$\begin{aligned} d_{\mathbb{L}^2}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(j)}) &\geq \beta \frac{\sqrt{2s}}{4} \quad \forall j \neq i, \\ d_{\text{TV}}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(j)}) &\geq \frac{\beta s}{16} \quad \forall j \neq i, \end{aligned}$$

as well as

$$\begin{aligned}
d_{\text{KL}}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(0)}) &\leq \sum_{j=1}^{\lfloor s/2 \rfloor} (\boldsymbol{\theta}_{2j-1}^{(0)} + \beta) \log \frac{\boldsymbol{\theta}_{2j-1}^{(0)} + \beta}{\boldsymbol{\theta}_{2j-1}^{(0)}} + (\boldsymbol{\theta}_{2j}^{(0)} - \beta) \log \frac{\boldsymbol{\theta}_{2j}^{(0)} - \beta}{\boldsymbol{\theta}_{2j}^{(0)}} \\
&\leq \sum_{j=1}^{\lfloor k/2 \rfloor} \frac{\beta}{\boldsymbol{\theta}_{2j-1}^{(0)}} (\boldsymbol{\theta}_{2j-1}^{(0)} + \beta) - \frac{\beta}{\boldsymbol{\theta}_{2j}^{(0)}} (\boldsymbol{\theta}_{2j}^{(0)} - \beta) \\
&= \beta^2 \sum_{j=1}^k \frac{1}{\boldsymbol{\theta}_j^{(0)}} = \beta^2 \sum_{j=1}^s s = \beta^2 s^2 \leq \frac{\alpha \log M}{n} \quad \forall i \in \{1, \dots, M\},
\end{aligned}$$

for $\beta = \sqrt{\alpha/(ns)}/4$, taking into account the fact that $2^{s/16} \leq M$. Now by applying the aforementioned theorem, we obtain for the non-contaminated setting ($\varepsilon = 0$)

$$\begin{aligned}
\inf_{\bar{\boldsymbol{\theta}}_n} \sup_{\mathcal{M}_n^{\text{HC}}(0, \Delta^{s-1})} \mathbf{P} \left(d_{\mathbb{L}^2}(\bar{\boldsymbol{\theta}}_n, \boldsymbol{\theta}^*) \geq \frac{\beta \sqrt{2s}}{2} \right) &\geq \frac{\sqrt{M}}{1 + \sqrt{M}} \left(1 - 2\alpha - \sqrt{\frac{2\alpha}{\log M}} \right), \\
\inf_{\bar{\boldsymbol{\theta}}_n} \sup_{\mathcal{M}_n^{\text{HC}}(0, \Delta^{s-1})} \mathbf{P} \left(d_{\text{TV}}(\bar{\boldsymbol{\theta}}_n, \boldsymbol{\theta}^*) \geq \frac{\beta s}{8} \right) &\geq \frac{\sqrt{M}}{1 + \sqrt{M}} \left(1 - 2\alpha - \sqrt{\frac{2\alpha}{\log M}} \right).
\end{aligned}$$

Setting $M = 2^{s/16}$ and $\alpha = 1/32$, by Markov's inequality, one concludes

$$\begin{aligned}
\inf_{\bar{\boldsymbol{\theta}}_n} R_{\mathbb{L}^2}^{\text{HC}}(n, 0, \Delta^{s-1}, \bar{\boldsymbol{\theta}}_n) &\geq c(1/n)^{1/2}, \\
\inf_{\bar{\boldsymbol{\theta}}_n} R_{\text{TV}}^{\text{HC}}(n, 0, \Delta^{s-1}, \bar{\boldsymbol{\theta}}_n) &\geq c(s/n)^{1/2},
\end{aligned}$$

where $c = 1/25600$. Since, $\sqrt{2} d_{\text{H}}(\boldsymbol{\theta}, \boldsymbol{\theta}') \geq d_{\text{TV}}(\boldsymbol{\theta}, \boldsymbol{\theta}')$ for any pair of points $(\boldsymbol{\theta}, \boldsymbol{\theta}')$ on the simplex, we have

$$\inf_{\bar{\boldsymbol{\theta}}_n} R_{\text{H}}^{\text{HC}}(n, 0, \Delta^{s-1}, \bar{\boldsymbol{\theta}}_n) \geq c(s/n)^{1/2}.$$

Finally, we apply (Chen et al., 2018, Theorem 5.1) stating in our case for any distance d

$$\inf_{\bar{\boldsymbol{\theta}}_n} R_d^{\text{HC}}(n, \varepsilon, \Delta^{s-1}, \bar{\boldsymbol{\theta}}_n) \geq c \left\{ \inf_{\bar{\boldsymbol{\theta}}_n} R_d^{\text{HC}}(n, 0, \Delta^{s-1}, \bar{\boldsymbol{\theta}}_n) + w_d(\varepsilon, \Delta^{s-1}) \right\},$$

for an universal constant c , which completes the proof of Theorem 8. \square

A.4 Proofs of bounds with high probability

Proof of Theorem 9 on page 48. Suppose $X_i = Y_i$ if $i \notin O$ where Y_1, \dots, Y_n are independently generated from the reference distribution \mathbf{P} so that $\mathbf{E}[Y_i] = \boldsymbol{\theta}^*$. For any Z_1, \dots, Z_n , let $\Phi_{\square}(Z_1, \dots, Z_n) := d_{\square}(\sum_{i=1}^n Z_i/n, \boldsymbol{\theta}^*)$, where \square refers here to the distances \mathbb{L}^2 or TV. Given

$\mathbf{Y}'_1, \dots, \mathbf{Y}'_n \in \Delta^{k-1}$ we have for every i

$$\Phi_{\square}(\mathbf{Y}_1, \dots, \mathbf{Y}_i, \dots, \mathbf{Y}_n) - \Phi_{\square}(\mathbf{Y}_1, \dots, \mathbf{Y}_{i-1}, \mathbf{Y}'_i, \mathbf{Y}_{i+1}, \dots, \mathbf{Y}_n) \leq \frac{1}{n} d_{\square}(\mathbf{Y}_i, \mathbf{Y}'_i).$$

Furthermore, it can easily be shown that the last term is bounded by $\sqrt{2}/n$ and $2/n$ for the distances \mathbb{L}^2 and TV, respectively. By bounded difference inequality (see for example Theorem 6.2 of [Boucheron et al. \(2013\)](#)) with probability at least $1 - \delta$

$$\begin{aligned} \Phi_{\mathbb{L}^2}(\mathbf{Y}_1, \dots, \mathbf{Y}_n) &\leq \mathbf{E}\Phi_{\mathbb{L}^2}(\mathbf{Y}_1, \dots, \mathbf{Y}_n) + (\log(1/\delta)/n)^{1/2} \\ &\leq (1/n)^{1/2} + (\log(1/\delta)/n)^{1/2}, \\ \Phi_{TV}(\mathbf{Y}_1, \dots, \mathbf{Y}_n) &\leq \mathbf{E}\Phi_{TV}(\mathbf{Y}_1, \dots, \mathbf{Y}_n) + (\log(2/\delta)/n)^{1/2} \\ &\leq (s/n)^{1/2} + (\log(2/\delta)/n)^{1/2}. \end{aligned}$$

Using $\Phi_{\square}(\mathbf{X}_1, \dots, \mathbf{X}_n) \leq \Phi_{\square}(\mathbf{Y}_1, \dots, \mathbf{Y}_n) + d_{\square}(\sum_{i \in O} \mathbf{X}_i/n, \sum_{i \in O} \mathbf{Y}_i/n)$, one can conclude the proof of the first two claims of the theorem.

For the Hellinger distance, the computations are more tedious. We have to separate the case of small θ_j^* . To this end, let $J = \{j : 0 < \theta_j^* < (1/n) \log(2s/\delta)\}$ and $J' = \{j : \theta_j^* \geq (1/n) \log(2s/\delta)\}$. We have

$$\begin{aligned} \sum_{j \in J} \left(\sqrt{\bar{\mathbf{Y}}_{n,j}} - \sqrt{\theta_j^*} \right)^2 &\leq \sum_{j \in J} (\bar{\mathbf{Y}}_{n,j} + \theta_j^*) \\ &\leq \frac{1}{n} \sum_{i=1}^n \left(\sum_{j \in J} \mathbf{Y}_{i,j} - \theta_j^* \right) + \frac{2s \log(2s/\delta)}{n}. \end{aligned}$$

Since the random variables $U_i := \left(\sum_{j \in J} \mathbf{Y}_{i,j} \right)$ are iid, positive, bounded by 1, the Bernstein inequality implies that

$$\frac{1}{n} \sum_{i=1}^n (U_i - \mathbf{E}[U_1]) \leq \sqrt{\frac{2\text{Var}(U_1) \log(2/\delta)}{n}} + \frac{\log(2/\delta)}{3n},$$

holds with probability at least $1 - \delta/2$ for $0 < \delta < 1$. One easily checks that $\sqrt{\text{Var}(U_1)} \leq \sum_{j \in J} \sqrt{\text{Var}(\mathbf{Y}_{1,j})} \leq s \sqrt{\log(2s/\delta)/n}$. Therefore, with probability at least $1 - \delta/2$, we have

$$\frac{1}{n} \sum_{i=1}^n \left(\sum_{j \in J} \mathbf{Y}_{i,j} - \theta_j^* \right) \leq \frac{\sqrt{2} s \log(2s/\delta)}{n} + \frac{\log(2/\delta)}{3n}.$$

This yields

$$\sum_{j \in J} \left(\sqrt{\bar{\mathbf{Y}}_{n,j}} - \sqrt{\theta_j^*} \right)^2 \leq \frac{3.5s \log(2s/\delta) + \log(2/\delta)}{n}, \quad (\text{A.2})$$

with probability at least $1 - \delta/2$.

On the other hand, we have

$$\sum_{j \in J'} \left(\sqrt{\bar{Y}_{n,j}} - \sqrt{\theta_j^*} \right)^2 \leq \sum_{j \in J'} \frac{(\bar{Y}_{n,j} - \theta_j^*)^2}{\theta_j^*} \leq s \max_{j \in J'} \frac{(\bar{Y}_{n,j} - \theta_j^*)^2}{\theta_j^*}. \quad (\text{A.3})$$

The Bernstein inequality and the union bound imply that, with probability at least $1 - \delta/2$, for all $j \in J'$,

$$\begin{aligned} |\bar{Y}_{n,j} - \theta_j^*| &\leq \sqrt{\frac{2\text{Var}(\mathbf{Y}_{1,j}) \log(2s/\delta)}{n}} + \frac{\log(2s/\delta)}{n} \\ &\leq \sqrt{\frac{2\theta_j^* \log(2s/\delta)}{n}} + \frac{\log(2s/\delta)}{n} \leq 2.5 \sqrt{\frac{\theta_j^* \log(2s/\delta)}{n}}. \end{aligned} \quad (\text{A.4})$$

Combining (A.3) and (A.4), we obtain

$$\sum_{j \in J'} \left(\sqrt{\bar{Y}_{n,j}} - \sqrt{\theta_j^*} \right)^2 \leq 2.5^2 \frac{s \log(2s/\delta)}{n}, \quad (\text{A.5})$$

with probability at least $1 - \delta/2$. Finally, inequalities (A.2) and (A.5) together lead to

$$d_{\text{H}}^2(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*) = \frac{1}{2} \sum_{j=1}^n \left(\sqrt{\bar{Y}_{n,j}} - \sqrt{\theta_j^*} \right)^2 \leq \frac{5s \log(2s/\delta)}{n} + \frac{\log(2/\delta)}{2n},$$

which is true with probability at least $1 - \delta$. Using the triangle inequality and the fact that the Hellinger distance is smaller than the square root of the TV-distance, we get

$$\begin{aligned} d_{\text{H}}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}^*) &\leq d_{\text{H}}(\bar{\mathbf{X}}_n, \bar{\mathbf{Y}}_n) + d_{\text{H}}(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*) \\ &\leq \sqrt{d_{\text{TV}}(\bar{\mathbf{X}}_n, \bar{\mathbf{Y}}_n)} + d_{\text{H}}(\bar{\mathbf{Y}}_n, \boldsymbol{\theta}^*) \\ &\leq \varepsilon^{1/2} + \sqrt{\frac{5s \log(2s/\delta)}{n}} + \sqrt{\frac{\log(2/\delta)}{2n}}, \end{aligned}$$

with probability at least $1 - \delta$. This completes the proof of the theorem. \square

A.5 Proofs of instance based bounds

Proof of Proposition 5 on page 49. In the adversarial model, we have $\mathbf{X}_i = \mathbf{Y}_i$ if $i \notin O$ where $\mathbf{Y}_1, \dots, \mathbf{Y}_n$ are generated from the reference distribution $\boldsymbol{\theta}^*$. By Proposition 3 and Lemma 6 in Berend and Kontorovich (2013) and using the triangle inequality

$$\begin{aligned} \|\bar{\mathbf{Y}}_n - \boldsymbol{\theta}^*\|_1 - \frac{|O|}{n} \sup_{\mathbf{x}, \mathbf{y} \in \Delta^\infty} \|\mathbf{x} - \mathbf{y}\|_1 &\leq d_{\text{TV}}(\bar{\mathbf{X}}_n, \boldsymbol{\theta}^*) \\ &\leq \|\bar{\mathbf{Y}}_n - \boldsymbol{\theta}^*\|_1 + \frac{|O|}{n} \sup_{\mathbf{x}, \mathbf{y} \in \Delta^\infty} \|\mathbf{x} - \mathbf{y}\|_1, \end{aligned}$$

the proposition is proved. □

Proof of Theorem 10 on page 49. Using Theorem 2.1 in [Cohen et al. \(2020\)](#) and applying the same triangle inequality as in the proof of Proposition 5, we conclude the theorem. □

Appendix B

Proofs for Chapter 3

B.1	Proof of Theorem 11	92
B.1.1	Bounding the projected error of the average of filtered observations	92
B.1.2	Bounding the error of the geometric median of projected observations	94
B.1.3	Bounding the number of filtered out observations	95
B.1.4	Estimating the mean from a low-dimensional adversarial projection	97
B.1.5	Bounding stochastic errors	98
B.1.6	Putting all the pieces together	102
B.2	Proof of Theorem 13	106
B.3	Extension to Sub-Gaussian distributions	111
B.3.1	Proof of Theorem 12	113

B.1 Proof of Theorem 11

Before proving Theorem 11, we provide some auxiliary lemmas and propositions.

B.1.1 Bounding the projected error of the average of filtered observations

For any $J \subset [n]$, we define $\bar{\mathbf{Z}}_J$ and $\hat{\Sigma}_J^Z$ as the sample average and sample covariance matrix of the subsample $\{\mathbf{Z}_i : i \in J\}$, that is

$$\bar{\mathbf{Z}}_J = \frac{1}{|J|} \sum_{i \in J} \mathbf{Z}_i, \quad \hat{\Sigma}_J^Z = \frac{1}{|J|} \sum_{i \in J} \mathbf{Z}_i \mathbf{Z}_i^\top - \bar{\mathbf{Z}}_J \bar{\mathbf{Z}}_J^\top.$$

The main building block of the proof is the following result.

Proposition 6. *Let $\mathcal{S} \subset [n]$ be an arbitrary set. We define its subsets $\mathcal{S}_{\mathcal{I}} = \mathcal{S} \cap \mathcal{I}$ and $\mathcal{S}_{\mathcal{O}} = \mathcal{S} \cap \mathcal{O}$. Let $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ and $\boldsymbol{\mu}^Z$ be arbitrary points in \mathbb{R}^q with $q \geq 2$. Let Σ^Z be an arbitrary $q \times q$*

covariance matrix and let \mathbf{P}_k be the projection matrix projecting onto the subspace spanned by the bottom k eigenvectors of $\hat{\Sigma}_S^Z - \Sigma^Z$, for $k = 1, \dots, q-1$. We have

$$\|\mathbf{P}_k(\bar{\mathbf{Z}}_S - \boldsymbol{\mu}^Z)\|_2 \leq \left\{ 2\omega_{\mathcal{O}} \|\hat{\Sigma}_{S_I}^Z - \Sigma^Z\|_{\text{op}} + \frac{\omega_{\mathcal{O}}^2}{\omega_{\mathcal{I}}} \left((\lambda_q - \lambda_1)(\Sigma^Z) + \frac{\delta_Z^2}{q-k} \right) \right\}^{1/2} + \|\mathbf{P}_k \bar{\boldsymbol{\xi}}_{S_I}^Z\|_2,$$

where $\omega_{\mathcal{O}} = |\mathcal{S}_{\mathcal{O}}|/|S|$, $\omega_{\mathcal{I}} = 1 - \omega_{\mathcal{O}}$, $\boldsymbol{\xi}_i^Z = \mathbf{Z}_i - \boldsymbol{\mu}^Z$ and $\delta_Z = \inf_{\boldsymbol{\mu}} \max_{i \in S} \|\mathbf{Z}_i - \boldsymbol{\mu}\|_2$. Furthermore, if $|S| \leq q-k$, then

$$\|\mathbf{P}_k(\bar{\mathbf{Z}}_S - \boldsymbol{\mu}^Z)\|_2 \leq \left\{ 2\omega_{\mathcal{O}} \|\hat{\Sigma}_{S_I}^Z - \Sigma^Z\|_{\text{op}} + \frac{\omega_{\mathcal{O}}^2}{\omega_{\mathcal{I}}} (\lambda_q - \lambda_1)(\Sigma^Z) \right\}^{1/2} + \|\mathbf{P}_k \bar{\boldsymbol{\xi}}_{S_I}^Z\|_2.$$

Proof. Since there is no risk of confusion, we remove the superscript Z from Σ^Z , $\hat{\Sigma}_S^Z$ and so on. Since $\bar{\mathbf{Z}}_S = \omega_{\mathcal{I}} \bar{\mathbf{Z}}_{S_I} + \omega_{\mathcal{O}} \bar{\mathbf{Z}}_{S_{\mathcal{O}}}$ yields $\bar{\mathbf{Z}}_S - \bar{\mathbf{Z}}_{S_I} = \omega_{\mathcal{O}}(\bar{\mathbf{Z}}_{S_{\mathcal{O}}} - \bar{\mathbf{Z}}_{S_I})$, the triangle inequality implies that

$$\begin{aligned} \|\mathbf{P}_k(\bar{\mathbf{Z}}_S - \boldsymbol{\mu})\|_2 &\leq \|\mathbf{P}_k(\bar{\mathbf{Z}}_S - \bar{\mathbf{Z}}_{S_I})\|_2 + \|\mathbf{P}_k(\bar{\mathbf{Z}}_{S_I} - \boldsymbol{\mu})\|_2 \\ &\leq \omega_{\mathcal{O}} \|\mathbf{P}_k(\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})\|_2 + \|\mathbf{P}_k \bar{\boldsymbol{\xi}}_{S_I}\|_2. \end{aligned} \quad (\text{B.1})$$

Moreover, one can check that

$$\begin{aligned} \hat{\Sigma}_S &= \omega_{\mathcal{I}} \hat{\Sigma}_{S_I} + \omega_{\mathcal{O}} \hat{\Sigma}_{S_{\mathcal{O}}} + \omega_{\mathcal{I}} \omega_{\mathcal{O}} (\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})^{\otimes 2} \\ &\succeq \omega_{\mathcal{I}} \hat{\Sigma}_{S_I} + \omega_{\mathcal{I}} \omega_{\mathcal{O}} (\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})^{\otimes 2}. \end{aligned} \quad (\text{B.2})$$

Hence, multiplying from left and right by \mathbf{P}_k and computing the largest eigenvalue of both sides, we get

$$\begin{aligned} \omega_{\mathcal{I}} \omega_{\mathcal{O}} \|\mathbf{P}_k(\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})\|_2^2 &\leq \lambda_q(\mathbf{P}_k(\hat{\Sigma}_S - \omega_{\mathcal{I}} \hat{\Sigma}_{S_I}) \mathbf{P}_k^{\top}) \\ &\leq \lambda_q(\mathbf{P}_k(\hat{\Sigma}_S - \Sigma) \mathbf{P}_k^{\top}) + \lambda_q(\Sigma - \omega_{\mathcal{I}} \hat{\Sigma}_{S_I}) \\ &\leq \lambda_k(\hat{\Sigma}_S - \Sigma) + \omega_{\mathcal{I}} \lambda_q(\Sigma - \hat{\Sigma}_{S_I}) + \omega_{\mathcal{O}} \lambda_q(\Sigma). \end{aligned} \quad (\text{B.3})$$

On the other hand, using the Weyl inequality (several times) and the identity (B.2), we get

$$\begin{aligned} \lambda_k(\hat{\Sigma}_S - \Sigma) &\leq \omega_{\mathcal{I}} \lambda_q(\hat{\Sigma}_{S_I} - \Sigma) + \omega_{\mathcal{O}} \lambda_k(\hat{\Sigma}_{S_{\mathcal{O}}} - \Sigma + \omega_{\mathcal{I}}(\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})^{\otimes 2}) \\ &\leq \omega_{\mathcal{I}} \lambda_q(\hat{\Sigma}_{S_I} - \Sigma) + \omega_{\mathcal{O}} \lambda_{k+1}(\hat{\Sigma}_{S_{\mathcal{O}}} - \Sigma) + \omega_{\mathcal{I}} \lambda_{q-1}((\bar{\mathbf{Z}}_{S_I} - \bar{\mathbf{Z}}_{S_{\mathcal{O}}})^{\otimes 2}) \\ &\leq \omega_{\mathcal{I}} \lambda_q(\hat{\Sigma}_{S_I} - \Sigma) + \omega_{\mathcal{O}} \lambda_{k+1}(\hat{\Sigma}_{S_{\mathcal{O}}}) - \omega_{\mathcal{O}} \lambda_1(\Sigma). \end{aligned} \quad (\text{B.4})$$

For the middle term of the right hand side, we can use the following upper bound

$$\begin{aligned}
\lambda_{k+1}(\widehat{\Sigma}_{\mathcal{S}_O}) &\leq \frac{\lambda_{k+1}(\widehat{\Sigma}_{\mathcal{S}_O}) + \dots + \lambda_q(\widehat{\Sigma}_{\mathcal{S}_O})}{q-k} \\
&\leq \frac{\text{Tr}(\widehat{\Sigma}_{\mathcal{S}_O})}{q-k} = \frac{1}{q-k} \text{Tr} \left(\frac{1}{|\mathcal{S}_O|} \sum_{i \in \mathcal{S}_O} (\mathbf{Z}_i - \bar{\mathbf{Z}}_{\mathcal{S}_O})^{\otimes 2} \right) \\
&= \frac{1}{q-k} \inf_{\boldsymbol{\mu}} \text{Tr} \left(\frac{1}{|\mathcal{S}_O|} \sum_{i \in \mathcal{S}_O} (\mathbf{Z}_i - \boldsymbol{\mu})^{\otimes 2} \right) \\
&= \inf_{\boldsymbol{\mu}} \frac{1}{(q-k)|\mathcal{S}_O|} \sum_{i \in \mathcal{S}_O} \|\mathbf{Z}_i - \boldsymbol{\mu}\|_2^2 \leq \frac{\delta_Z^2}{q-k}.
\end{aligned}$$

Combining (B.3), (B.4) and the last display, we get

$$\omega_{\mathcal{I}} \omega_{\mathcal{O}} \|\mathbf{P}_k(\bar{\mathbf{Z}}_{\mathcal{S}_I} - \bar{\mathbf{Z}}_{\mathcal{S}_O})\|_2^2 \leq 2\omega_{\mathcal{I}} \|\widehat{\Sigma}_{\mathcal{S}_I} - \Sigma\|_{\text{op}} + \omega_{\mathcal{O}} (\lambda_q(\Sigma) - \lambda_1(\Sigma)) + \frac{\omega_{\mathcal{O}} \delta_Z^2}{q-k}.$$

Dividing both sides by $\omega_{\mathcal{I}} \omega_{\mathcal{O}}$, we arrive at

$$\|\mathbf{P}_k(\bar{\mathbf{Z}}_{\mathcal{S}_I} - \bar{\mathbf{Z}}_{\mathcal{S}_O})\|_2^2 \leq \frac{2}{\omega_{\mathcal{O}}} \|\widehat{\Sigma}_{\mathcal{S}_I} - \Sigma\|_{\text{op}} + \frac{1}{\omega_{\mathcal{I}}} \left(\lambda_q(\Sigma) - \lambda_1(\Sigma) + \frac{\delta_Z^2}{q-k} \right).$$

Combining this inequality with (B.1), we obtain the first claim of the proposition. To get the second claim, we simply remark that $\lambda_{k+1}(\widehat{\Sigma}_{\mathcal{S}_O}) = 0$ since the rank of the matrix $\widehat{\Sigma}_{\mathcal{S}_O}$ is less than $|\mathcal{S}_O| \leq |\mathcal{S}| \leq q-k$. \square

B.1.2 Bounding the error of the geometric median of projected observations

We assume in this section that V is a linear subspace of \mathbb{R}^p of dimension k and consider the geometric median $\widehat{\boldsymbol{\mu}}_V^{\text{GM}}$ of the projected vectors $\mathbf{P}_V \mathbf{X}_1, \dots, \mathbf{P}_V \mathbf{X}_n$ that is

$$\widehat{\boldsymbol{\mu}}_V^{\text{GM}} \in \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^p} \sum_{i=1}^n \|\mathbf{P}_V \mathbf{X}_i - \boldsymbol{\mu}\|_2.$$

Lemma 2. *With probability at least $1 - \delta$, for all linear subspaces $V \subset \mathbb{R}^p$, we have*

$$\frac{\|\widehat{\boldsymbol{\mu}}_V^{\text{GM}} - \mathbf{P}_V \boldsymbol{\mu}^*\|_2}{\sqrt{\dim(V)}} \leq \frac{2\sqrt{\|\Sigma\|_{\text{op}}}}{1-2\varepsilon} \left(1 + \frac{\sqrt{\mathbf{r}_{\Sigma}} + \sqrt{2 \log(1/\delta)}}{\sqrt{n}} \right).$$

Proof. It follows from (Dalalyan and Minasyan, 2020, Lemma 2) that

$$\|\widehat{\boldsymbol{\mu}}_V^{\text{GM}} - \mathbf{P}_V \boldsymbol{\mu}^*\|_2 \leq \frac{2}{n(1-2\varepsilon)} \sum_{i=1}^n \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2.$$

Upper bounding the last term using Cauchy-Schwartz's inequality, one obtains

$$\|\hat{\mu}_V^{\text{GM}} - \mathbb{P}_V \mu^*\|_2 \leq \frac{2}{\sqrt{n}(1-2\varepsilon)} \left(\sum_{i=1}^n \|\mathbb{P}_V \xi_i\|_2^2 \right)^{1/2}.$$

Let e_1, \dots, e_k be any orthonormal basis of V , with $k = \dim(V)$. We have

$$\begin{aligned} \|\hat{\mu}_V^{\text{GM}} - \mathbb{P}_V \mu^*\|_2 &\leq \frac{2}{\sqrt{n}(1-2\varepsilon)} \left(\sum_{i=1}^n \sum_{\ell=1}^k |e_\ell^\top \xi_i|^2 \right)^{1/2} \\ &\leq \frac{2}{\sqrt{n}(1-2\varepsilon)} \left(k \sup_{\|e\|_2=1} \sum_{i=1}^n |e^\top \xi_i|^2 \right)^{1/2} \\ &= \frac{2}{\sqrt{n}(1-2\varepsilon)} \left(k \sup_{\|e\|_2=1} \|[\xi_1, \dots, \xi_n]^\top e\|_2^2 \right)^{1/2} \\ &= \frac{2\sqrt{k}}{\sqrt{n}(1-2\varepsilon)} \|\xi_1, \dots, \xi_n\|_{\text{op}}. \end{aligned}$$

By Corollary 5.35 in (Vershynin, 2012), the inequality

$$\|\xi_1, \dots, \xi_n\|_{\text{op}} \leq \|\Sigma\|_{\text{op}}^{1/2} (\sqrt{n} + \sqrt{r_\Sigma} + \sqrt{2 \log(1/\delta)})$$

holds with probability at least $1 - \delta$. This yields the desired inequality. \square

B.1.3 Bounding the number of filtered out observations

Throughout this section, we assume without loss of generality that $\|\Sigma\|_{\text{op}} = 1$.

Lemma 3. *Let τ and δ be two numbers from $(0, 1)$. Define*

$$z = 1 + \frac{\sqrt{r_\Sigma} + \sqrt{2 \log(1/\delta)}}{\sqrt{n\tau}} + \sqrt{2 + 2 \log(1/\tau)}.$$

With probability at least $1 - \delta$, we have

$$\sup_V \sum_{i=1}^n \mathbb{1}(\|\mathbb{P}_V \xi_i\|_2^2 > z^2 \dim(V)) \leq n\tau,$$

where the supremum is over all linear subspaces V of \mathbb{R}^p .

Proof. Let us define the random variable

$$T_n = \sup_V \sum_{i=1}^n \mathbb{1}(\|\mathbb{P}_V \xi_i\|_2^2 > z^2 \dim(V)).$$

In what follows, we write d_V for the dimension of the subspace V . We check that

$$\begin{aligned}
\mathbf{P}(T_n > n\tau) &\leq \mathbf{P}\left(\exists V \subset \mathbb{R}^p, \exists J \subset [n] \text{ with } |J| = n\tau, \text{ s.t. } \min_{i \in J} \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2^2 \geq z^2 \dim(V)\right) \\
&\leq \sum_{\substack{J \subseteq [n] \\ |J|=n\tau}} \mathbf{P}\left(\sup_V \min_{i \in J} \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2^2 / d_V \geq z^2\right) \\
&= \binom{n}{n\tau} \mathbf{P}\left(\sup_V \min_{i \in \{1, \dots, n\tau\}} \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2^2 / d_V \geq z^2\right) \\
&\leq \binom{n}{n\tau} \mathbf{P}\left(\sup_V \frac{1}{n\tau d_V} \sum_{i=1}^{n\tau} \|\mathbf{P}_V \boldsymbol{\xi}_i\|_2^2 \geq z^2\right). \tag{B.5}
\end{aligned}$$

Given a linear subspace $V \subset \mathbb{R}^p$ of dimension d_V , let $\mathbf{e}_1^V, \dots, \mathbf{e}_{d_V}^V$ be an orthonormal basis of V . Using (B.5), we get

$$\begin{aligned}
\mathbf{P}(T_n \geq n\tau) &\leq \binom{n}{n\tau} \mathbf{P}\left(\sup_V \frac{1}{n\tau d_V} \sum_{l=1}^{d_V} \sum_{i=1}^{n\tau} |\boldsymbol{\xi}_i^\top \mathbf{e}_l^V|^2 \geq z^2\right) \\
&\leq \binom{n}{n\tau} \mathbf{P}\left(\sup_{\|\mathbf{e}\|=1} \frac{1}{n\tau} \sum_{i=1}^{n\tau} |\boldsymbol{\xi}_i^\top \mathbf{e}|^2 \geq z^2\right) \\
&= \binom{n}{n\tau} \mathbf{P}\left(\|[\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{n\tau}]\|_{\text{op}} \geq z\sqrt{n\tau}\right).
\end{aligned}$$

By (Vershynin, 2012, Corollary 5.35), the inequality

$$\|[\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{n\tau}]\|_{\text{op}} \leq (\sqrt{n\tau} + \sqrt{\mathbf{r}_\Sigma} + \sqrt{2 \log(1/\delta_0)})$$

holds with probability at least $1 - \delta_0$. Taking $\delta_0 = \delta / \binom{n}{n\tau}$ and using the inequality $\log \binom{n}{n\tau} \leq n\tau(1 + \log(1/\tau))$, we can conclude that $\mathbf{P}(T_n \geq n\tau) \leq \delta$. This proves the lemma. \square

Lemma 4. Let $\tau \in (0, 1/2)$ and $\delta \in (0, 1/2)$ be arbitrary. Set

$$t = \frac{3 - 2\varepsilon^*}{1 - 2\varepsilon^*} \left(1 + \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n\tau}}\right) + \sqrt{2 + 2 \log(1/\tau)},$$

where $\varepsilon^* < 1/2$. Assume that $\{\mathbf{X}_i\}$ are drawn from $\text{GAC}(\boldsymbol{\mu}^*, \Sigma, \varepsilon)$ model with $\varepsilon \leq \varepsilon^*$. Then, with probability at least $1 - \delta$, for any linear subspace V of \mathbb{R}^p , the inequalities

$$\begin{aligned}
N_V &= \sum_{i \in \mathcal{I}} \mathbf{1}\left(\frac{\|\mathbf{P}_V \mathbf{X}_i - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2}{\sqrt{\dim(V)}} \leq t\right) \geq n(1 - \varepsilon - \tau), \\
\frac{\|\mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2}{\sqrt{\dim(V)}} &\leq \frac{2}{1 - 2\varepsilon} \left(1 + \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2 \log(2/\delta)}}{\sqrt{n}}\right),
\end{aligned}$$

where $\hat{\boldsymbol{\mu}}_V^{\text{GM}}$ is the geometric median of $\mathbf{P}_V \mathbf{X}_1, \dots, \mathbf{P}_V \mathbf{X}_n$.

Proof. To avoid unnecessary technicalities, we assume in this proof that $n\tau$ is an integer. We

also write

$$t_1 = \frac{2}{1-2\varepsilon} \left(1 + \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2\log(2/\delta)}}{\sqrt{n}} \right)$$

$$t_2 = 1 + \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2\log(2/\delta)}}{\sqrt{n\tau}} + \sqrt{2 + 2\log(1/\tau)},$$

so that $t \geq t_1 + t_2$. Simple algebra yields

$$\begin{aligned} N_V &= \sum_{i \in \mathcal{I}} \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i + \mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 \leq t\sqrt{\dim(V)}) \\ &\geq \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i + \mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 \leq t\sqrt{\dim(V)}) - n\varepsilon \\ &= n - n\varepsilon - \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i + \mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 > t\sqrt{\dim(V)}) \\ &\geq n - n\varepsilon - \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i\|_2 + \|\mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 > t\sqrt{\dim(V)}), \end{aligned}$$

where in the last step, we used the triangle inequality. According to Lemma 2, on an event Ω_1 of probability at least $1 - \delta/2$, we have $\|\mathbf{P}_V \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 \leq t_1 \sqrt{\dim(V)}$ for every V . This implies that on this event,

$$N_V \geq n - n\varepsilon - \sum_{i=1}^n \mathbb{1}(\|\mathbf{P}_V \boldsymbol{\xi}_i\|_2 > t_2 \sqrt{\dim(V)}), \quad \forall V \subset \mathbb{R}^p.$$

Using Lemma 3, we get that on an event Ω_2 of probability at least $1 - \delta/2$, the sum on the right hand side of the last display is less than $n\tau$. Therefore, on the intersection of the events Ω_1 and Ω_2 , we have $N_V \geq n(1 - \varepsilon - \tau)$ for every linear subspace V of \mathbb{R}^p . \square

B.1.4 Estimating the mean from a low-dimensional adversarial projection

In this section, we consider the following problem. We assume that for a q dimensional linear subspace V of \mathbb{R}^p , which can depend on the sample $\mathbf{X}_1, \dots, \mathbf{X}_n$, we observe the projected data $\mathbf{P}_V \mathbf{X}_1, \dots, \mathbf{P}_V \mathbf{X}_n$. The goal is to estimate the projected mean $\boldsymbol{\mu}_V^* = \mathbf{P}_V \boldsymbol{\mu}^* \in \mathbb{R}^q$. We will estimate $\boldsymbol{\mu}_V^*$ by the mean of filtered data points. More precisely, let $\hat{\boldsymbol{\mu}}_V^{\text{GM}}$ be the geometric median of $\mathbf{P}_V \mathbf{X}_1, \dots, \mathbf{P}_V \mathbf{X}_n$. We set

$$\mathcal{S}_V = \{i \in [n] : \|\mathbf{P}_V \mathbf{X}_i - \hat{\boldsymbol{\mu}}_V^{\text{GM}}\|_2 \leq t\sqrt{q}\}, \quad N_V = |\mathcal{S}_V|,$$

where t is a real number. We estimate $\boldsymbol{\mu}_V^*$ by

$$\hat{\boldsymbol{\mu}}_V = \mathbf{P}_V \bar{\mathbf{X}}_{\mathcal{S}_V} = \frac{1}{N_V} \sum_{i \in \mathcal{S}_V} \mathbf{P}_V \mathbf{X}_i.$$

Lemma 5. For every positive threshold $t > 0$, we have

$$\|\hat{\mu}_V - \mu_V^*\|_2 \leq \frac{\|P_V \bar{\xi}\|_2}{N_V} + \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V^c \cup \mathcal{O}} P_V \xi_i \right\|_2 + \frac{n\varepsilon(t\sqrt{q} + \|\hat{\mu}_V^{\text{GM}} - \mu_V^*\|_2)}{N_V}.$$

Proof. For this estimator, using the triangle inequality and the fact that $X_i = \mu^* + \xi_i$ for every $i \in \mathcal{I}$, we have

$$\begin{aligned} \|\hat{\mu}_V - \mu_V^*\|_2 &= \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V} (P_V X_i - P_V \mu^*) \right\|_2 \\ &\leq \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V \cap \mathcal{I}} P_V \xi_i \right\|_2 + \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V \cap \mathcal{O}} P_V (X_i - \mu^*) \right\|_2 \\ &\leq \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V \cap \mathcal{I}} P_V \xi_i \right\|_2 + \frac{n\varepsilon(t\sqrt{q} + \|\hat{\mu}_V^{\text{GM}} - \mu_V^*\|_2)}{N_V}. \end{aligned}$$

Using once again the triangle inequality, we arrive at

$$\begin{aligned} \|\hat{\mu}_V - \mu_V^*\|_2 &\leq \frac{\|P_V \bar{\xi}\|_2}{N_V} + \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V^c \cup \mathcal{O}} P_V \xi_i \right\|_2 + \frac{n\varepsilon(t\sqrt{q} + \|\hat{\mu}_V^{\text{GM}} - \mu_V^*\|_2)}{N_V} \\ &\leq \frac{\|P_V \bar{\xi}\|_2}{N_V} + \frac{1}{N_V} \left\| \sum_{i \in \mathcal{S}_V^c \cup \mathcal{O}} P_V \xi_i \right\|_2 + \frac{n\varepsilon(t\sqrt{q} + \|\hat{\mu}_V^{\text{GM}} - \mu_V^*\|_2)}{N_V}. \end{aligned}$$

This completes the proof. \square

B.1.5 Bounding stochastic errors

Throughout this section, without loss of generality, we assume that $\|\Sigma\|_{\text{op}} = 1$.

Lemma 6. For any positive integer $m \leq n$ and any $t > 0$, we have

$$\mathbf{P} \left(\max_{|S| \geq n-m} \left\| \frac{1}{|S|} \sum_{i \in S} P \xi_i \right\|_2 \leq \frac{n\|P \bar{\xi}_n\|_2}{n-m} + \frac{\sqrt{m(2\mathbf{r}_\Sigma + 3t)} + m\sqrt{3 \log(2ne/m)}}{n-m} \right) \geq 1 - e^{-t}.$$

Proof. Using the triangle inequality, one has

$$\begin{aligned} \frac{1}{|S|} \left\| \sum_{i \in S} P \xi_i \right\|_2 &\leq \frac{1}{n-m} \left\| \sum_{i \in S} P \xi_i \right\|_2 = \frac{1}{n-m} \left\| \sum_{i=1}^n P \xi_i - \sum_{i \in S^c} P \xi_i \right\|_2 \\ &\leq \frac{1}{n-m} \left\| \sum_{i=1}^n P \xi_i \right\|_2 + \frac{1}{n-m} \left\| \sum_{i \in S^c} P \xi_i \right\|_2 \\ &\leq \frac{n\|P \bar{\xi}_n\|_2}{n-m} + \frac{1}{n-m} \max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2. \end{aligned} \tag{B.6}$$

For $s \in [1, m]$, we choose t_s by

$$t_s = 3s \log \left(\frac{2ne}{s} \right) + 3t,$$

so that $t_s \leq t_m = 3m \log \left(\frac{2ne}{m} \right) + 3t$ and

$$\left(\frac{ne}{s} \right)^s e^{-t_s/3} \leq 2^{-s} e^{-t}.$$

For every J of cardinality m , the random variable $\|\sum_{j \in J} \xi_j\|_2^2$ has the same distribution as $m \sum_{j=1}^p \lambda_j(\Sigma) \alpha_j^2$, where $\alpha_1, \dots, \alpha_p$ are i.i.d. standard Gaussian. Hence, using the union bound, the well-known upper bound on the binomial coefficients and (Comminges and Dalalyan, 2012, Lemma 8), we have

$$\begin{aligned} \mathbf{P} \left(\max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2^2 \geq m(2\mathbf{r}_\Sigma + t_m) \right) &\leq \sum_{s=1}^m \binom{n}{s} \mathbf{P} \left(\left\| \sum_{i=1}^s \xi_i \right\|_2^2 \geq m(2\mathbf{r}_\Sigma + t_m) \right) \\ &\leq \sum_{s=1}^m \left(\frac{ne}{s} \right)^s \mathbf{P} \left(\left\| \sum_{i=1}^s \xi_i \right\|_2^2 \geq s(2\mathbf{r}_\Sigma + t_s) \right) \\ &\leq \sum_{s=1}^m \left(\frac{ne}{s} \right)^s e^{-t_s/3} \leq e^{-t}. \end{aligned}$$

This entails that with probability at least $1 - e^{-t}$, we have

$$\max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2 \leq \sqrt{m(2\mathbf{r}_\Sigma + t_m)} \leq \sqrt{m(2\mathbf{r}_\Sigma + 3t)} + m\sqrt{3 \log(2ne/m)}.$$

Combining this inequality with (B.6), we get the claim of the lemma. \square

In the two next lemmas, given a set $\mathcal{S} \subset [n]$, we look at the sample average and sample covariance matrix of the subsample $\{\mathbf{X}_i : i \in \mathcal{S}\}$,

$$\bar{\mathbf{X}}_{\mathcal{S}} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \mathbf{X}_i, \quad \hat{\Sigma}_{\mathcal{S}} = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \mathbf{X}_i \mathbf{X}_i^\top - \bar{\mathbf{X}}_{\mathcal{S}} \bar{\mathbf{X}}_{\mathcal{S}}^\top.$$

Lemma 7. *There exists a positive constant A such that, for any positive integer $m \leq n$ and any $t \geq 1$, with probability at least $1 - 2e^{-t}$, the inequality*

$$\|\hat{\Sigma}_{\mathcal{S}} - \Sigma\|_{\text{op}} \leq A \frac{\sqrt{n\mathbf{r}_\Sigma} + \mathbf{r}_\Sigma + m \log(2ne/m) + 2t}{n - m} + \|\bar{\xi}_{\mathcal{S}}\|_2^2$$

is satisfied for every $\mathcal{S} \subset [n]$ of cardinality $\geq n - m$.

Proof. The triangle inequality implies

$$\begin{aligned}\|\widehat{\Sigma}_{\mathcal{S}} - \Sigma\|_{\text{op}} &\leq \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} (\mathbf{X}_i - \boldsymbol{\mu})^{\otimes 2} - \Sigma \right\|_{\text{op}} + \|(\boldsymbol{\mu} - \bar{\mathbf{X}}_{\mathcal{S}})^{\otimes 2}\|_{\text{op}} \\ &= \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma \right\|_{\text{op}} + \|\bar{\boldsymbol{\xi}}_{\mathcal{S}}\|_2^2.\end{aligned}\tag{B.7}$$

Using again the triangle inequality, one gets

$$\begin{aligned}\left\| \sum_{i \in \mathcal{S}} \boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma \right\|_{\text{op}} &\leq \left\| \sum_{i \in \mathcal{S}} (\boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma) \right\|_{\text{op}} \\ &\leq \left\| \sum_{i=1}^n (\boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma) \right\|_{\text{op}} + \left\| \sum_{i \in \mathcal{S}^C} (\boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma) \right\|_{\text{op}} \\ &\leq \left\| \sum_{i=1}^n (\boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma) \right\|_{\text{op}} + \max_{|J| \leq m} \left\| \sum_{i \in J} (\boldsymbol{\xi}_i \boldsymbol{\xi}_i^{\top} - \Sigma) \right\|_{\text{op}}.\end{aligned}\tag{B.8}$$

In view of (Koltchinskii and Lounici, 2017, Theorems 4 and 5), one can show that there exists a positive universal constant A_1 such that for every $t \geq 1$ and every set J of cardinality s , the inequality

$$\left\| \sum_{j \in J} \boldsymbol{\xi}_j \boldsymbol{\xi}_j^{\top} - s\Sigma \right\|_{\text{op}} \leq A_1(\sqrt{m\mathbf{r}_{\Sigma}} + \mathbf{r}_{\Sigma} + t)$$

is satisfied with probability at least $1 - e^{-t}$. We define t_s by

$$t_s = s \log \left(\frac{2ne}{s} \right) + t,$$

so that $t_s \leq t_m = m \log \left(\frac{2ne}{m} \right) + t$ and

$$\left(\frac{ne}{s} \right)^s e^{-t_s} \leq 2^{-s} e^{-t}.$$

Applying the union bound and the well-known upper bound on the binomial coefficients, this yields

$$\begin{aligned}\mathbf{P} \left(\max_{|J| \leq m} \left\| \sum_{j \in J} (\boldsymbol{\xi}_j \boldsymbol{\xi}_j^{\top} - \Sigma) \right\|_{\text{op}} \geq A_1(\sqrt{m\mathbf{r}_{\Sigma}} + \mathbf{r}_{\Sigma} + t_m) \right) \\ \leq \sum_{s=1}^m \binom{n}{s} \mathbf{P} \left(\left\| \sum_{j=1}^s \boldsymbol{\xi}_j \boldsymbol{\xi}_j^{\top} - s\Sigma \right\|_{\text{op}} \geq A_1(\sqrt{m\mathbf{r}_{\Sigma}} + \mathbf{r}_{\Sigma} + t_s) \right) \\ \leq \sum_{s=1}^m \left(\frac{ne}{s} \right)^s e^{-t_s} \leq e^{-t}.\end{aligned}$$

One deduces from (B.8) that, with probability at least $1 - 2e^{-t}$,

$$\begin{aligned} \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} \xi_i \xi_i^\top - \Sigma \right\|_{\text{op}} &\leq A_1 \frac{(\sqrt{n} + \sqrt{m})\sqrt{\mathbf{r}_\Sigma} + 2\mathbf{r}_\Sigma + m \log(2ne/m) + 2t}{n - m} \\ &\leq A \frac{\sqrt{n\mathbf{r}_\Sigma} + \mathbf{r}_\Sigma + m \log(2ne/m) + 2t}{n - m}. \end{aligned}$$

Combining this with (B.7), one gets the claim of the lemma. \square

Lemma 8. *For any positive integer $m \leq n$ and any $t > 0$, with probability at least $1 - 4e^{-t}$, the inequality*

$$\|\hat{\Sigma}_S - \Sigma\|_{\text{op}} \leq \frac{5p + (8 \log(2ne/m) + 2)m + 7t}{n - m} + 2 \frac{\sqrt{p} + \sqrt{t}}{\sqrt{n} - \sqrt{m}} + \|\bar{\xi}_S\|_2^2$$

is satisfied for every $S \subset [n]$ of cardinality $\geq n - m$.

Proof. In this proof, without loss of generality, we assume that the matrix Σ is invertible. In view of (B.7) and (B.8), we have

$$\|\hat{\Sigma}_S - \Sigma\|_{\text{op}} \leq \left\| \frac{1}{|\mathcal{S}|} \sum_{i=1}^n (\xi_i \xi_i^\top - \Sigma) \right\|_{\text{op}} + \max_{|J| \leq m} \left\| \frac{1}{|\mathcal{S}|} \sum_{i \in J} (\xi_i \xi_i^\top - \Sigma) \right\|_{\text{op}} + \|\bar{\xi}_S\|_2^2. \quad (\text{B.9})$$

Let us define $\zeta_i := \Sigma^{-1/2} \xi_i$ for all $i \in [n]$. For every set J of cardinality s , it holds that

$$\begin{aligned} \left\| \sum_{i \in J} (\xi_i \xi_i^\top - \Sigma) \right\|_{\text{op}} &\leq \|\Sigma\|_{\text{op}} \left\| \sum_{i \in J} (\zeta_i \zeta_i^\top - \mathbf{I}_p) \right\|_{\text{op}} \\ &= \max \left(\lambda_{\max} \left(\sum_{i \in J} \zeta_i \zeta_i^\top \right) - s, s - \lambda_{\min} \left(\sum_{i \in J} \zeta_i \zeta_i^\top \right) \right) \\ &= \max \left(\sigma_{\max}^2(\zeta_J) - s, s - \sigma_{\min}^2(\zeta_J) \right), \end{aligned}$$

where ζ_J is the $s \times n$ random matrix obtained by concatenating the vectors ζ_i with $i \in J$. By (Vershynin, 2012, Corollary 5.35), we know that for every $x > 0$

$$\sqrt{s} - \sqrt{p} - x \leq \sigma_{\min}(\zeta_J) \leq \sigma_{\max}(\zeta_J) \leq \sqrt{s} + \sqrt{p} + x$$

is true with probability at least $1 - 2e^{-x^2/2}$. This yields¹

$$\begin{aligned} \left\| \sum_{i \in J} (\xi_i \xi_i^\top - \Sigma) \right\|_{\text{op}} &\leq \max \left((\sqrt{p} + x)(2\sqrt{s} + \sqrt{p} + x), (\sqrt{p} + x)(2\sqrt{s} - \sqrt{p} - x) \right) \\ &\leq p + x^2 + 2\sqrt{ps} + 2x\sqrt{p} + 2x\sqrt{s} \end{aligned}$$

with probability at least $1 - 2e^{-x^2/2}$. By applying the same technique as in the proof of

¹We provide the argument only in the case $\sqrt{s} \geq \sqrt{p} + x$, but the conclusion is true for every value s .

Lemma 7, we can set

$$t_s = 2\sqrt{s \log\left(\frac{2ne}{s}\right)} + t,$$

and obtain

$$\mathbf{P}\left(\max_{|J| \leq m} \left\| \sum_{j \in J} (\xi_j \xi_j^\top - \Sigma) \right\|_{\text{op}} \geq p + t_m^2 + 2\sqrt{pm} + 2t_m\sqrt{p} + 2t_m\sqrt{m}\right) \leq 2e^{-t}.$$

Hence, going back to (B.9), we can show that the inequalities

$$\begin{aligned} \|\widehat{\Sigma}_S - \Sigma\|_{\text{op}} &\leq \frac{p + t + 2\sqrt{pn} + 2\sqrt{tp} + 2\sqrt{tn}}{n - m} + \frac{p + 4t + 4m \log(2ne/m) + 2\sqrt{pm}}{n - m} \\ &\quad + \frac{4(\sqrt{p} + \sqrt{m})\sqrt{m \log(2ne/m)} + t}{n - m} + \|\bar{\xi}_S\|_2^2 \\ &\leq \frac{5p + 8m \log(2ne/m) + 2m + 7t}{n - m} + \frac{2(\sqrt{p} + \sqrt{t})(\sqrt{n} + \sqrt{m})}{n - m} + \|\bar{\xi}_S\|_2^2 \end{aligned}$$

hold with probability at least $1 - 4e^{-t}$, and this proves the lemma. \square

B.1.6 Putting all the pieces together

All the ingredients provided, we can now compile the complete proof of Theorem 11.

Taking $\mathbf{U}_L := \mathbf{V}_L$, the algorithm detailed in (3.2) returns $\widehat{\mu}^{\text{SDR}} = \sum_{\ell=0}^L \widehat{\mu}^{(\ell)}$ with $\widehat{\mu}^{(\ell)} \in \mathcal{U}_\ell = \text{Im}(\mathbf{V}_\ell \mathbf{U}_\ell^\top)$ for every $\ell \in \{0, \dots, L\}$ where the two-by-two orthogonal subspaces $\mathcal{U}_0, \dots, \mathcal{U}_L$ span the whole space \mathbb{R}^p . This allows us to decompose the risk:

$$\begin{aligned} \|\widehat{\mu}^{\text{SDR}} - \mu^*\|_2^2 &= \sum_{\ell=0}^L \|\widehat{\mu}^{(\ell)} - \mathbf{P}_{\mathcal{U}_\ell} \mu^*\|_2^2 = \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell}(\bar{\mathbf{X}}_\ell - \mu^*)\|_2^2 \\ &= \sum_{\ell=0}^L \|\mathbf{P}_\ell(\mathbf{V}_\ell^\top \bar{\mathbf{X}}_\ell - \mathbf{V}_\ell^\top \mu^*)\|_2^2, \end{aligned}$$

where $\mathbf{P}_\ell := \mathbf{U}_\ell^\top \mathbf{U}_\ell$ is the projection matrix projecting onto the subspace of \mathbb{R}^{p_ℓ} spanned by the bottom $p_\ell - p_{\ell+1}$ eigenvectors of $\mathbf{V}_\ell^\top (\widehat{\Sigma}^{(\ell)} - \Sigma) \mathbf{V}_\ell$ for $\ell = 0, \dots, L$ with the convention that $p_{L+1} = 0$.

For $\ell \in \{0, \dots, L-1\}$, we intend to apply Proposition 6 to $\mathbf{Z}_i = \mathbf{V}_\ell^\top \mathbf{X}_i$ and $\mu^Z = \mathbf{V}_\ell^\top \mu^*$ in order to upper bound the error term $\text{Err}_\ell := \|\mathbf{P}_\ell(\mathbf{V}_\ell^\top \bar{\mathbf{X}}_\ell - \mathbf{V}_\ell^\top \mu^*)\|_2$. Using the inequalities

$$\|\mathbf{V}^\top (\widehat{\Sigma}^{(\ell)} - \Sigma) \mathbf{V}\|_{\text{op}} \leq \|\widehat{\Sigma}^{(\ell)} - \Sigma\|_{\text{op}}, \quad \lambda_{p_\ell}(\mathbf{V}^\top \Sigma \mathbf{V}) \leq \lambda_p(\Sigma), \quad \lambda_1(\mathbf{V}^\top \Sigma \mathbf{V}) \geq \lambda_1(\Sigma)$$

that are true for every orthogonal matrix \mathbf{V} , and keeping in mind the definition of P_ℓ , we get

$$\text{Err}_\ell \leq \left\{ 2\omega_{\mathcal{O}} \|\widehat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + \frac{\omega_{\mathcal{O}}^2}{1 - \omega_{\mathcal{O}}} \left((\lambda_p - \lambda_1)(\boldsymbol{\Sigma}) + \frac{\delta_\ell^2}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_\ell \mathbf{V}_\ell^\top \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2,$$

where we have used the notation

$$\omega_{\mathcal{O}} = \max_{\ell} \frac{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|}{|\mathcal{S}^{(\ell)}|}, \quad \boldsymbol{\xi}_i = \mathbf{X}_i - \boldsymbol{\mu}^*$$

and $\delta_\ell = \inf_{\boldsymbol{\mu}} \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top (\mathbf{X}_i - \boldsymbol{\mu})\|_2$. Note that when \mathcal{O} and $(\mathcal{S}^{(\ell)} \cap \mathcal{I})^c$ are of cardinality less than $n\varepsilon$ and $n(\varepsilon + \tau)$, respectively, we have

$$\frac{|\mathcal{S}^{(\ell)}|}{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|} = \frac{|\mathcal{S}^{(\ell)} \cap \mathcal{I}| + |\mathcal{S}^{(\ell)} \cap \mathcal{O}|}{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|} = \frac{|\mathcal{S}^{(\ell)} \cap \mathcal{I}|}{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|} + 1 \geq \frac{n(1 - \varepsilon - \tau)}{n\varepsilon} + 1 = \frac{1 - \tau}{\varepsilon}$$

and, therefore, $\omega_{\mathcal{O}} \leq \varepsilon/(1 - \tau)$. We set $\eta := \varepsilon + \tau \leq 3/4$ and apply Lemma 4 to infer that $\omega_{\mathcal{O}} \leq \omega_{\mathcal{O}}/(1 - \omega_{\mathcal{O}}) \leq \varepsilon/(1 - \eta) \leq 4\varepsilon$ is true with probability at least $1 - \delta$. Furthermore, we know that $\delta_\ell \leq \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(\ell)}\|_2 \leq t\sqrt{p_\ell}$. This yields

$$\text{Err}_\ell \leq \left\{ 8\varepsilon \|\widehat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + 16\varepsilon^2 \left((\lambda_p - \lambda_1)(\boldsymbol{\Sigma}) + \frac{t^2 p_\ell}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2.$$

Let us introduce the shorthand

$$T_1 = \max_{\ell \in [L]} \|\widehat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + \varepsilon(\lambda_p - \lambda_1)(\boldsymbol{\Sigma}).$$

This leads to

$$\text{Err}_\ell \leq \left\{ 8\varepsilon T_1 + \frac{16\varepsilon^2 t^2 p_\ell}{p_{\ell+1}} \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2. \quad (\text{B.10})$$

For the last error term, since $p_L = 1$ we have by Lemma 5

$$\begin{aligned} \text{Err}_L &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(L)}}\|_2 + \frac{n\varepsilon(t\sqrt{p_L} + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)}{|\mathcal{S}^{(L)}|} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(L)}}\|_2 + \frac{\varepsilon t + \varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2}{1 - \eta} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(L)}}\|_2 + 4\varepsilon t + 4\varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2. \end{aligned} \quad (\text{B.11})$$

Combining (B.10), (B.11), inequality $p_\ell \leq ep_{\ell+1}$, as well as the Minkowski inequality, we get

$$\begin{aligned}
\|\boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}^{\text{SDR}}\|_2 &= \left\{ \sum_{\ell=0}^L \text{Err}_\ell^2 \right\}^{1/2} \\
&\leq \left\{ 8\varepsilon L(T_1 + e\varepsilon t^2) + 16\varepsilon^2(t + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)^2 \right\}^{1/2} + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} \\
&\leq 2\sqrt{2\varepsilon L T_1} + 9\varepsilon t\sqrt{L} + 4\varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2}. \tag{B.12}
\end{aligned}$$

To ease notation, let us set

$$r_n = \left(\frac{2\mathbf{r}_\Sigma + 3\log(2/\delta)}{n} \right)^{1/2}.$$

In view of Lemma 6, with probability at least $1 - \delta$, we have

$$\begin{aligned}
\left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} &\leq \left\{ \sum_{\ell=0}^L \left(\frac{\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2}{1-\eta} + \frac{r_n\sqrt{\eta} + \eta\sqrt{3\log(2e/\eta)}}{1-\eta} \right)^2 \right\}^{1/2} \\
&\leq \left\{ \sum_{\ell=0}^L (4\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2 + 4r_n\sqrt{\eta} + 10\eta\sqrt{\log(2/\eta)})^2 \right\}^{1/2} \\
&\leq 4\|\bar{\boldsymbol{\xi}}_n\|_2 + 4r_n\sqrt{\eta L} + 10\eta\sqrt{L\log(2/\eta)}.
\end{aligned}$$

Since the random variable $\|\bar{\boldsymbol{\xi}}_n\|_2^2$ has the same distribution as $\frac{1}{n} \sum_{j=1}^p \lambda_j(\Sigma) \gamma_j^2$, where $\gamma_1, \dots, \gamma_p$ are i.i.d. standard Gaussian, by (Comminges and Dalalyan, 2012, Lemma 8) we have

$$\|\bar{\boldsymbol{\xi}}_n\|_2^2 \leq \frac{2\mathbf{r}_\Sigma + 3\log(2/\delta)}{n} = r_n^2$$

with probability at least $1 - \delta$. Therefore, with probability at least $1 - 2\delta$,

$$\left(\sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right)^{1/2} \leq 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}. \tag{B.13}$$

Next, Lemma 4 and the fact that $p_L = \dim(\mathcal{U}_L) = 1$ imply that with probability at least $1 - \delta$

$$\|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 \leq \frac{2(1 + \sqrt{2}r_n)}{1 - 2\varepsilon}. \tag{B.14}$$

Recall that we have chosen t in such a way that

$$t \leq \frac{3(1 + \sqrt{2}r_n/\sqrt{\tau})}{1 - 2\varepsilon^*} + 1.6\sqrt{\log(2/\tau)}. \tag{B.15}$$

Combining (B.12), (B.13), (B.14) and (B.15), we arrive at the inequality

$$\begin{aligned}\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 &\leq 2\sqrt{2\varepsilon LT_1} + 9\varepsilon t\sqrt{L} + \frac{8\varepsilon(1 + \sqrt{2}r_n)}{1 - 2\varepsilon} + 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)} \\ &\leq 2\sqrt{2\varepsilon LT_1} + \frac{27\varepsilon\sqrt{L}(1 + \sqrt{2}r_n/\sqrt{\tau})}{1 - 2\varepsilon^*} + 14.4\varepsilon\sqrt{L\log(2/\tau)} \\ &\quad + \frac{8\varepsilon(1 + \sqrt{2}r_n)}{1 - 2\varepsilon} + 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}\end{aligned}$$

that holds with probability at least $1 - 3\delta$. In the upper bound obtained above, only the term T_1 remains random. We can upper bound this term using Lemma 7. It implies that with probability at least $1 - 2\delta$, we have

$$\begin{aligned}T_1 &\leq A \frac{\sqrt{n\mathbf{r}_\Sigma} + \mathbf{r}_\Sigma + n\eta\log(2e/\eta) + 2\log(1/\delta)}{n(1 - \eta)} + (4r_n(1 + \sqrt{\eta}) + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon \\ &\leq 2A(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta)) + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon.\end{aligned}$$

Consequently,

$$\begin{aligned}\sqrt{\varepsilon T_1} &\leq \{2A\varepsilon(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta))\}^{1/2} + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})\sqrt{\varepsilon} + \varepsilon \\ &\leq \{2A\varepsilon(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta))\}^{1/2} + 5.4r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon \\ &\leq \{2A\varepsilon(\sqrt{2}r_n + 4\eta\log(2/\eta))\}^{1/2} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon \\ &\leq \sqrt{2\sqrt{2}Ar_n\varepsilon} + 2\sqrt{2A\varepsilon\eta\log(2/\eta)} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon \\ &\leq \varepsilon + Ar_n/\sqrt{2} + 2\eta\sqrt{2A\log(2/\eta)} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon \\ &\leq (A/\sqrt{2} + \sqrt{A} + 5.4)r_n + (7.1 + 2\sqrt{2A})\tau\sqrt{\log(2/\tau)} + (9.1 + 2\sqrt{2A})\varepsilon\sqrt{\log(2/\varepsilon)}.\end{aligned}$$

These inequalities imply that there exists a universal constant C such that

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C(r_n + \tau\sqrt{\log(2/\tau)} + \varepsilon\sqrt{\log(2/\varepsilon)} + r_n\varepsilon/\sqrt{\tau})\sqrt{L}}{1 - 2\varepsilon^*}. \quad (\text{B.16})$$

We choose

$$\tau = \frac{1}{4} \bigwedge \frac{\bar{r}_n}{\sqrt{\log_+(2/\bar{r}_n)}}, \quad \text{with} \quad \bar{r}_n = \frac{\sqrt{\mathbf{r}_\Sigma} + \sqrt{2\log(2/\delta)}}{\sqrt{n}}.$$

Note that $r_n \leq \sqrt{2}\bar{r}_n$ and, furthermore, $\tau = 1/4$ whenever $\bar{r}_n \geq 1/2$. Therefore, $r_n\varepsilon/\sqrt{\tau} \leq r_n + \varepsilon$. Inserting this value of τ in (B.16) leads to

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C(r_n + \varepsilon\sqrt{\log(2/\varepsilon)})\sqrt{L}}{1 - 2\varepsilon^*}.$$

where C is a universal constant, the value of which is not necessarily the same in different places where it appears. Replacing r_n by its expression, and upper bounding L by $2\log p$, we

arrive at

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C\sqrt{\log p}}{1-2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_{\Sigma}}{n}} + \varepsilon\sqrt{\log(2/\varepsilon)} + \sqrt{\frac{\log(1/\delta)}{n}} \right).$$

Note that this inequality holds true on an event of probability at least $1 - 5\delta$.

To prove the second part of the theorem, we use Lemma 8 instead of Lemma 7 in order to bound the term T_1 . Moreover, in the definitions of r_n and \bar{r}_n the effective rank \mathbf{r}_{Σ} is replaced by the dimension p . Then, with probability at least $1 - 2\delta$, we have

$$\begin{aligned} T_1 &\leq \frac{5p + n\eta(8\log(2e/\eta) + 2) + 7\log(2/\delta)}{n(1-\eta)} + 2\frac{\sqrt{p} + \sqrt{\log(2/\delta)}}{\sqrt{n}(1-\sqrt{\eta})} + \\ &\quad + (4r_n(1+\sqrt{\eta}) + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon \\ &\leq (10r_n^2 + 15r_n + 72\eta\log(2/\eta)) + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon. \end{aligned}$$

Then, repeating the same steps as for the previous case where the effective rank is used instead of dimension we arrive at the following inequality

$$\begin{aligned} \sqrt{\varepsilon T_1} &\leq \{\varepsilon(10r_n^2 + 15r_n + 72\eta\log(2/\eta))\}^{1/2} + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})\sqrt{\varepsilon} + \varepsilon \\ &\leq 12r_n + 16\tau\sqrt{\log(2/\tau)} + 18\varepsilon\sqrt{\log(2/\varepsilon)}. \end{aligned}$$

Combining the obtained inequalities, plugging in the values of τ and r_n and bounding L by $2\log p$ we arrive at a final bound which reads as

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{156\sqrt{2\log p}}{1-2\varepsilon^*} \left(\sqrt{\frac{2p}{n}} + \varepsilon\sqrt{\log(2/\varepsilon)} + \sqrt{\frac{3\log(2/\delta)}{n}} \right),$$

which concludes the proof.

B.2 Proof of Theorem 13

The proof follows the same steps as that of Theorem 11. The assumption $\|\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma$ gives upper and lower bounds on the effective rank of $\tilde{\Sigma}$ using that of Σ , which we formulate as a separate lemma. Therefore, the choice of the threshold parameter \tilde{t}_{γ} stated in Theorem 13 makes Lemmas 3 and 4 applicable to this case as well. To bound the operator norm of $\|\hat{\Sigma}^{(\ell)} - \tilde{\Sigma}\|_{\text{op}}$ we make use of the assumption $\|\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma$ and Lemma 7 using triangle inequality. We provide the full proof for reader's convenience.

Before proceeding with the proof we first formulate and prove an auxiliary lemma for bounding the effective rank of $\tilde{\Sigma}$ using that of Σ .

Lemma 9. *Let Σ and $\tilde{\Sigma}$ be symmetric positive definite matrices such that $\|\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma$*

$\mathbf{I}_p\|_{\text{op}} \leq \gamma$. Then,

$$\mathbf{r}_{\Sigma} \cdot \frac{1-\gamma}{1+\gamma} \leq \mathbf{r}_{\tilde{\Sigma}} \leq \mathbf{r}_{\Sigma} \cdot \frac{1+\gamma}{1-\gamma}.$$

Proof of Lemma 9. We start with upper- and lower-bounding the operator norm of $\tilde{\Sigma}$. Using triangle inequality we have

$$|\|\tilde{\Sigma}\|_{\text{op}} - \|\Sigma\|_{\text{op}}| \leq \|\tilde{\Sigma} - \Sigma\|_{\text{op}} \leq \|\Sigma\|_{\text{op}} \cdot \|\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2} - \mathbf{I}_p\|_{\text{op}} \leq \gamma\|\Sigma\|_{\text{op}}.$$

This readily yields

$$(1-\gamma)\|\Sigma\|_{\text{op}} \leq \|\tilde{\Sigma}\|_{\text{op}} \leq (1+\gamma)\|\Sigma\|_{\text{op}}. \quad (\text{B.17})$$

Moreover, for any pair of positive definite matrices $\mathbf{A}, \mathbf{B} \succeq 0$ the following holds $\text{Tr}(\mathbf{A})\lambda_1(\mathbf{B}) \leq \text{Tr}(\mathbf{A}\mathbf{B}) \leq \text{Tr}(\mathbf{A}) \cdot \|\mathbf{B}\|_{\text{op}}$. Hence, combining the cyclic property of trace, the trace inequality and the fact that the spectrum of the matrix $\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2}$ is between $1-\gamma$ and $1+\gamma$, we get both the upper and the lower bounds for $\text{Tr}(\tilde{\Sigma})$. The upper bound reads as

$$\begin{aligned} \text{Tr}(\tilde{\Sigma}) &= \text{Tr}((\Sigma^{1/2}\Sigma^{-1/2})\tilde{\Sigma}(\Sigma^{-1/2}\Sigma^{1/2})) = \text{Tr}(\Sigma(\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2})) \\ &\leq \|\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2}\|_{\text{op}} \text{Tr}(\Sigma) \leq (1+\gamma) \text{Tr}(\Sigma). \end{aligned} \quad (\text{B.18})$$

Similarly, the lower bound can be obtained as follows

$$\text{Tr}(\tilde{\Sigma}) \geq \lambda_1(\Sigma^{-1/2}\tilde{\Sigma}\Sigma^{-1/2}) \text{Tr}(\Sigma) \geq (1-\gamma) \text{Tr}(\Sigma). \quad (\text{B.19})$$

Therefore, combining (B.17) and (B.19) we get the lower bound for $\mathbf{r}_{\tilde{\Sigma}}$, while combining (B.17) and (B.18) yields the upper bound, concluding the proof of the lemma. \square

Taking $\mathbf{U}_L := \mathbf{V}_L$, the Algorithm 2 returns $\hat{\boldsymbol{\mu}}^{\text{SDR}} = \sum_{\ell=0}^L \hat{\boldsymbol{\mu}}^{(\ell)}$ with $\hat{\boldsymbol{\mu}}^{(\ell)} \in \mathcal{U}_{\ell} = \text{Im}(\mathbf{V}_{\ell}\mathbf{U}_{\ell}^{\top})$ for every $\ell \in \{0, \dots, L\}$ where the two-by-two orthogonal subspaces $\mathcal{U}_0, \dots, \mathcal{U}_L$ span the whole space \mathbb{R}^p . This allows us to decompose the risk:

$$\begin{aligned} \|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2^2 &= \sum_{\ell=0}^L \|\hat{\boldsymbol{\mu}}^{(\ell)} - \mathbf{P}_{\mathcal{U}_{\ell}}\boldsymbol{\mu}^*\|_2^2 = \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_{\ell}}(\bar{\mathbf{X}}_{\ell} - \boldsymbol{\mu}^*)\|_2^2 \\ &= \sum_{\ell=0}^L \|\mathbf{P}_{\ell}(\mathbf{V}_{\ell}^{\top}\bar{\mathbf{X}}_{\ell} - \mathbf{V}_{\ell}^{\top}\boldsymbol{\mu}^*)\|_2^2, \end{aligned}$$

where $\mathbf{P}_{\ell} := \mathbf{U}_{\ell}^{\top}\mathbf{U}_{\ell}$ is the projection matrix projecting onto the subspace of $\mathbb{R}^{p_{\ell}}$ spanned by the bottom $p_{\ell} - p_{\ell+1}$ eigenvectors of $\mathbf{V}_{\ell}^{\top}(\hat{\Sigma}^{(\ell)} - \tilde{\Sigma})\mathbf{V}_{\ell}$ for $\ell = 0, \dots, L$ with the convention that $p_{L+1} = 0$.

For $\ell \in \{0, \dots, L-1\}$, we intend to apply Proposition 6 to $\mathbf{Z}_i = \mathbf{V}_{\ell}^{\top}\mathbf{X}_i$ and $\boldsymbol{\mu}^Z = \mathbf{V}_{\ell}^{\top}\boldsymbol{\mu}^*$ in

order to upper bound the error term $\text{Err}_\ell := \|\mathbf{P}_\ell(\mathbf{V}_\ell^\top \bar{\mathbf{X}}_\ell - \mathbf{V}_\ell^\top \boldsymbol{\mu}^*)\|_2$. Using the inequalities

$$\|\mathbf{V}^\top(\hat{\boldsymbol{\Sigma}}^{(\ell)} - \tilde{\boldsymbol{\Sigma}})\mathbf{V}\|_{\text{op}} \leq \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \tilde{\boldsymbol{\Sigma}}\|_{\text{op}}, \quad \lambda_{p_\ell}(\mathbf{V}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{V}) \leq \lambda_p(\tilde{\boldsymbol{\Sigma}}), \quad \lambda_1(\mathbf{V}^\top \tilde{\boldsymbol{\Sigma}} \mathbf{V}) \geq \lambda_1(\tilde{\boldsymbol{\Sigma}})$$

that are true for every orthogonal matrix \mathbf{V} , and keeping in mind the definition of \mathbf{P}_ℓ , we get

$$\text{Err}_\ell \leq \left\{ 2\omega_{\mathcal{O}} \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \tilde{\boldsymbol{\Sigma}}\|_{\text{op}} + \frac{\omega_{\mathcal{O}}^2}{1 - \omega_{\mathcal{O}}} \left((\lambda_p - \lambda_1)(\tilde{\boldsymbol{\Sigma}}) + \frac{\delta_\ell^2}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_\ell \mathbf{V}_\ell^\top \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(\ell)}}\|_2,$$

where we have used the notation

$$\omega_{\mathcal{O}} = \max_{\ell} \frac{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|}{|\mathcal{S}^{(\ell)}|}, \quad \boldsymbol{\xi}_i = \mathbf{X}_i - \boldsymbol{\mu}^*$$

and $\delta_\ell = \inf_{\boldsymbol{\mu}} \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top (\mathbf{X}_i - \boldsymbol{\mu})\|_2$. Note that when \mathcal{O} and $(\mathcal{S}_T^{(\ell)})^c$ are of cardinality less than $n\varepsilon$ and $n(\varepsilon + \tau)$, respectively, we have $\omega_{\mathcal{O}} \leq \varepsilon/(1 - \tau)$ and

$$\frac{\omega_{\mathcal{O}}}{1 - \omega_{\mathcal{O}}} \leq \frac{\varepsilon}{1 - \varepsilon - \tau}.$$

Since $\mathbf{C}_\gamma \mathbf{r}_{\tilde{\Sigma}} \geq \mathbf{r}_{\Sigma}$ (by Lemma 9) then Lemma 4 holds for the new threshold \tilde{t}_γ as well. We set $\eta := \varepsilon + \tau \leq 3/4$ and apply Lemma 4 to infer that $\omega_{\mathcal{O}} \leq \varepsilon/(1 - \eta) \leq 4\varepsilon$ is true with probability at least $1 - \delta$. Furthermore, we know that $\delta_\ell \leq \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(\ell)}\|_2 \leq \tilde{t}_\gamma \sqrt{p_\ell}$. This yields

$$\text{Err}_\ell \leq \left\{ 8\varepsilon \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \tilde{\boldsymbol{\Sigma}}\|_{\text{op}} + 16\varepsilon^2 \left((\lambda_p - \lambda_1)(\tilde{\boldsymbol{\Sigma}}) + \frac{\tilde{t}_\gamma^2 p_\ell}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(\ell)}}\|_2.$$

Let us introduce the shorthand $\tilde{T}_1 = \max_{\ell \in [L]} \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \tilde{\boldsymbol{\Sigma}}\|_{\text{op}} + \varepsilon(\lambda_p - \lambda_1)(\tilde{\boldsymbol{\Sigma}})$. This leads to

$$\text{Err}_\ell \leq \left\{ 8\varepsilon \tilde{T}_1 + \frac{16\varepsilon^2 \tilde{t}_\gamma^2 p_\ell}{p_{\ell+1}} \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(\ell)}}\|_2. \quad (\text{B.20})$$

For the last error term, since $p_L = 1$ we have by Lemma 5

$$\begin{aligned} \text{Err}_L &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(L)}}\|_2 + \frac{n\varepsilon(\tilde{t}_\gamma \sqrt{p_L} + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)}{|\mathcal{S}^{(L)}|} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(L)}}\|_2 + \frac{\varepsilon \tilde{t}_\gamma + \varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2}{1 - \eta} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{\mathcal{S}_T^{(L)}}\|_2 + 4\varepsilon \tilde{t}_\gamma + 4\varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2. \end{aligned} \quad (\text{B.21})$$

Combining (B.20), (B.21), inequality $p_\ell \leq ep_{\ell+1}$, as well as the Minkowski inequality, we get

$$\begin{aligned}
\|\boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}^{\text{SDR}}\|_2 &= \left\{ \sum_{\ell=0}^L \text{Err}_\ell^2 \right\}^{1/2} \\
&\leq \left\{ 8\varepsilon L(\tilde{T}_1 + e\varepsilon t^2) + 16\varepsilon^2(\tilde{t}_\gamma + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)^2 \right\}^{1/2} + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} \\
&\leq 2\sqrt{2\varepsilon L\tilde{T}_1} + 9\varepsilon\tilde{t}_\gamma\sqrt{L} + 4\varepsilon\|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2}. \tag{B.22}
\end{aligned}$$

To ease notation, let us set

$$r_n = \left(\frac{2\mathbf{r}_\Sigma + 3\log(2/\delta)}{n} \right)^{1/2}.$$

In view of Lemma 6, with probability at least $1 - \delta$, we have

$$\begin{aligned}
\left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} &\leq \left\{ \sum_{\ell=0}^L \left(\frac{\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2}{1-\eta} + \frac{r_n\sqrt{\eta} + \eta\sqrt{3\log(2e/\eta)}}{1-\eta} \right)^2 \right\}^{1/2} \\
&\leq \left\{ \sum_{\ell=0}^L (4\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2 + 4r_n\sqrt{\eta} + 10\eta\sqrt{\log(2/\eta)})^2 \right\}^{1/2} \\
&\leq 4\|\bar{\boldsymbol{\xi}}_n\|_2 + 4r_n\sqrt{\eta L} + 10\eta\sqrt{L\log(2/\eta)}.
\end{aligned}$$

Since $\|\bar{\boldsymbol{\xi}}_n\|_2^2$ has the same distribution as $\frac{1}{n} \sum_{j=1}^p \lambda_j(\Sigma) \gamma_j^2$, where $\gamma_1, \dots, \gamma_p$ are i.i.d. standard Gaussian, by (Comminges and Dalalyan, 2012, Lemma 8) we have

$$\|\bar{\boldsymbol{\xi}}_n\|_2^2 \leq \frac{2\mathbf{r}_\Sigma + 3\log(2/\delta)}{n} = r_n^2$$

with probability at least $1 - \delta$. Therefore, with probability at least $1 - 2\delta$,

$$\left(\sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right)^{1/2} \leq 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}. \tag{B.23}$$

Next, Lemma 4 and the fact that $p_L = \dim(\mathcal{U}_L) = 1$ imply that with probability at least $1 - \delta$

$$\|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 \leq \frac{2(1 + \sqrt{2}r_n)}{1 - 2\varepsilon}. \tag{B.24}$$

Recall that we have chosen \tilde{t}_γ in such a way that

$$\tilde{t}_\gamma \leq \frac{3(1 + C_\gamma\sqrt{2}r_n/\sqrt{\tau})}{1 - 2\varepsilon^*} + 1.6\sqrt{\log(2/\tau)}. \tag{B.25}$$

Combining (B.22), (B.23), (B.24) and (B.25), we arrive at the inequality

$$\begin{aligned}\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 &\leq 2\sqrt{2\varepsilon L\tilde{T}_1} + 9\varepsilon\tilde{t}_\gamma\sqrt{L} + \frac{8\varepsilon(1 + C_\gamma\sqrt{2}r_n)}{1 - 2\varepsilon} + 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)} \\ &\leq 2\sqrt{2\varepsilon L\tilde{T}_1} + \frac{27\varepsilon\sqrt{L}(1 + C_\gamma\sqrt{2}r_n/\sqrt{\tau})}{1 - 2\varepsilon^*} + 14.4\varepsilon\sqrt{L\log(2/\tau)} \\ &\quad + \frac{8\varepsilon(1 + C_\gamma\sqrt{2}r_n)}{1 - 2\varepsilon} + 4r_n(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}\end{aligned}$$

that holds with probability at least $1 - 3\delta$. In the upper bound obtained above, only the term \tilde{T}_1 remains random. To bound \tilde{T}_1 we first apply a triangle inequality then use Lemma 7. It implies that with probability at least $1 - 2\delta$, we have

$$\begin{aligned}\tilde{T}_1 &\leq A \frac{\sqrt{n\mathbf{r}_\Sigma} + \mathbf{r}_\Sigma + n\eta\log(2e/\eta) + 2\log(1/\delta)}{n(1 - \eta)} + (4r_n(1 + \sqrt{\eta}) + 10\eta\sqrt{\log(2/\eta)})^2 + (1 + \gamma)\varepsilon + \gamma \\ &\leq 2A(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta)) + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon + 2\gamma.\end{aligned}$$

Consequently,

$$\begin{aligned}\sqrt{\varepsilon\tilde{T}_1} &\leq \{2A\varepsilon(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta))\}^{1/2} + (7.5r_n + 10\eta\sqrt{\log(2/\eta)})\sqrt{\varepsilon} + \varepsilon + \sqrt{2\varepsilon\gamma} \\ &\leq \{2A\varepsilon(\sqrt{2}r_n + r_n^2 + 4\eta\log(2/\eta))\}^{1/2} + 5.4r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon + \sqrt{2\varepsilon\gamma} \\ &\leq \{2A\varepsilon(\sqrt{2}r_n + 4\eta\log(2/\eta))\}^{1/2} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon + \sqrt{2\varepsilon\gamma} \\ &\leq \sqrt{2\sqrt{2}Ar_n\varepsilon} + 2\sqrt{2A\varepsilon\eta\log(2/\eta)} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon + \sqrt{2\varepsilon\gamma} \\ &\leq \varepsilon + Ar_n/\sqrt{2} + 2\eta\sqrt{2A\log(2/\eta)} + (\sqrt{A} + 5.4)r_n + 7.1\eta\sqrt{\log(2/\eta)} + \varepsilon + \sqrt{2\varepsilon\gamma} \\ &\leq (A/\sqrt{2} + \sqrt{A} + 5.4)r_n + (7.1 + 2\sqrt{2A})\eta\sqrt{\log(2/\eta)} + (9.1 + 2\sqrt{2A})\varepsilon\sqrt{\log(2/\eta)} + \sqrt{2\varepsilon\gamma}.\end{aligned}$$

These inequalities imply that there exists a universal constant C such that

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C(C_\gamma r_n + \tau\sqrt{\log(2/\tau)} + \varepsilon\sqrt{\log(2/\varepsilon)} + r_n\varepsilon/\sqrt{\tau} + \sqrt{\varepsilon\gamma})\sqrt{L}}{1 - 2\varepsilon^*}. \quad (\text{B.26})$$

Let us denote $\log_+(x) = \max\{0, \log(x)\}$ the positive part of logarithm, then we choose

$$\tau = \frac{1}{4} \bigwedge \frac{\tilde{r}_n}{\sqrt{\log_+(2/\tilde{r}_n)}}, \quad \text{with} \quad \tilde{r}_n = \frac{\sqrt{C_\gamma\mathbf{r}_\Sigma} + \sqrt{2\log(2/\delta)}}{\sqrt{n}}.$$

Note that $r_n \leq \sqrt{2}\tilde{r}_n$ and, furthermore, $\tau = 1/4$ whenever $\tilde{r}_n \geq 1/2$. Therefore, $r_n\varepsilon/\sqrt{\tau} \leq r_n + \varepsilon$. Inserting this value of τ in (B.26) leads to

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2 \leq \frac{C(C_\gamma r_n + \varepsilon\sqrt{\log(2/\varepsilon)} + \sqrt{\varepsilon\gamma})\sqrt{L}}{1 - 2\varepsilon^*}.$$

where C is a universal constant, the value of which is not necessarily the same in different places where it appears. Replacing r_n by its expression, upper bounding L by $2\log p$, and

using the fact that $C_\gamma \leq 3$ for $\gamma \in (0, 1/2]$ we arrive at

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C \sqrt{\log p}}{1 - 2\varepsilon^*} \left(\sqrt{\frac{\mathbf{r}_\Sigma}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} + \varepsilon \sqrt{\log(2/\varepsilon)} + \sqrt{\varepsilon\gamma} \right).$$

Note that this inequality holds true on an event of probability at least $1 - 5\delta$.

B.3 Extension to Sub-Gaussian distributions

This section is devoted to the proof of Theorem 12, which is an extension of Theorem 11 to the case when the $1 - \varepsilon$ portion of observations are sub-Gaussian. First, we formulate some technical lemmas necessary for the proof of Theorem 12 postponing the full proof to the end of the present section.

Recall that a random vector ζ with zero mean and identity covariance matrix is sub-Gaussian with variance proxy $\mathfrak{s} > 0$, $\zeta \sim \text{SG}_p(\mathfrak{s})$, if

$$\mathbb{E}[e^{\mathbf{v}^\top \zeta}] \leq \exp \left\{ \frac{\mathfrak{s}}{2} \|\mathbf{v}\|^2 \right\}, \quad \forall \mathbf{v} \in \mathbb{R}^p.$$

The concentration inequality for sub-Gaussian vectors is a well-known fact (see, e.g. (Rigollet and Hütter, 2019), Theorem 1.19) that if $\zeta \sim \text{SG}_p(\mathfrak{s})$ then for all $\delta \in (0, 1)$, it holds

$$\mathbf{P}(\|\zeta\|_2 \leq 4\sqrt{p\mathfrak{s}} + \sqrt{8\mathfrak{s} \log(1/\delta)}) \geq 1 - \delta. \quad (\text{B.27})$$

In the definition of $\text{SGAC}(\mu^*, \Sigma, \mathfrak{s}, \varepsilon)$ we assume that the $1 - \varepsilon$ portion of the data $\{\mathbf{X}_i\}_{i=1}^n$ are sub-Gaussian, that is $\mathbf{X}_i = \mu^* + \Sigma^{1/2} \zeta_i$ for all $i \in \mathcal{I}$, where the set \mathcal{I} is of cardinality at least $(1 - \varepsilon)n$. Denote $\xi_i = \Sigma^{1/2} \zeta_i$ for all $i = 1, \dots, n$ and assume that $\|\Sigma\|_{\text{op}} = 1$.

First, we show that with the choice of threshold parameter the analogous to Lemma 2, Lemma 3, Lemma 4 lemmas hold true. Notice that all three lemmas are using the concentration bound for the operator norm of (sub)-Gaussian vectors. In the case of Gaussian vectors we make use of (Vershynin, 2012, Corollary 5.35), while the analogous result for the sub-Gaussian distributions is also known (Vershynin, 2012, Theorem 5.39). For the readers convenience the latter is formulated in Lemma 1.

Lemma 10. *Let $J \subset \{1, \dots, n\}$ be a subset of cardinality m . For every $\delta \in (0, 1)$, it holds that*

$$\mathbf{P} \left(\left\| \sum_{j \in J} \xi_j \right\|_2 \leq 4\sqrt{pm\mathfrak{s}} + \sqrt{8m\mathfrak{s} \log(1/\delta)} \right) \geq 1 - \delta.$$

Proof of Lemma 10. Without loss of generality, we assume that $J = \{1, \dots, m\}$. On the one hand, $\|\Sigma\|_{\text{op}} = 1$ implies that

$$\left\| \sum_{i=1}^m \xi_i \right\|_2 \leq \left\| \sum_{i=1}^m \zeta_i \right\|_2.$$

On the other hand, $\zeta_1 + \dots + \zeta_m \sim \text{SG}_p(m\mathfrak{s})$. Applying inequality (B.27) to this random variable yields the desired result. \square

We now state the versions of Lemma 2 and Lemma 4 that are valid in the setting of sub-Gaussian vectors. Notice that the only difference is in the choice of the threshold; thanks to Lemma 1, the threshold now includes the universal constant C_0 and the variance proxy \mathfrak{s} . The proofs of these lemmas are omitted, since they are the same as in the Gaussian case presented in Appendix B.1.3 (except for bounding the operator norm of a matrix with sub-Gaussian columns we use Lemma 1).

Lemma 11. *With probability at least $1 - \delta$, for all linear subspaces $V \subset \mathbb{R}^p$, we have*

$$\frac{\|\hat{\mu}_V^{\text{GM}} - P_V \mu^*\|_2}{\sqrt{\dim(V)}} \leq \frac{2\sqrt{\|\Sigma\|_{\text{op}}}}{1 - 2\varepsilon} \left(1 + \frac{C_0(\mathfrak{s}\sqrt{p} + 2\mathfrak{s}\sqrt{\log(1/\delta)})}{\sqrt{n}} \right),$$

where the constant C_0 is the same constant as in Lemma 1.

Lemma 12. *Let τ and δ be two numbers from $(0, 1)$. Define*

$$z = 1 + \frac{C_0\mathfrak{s}(\sqrt{p} + \sqrt{2\log(1/\delta)})}{\sqrt{n\tau}} + C_0\mathfrak{s}\sqrt{2 + 2\log(1/\tau)}$$

with the same constant C_0 as in Lemma 1. Then, with probability at least $1 - \delta$, we have

$$\sup_V \sum_{i=1}^n \mathbb{1}(\|P_V \xi_i\|_2^2 > z^2 \dim(V)) \leq n\tau,$$

where the supremum is over all linear subspaces V of \mathbb{R}^p .

Lemma 13 (Koltchinskii and Lounici (2017), Theorem 9). *There is a constant $A_3 > 0$ depending only on the variance proxy τ such that for every pair of integers $n \geq 1$ and $p \geq 1$, we have*

$$\mathbf{P}\left(\|\zeta_{1:n} \zeta_{1:n}^\top - n\Sigma\|_{\text{op}} \geq A_3\left(\sqrt{(p+t)n} + p + t\right)\right) \leq e^{-t}, \quad \forall t \geq 1.$$

Lemma 14. *There exists a positive constant A such that, for any positive integer $m \leq n$ and any $t \geq 1$, with probability at least $1 - 2e^{-t}$, the inequality*

$$\|\hat{\Sigma}_S - \Sigma\|_{\text{op}} \leq A \frac{\sqrt{np} + p + m \log(2ne/m) + 2t}{n - m} + \|\bar{\xi}_S\|_2^2$$

is satisfied for every $S \subset [n]$ of cardinality $\geq n - m$.

Proof. The proof of this lemma is similar to the proof of Lemma 7 with the only difference that instead of Theorems 4 and 5 from (Koltchinskii and Lounici, 2017) we now use Lemma 13. \square

Lemma 15. For any positive integer $m \leq n$ and any $t > 0$, with probability at least $1 - e^{-t}$, we have

$$\max_{|S| \geq n-m} \left\| \frac{1}{|S|} \sum_{i \in S} \mathbf{P} \xi_i \right\|_2 \leq \frac{n \|\mathbf{P} \bar{\xi}_n\|_2}{n-m} + \frac{\sqrt{m\mathfrak{s}}(4\sqrt{p} + 2\sqrt{2t}) + 2m\sqrt{2\mathfrak{s} \log(2ne/m)}}{n-m}.$$

Proof. The proof of this theorem is similar to that of Lemma 6, with the only exception that now we need to bound the maximum of a norm of a sum of at most m sub-Gaussian vectors, where the maximum is taken over all subsets of $[n]$ of size at most m . Since, each sub-Gaussian vector has a variance proxy \mathfrak{s} then using Lemma 10 along with union bound, we have

$$\begin{aligned} \mathbf{P} \left(\max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2 \geq \sqrt{m\mathfrak{s}}(4\sqrt{p} + t_m) \right) &\leq \sum_{l=1}^m \binom{n}{l} \mathbf{P} \left(\left\| \sum_{i=1}^l \xi_i \right\|_2 \geq \sqrt{l\mathfrak{s}}(4\sqrt{p} + t_l) \right) \\ &\leq \sum_{l=1}^m \left(\frac{ne}{l} \right)^l \mathbf{P} \left(\left\| \sum_{i=1}^l \xi_i \right\|_2 \geq \sqrt{l\mathfrak{s}}(4\sqrt{p} + t_l) \right) \\ &\leq \sum_{l=1}^m \left(\frac{ne}{l} \right)^l e^{-t_l/3} \leq e^{-t}. \end{aligned}$$

Therefore, we obtain that with probability at least $1 - e^{-t}$ the inequality

$$\max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2 \leq \sqrt{m\mathfrak{s}}(4\sqrt{p} + 2\sqrt{2t}) + 2m\sqrt{2\mathfrak{s} \log(2ne/m)}$$

holds. Then, combining with

$$\frac{1}{|S|} \left\| \sum_{i \in S} \mathbf{P} \xi_i \right\|_2 \leq \frac{n \|\mathbf{P} \bar{\xi}_n\|_2}{n-m} + \frac{1}{n-m} \max_{|J| \leq m} \left\| \sum_{i \in J} \xi_i \right\|_2.$$

yields the desired result. □

B.3.1 Proof of Theorem 12

All the ingredients provided, we can now compile the complete proof of Theorem 12.

Taking $\mathbf{U}_L := \mathbf{V}_L$, the algorithm detailed in (3.2) returns $\hat{\boldsymbol{\mu}}^{\text{SDR}} = \sum_{\ell=0}^L \hat{\boldsymbol{\mu}}^{(\ell)}$ with $\hat{\boldsymbol{\mu}}^{(\ell)} \in \mathcal{U}_\ell = \text{Im}(\mathbf{V}_\ell \mathbf{U}_\ell^\top)$ for every $\ell \in \{0, \dots, L\}$ where the two-by-two orthogonal subspaces $\mathcal{U}_0, \dots, \mathcal{U}_L$ span the whole space \mathbb{R}^p . This allows us to decompose the risk:

$$\|\hat{\boldsymbol{\mu}}^{\text{SDR}} - \boldsymbol{\mu}^*\|_2^2 = \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell}(\bar{\mathbf{X}}_\ell - \boldsymbol{\mu}^*)\|_2^2 = \sum_{\ell=0}^L \|\mathbf{P}_\ell(\mathbf{V}_\ell^\top \bar{\mathbf{X}}_\ell - \mathbf{V}_\ell^\top \boldsymbol{\mu}^*)\|_2^2,$$

where $\mathbf{P}_\ell := \mathbf{U}_\ell^\top \mathbf{U}_\ell$ is the projection matrix projecting onto the subspace of \mathbb{R}^{p_ℓ} spanned by the bottom $p_\ell - p_{\ell+1}$ eigenvectors of $\mathbf{V}_\ell^\top (\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}) \mathbf{V}_\ell$ for $\ell = 0, \dots, L$ with the convention that

$p_{L+1} = 0$.

For $\ell \in \{0, \dots, L-1\}$, we intend to apply Proposition 6 to $\mathbf{Z}_i = \mathbf{V}_\ell^\top \mathbf{X}_i$ and $\boldsymbol{\mu}^Z = \mathbf{V}_\ell^\top \boldsymbol{\mu}^*$ in order to upper bound the error term $\text{Err}_\ell := \|\mathbf{P}_\ell(\mathbf{V}_\ell^\top \bar{\mathbf{X}}_\ell - \mathbf{V}_\ell^\top \boldsymbol{\mu}^*)\|_2$. Using the inequalities

$$\|\mathbf{V}^\top (\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}) \mathbf{V}\|_{\text{op}} \leq \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}}, \quad \lambda_{p_\ell}(\mathbf{V}^\top \boldsymbol{\Sigma} \mathbf{V}) \leq \lambda_p(\boldsymbol{\Sigma}), \quad \lambda_1(\mathbf{V}^\top \boldsymbol{\Sigma} \mathbf{V}) \geq \lambda_1(\boldsymbol{\Sigma})$$

that are true for every orthogonal matrix \mathbf{V} , and keeping in mind the definition of \mathbf{P}_ℓ , we get

$$\text{Err}_\ell \leq \left\{ 2\omega_{\mathcal{O}} \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + \frac{\omega_{\mathcal{O}}^2}{1 - \omega_{\mathcal{O}}} \left((\lambda_p - \lambda_1)(\boldsymbol{\Sigma}) + \frac{\delta_\ell^2}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_\ell \mathbf{V}_\ell^\top \bar{\boldsymbol{\xi}}_{S_I^{(\ell)}}\|_2,$$

where we have used the notation

$$\omega_{\mathcal{O}} = \max_{\ell} \frac{|\mathcal{S}^{(\ell)} \cap \mathcal{O}|}{|\mathcal{S}^{(\ell)}|}, \quad \boldsymbol{\xi}_i = \mathbf{X}_i - \boldsymbol{\mu}^*$$

and $\delta_\ell = \inf_{\boldsymbol{\mu}} \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top (\mathbf{X}_i - \boldsymbol{\mu})\|_2$. Note that when \mathcal{O} and $(\mathcal{S}_I^{(\ell)})^c$ are of cardinality less than $n\varepsilon$ and $n(\varepsilon + \tau)$, respectively, we have $\omega_{\mathcal{O}} \leq \varepsilon/(1 - \tau)$ and $\frac{\omega_{\mathcal{O}}}{1 - \omega_{\mathcal{O}}} \leq \frac{\varepsilon}{1 - \varepsilon - \tau}$.

We set $\eta := \varepsilon + \tau \leq 3/4$ and apply Lemma 12 to infer that $\omega_{\mathcal{O}} \leq \varepsilon/(1 - \eta) \leq 4\varepsilon$ is true with probability at least $1 - \delta$. Furthermore, we know that $\delta_\ell \leq \max_{i \in \mathcal{S}^{(\ell)}} \|\mathbf{V}_\ell^\top \mathbf{X}_i - \bar{\boldsymbol{\mu}}^{(\ell)}\|_2 \leq t\sqrt{p_\ell}$. This yields

$$\text{Err}_\ell \leq \left\{ 8\varepsilon \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + 16\varepsilon^2 \left((\lambda_p - \lambda_1)(\boldsymbol{\Sigma}) + \frac{t^2 p_\ell}{p_{\ell+1}} \right) \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{S_I^{(\ell)}}\|_2.$$

Let us introduce the shorthand $T_1 = \max_{\ell \in [L]} \|\hat{\boldsymbol{\Sigma}}^{(\ell)} - \boldsymbol{\Sigma}\|_{\text{op}} + \varepsilon(\lambda_p - \lambda_1)(\boldsymbol{\Sigma})$. This leads to

$$\text{Err}_\ell \leq \left\{ 8\varepsilon T_1 + \frac{16\varepsilon^2 t^2 p_\ell}{p_{\ell+1}} \right\}^{1/2} + \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{S_I^{(\ell)}}\|_2. \quad (\text{B.28})$$

For the last error term, since $p_L = 1$ then, by the combination of Lemma 11 and Lemma 12, we have

$$\begin{aligned} \text{Err}_L &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{S_I^{(L)}}\|_2 + \frac{n\varepsilon(t\sqrt{p_L} + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)}{|\mathcal{S}^{(L)}|} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{S_I^{(L)}}\|_2 + \frac{\varepsilon t + \varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2}{1 - \eta} \\ &\leq \|\mathbf{P}_{\mathcal{U}_L} \bar{\boldsymbol{\xi}}_{S_I^{(L)}}\|_2 + 4\varepsilon t + 4\varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2. \end{aligned} \quad (\text{B.29})$$

Combining (B.28), (B.29), inequality $p_\ell \leq ep_{\ell+1}$, as well as the Minkowski inequality, we get

$$\begin{aligned}
\|\boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}^{\text{SDR}}\|_2 &= \left\{ \sum_{\ell=0}^L \text{Err}_\ell^2 \right\}^{1/2} \\
&\leq \left\{ 8\varepsilon L(T_1 + e\varepsilon t^2) + 16\varepsilon^2(t + \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2)^2 \right\}^{1/2} + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} \\
&\leq 2\sqrt{2\varepsilon LT_1} + 9\varepsilon t\sqrt{L} + 4\varepsilon \|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 + \left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2}. \tag{B.30}
\end{aligned}$$

To ease notation, let us set

$$r_{n,s} = \frac{4\sqrt{s}(\sqrt{p} + 2\sqrt{\log(2/\delta)})}{\sqrt{n}}.$$

In view of Lemma 15, with probability at least $1 - \delta$, we have

$$\begin{aligned}
\left\{ \sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right\}^{1/2} &\leq \left\{ \sum_{\ell=0}^L \left(\frac{\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2}{1-\eta} + \frac{r_{n,s}\sqrt{\eta} + 2\eta\sqrt{2\log(2e/\eta)}}{1-\eta} \right)^2 \right\}^{1/2} \\
&\leq \left\{ \sum_{\ell=0}^L (4\|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_n\|_2 + 4r_{n,s}\sqrt{\eta} + 10\eta\sqrt{\log(2/\eta)})^2 \right\}^{1/2} \\
&\leq 4\|\bar{\boldsymbol{\xi}}_n\|_2 + 4r_{n,s}\sqrt{\eta L} + 10\eta\sqrt{L\log(2/\eta)}.
\end{aligned}$$

Moreover, since the random variable $\bar{\boldsymbol{\xi}}_n$ is sub-Gaussian with variance proxy s/n , hence by Lemma 8 we have

$$\|\bar{\boldsymbol{\xi}}_n\|_2^2 \leq \frac{16s(\sqrt{p} + 2\sqrt{\log(2/\delta)})^2}{n} = r_{n,s}^2$$

with probability at least $1 - \delta$. Therefore, with probability at least $1 - 2\delta$,

$$\left(\sum_{\ell=0}^L \|\mathbf{P}_{\mathcal{U}_\ell} \bar{\boldsymbol{\xi}}_{\mathcal{S}_\mathcal{I}^{(\ell)}}\|_2^2 \right)^{1/2} \leq 4r_{n,s}(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}. \tag{B.31}$$

Next, the combination of Lemma 11 and Lemma 12 and the fact that $p_L = \dim(\mathcal{U}_L) = 1$ imply that with probability at least $1 - \delta$

$$\|\mathbf{P}_{\mathcal{U}_L} \boldsymbol{\mu}^* - \hat{\boldsymbol{\mu}}_{\mathcal{U}_L}^{\text{GM}}\|_2 \leq \frac{2(1 + Cr_{n,s}\sqrt{s}/4)}{1 - 2\varepsilon}, \tag{B.32}$$

where C is the same universal constant as in Lemma 1. Recall that we have chosen t in such a way that

$$t \leq \frac{3(1 + C_0 r_{n,s} \sqrt{s}/4\sqrt{\tau})}{1 - 2\varepsilon^*} + 1.6C_0 s \sqrt{\log(2/\tau)}. \tag{B.33}$$

Combining (B.30), (B.31), (B.32) and (B.33), we arrive at the inequality

$$\begin{aligned}\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 &\leq 2\sqrt{2\varepsilon LT_1} + 9\varepsilon t\sqrt{L} + \frac{8\varepsilon(1 + Cr_{n,s}\sqrt{s}/4)}{1 - 2\varepsilon} + 4r_{n,s}(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)} \\ &\leq 2\sqrt{2\varepsilon LT_1} + \frac{27\varepsilon\sqrt{L}(1 + Cr_{n,s}\sqrt{s}/4\sqrt{\tau})}{1 - 2\varepsilon^*} + 14.4C_s\varepsilon\sqrt{L\log(2/\tau)} \\ &\quad + \frac{8\varepsilon(1 + Cr_{n,s}\sqrt{s}/4)}{1 - 2\varepsilon} + 4r_{n,s}(1 + \sqrt{L\eta}) + 10\eta\sqrt{L\log(2/\eta)}\end{aligned}$$

that holds with probability at least $1 - 3\delta$. In the upper bound obtained above, only the term T_1 remains random. We can upper bound this term using Lemma 14. It implies that with probability at least $1 - 2\delta$, we have

$$\begin{aligned}T_1 &\leq A \frac{\sqrt{np} + p + n\eta\log(2e/\eta) + 2\log(1/\delta)}{n(1 - \eta)} + (4r_{n,s}(1 + \sqrt{\eta}) + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon \\ &\leq A_s(r_{n,s} + r_{n,s}^2 + 8\eta\log(2/\eta)) + (7.5r_{n,s} + 10\eta\sqrt{\log(2/\eta)})^2 + \varepsilon,\end{aligned}$$

where A_s is a constant depending only on the variance proxy s , the value of which is not necessarily the same in further simplifications of the expression from the last display. Then, using the triangle inequality several times we arrive at the following expression

$$\begin{aligned}\sqrt{\varepsilon T_1} &\leq \{A_s\varepsilon(r_{n,s} + r_{n,s}^2 + 8\eta\log(2/\eta))\}^{1/2} + (7.5r_{n,s} + 10\eta\sqrt{\log(2/\eta)})\sqrt{\varepsilon} + \varepsilon \\ &\leq (A_s + \sqrt{A_s/2} + 5.4)r_{n,s} + (7.1 + 2\sqrt{2A_s})\tau\sqrt{\log(2/\tau)} + (9.1 + 2\sqrt{2A_s})\varepsilon\sqrt{\log(2/\varepsilon)}.\end{aligned}$$

These inequalities imply that there is a universal constant C such that

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C_s(A_s r_{n,s}/\sqrt{s} + \tau\sqrt{\log(2/\tau)} + \varepsilon\sqrt{\log(2/\varepsilon)} + r_{n,s}\varepsilon/\sqrt{s\tau})\sqrt{L}}{1 - 2\varepsilon^*}. \quad (\text{B.34})$$

Let us denote $\log_+(x) = \max\{0, \log(x)\}$ the positive part of logarithm, then we choose

$$\tau = \frac{1}{4} \bigwedge \frac{\bar{r}_{n,s}}{\sqrt{\log_+(2/\bar{r}_{n,s})}}, \quad \text{with} \quad \bar{r}_{n,s} = \frac{3\sqrt{s}(\sqrt{p} + 2\sqrt{\log(2/\delta)})}{\sqrt{n}}.$$

Note that $r_{n,s} \leq \sqrt{2}\bar{r}_{n,s}$ and, furthermore, $\tau = 1/4$ whenever $\bar{r}_{n,s} \geq 1/2$. Therefore, $r_{n,s}\varepsilon/\sqrt{\tau} \leq r_{n,s} + \varepsilon$. Inserting this value of τ in (B.34) leads to

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C_s(A_s r_{n,s}/\sqrt{s} + \varepsilon\sqrt{\log(2/\varepsilon)})\sqrt{L}}{1 - 2\varepsilon^*}.$$

where C is a universal constant. Replacing $r_{n,s}$ by its expression, and upper bounding L by $2\log p$, we arrive at

$$\|\hat{\mu}^{\text{SDR}} - \mu^*\|_2 \leq \frac{C_s\sqrt{\log p}}{1 - 2\varepsilon^*} \left(A_s\sqrt{\frac{p}{n}} + \varepsilon\sqrt{\log(2/\varepsilon)} + A_s\sqrt{\frac{\log(1/\delta)}{n}} \right).$$

Note that this inequality holds true on an event of probability at least $1 - 5\delta$.

Bibliography

- Alon, N., Matias, Y., and Szegedy, M. (1996). The space complexity of approximating the frequency moments. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing*, STOC '96, page 20–29, New York, NY, USA. Association for Computing Machinery.
- Alquier, P. and Gerber, M. (2020). Universal robust regression via maximum mean discrepancy. *arXiv: Statistics Theory*.
- Audibert, J.-Y. and Catoni, O. (2011). Robust linear least squares regression. *Ann. Statist.*, 39(5):2766–2794.
- Bahmani, S. (2021). Nearly optimal robust mean estimation via empirical characteristic function. *Bernoulli*, 27(3):2139 – 2158.
- Bakshi, A. and Prasad, A. (2021). *Robust Linear Regression: Optimal Rates in Polynomial Time*, page 102–115. Association for Computing Machinery, New York, NY, USA.
- Balakrishnan, S., Du, S. S., Li, J., and Singh, A. (2017). Computationally efficient robust sparse estimation in high dimensions. In *COLT 2017*, pages 169–212.
- Batani, A.-H. and Dalalyan, A. S. (2020). Confidence regions and minimax rates in outlier-robust estimation on the probability simplex. *Electron. J. Statist.*, 14(2):2653–2677.
- Batani, A.-H., Minasyan, A., and Dalalyan, A. S. (2022). Nearly minimax robust estimator of the mean vector by iterative spectral dimension reduction.
- Bellec, P. C. (2016). Adaptive confidence sets in shape restricted regression. *arXiv preprint arXiv:1601.05766*.
- Bellec, P. C. (2018). Sharp oracle inequalities for least squares estimators in shape restricted regression. *Ann. Statist.*, 46(2):745–780.
- Berend, D. and Kontorovich, A. (2013). A sharp estimate of the binomial mean absolute deviation with applications. *Statistics & Probability Letters*, 83(4):1254 – 1259.
- Bhatia, K., Jain, P., Kamalaruban, P., and Kar, P. (2017). Consistent robust regression. In *Advances in Neural Information Processing Systems*, pages 2110–2119.

- Boucheron, S., Lugosi, G., and Massart, P. (2013). *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford.
- Braess, D. and Sauer, T. (2004). Bernstein polynomials and learning theory. *Journal of Approximation Theory*, 128(2):187 – 206.
- Brunel, V.-E. (2019). Concentration of the empirical level sets of Tukey’s halfspace depth. *Probability Theory and Related Fields*, 173(3-4):1165–1196.
- Cai, T. T. and Low, M. G. (2006). Adaptive confidence balls. *Ann. Statist.*, 34(1):202–228.
- Cardot, H., Cenac, P., and Zitt, P.-A. (2013). Efficient and fast estimation of the geometric median in hilbert spaces with an averaged stochastic gradient algorithm. *Bernoulli*, 19:18–43.
- Carpentier, A., Delattre, S., Roquain, E., and Verzelen, N. (2018). Estimating minimum effect with outlier selection. *arXiv e-prints*, page arXiv:1809.08330.
- Catoni, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 48(4):1148 – 1185.
- Catoni, O. and Giulini, I. (2017). Dimension-free pac-bayesian bounds for matrices, vectors, and linear least squares regression. *arXiv: Statistics Theory*.
- Catoni, O. and Giulini, I. (2018). Dimension-free pac-bayesian bounds for the estimation of the mean of a random vector. *arXiv: Statistics Theory*.
- Chen, M., Gao, C., and Ren, Z. (2016). A general decision theory for Huber’s ε -contamination model. *Electronic Journal of Statistics*, 10(2):3752–3774.
- Chen, M., Gao, C., and Ren, Z. (2018). Robust covariance and scatter matrix estimation under Huber’s contamination model. *Ann. Statist.*, 46(5):1932–1960.
- Chen, S., Li, J., and Moitra, A. (2020). Efficiently learning structured distributions from untrusted batches. In *Proceedings of STOC 2020, June 22-26*, pages 960–973. ACM.
- Chen, Y., Caramanis, C., and Mannor, S. (2013). Robust sparse regression under adversarial corruption. In *International Conference on Machine Learning*, pages 774–782.
- Cheng, Y., Diakonikolas, I., and Ge, R. (2019a). High-dimensional robust mean estimation in nearly-linear time. In *Proceedings of SODA 2019*, pages 2755–2771.
- Cheng, Y., Diakonikolas, I., Ge, R., and Soltanolkotabi, M. (2020). High-dimensional robust mean estimation via gradient descent. *CoRR*, abs/2005.01378.
- Cheng, Y., Diakonikolas, I., Ge, R., and Woodruff, D. P. (2019b). Faster algorithms for high-dimensional robust covariance estimation. In Beygelzimer, A. and Hsu, D., editors, *COLT 2019*, volume 99 of *Proceedings of Machine Learning Research*, pages 727–757. PMLR.

- Cheng, Y., Diakonikolas, I., Kane, D. M., Ge, R., Gupta, S., and Soltanolkotabi, M. (2021). Outlier-robust sparse estimation via non-convex optimization. *CoRR*, abs/2109.11515.
- Cherapanamjeri, Y., Aras, E., Tripuraneni, N., Jordan, M. I., Flammarion, N., and Bartlett, P. L. (2020). Optimal robust linear regression in nearly linear time. *ArXiv*, abs/2007.08137.
- Cherapanamjeri, Y., Flammarion, N., and Bartlett, P. L. (2019). Fast mean estimation with sub-Gaussian rates. In *COLT 2019*, volume 99 of *Proceedings of Machine Learning Research*, pages 786–806.
- Chinot, G. (2020). Erm and rerm are optimal estimators for regression problems when malicious outliers corrupt the labels. *Electron. J. Statist.*, 14(2):3563–3605.
- Chinot, G., Lecué, G., and Lerasle, M. (2018). Statistical learning with Lipschitz and convex loss functions. *arXiv e-prints*, page arXiv:1810.01090.
- Chinot, G., Lecué, G., and Lerasle, M. (2020). Robust high dimensional learning for Lipschitz and convex losses. *J. Mach. Learn. Res.*, 21:Paper No. 233, 47.
- Cohen, D., Kontorovich, A., and Wolfer, G. (2020). Learning discrete distributions with infinite support.
- Collier, O. and Dalalyan, A. S. (2019). Multidimensional linear functional estimation in sparse gaussian models and robust estimation of the mean. *Electron. J. Statist.*, 13(2):2830–2864.
- Comminges, L., Collier, O., Ndaoud, M., and Tsybakov, A. B. (2021). Adaptive robust estimation in sparse vector model. *Ann. Statist.*, 49(3):1347–1377.
- Comminges, L. and Dalalyan, A. S. (2012). Tight conditions for consistency of variable selection in the context of high dimensionality. *The Annals of Statistics*, 40(5):2667 – 2696.
- Dalalyan, A. and Chen, Y. (2012). Fused sparsity and robust estimation for linear models with unknown variance. In *Advances in Neural Information Processing Systems 25*, pages 1259–1267. Curran Associates, Inc.
- Dalalyan, A. and Thompson, P. (2019). Outlier-robust estimation of a sparse linear model using ℓ_1 -penalized Huber’s M-estimator. In *NeurIPS 32*, pages 13188–13198.
- Dalalyan, A. S. and Minasyan, A. (2020). All-in-one robust estimator of the gaussian mean. *math.ST*, arXiv:2002.01432.
- Dalalyan, A. S. and Sebbar, M. (2018). Optimal Kullback-Leibler aggregation in mixture density estimation by maximum likelihood. *Math. Stat. Learn.*, 1(1):1–35.
- Depersin, J. (2020). A spectral algorithm for robust regression with subgaussian rates. *ArXiv*, abs/2007.06072.

- Depersin, J. and Lecué, G. (2021). On the robustness to adversarial corruption and to heavy-tailed data of the stahel-donoho median of means. *arXiv preprint arXiv:2101.09117*.
- Depersin, J. and Lecué, G. (2022). Robust sub-Gaussian estimation of a mean vector in nearly linear time. *Ann. Statist.*, 50(1):511–536.
- Devroye, L., Lerasle, M., Lugosi, G., and Oliveira, R. I. (2016). Sub-gaussian mean estimators. *Ann. Statist.*, 44(6):2695–2725.
- Devroye, L. and Lugosi, G. (2000). *Combinatorial methods in density estimation*. Springer Series in Statistics.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. (2016a). Robust estimators in high dimensions without the computational intractability. In *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016*, pages 655–664.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. (2017). Being robust (in high dimensions) can be practical. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 999–1008.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Moitra, A., and Stewart, A. (2018a). Robustly learning a gaussian: Getting optimal error, efficiently. In Czumaj, A., editor, *Proceedings of SODA 2018*, pages 2683–2702. SIAM.
- Diakonikolas, I., Kamath, G., Kane, D. M., Li, J., Steinhardt, J., and Stewart, A. (2019a). Sever: A robust meta-algorithm for stochastic optimization. *ArXiv*, abs/1803.02815.
- Diakonikolas, I. and Kane, D. M. (2019). Recent advances in algorithmic high-dimensional robust statistics. *CoRR*, abs/1911.05911.
- Diakonikolas, I., Kane, D. M., and Pensia, A. (2020). Outlier robust mean estimation with subgaussian rates via stability. In Larochele, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1830–1840. Curran Associates, Inc.
- Diakonikolas, I., Kane, D. M., and Stewart, A. (2016b). Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. *Electron. Colloquium Comput. Complex.*, 23:177.
- Diakonikolas, I., Kong, W., and Stewart, A. (2019b). Efficient algorithms and lower bounds for robust linear regression. *SODA '19*, pages 2745–2754, USA. Society for Industrial and Applied Mathematics.

- Diakonikolas, I., Li, J., and Schmidt, L. (2018b). Fast and sample near-optimal algorithms for learning multidimensional histograms. In *COLT 2018, Stockholm, Sweden, 6-9 July 2018.*, pages 819–842.
- Dong, Y., Hopkins, S. B., and Li, J. (2019). Quantum entropy scoring for fast robust mean estimation and improved outlier detection. In *NeurIPS 2019*, pages 6065–6075.
- Donoho, D. L. and Gasko, M. (1992). Breakdown Properties of Location Estimates Based on Halfspace Depth and Projected Outlyingness. *The Annals of Statistics*, 20(4):1803 – 1827.
- Donoho, D. L. and Huber, P. J. (1983). The notion of breakdown point. *Festschr. for Erich L. Lehmann*, 157-184 (1983).
- Donoho, D. L. and Montanari, A. (2016). High dimensional robust m-estimation: asymptotic variance via approximate message passing. *Probability Theory and Related Fields*, 166(3):935–969.
- Feng, J., Xu, H., Mannor, S., and Yan, S. (2014). Robust logistic regression and classification. In *Advances in Neural Information Processing Systems 27*, pages 253–261. Curran Associates, Inc.
- Gao, C. (2020). Robust regression via multivariate regression depth. *Bernoulli*, 26(2):1139 – 1170.
- Gao, C., Liu, J., Yao, Y., and Zhu, W. (2018). Robust estimation and generative adversarial nets.
- Gao, C., Yao, Y., and Zhu, W. (2020). Generative adversarial nets for robust scatter estimation: A proper scoring rule perspective. *J. Mach. Learn. Res.*, 21:160:1–160:48.
- Goes, J., Lerman, G., and Nadler, B. (2020). Robust sparse covariance estimation by thresholding Tyler’s M-estimator. *Ann. Statist.*, 48(1):86–110.
- Guntuboyina, A. and Sen, B. (2018). Nonparametric shape-restricted regression. *Statist. Sci.*, 33(4):568–594.
- Halko, N., Martinsson, P. G., and Tropp, J. A. (2011). Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, 53(2):217–288.
- Hoffmann, M. and Nickl, R. (2011). On adaptive inference and confidence bands. *Ann. Statist.*, 39(5):2383–2409.
- Hopkins, S. B. (2020). Mean estimation with sub-Gaussian rates in polynomial time. *The Annals of Statistics*, 48(2):1193 – 1213.

- Hsu, D. and Sabato, S. (2016). Loss minimization and parameter estimation with heavy tails. *J. Mach. Learn. Res.*, 17(1):543–582.
- Huber, P. and Ronchetti, E. (2011). *Robust Statistics*. Wiley Series in Probability and Statistics. Wiley.
- Huber, P. J. (1964). Robust estimation of a location parameter. *Ann. Math. Statist.*, 35(1):73–101.
- Jain, A. and Orlitsky, A. (2019). Robust learning of discrete distributions from batches. *CoRR*, arXiv:1911.08532.
- Jain, A. and Orlitsky, A. (2021). Robust density estimation from batches: The best things in life are (nearly) free. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4698–4708. PMLR.
- Jerrum, M. R., Valiant, L. G., and Vazirani, V. V. (1986). Random generation of combinatorial structures from a uniform distribution. *Theoretical Computer Science*, 43:169–188.
- Joly, E., Lugosi, G., and Imbuzeiro Oliveira, R. (2017). On the estimation of the mean of a random vector. *Electron. J. Statist.*, 11(1):440 – 451.
- Kamath, S., Orlitsky, A., Pichapati, D., and Suresh, A. T. (2015). On learning distributions from their samples. In *Proceedings of The 28th Conference on Learning Theory, COLT 2015, Paris, France, July 3-6, 2015*, volume 40 of *JMLR Workshop and Conference Proceedings*, pages 1066–1100. JMLR.org.
- Kearns, M. (1998). Efficient noise-tolerant learning from statistical queries. *J. ACM*, 45(6):983–1006.
- Koltchinskii, V. and Lounici, K. (2017). Concentration inequalities and moment bounds for sample covariance operators. *Bernoulli*, 23(1):110 – 133.
- Lai, K. A., Rao, A. B., and Vempala, S. S. (2016). Agnostic estimation of mean and covariance. In *IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016*, pages 665–674.
- Lecué, G. and Lerasle, M. (2017). Learning from MOM’s principles: Le Cam’s approach. *Stoch. Proc. App.*, to appear.
- Lecué, G., Lerasle, M., and Mathieu, T. (2020). Robust classification via mom minimization. *Machine Learning*, pages 1–31.
- Lei, Z., Luh, K., Venkat, P., and Zhang, F. (2019). A fast spectral algorithm for mean estimation with sub-gaussian rates. *ArXiv*, abs/1908.04468.

- Liu, H. and Gao, C. (2019). Density estimation with contamination: minimax rates and theory of adaptation. *Electron. J. Stat.*, 13(2):3613–3653.
- Liu, J., Cosman, P. C., and Rao, B. D. (2018). Robust linear regression via ℓ_0 regularization. *IEEE Transactions on Signal Processing*, 66(3):698–713.
- Liu, J., Deshmukh, A., and Veeravalli, V. V. (2020a). Robust mean estimation in high dimensions via ℓ_0 minimization.
- Liu, L., Shen, Y., Li, T., and Caramanis, C. (2020b). High dimensional robust sparse regression. In Chiappa, S. and Calandra, R., editors, *AISTATS 2020*, volume 108 of *Proceedings of Machine Learning Research*, pages 411–421. PMLR.
- Liu, X., Mosler, K., and Mozharovskyi, P. (2019). Fast computation of tukey trimmed regions and median in dimension $p > 2$. *Journal of Computational and Graphical Statistics*, 28:682–697.
- Liu, X., Zuo, Y., and Wang, Q. (2017). Finite sample breakdown point of tukey’s halfspace median. *Sci. China Math*, 60:861–874.
- Lopuhaa, H. P. and Rousseeuw, P. J. (1991). Breakdown points of affine equivariant estimators of multivariate location and covariance matrices. *The Annals of Statistics*, 19(1):229–248.
- Lugosi, G. and Mendelson, S. (2019a). Mean estimation and regression under heavy-tailed distributions: A survey. *Found Comput Math*, 19:1145–1190.
- Lugosi, G. and Mendelson, S. (2019b). Sub-gaussian estimators of the mean of a random vector. *Ann. Statist.*, 47(2):783–794.
- Lugosi, G. and Mendelson, S. (2021). Robust multivariate mean estimation: The optimality of trimmed mean. *The Annals of Statistics*, 49(1):393 – 410.
- Maronna, R., Martin, D., and Yohai, V. (2006). *Robust Statistics: Theory and Methods*. Wiley Series in Probability and Statistics. Wiley.
- Minasyan, A. G. (2020). Excess-risk consistency of group-hard thresholding estimator in robust estimation of Gaussian mean. *Journal of Contemporary Mathematical Analysis*, 55(3):208–212.
- Minsker, S. (2015). Geometric median and robust estimation in banach spaces. *Bernoulli*, 21(4):2308–2335.
- Minsker, S. (2018a). Sub-gaussian estimators of the mean of a random matrix with heavy-tailed entries. *Ann. Statist.*, 46(6A):2871–2903.
- Minsker, S. (2018b). Uniform bounds for robust mean estimators.

- Minsker, S. and Ndaoud, M. (2021). Robust and efficient mean estimation: an approach based on the properties of self-normalized sums. *Electronic Journal of Statistics*, 15(2):6036 – 6070.
- Nemirovsky, A. and Yudin, D. (1983). Problem complexity and method efficiency in optimization.
- Nguyen, N. H. and Tran, T. D. (2013). Robust lasso with missing and grossly corrupted observations. *IEEE Trans. Inform. Theory*, 59(4):2036–2058.
- Pensia, A., Jog, V., and Loh, P.-L. (2020). Robust regression with covariate filtering: Heavy tails and adversarial contamination. *preprint arXiv:2009.12976*.
- Prasad, A., Balakrishnan, S., and Ravikumar, P. (2019). A unified approach to robust mean estimation.
- Prasad, A., Balakrishnan, S., and Ravikumar, P. (2020). A robust univariate mean estimator is all you need. In *AISTATS*.
- Qiao, M. and Valiant, G. (2018). Learning discrete distributions from untrusted batches. In *9th Innovations in Theoretical Computer Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA*, volume 94 of *LIPIcs*, pages 47:1–47:20. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- Rigollet, P. and Hütter, J.-C. (2019). High dimensional statistics. *Lecture notes (MIT)*.
- Rousseeuw, P., Hampel, F., Ronchetti, E., and Stahel, W. (2011). *Robust Statistics: The Approach Based on Influence Functions*. Wiley Series in Probability and Statistics. Wiley.
- Rousseeuw, P. J. and Hubert, M. (1999). Regression depth. *Journal of the American Statistical Association*, 94(446):388–402.
- Sasai, T. and Fujisawa, H. (2021). Adversarial robust weighted huber regression. *arXiv preprint arXiv:2102.11120*.
- Steinhardt, J., Charikar, M., and Valiant, G. (2017). Resilience: A criterion for learning in the presence of arbitrary outliers.
- Tsybakov, A. B. (2009). *Introduction to Nonparametric Estimation*. Springer, New York.
- Tukey, J. W. (1975). Mathematics and the picturing of data. *Proc. int. Congr. Math.*, Vancouver 1974, Vol. 2, 523-531.
- Vershynin, R. (2012). *Introduction to the non-asymptotic analysis of random matrices*, page 210–268. Cambridge University Press.

- Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Wang, Z. and Tan, Z. (2021). Tractable and near-optimal adversarial algorithms for robust estimation in contaminated gaussian models.
- Xia, D. and Koltchinskii, V. (2016). Estimation of low rank density matrices: Bounds in Schatten norms and other distances. *Electron. J. Statist.*, 10(2):2717–2745.
- Zhu, B., Jiao, J., and Jordan, M. I. (2022). Robust estimation for nonparametric families via generative adversarial networks. *CoRR*, abs/2202.01269.
- Zhu, B., Jiao, J., and Steinhardt, J. (2019). Generalized resilience and robust statistics. *ArXiv*, abs/1909.08755.
- Zhu, B., Jiao, J., and Steinhardt, J. (2020). When does the Tukey median work? *2020 IEEE International Symposium on Information Theory (ISIT)*, pages 1201–1206.
- Zuo, Y. and Serfling, R. (2000). General notions of statistical depth function. *Ann. Statist.*, 28:461–482.

Titre: Contributions à l'estimation robuste : optimalité minimax vs. efficacité calculatoire

Mots clés: robustesse, approche minimax, efficacité calculatoire, statistique en grande dimension

Résumé: En statistique et en théorie de l'apprentissage statistique, on suppose souvent que les échantillons sont distribués indépendamment et identiquement selon une distribution de probabilité de référence. Une approche plus réaliste pourrait consister à relaxer cette hypothèse en permettant à une fraction des échantillons de ne pas nécessairement suivre la distribution de référence. Ces échantillons désobéissants, appelés données aberrantes, peuvent considérablement détériorer la performance des estimateurs classiques. Dans ce travail, nous cherchons à estimer la moyenne des distributions de référence par des estimateurs robustes aux données aberrantes. Nous nous intéressons au comportement non-asymptotique des estimateurs.

Dans un premier temps, nous décrivons divers modèles de contamination qui déterminent la nature des données aberrantes

parmi nos observations. Puis, nous considérons le problème de l'estimation de la moyenne d'une distribution dont le support est le simplexe de probabilité de dimension k dans le cas où une fraction ε d'observations sont des données aberrantes générées par un adversaire. Un exemple particulier simple est le problème de l'estimation de la distribution d'une variable aléatoire discrète. Dans un deuxième temps, nous étudions le problème de l'estimation robuste de la moyenne d'une distribution gaussienne. Les estimateurs minimax-optimaux connus pour ce problème ne sont pas calculables en temps polynomial. Nous introduisons un estimateur efficace basé sur la réduction spectrale de dimension et établissons une borne supérieure sur son erreur qui est minimax-optimale modulo des facteurs logarithmiques.

Title: Contributions to robust estimation: minimax optimality vs. computational tractability

Keywords: robustness, minimax approach, computational efficiency, high-dimensional data

Abstract: In statistics and learning theory, it is common to assume that samples are independently and identically distributed according to a reference probability distribution. A more realistic approach could be to relax this assumption by allowing a fraction of samples to not necessarily follow the reference distribution. These disobeying samples, called outliers, may drastically skew the classical estimators. In this work, we aim to estimate the mean of reference distributions by estimators robust to outliers. We are interested in the non-asymptotic behavior of the estimators.

In the first stage, we describe various contamination models which determine the nature of the outliers among our observations.

Then, we consider the problem of estimating the mean of a distribution supported by the k -dimensional probability simplex in the setting where an ε fraction of observations are outliers generated by an adversary. A simple particular example is the problem of estimating the distribution of a discrete random variable. In the second stage, we study the problem of robust estimation of the mean of a Gaussian distribution. The known minimax-optimal estimators for this problem are not computationally tractable. We introduce a computationally efficient estimator based on spectral dimension reduction and establish a finite sample upper bound on its error that is minimax-optimal up to logarithmic factors.

