



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

MASTER THESIS

Filip Bočinec

**Halfspace depth for location and
scatter: robustness and minimax
optimality**

Department of Probability and Mathematical Statistics

Supervisor of the master thesis: doc. Mgr. Stanislav Nagy, Ph.D.

Study programme: Probability, Mathematical Statistics
and Econometrics

Prague 2024

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date
Author's signature

I want to express my deepest gratitude to my supervisor, doc. Mgr. Stanislav Nagy, Ph.D., for his invaluable mentoring and guidance. His expertise, patience, and encouragement have been instrumental in completing this thesis. I also extend my heartfelt thanks to my family and friends for their unwavering support and encouragement. Their love and belief in me have been a constant source of strength and motivation.

Title: Halfspace depth for location and scatter: robustness and minimax optimality

Author: Filip Bočinec

Department: Department of Probability and Mathematical Statistics

Supervisor: doc. Mgr. Stanislav Nagy, Ph.D., Department of Probability and Mathematical Statistics

Abstract: This thesis explores the concepts of location and scatter halfspace depth. Location halfspace depth is a well-established tool in nonparametric statistics, while scatter halfspace depth represents a newer concept that is currently undergoing active research. The primary goal of this work is to present the fundamental properties of halfspace depth for both location and scatter, with a special emphasis on the robustness of the corresponding medians. A significant portion of the thesis is dedicated to examining the minimax optimality of the location and scatter halfspace median. It provides a detailed framework concerning the rates of convergence and minimax optimal estimators. By employing this framework, the thesis demonstrates that both the location halfspace median and the scatter halfspace median achieve minimax optimality within Huber's contamination model. This finding underscores both the robustness and the rate optimality of these estimators.

Keywords: multivariate analysis, halfspace depth, scatter halfspace depth, robustness, minimax optimality

Contents

Introduction	3
Notation	5
1 Location halfspace depth	7
1.1 General notion of location statistical depth	7
1.2 Definition of location halfspace depth	8
1.3 Basic properties of location halfspace depth	11
1.4 Sample location halfspace depth and halfspace median	14
1.5 Robustness of halfspace median	16
2 Scatter halfspace depth	25
2.1 Axiomatic approach to scatter statistical depth	25
2.2 Definition of scatter halfspace depth	27
2.3 Basic properties of scatter halfspace depth	30
2.4 Sample scatter halfspace depth and scatter halfspace median	32
2.5 Robustness of scatter halfspace median	33
3 Minimax optimality	35
3.1 Convergence rate	35
3.2 Distance of probability measures	38
3.3 Huber's contamination model	41
3.4 Optimality of location halfspace median	44
3.5 Optimality of scatter halfspace median	53
Conclusion	65
Bibliography	67
Auxiliary theorems	69

Introduction

Since the pioneering work of Tukey [1975], the concept of statistical depth has become an invaluable tool in multivariate and nonparametric statistics. Statistical depth functions provide robust methods for understanding and summarizing multivariate data, making them essential in modern data analysis. Notable applications include outlier detection [Serfling, 2006], depth-based classification and clustering [Jörnsten, 2004], and rank and sign tests [Hettmansperger and Oja, 1994].

The classical location halfspace depth, rooted in the notion of centrality, measures how central a point is with respect to a distribution. It generalizes the univariate concept of the median to higher dimensions, offering a robust location functional. This is particularly important in multivariate analysis, where outliers or skewed distributions can heavily influence traditional measures such as the mean.

Scatter halfspace depth, on the other hand, extends the idea of data depth to the scatter of the data. Traditional measures of scatter, such as covariance matrices, assume specific distributional forms and are sensitive to outliers, which can distort the analysis. Scatter halfspace depth offers a nonparametric alternative that does not rely on these assumptions. It measures the “centrality” of scatter matrices within the space of all positive definite matrices with respect to a distribution. This thesis explores both location halfspace depth and scatter halfspace depth, delving into their theoretical foundations. It is divided into three comprehensive chapters, each addressing different aspects of these concepts.

In the first chapter, we introduce location halfspace depth and discuss its main properties in detail. We further examine its robust characteristics, such as the influence function, gross error sensitivity, and breakdown point. Its robust properties make location halfspace depth a powerful tool in multivariate analysis.

In the second chapter, we turn our attention to scatter halfspace depth. We introduce its definition and discuss its essential properties, drawing analogies to location halfspace depth where appropriate. Notably, because the third chapter forms the main part of the thesis, several theorems in the first and second chapters only reference proofs from the literature to maintain a manageable scope for this work.

The third chapter constitutes the core of this thesis. We begin by introducing key concepts of minimax estimation theory, such as convergence rate and minimax optimal convergence rate. This framework is then used in Sections 3.4 and 3.5, where we demonstrate in detail that the location halfspace median and scatter halfspace median are minimax optimal estimators in Huber’s contamination model. Huber’s contamination model is a robust statistical framework that accounts for data contamination, making it highly relevant for real-world applications where data imperfections are common.

These findings indicate that the above-mentioned estimators are not only robust but also achieve an optimal convergence rate under contamination. This means that they provide relatively reliable estimates even when the data contains a significant amount of noise or outliers. These results have noteworthy implications for various applications, including financial data analysis, biomed-

cal research, and any field that relies on robust multivariate statistical methods.

In summary, this thesis provides a thorough examination of location and scatter halfspace depths, highlighting their theoretical properties. It underscores the importance of statistical depth functions in modern data analysis and their potential to provide more accurate and insightful results in the presence of complex, multivariate data.

Notation

Symbol	Description
\mathbb{R}^d	Euclidean space of dimension d
GL_d	Set of all non-singular $d \times d$ real matrices
PD_d	Set of all symmetric positive definite $d \times d$ real matrices
\mathbf{I}_d	$d \times d$ identity matrix
\mathbf{x}, \mathbf{A}	Vector \mathbf{x} , matrix \mathbf{A} (bold font)
$(a, b), [a, b]$	Open interval, closed interval
\subset	Subset, proper subset
$\langle \cdot, \cdot \rangle, \ \cdot\ $	Standard inner product, Euclidean norm
$\text{int}(\cdot), \text{cl}(\cdot), \text{bd}(\cdot)$	Interior, closure and boundary
$B_d(\mathbf{x}, r)$	Closed ball in \mathbb{R}^d with center \mathbf{x} and radius r
$\mathbf{1}_A(\cdot)$	Indicator function of set A
S^{d-1}	Unit sphere in \mathbb{R}^d
$\ \cdot\ _F$	Frobenius matrix norm
$\ \cdot\ _{\text{op}}$	Operator matrix norm
λ_d	Lebesgue measure on \mathbb{R}^d
$(\Omega, \mathcal{F}, \mathbb{P})$	Probability space considered in the following text
\mathcal{B}^d	Borel σ -algebra on \mathbb{R}^d
$\mathcal{P}(\mathbb{R}^d)$	Space of all Borel probability measures on \mathbb{R}^d
$P_{\mathbf{X}}$	Distribution of random vector \mathbf{X}
$\mathbf{X} \sim P$	Random vector \mathbf{X} with distribution P
$\mathbf{X} \stackrel{D}{=} \mathbf{Y}$	Random vectors \mathbf{X}, \mathbf{Y} with equal distributions
$\delta_{\mathbf{x}}$	Dirac measure concentrated in a point \mathbf{x}
$N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$	d -variate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$
Φ, φ	Distribution function and density of $N_1(0, 1)$
$P \ll Q$	P absolutely continuous w.r.t. Q
dP/dQ	Radon–Nikodym derivative of P w.r.t. Q
$ \mathbf{A} $	Determinant of matrix \mathbf{A}
$\text{tr}(\mathbf{A})$	Trace of matrix \mathbf{A}
$\sigma_{\max}(\mathbf{S}), \sigma_{\min}(\mathbf{S})$	Maximal and minimal eigenvalue of matrix $\mathbf{S} \in \text{PD}_d$
$\mathbf{S}^{1/2}$	Unique matrix $\mathbf{S}^{1/2} \in \text{PD}_d$ such that $\mathbf{S}^{1/2}\mathbf{S}^{1/2} = \mathbf{S} \in \text{PD}_d$
$a \vee b, a \wedge b$	Maximum of $\{a, b\}$, minimum of $\{a, b\}$
a.e., a.s., w.r.t.	Almost everywhere, almost surely, with respect to

1. Location halfspace depth

We begin this thesis by defining a location parameter of a probability distribution.

Definition 1 (Location parameter) *Let $P \in \mathcal{P}(\mathbb{R}^d)$. Consider a functional $\Xi: \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}^d$ satisfying $\Xi(P_{\mathbf{A}\mathbf{X}+\mathbf{b}}) = \mathbf{A}\Xi(P_{\mathbf{X}}) + \mathbf{b}$ for any $\mathbf{X} \in \mathcal{P}_{\mathbf{X}} \subseteq \mathcal{P}$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{b} \in \mathbb{R}^d$ such that $P_{\mathbf{A}\mathbf{X}+\mathbf{b}} \in \mathcal{P}$. Then Ξ is called a location functional and $\Xi(P)$ is called a location parameter of $P \in \mathcal{P}$.*

For the set \mathcal{P} of all integrable d -variate random vectors, the mean value functional $\Xi(P) := \int_{\mathbb{R}^d} \mathbf{x} dP(\mathbf{x})$ is an example of a location functional. Another example is the univariate median functional defined on $\mathcal{P} = \mathcal{P}(\mathbb{R})$, i.e. $\Xi(P) = \arg \max_{x \in \mathbb{R}} \min \{P((-\infty, x]), P([x, \infty))\}$.

Considering a univariate probability distribution $P \in \mathcal{P}(\mathbb{R})$, one is able to tell how central any point $x \in \mathbb{R}$ is. This can be done by comparing the distribution function of P at x to $1/2$. That is, we can consider the median as the most central point, and points far from the median can be seen as the least central points. However, this approach strongly relies on the ordering of \mathbb{R} and therefore cannot be straightforwardly generalized for multivariate distributions. The concept of location statistical depth offers a way to organize data points in a center-outward manner with respect to a given multivariate distribution. Therefore, it generalizes univariate rank and order statistics to multivariate data.

1.1 General notion of location statistical depth

As stated above, location statistical depth aims to quantify centrality. To do so, we need to define some notions of symmetry for d -variate random vectors (or equivalently for Borel probability measures on \mathbb{R}^d). In the following text, $d \in \mathbb{N}$ will always represent the dimension of the data points.

Definition 2 (Symmetry of random vector) *Let $\mathbf{X} \in \mathcal{P}(\mathbb{R}^d)$ be a random vector and consider a fixed point $\mathbf{x}_0 \in \mathbb{R}^d$. We say that \mathbf{X} (or P) is centrally symmetric about its center \mathbf{x}_0 if*

$$\mathbf{X} - \mathbf{x}_0 \stackrel{D}{=} \mathbf{x}_0 - \mathbf{X}.$$

More generally, a random vector \mathbf{X} (or its distribution P) is said to be angularly symmetric about its center \mathbf{x}_0 if $(\mathbf{X} - \mathbf{x}_0) / \|\mathbf{X} - \mathbf{x}_0\|$ is centrally symmetric about $\mathbf{0}$. As an even more general notion, we say that \mathbf{X} (or P) is halfspace symmetric about its center \mathbf{x}_0 if $P(\mathbf{X} \in H) \geq 1/2$ for every closed halfspace² H containing \mathbf{x}_0 .

Note that every centrally symmetric distribution is also angularly symmetric, and every angularly symmetric distribution is halfspace symmetric. Based on the desirable properties a location statistical depth function should possess, Zuo and Serfling [2000a] introduced the subsequent definition.

¹For $\omega \in \mathbb{R}^d$ such that $\mathbf{X}(\omega) = \mathbf{x}_0$, this is considered to equal $\mathbf{0}$.

²By halfspace we mean any set in the form $\{\mathbf{x} \mid \mathbf{x} \cdot \mathbf{u} \leq c\}$ for some $\mathbf{u} \in \mathbb{S}^{d-1}$ and $c \in \mathbb{R}$.

Definition 3 (Location statistical depth) *A mapping $D: \mathbb{R}^d \times P(\mathbb{R}^d) \rightarrow [0, \infty)$ is called a location statistical depth (or simply depth) if it is bounded and satisfies the following properties:*

(A1) *For any random vector $\mathbf{X} \sim P_{\mathbf{X}} \in P(\mathbb{R}^d)$, $\mathbf{A} \in \text{GL}_d$, and $\mathbf{x}, \mathbf{b} \in \mathbb{R}^d$ it holds that $D(\mathbf{A}\mathbf{x} + \mathbf{b}; P_{\mathbf{A}\mathbf{X} + \mathbf{b}}) = D(\mathbf{x}; P_{\mathbf{X}})$.*

(A2) *$D(\mathbf{x}_0; P) = \max_{\mathbf{x} \in \mathbb{R}^d} D(\mathbf{x}; P)$ holds for any symmetric $P \in P(\mathbb{R}^d)$ with center \mathbf{x}_0 .*

(A3) *For any $P \in P(\mathbb{R}^d)$ with a deepest point \mathbf{x}_0 and any $\mathbf{x} \in \mathbb{R}^d$, the mapping $\alpha \mapsto D(\mathbf{x}_0 + \alpha(\mathbf{x} - \mathbf{x}_0); P)$ is non-increasing on $[0, 1]$.*

(A4) *For any fixed $P \in P(\mathbb{R}^d)$, $D(\mathbf{x}; P) \rightarrow 0$ as $\|\mathbf{x}\| \rightarrow \infty$.*

By center, we mean centrality as described by one of the three types introduced in Definition 2.

The above definition should not be taken dogmatically. Numerous reasonable depth functions proposed in the literature may not satisfy all of the properties (A1) – (A4) (for a detailed analysis, see [Mosler and Mozharovskyi, 2022, Section 3]). Instead, these properties serve as guidelines. Property (A1) is often referred to as *a ne invariance*. It emphasizes that statistical procedures based on depth should not depend on the choice of the coordinate system. Properties (A2) and (A3) summarize the idea that the depth function should quantify the degree of centrality of the point of interest. Property (A4) indicates that the measure of centrality of points far from the origin should be close to zero.

Definition 3 defines a population version of depth. In order to define the sample version, consider a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ generated by a distribution $P \in P(\mathbb{R}^d)$. We define the empirical probability measure

$$\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}.$$

Note that this is a random measure depending on the random sample. A sample version of depth is then defined as $D(\cdot; \hat{P}_n)$. The point that maximizes depth can be thought of as a generalization of the univariate median. That is,

$$\arg \max_{\mathbf{x} \in \mathbb{R}^d} D(\mathbf{x}; \hat{P}_n)$$

could serve as an estimator of the location parameter for d -variate random samples.

1.2 Definition of location halfspace depth

The location halfspace depth (introduced by Tukey [1975], therefore sometimes referred to as the Tukey depth) is defined as follows.

Definition 4 (Location halfspace depth) *Let $\mathbf{X} \sim P \in P(\mathbb{R}^d)$ and $\mathbf{x} \in \mathbb{R}^d$. The location halfspace depth (or simply halfspace depth) of \mathbf{x} w.r.t. P is defined as*

$$hD(\mathbf{x}; P) := \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} P(\langle \mathbf{X}, \mathbf{u} \rangle \leq \langle \mathbf{x}, \mathbf{u} \rangle).$$

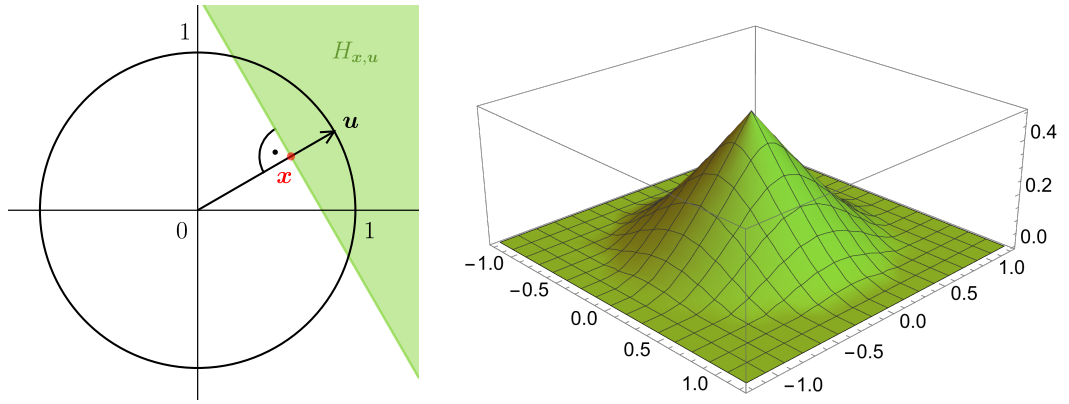
Note that for any fixed $P = P(\mathbb{R}^d)$, the function $hD(\cdot; P)$ is defined everywhere in \mathbb{R}^d and it maps \mathbb{R}^d to a subset of $[0, 1]$. For a point $\mathbf{x} \in \mathbb{R}^d$ and a direction $\mathbf{u} \in S^{d-1}$ denote

$$H_{\mathbf{x}, \mathbf{u}} := \{\mathbf{y} \mid \mathbf{y}, \mathbf{u} \cdot \mathbf{x}, \mathbf{u}\} \subset \mathbb{R}^d$$

the closed halfspace with \mathbf{x} lying on its boundary and inner normal \mathbf{u} . Then, the definition of halfspace depth can be rewritten as

$$hD(\mathbf{x}; P) = \inf_{\mathbf{u} \in S^{d-1}} P(H_{\mathbf{x}, \mathbf{u}}).$$

When calculating the depth of some point $\mathbf{x} \in \mathbb{R}^d$, one is looking for a halfspace $H_{\mathbf{x}, \mathbf{u}}$ containing this point on its boundary with minimal probability. Every such halfspace is called a *minimizing halfspace*, and the corresponding direction \mathbf{u} is called a *minimizing direction*.



(a) Minimizing halfspace $H_{\mathbf{x}, \mathbf{u}}$ at point \mathbf{x} . Its boundary is perpendicular to the line connecting \mathbf{x} and $\mathbf{0}$. (b) 3D plot of the resulting halfspace depth.

Figure 1: The situation from Example 1.

Example 1 Let P be the uniform distribution on the centered unit disc in \mathbb{R}^2 , i.e. $\frac{dP}{d\lambda_2} = \frac{1}{\pi} \mathbf{1}_{B_2(\mathbf{0}, 1)}$. Fix a point $\mathbf{x} \in \mathbb{R}^2$, $\|\mathbf{x}\| < 1$, and consider a direction $\mathbf{u} \in S^1$. Then, the probability $P(H_{\mathbf{x}, \mathbf{u}})$ equals the area of the disk segment, which is cut off from $B_2(\mathbf{0}, 1)$ by $H_{\mathbf{x}, \mathbf{u}}$, divided by π . See Figure 1a for reference. By basic tools of planimetry, we get

$$P(H_{\mathbf{x}, \mathbf{u}}) = \frac{1}{\pi} \left(\arccos(\|\mathbf{x}, \mathbf{u}\|) - \|\mathbf{x}, \mathbf{u}\| \sqrt{1 - \|\mathbf{x}, \mathbf{u}\|^2} \right). \quad (1.1)$$

Note that (1.1) is continuous in $\|\mathbf{x}, \mathbf{u}\|$. In order to minimize it over \mathbf{u} , we differentiate

$$\frac{d}{dt} \left(\arccos(t) - t \sqrt{1 - t^2} \right) = -2 \sqrt{1 - t^2} < 0 \quad \text{for } |t| < 1.$$

By the Cauchy-Schwarz inequality, $\|\mathbf{x}, \mathbf{u}\| \leq \|\mathbf{u}\| \|\mathbf{x}\| = \|\mathbf{x}\|$. If $\mathbf{x} = \mathbf{0}$, we obtain equality if and only if $\mathbf{u} = \mathbf{x} / \|\mathbf{x}\|$ and if $\mathbf{x} = \mathbf{0}$, then $\|\mathbf{x}, \mathbf{u}\| = \|\mathbf{x}\| = 0$ for all \mathbf{u} . Summarizing these observations, we have

$$hD(\mathbf{x}; P) = \inf_{\mathbf{u} \in S^1} P(H_{\mathbf{x}, \mathbf{u}}) = \frac{1}{\pi} \left(\arccos \|\mathbf{x}\| - \|\mathbf{x}\| \sqrt{1 - \|\mathbf{x}\|^2} \right) \quad \text{for } \|\mathbf{x}\| < 1.$$

For a point \mathbf{x} such that $\|\mathbf{x}\| = 1$, its halfspace depth is obviously 0. The resulting halfspace depth is illustrated in Figure 1b.

A minimizing halfspace does not have to be unique. For instance, consider Example 1 and $\mathbf{x} = \mathbf{0}$. Then, every direction $\mathbf{u} \in S^1$ is minimizing. It is also possible that the value of halfspace depth is not attained in any direction, as shown in the following example.

Example 2 Let P be the mixture of the uniform distribution on the centered unit disc and the Dirac measure concentrated in the point $\mathbf{z} = (1, -1)^\top$ with equal weights. See Figure 2 for reference. Consider the point $\mathbf{x} = (1, 0)^\top$. It is easy to see that $hD(\mathbf{x}; P) = 0$. Indeed, let $\mathbf{u}(\theta) = (\cos(\theta), \sin(\theta))^\top$. Then $P(H_{\mathbf{x}, \mathbf{u}(\theta)}) = 0$ as $\theta \rightarrow 0^+$, but $P(H_{\mathbf{x}, \mathbf{u}(\theta)}) > 0$ for all θ .

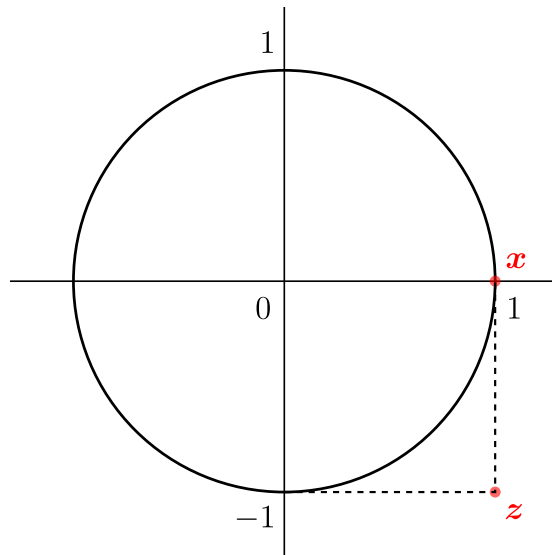


Figure 2: The situation from Example 2: For a mixture of the uniform distribution on the centered unit disc and the Dirac measure concentrated in the point \mathbf{z} , there is no minimizing direction at point \mathbf{x} .

For so-called *smooth* probability distributions, this situation can not happen.

Definition 5 (Smooth distribution) Consider $P \in P(\mathbb{R}^d)$ and $\mathbf{x} \in \mathbb{R}^d$. Then, P is said to be smooth at \mathbf{x} if $P(\text{bd}(H)) = 0$ for every halfspace $H \in \mathbb{R}^d$ with \mathbf{x} on its boundary. More generally, P is said to be smooth if it is smooth at every point $\mathbf{x} \in \mathbb{R}^d$.

If $P \in P(\mathbb{R}^d)$ is smooth at $\mathbf{x} \in \mathbb{R}^d$, then the minimizing halfspace at \mathbf{x} always exists. This is an easy application of the Dominated convergence theorem (Theorem A1 in the Appendix). Indeed, fix a point \mathbf{x} and consider a sequence $\{\mathbf{u}_i\}_{i=1}^\infty \subset S^{d-1}$ such that $P(H_{\mathbf{x}, \mathbf{u}_i}) \rightarrow hD(\mathbf{x}; P)$ as $i \rightarrow \infty$. As S^{d-1} is a compact set, the sequence $\{\mathbf{u}_i\}_{i=1}^\infty$ has a cluster point. Without loss of generality, let $\mathbf{u}_i \rightarrow \mathbf{u} \in S^{d-1}$ as $i \rightarrow \infty$. Then,

$$P(H_{\mathbf{x}, \mathbf{u}}) - P(H_{\mathbf{x}, \mathbf{u}_i}) = \int_{\mathbb{R}^d} (\mathbf{1}_{H_{\mathbf{x}, \mathbf{u}}}(\mathbf{y}) - \mathbf{1}_{H_{\mathbf{x}, \mathbf{u}_i}}(\mathbf{y})) dP(\mathbf{y}) \xrightarrow{i \rightarrow \infty} 0,$$

that is

$$hD(\mathbf{x}; P) = P(H_{\mathbf{x}, \mathbf{u}}).$$

Regarding the assumptions of the Dominated convergence theorem, the difference in indicator functions is obviously bounded, and as $i \rightarrow \infty$, we have $\mathbf{u}_i \rightarrow \mathbf{u}$. Therefore, $\mathbf{1}_{H_{\mathbf{x}, \mathbf{u}}}(\mathbf{y}) - \mathbf{1}_{H_{\mathbf{x}, \mathbf{u}_i}}(\mathbf{y}) \rightarrow 0$ for all $\mathbf{y} \in \mathbb{R}^d \setminus \text{bd}(H_{\mathbf{x}, \mathbf{u}})$, that is P -a.e. because P is smooth at \mathbf{x} .

It is also easy to see that all probability distributions $P \in \mathcal{P}(\mathbb{R}^d)$ which are absolutely continuous w.r.t. λ_d are smooth. But these two conditions are not equivalent. For example, consider the uniform distribution on the unit circle in \mathbb{R}^2 . In the following section, we show that the location halfspace depth possesses more convenient properties for *smooth* distributions.

1.3 Basic properties of location halfspace depth

This section summarizes the well-known properties of the halfspace depth function from the literature. The goal of this work is not to establish fundamental properties but rather to emphasize robustness and minimax optimality of the location halfspace depth. Therefore, we do not provide proof for most theorems in this section.

The most important property of the location halfspace depth is that under certain conditions, it is a location statistical depth in the sense of Definition 3, as stated in the following theorem.

Theorem 6 *The location halfspace depth hD always satisfies properties (A1), (A3) and (A4) from Definition 3. For a smooth P , it also satisfies property (A2) for the halfspace symmetry (therefore also for the central and angular symmetry).*

This claim follows from several other fundamental properties of the halfspace depth that we are going to introduce before proving Theorem 6. The first easy consequence of its definition is quasiconcavity. A function $f: \mathbb{R}^d \rightarrow [0, \infty)$ is said to be quasiconcave if for any $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d$ and $\alpha \in [0, 1]$ it holds that

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \geq \min \{f(\mathbf{x}_1), f(\mathbf{x}_2)\}.$$

Theorem 7 *For any fixed $P \in \mathcal{P}(\mathbb{R}^d)$, the halfspace depth $hD(\cdot; P)$ is quasiconcave.*

Proof. Refer to [Rousseeuw and Ruts, 1999, Proposition 1].

This property illuminates the geometry of the halfspace depth function. To be specific, we define α -depth regions as the upper-level sets of the halfspace depth function.

Definition 8 (α -depth region) *For any fixed $P \in \mathcal{P}(\mathbb{R}^d)$ and $\alpha \in [0, 1]$, the α -depth region of P is defined as*

$$D(\alpha; P) := \left\{ \mathbf{x} \in \mathbb{R}^d \mid hD(\mathbf{x}; P) \geq \alpha \right\} \subseteq \mathbb{R}^d.$$

With this notion, the quasiconcavity of $hD(\cdot; P)$ is equivalent to the convexity of all α -depth regions. Indeed, fix $\alpha \in [0, 1]$ and $\mathbf{x}_1, \mathbf{x}_2 \in D(\alpha; P)$. Then, for any $\beta \in [0, 1]$ we have

$$hD(\beta\mathbf{x}_1 + (1 - \beta)\mathbf{x}_2; P) \geq \min \{hD(\mathbf{x}_1; P), hD(\mathbf{x}_2; P)\} \geq \alpha \quad (1.2)$$

by the quasiconcavity of the halfspace depth. Hence $\beta\mathbf{x}_1 + (1 - \beta)\mathbf{x}_2 \in D(\alpha; P)$ and $D(\alpha; P)$ is a convex set. The other implication can be shown in a similar way.

Another property of interest is continuity. A function $f: \mathbb{R}^d \rightarrow [0, \infty)$ is defined to be upper semicontinuous if for any $t \geq 0$ the sets $\{\mathbf{x} \in \mathbb{R}^d / f(\mathbf{x}) \leq t\}$ are closed. Every continuous function is obviously upper semicontinuous. The halfspace depth function is always upper semicontinuous. Further, it is continuous under certain conditions. Specifically, for a smooth distribution P , the function $hD(\cdot; P)$ is continuous, as demonstrated in the following theorem.

Theorem 9 *For any fixed $P \in \mathcal{P}(\mathbb{R}^d)$, the function $hD(\cdot; P)$ is upper semicontinuous. If P is smooth, then $hD(\cdot; P)$ is continuous.*

Proof. Refer to [Donoho and Gasko, 1992, Lemma 6.1]. This lemma proves the continuity of the halfspace depth for absolutely continuous distributions. However, it is easy to see that its proof works for any smooth distribution.

The upper semi-continuity, together with Theorem 7, implies that the α -depth regions possess very plausible qualities to work with.

Theorem 10 *For any fixed $P \in \mathcal{P}(\mathbb{R}^d)$ and $\alpha \in [0, 1]$, the α -depth region $D(\alpha; P)$ is a closed convex set. If $\alpha = 0$, $D(\alpha; P)$ is also bounded, hence compact.*

Proof. This is a consequence of our Theorem 7, Theorem 9 and [Rousseeuw and Ruts, 1999, Proposition 5].

In Theorem 9, we noticed that the halfspace depth behaves especially nicely when the underlying probability distribution is smooth. This condition also implies that the halfspace depth function can be bounded from above by $1/2$, as formulated in the following theorem.

Theorem 11 *Let $P \in \mathcal{P}(\mathbb{R}^d)$ be smooth. Then $hD(\mathbf{x}; P) \leq 1/2$ holds for all $\mathbf{x} \in \mathbb{R}^d$.*

Proof. Refer to [Rousseeuw and Ruts, 1999, Proposition 10]. This proposition establishes the upper bound for absolutely continuous distributions. However, it is evident that the proof applies to any smooth distribution.

Theorem 11 holds only for smooth probability distributions. Otherwise, there is no general upper bound smaller than 1. For example, consider $hD(\cdot; \delta_{\mathbf{x}})$ for some $\mathbf{x} \in \mathbb{R}^d$. Then, it follows from the definition that $hD(\mathbf{x}; \delta_{\mathbf{x}}) = 1$. Now, we prove Theorem 6 to demonstrate the use of some of its basic properties mentioned above.

Proof of Theorem 6.

(A1) Let $\mathbf{X} \sim P_{\mathbf{X}} \in P(\mathbb{R}^d)$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{x}, \mathbf{b} \in \mathbb{R}^d$. Then

$$\begin{aligned} hD(\mathbf{x}; P_{\mathbf{X}}) &= \inf_{\mathbf{u} \in S^{d-1}} P(\langle \mathbf{X}, \mathbf{u} \rangle \leq \langle \mathbf{x}, \mathbf{u} \rangle) \\ &= \inf_{\mathbf{u} \in S^{d-1}} P\left(\frac{\langle \mathbf{X}, \mathbf{A}^T \mathbf{u} \rangle}{\|\mathbf{A}^T \mathbf{u}\|} \leq \frac{\langle \mathbf{x}, \mathbf{A}^T \mathbf{u} \rangle}{\|\mathbf{A}^T \mathbf{u}\|}\right) \\ &= \inf_{\mathbf{u} \in S^{d-1}} P(\langle \mathbf{A}\mathbf{X}, \mathbf{u} \rangle \leq \langle \mathbf{A}\mathbf{x}, \mathbf{u} \rangle) \\ &= \inf_{\mathbf{u} \in S^{d-1}} P(\langle \mathbf{A}\mathbf{X} + \mathbf{b}, \mathbf{u} \rangle \leq \langle \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{u} \rangle) \\ &= hD(\mathbf{A}\mathbf{x} + \mathbf{b}; P_{\mathbf{A}\mathbf{X} + \mathbf{b}}). \end{aligned}$$

In the second equality we used that \mathbf{A}^T is non-singular, hence the mapping $\mathbf{u} \mapsto (\mathbf{A}^T \mathbf{u}) / \|\mathbf{A}^T \mathbf{u}\|$ maps S^{d-1} onto S^{d-1} .

(A2) Consider $\mathbf{X} \sim P_{\mathbf{X}}$ for a smooth $P_{\mathbf{X}}$ and suppose that \mathbf{X} is halfspace symmetric about $\mathbf{x}_0 \in \mathbb{R}^d$. Using the affine invariance, we can assume that $\mathbf{x}_0 = \mathbf{0}$. By Theorem 11, it suffices to show that $hD(\mathbf{0}; P_{\mathbf{X}}) = 1/2$. For any $\mathbf{u} \in S^{d-1}$, we have (using the smoothness of P) $P(\langle \mathbf{X}, \mathbf{u} \rangle = 0) = 0$. Therefore, we obtain

$$P(\langle \mathbf{X}, \mathbf{u} \rangle \leq 0) = 1 - P(\langle \mathbf{X}, \mathbf{u} \rangle > 0)$$

which, using the halfspace symmetry, implies $P(\langle \mathbf{X}, \mathbf{u} \rangle \leq 0) = 1/2$ for all $\mathbf{u} \in S^{d-1}$, hence $hD(\mathbf{0}; P_{\mathbf{X}}) = 1/2$.

(A3) Consider $\mathbf{X} \sim P_{\mathbf{X}}$ for a smooth $P_{\mathbf{X}}$ and let \mathbf{x}_0 be a deepest point, i.e., let $\mathbf{x}_0 = \arg \max_{\mathbf{x} \in \mathbb{R}^d} hD(\mathbf{x}; P_{\mathbf{X}})$. The function $hD(\cdot; P_{\mathbf{X}})$ is quasiconcave by Theorem 7. For a contradiction, suppose that there are $\mathbf{x} \in \mathbb{R}^d$ and $0 < \lambda_1 < \lambda_2 < 1$ such that

$$hD(\mathbf{x}_0 + \lambda_1(\mathbf{x} - \mathbf{x}_0); P_{\mathbf{X}}) < hD(\mathbf{x}_0 + \lambda_2(\mathbf{x} - \mathbf{x}_0); P_{\mathbf{X}}).$$

But due to the fact that $\mathbf{x}_0 + \lambda_1(\mathbf{x} - \mathbf{x}_0)$ lies in the line segment between \mathbf{x}_0 and $\mathbf{x}_0 + \lambda_2(\mathbf{x} - \mathbf{x}_0)$, this contradicts the quasiconcavity.

(A4) Fix $P \in P(\mathbb{R}^d)$. Consider $\varepsilon > 0$. Then by Theorem 10, the ε -depth region $D(\varepsilon; P)$ is bounded. This implies that $hD(\mathbf{x}; P) < \varepsilon$ for all \mathbf{x} such that $\|\mathbf{x}\|$ is large enough. Therefore, $hD(\mathbf{x}; P) \rightarrow 0$ as $\|\mathbf{x}\| \rightarrow \infty$.

Let us also show what the halfspace depth looks like for Gaussian distributions to demonstrate how it works.

Example 3 Consider $\mathbf{X} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for some $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} \in \text{PD}_d$. Recall that the distribution function of the univariate standard Gaussian distribution is denoted by Φ . For all $\mathbf{x} \in \mathbb{R}^d$ we have

$$hD(\mathbf{x}; N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = hD(\boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu}); N_d(\mathbf{0}, \mathbf{I}_d)) \quad (1.3)$$

by a line invariance. Therefore, we should first calculate $hD(\mathbf{y}; N_d(\mathbf{0}, \mathbf{I}_d))$ for all $\mathbf{y} \in \mathbb{R}^d$. For $\mathbf{Y} \sim N_d(\mathbf{0}, \mathbf{I}_d)$ and all $\mathbf{u} \in S^{d-1}$ we have $\mathbf{Y} \cdot \mathbf{u} \sim N_1(0, 1)$. Hence,

$$P(\mathbf{Y} \cdot \mathbf{u} \leq \mathbf{y} \cdot \mathbf{u}) = 1 - \Phi(\mathbf{y} \cdot \mathbf{u}) = 1 - \Phi(\|\mathbf{y}\|), \quad (1.4)$$

where the inequality is due to the Cauchy-Schwarz theorem and for $\mathbf{u} = \mathbf{y} / \|\mathbf{y}\|$, we obtain equality. Therefore, this choice of \mathbf{u} is the minimizing direction and $hD(\mathbf{y}; N_d(\mathbf{0}, \mathbf{I}_d)) = 1 - \Phi(\|\mathbf{y}\|)$. Returning back to the case of a general Gaussian distribution, by combining (1.3) and (1.4), we obtain that

$$hD(\mathbf{x}; N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = 1 - \Phi\left(\left\|\boldsymbol{\Sigma}^{-1/2}(\mathbf{x} - \boldsymbol{\mu})\right\|\right).$$

Subsequently, we can see that the halfspace depth of Gaussian distribution is uniquely maximized for the mean $\boldsymbol{\mu}$ with the maximum depth $1/2$. Moreover, for $\alpha \in (0, 1/2)$ we have

$$D(\alpha; P) = \left\{ \mathbf{x} \in \mathbb{R}^d \mid \sqrt{(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})} \leq \Phi^{-1}(1 - \alpha) \right\}.$$

These are the so-called Mahalanobis ellipsoids centered about $\boldsymbol{\mu}$. The $(1/2)$ -depth region $D(1/2; P)$ equals the singleton $\{\boldsymbol{\mu}\}$ and $D(\alpha; P) = \emptyset$ for $\alpha > 1/2$. This corresponds to the fact that the Gaussian distribution is centrally symmetric about its mean.

Calculating the depth function analytically is not an easy task in general. For some classes of distributions, it can be done relatively easily. Those are, for example, elliptically symmetric distributions and symmetric stable distributions, see [Chen and Tyler, 2004, Theorem 3.1]. For uniform distributions on convex bodies, this problem is closely related to the concept of floating bodies from geometry, as explored by Nagy et al. [2019].

1.4 Sample location halfspace depth and halfspace median

For halfspace depth to be useful from a statistical perspective, we must be able to estimate it based on a random sample from a distribution. Therefore, let us define the sample halfspace depth as follows.

Definition 12 (Sample halfspace depth) *Consider a random sample $\{\mathbf{X}_i\}_{i=1}^n$ from a distribution $P \in \mathcal{P}(\mathbb{R}^d)$ and $\mathbf{x} \in \mathbb{R}^d$. The sample halfspace depth of \mathbf{x} w.r.t. $\{\mathbf{X}_i\}_{i=1}^n$ is defined as*

$$\begin{aligned} hD(\mathbf{x}; \{\mathbf{X}_i\}_{i=1}^n) &:= hD(\mathbf{x}; \hat{P}_n) = \min_{\mathbf{u} \in S^{d-1}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{X}_i \cdot \mathbf{u} \leq \mathbf{x} \cdot \mathbf{u}\}} \\ &= \min_{\mathbf{u} \in S^{d-1}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{X}_i \in H_{\mathbf{x}, \mathbf{u}}\}}. \end{aligned}$$

This definition implies that the sample halfspace depth of \mathbf{x} is the smallest number of observations that can lie in a halfspace with \mathbf{x} on its boundary, divided by n . A fundamental quality of estimators lies in how they behave for large sample sizes n . It turns out that the sample halfspace depth is strongly uniformly consistent.

Theorem 13 For a random sample $\{\mathbf{X}_i\}_{i=1}^n \sim P \in \mathcal{P}(\mathbb{R}^d)$ we have

$$\sup_{\mathbf{x} \in \mathbb{R}^d} |hD(\mathbf{x}; \{\mathbf{X}_i\}_{i=1}^n) - hD(\mathbf{x}; P)| \xrightarrow{a.s.} 0.$$

Proof. Refer to [Donoho and Gasko, 1992, Section 6].

As mentioned above, the halfspace depth tells us how deep a given point is w.r.t. a distribution. It is, therefore, natural to define the deepest point as a generalization of the one-dimensional median. Such a point always exists, as stated in the following theorem.

Theorem 14 Consider a probability measure $P \in \mathcal{P}(\mathbb{R}^d)$. Then, there exists at least one point $\mathbf{y} \in \mathbb{R}^d$ such that $hD(\mathbf{y}; P) = \sup_{\mathbf{x} \in \mathbb{R}^d} hD(\mathbf{x}; P)$.

Proof. Refer to [Rousseeuw and Ruts, 1999, Proposition 7].

Therefore, we can define the following notions.

Definition 15 (Maximal depth) For $P \in \mathcal{P}(\mathbb{R}^d)$, we define the maximal depth

$$\Pi(P) := \max_{\mathbf{x} \in \mathbb{R}^d} hD(\mathbf{x}; P)$$

and the set of points of maximal halfspace depth

$$\text{Med}(P) := D(\Pi(P); P).$$

Using Theorems 10 and 14, the set $\text{Med}(P)$ is non-empty, convex, and compact. Therefore, the following definition is well grounded.

Definition 16 (Halfspace median) Define the functional $T: \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}^d$ so that it maps $P \in \mathcal{P}(\mathbb{R}^d)$ to the centroid³ of $\text{Med}(P)$, i.e. the centroid of

$$\arg \max_{\mathbf{x} \in \mathbb{R}^d} hD(\mathbf{x}; P).$$

The value $T(P)$ is called the halfspace median (or Tukey's median) of P . For a random sample $\{\mathbf{X}_i\}_{i=1}^n \sim P \in \mathcal{P}(\mathbb{R}^d)$, we define the sample halfspace median (also called sample Tukey's median) as $T(\hat{P}_n)$.

From Theorem 11, we know that the maximal depth is bounded from above by $1/2$ for smooth distributions. Conversely, it turns out that the lower bound of $1/(d+1)$ for the maximum depth holds for all distributions.

Theorem 17 For any $P \in \mathcal{P}(\mathbb{R}^d)$ we have $\Pi(P) \geq \frac{1}{d+1}$.

³Centroid of a bounded, non-empty convex set $A \subset \mathbb{R}^d$ is defined as the expectation of the uniform distribution on A .

Proof. This claim is a consequence of Helly’s theorem, which is a basic result in the field of discrete geometry. Proof can be found in [Rousseeuw and Ruts, 1999, Proposition 9].

The previous theorem implies that for any smooth $P \in \mathcal{P}(\mathbb{R}^d)$, we have $\Pi(P) = 1/2$. To conclude this section, we summarize the basic properties of the halfspace median in the following theorem.

Theorem 18 *Let $P \in \mathcal{P}(\mathbb{R}^d)$, then it holds that*

- $T(P)$ is a location parameter of P in terms of Definition 1.
- If $d = 1$, then $T(P)$ coincides with the standard univariate median⁴.
- The halfspace depth of $T(P)$ is at least $1/(d + 1)$. If P is smooth, then the halfspace depth of $T(P)$ is bounded from above by $1/2$.
- If \hat{P}_n is the empirical measure corresponding to P , then $T(\hat{P}_n) \xrightarrow[n]{a.s.} T(P)$.

Proof. The first property, affine equivariance, easily follows from the affine invariance of the halfspace depth (see Theorem 6). The second property arises from the fact that if $d = 1$, then $hD(x; P) = \min \{F(x), 1 - F(x^-)\}$, where F is the distribution function of P . The third property is derived from Theorem 11 and Theorem 17. The strong consistency follows from Lemma S.2.5 in the supplement to [Paindaveine and Van Bever, 2018].

1.5 Robustness of halfspace median

In this section, we discuss the robustness properties of the halfspace median. When inferring using parametric statistics, one assumes a parametric model. The properties of the estimators usually rely strongly on the fulfillment of the model assumptions. For example, the exact t-test assumes the normality of the underlying distribution. Without normality, the t-test works only asymptotically. The area of robust statistics explores how much statistical methods are affected by outliers or other small departures from model assumptions.

Take, for example, the problem of estimating the expectation (center) of a Gaussian distribution. The maximum likelihood theory yields the sample mean as an estimator. This estimator is unbiased and consistent. Furthermore, by the Cramér–Rao bound, we see that it is an efficient estimator. Thus, assuming a Gaussian model, it exhibits excellent properties. However, let us now assume that we are not entirely certain whether the data come from a Gaussian distribution. In the following text, we will demonstrate that in such situations, the sample halfspace median has good robustness properties against such violations of assumptions.

⁴This holds if we define the univariate median of $X \in \mathcal{P}(\mathbb{R})$ as the centerpoint of all points $x \in \mathbb{R}$ such that $P(X \leq x) \geq 1/2$ and $P(X \geq x) \geq 1/2$.

Assume that we have a statistical functional $\Xi: P(\mathbb{R}^d) \rightarrow \mathbb{R}^d$. In this section, we consider several measures of its robustness. The departure from model assumptions will be represented by contaminated distributions.

Definition 19 (Contaminated distribution) *Let $P, Q \in P(\mathbb{R}^d)$, $\varepsilon \in [0, 1]$ and $\mathbf{x} \in \mathbb{R}^d$. We consider the distribution $(1 - \varepsilon)P + \varepsilon Q \in P(\mathbb{R}^d)$ defined as*

$$((1 - \varepsilon)P + \varepsilon Q)(B) = (1 - \varepsilon)P(B) + \varepsilon Q(B) \quad \text{for } B \in \mathcal{B}^d.$$

Specifically, we denote

$$\begin{aligned} P_{(\varepsilon, Q)} &:= (1 - \varepsilon)P + \varepsilon Q, \\ P_{(\varepsilon, \mathbf{x})} &:= (1 - \varepsilon)P + \varepsilon \delta_{\mathbf{x}}. \end{aligned}$$

One can define the influence function and gross error sensitivity using the notion of contaminated distributions.

Definition 20 (Influence function, gross error sensitivity) *Consider a location statistical functional $\Xi: P(\mathbb{R}^d) \rightarrow \mathbb{R}^d$ and a probability measure $P \in P(\mathbb{R}^d)$. The influence function of Ξ at P is a function of $\mathbf{x} \in \mathbb{R}^d$ defined as*

$$IF(\mathbf{x}; \Xi, P) := \lim_{\varepsilon \rightarrow 0^+} \frac{\Xi(P_{(\varepsilon, \mathbf{x})}) - \Xi(P)}{\varepsilon},$$

provided that the limit exists for all $\mathbf{x} \in \mathbb{R}^d$. The gross error sensitivity of Ξ at P is defined as

$$\gamma(\Xi, P) := \sup_{\mathbf{x} \in \mathbb{R}^d} IF(\mathbf{x}; \Xi, P).$$

The influence function (being one-sided Gateaux derivative of Ξ) describes how the contamination by a point \mathbf{x} with infinitesimally small mass ε changes the value of the statistical functional, normalized by ε . In general, one seeks functionals with bounded influence function to get desirable robustness properties. On the other hand, gross error sensitivity evaluates the maximum effect that an infinitesimally small point-mass contamination can have, indicating its reliability under extreme conditions.

Both previous measures indicate robustness only w.r.t. a point-mass contamination. A stronger approach is to measure the maximum impact that any form of contamination can exert on a functional. This leads to a definition of maximum bias and breakdown point.

Definition 21 (Maximum bias, breakdown point) *Consider a statistical functional $\Xi: P(\mathbb{R}^d) \rightarrow \mathbb{R}^d$, a probability measure $P \in P(\mathbb{R}^d)$ and $\varepsilon \in [0, 1]$. The maximum bias function of Ξ at P is defined as*

$$B(\varepsilon; \Xi, P) := \sup_{Q \in P(\mathbb{R}^d)} \left\| \Xi(P_{(\varepsilon, Q)}) - \Xi(P) \right\|.$$

The breakdown point of Ξ at P is defined as

$$\varepsilon(\Xi, P) := \inf \{ \varepsilon \in (0, 1] \mid B(\varepsilon; \Xi, P) = \infty \}.$$

The breakdown point is the maximal proportion ε of contamination the statistical functional can handle before producing an arbitrarily large bias.

Recall that the functional corresponding to the halfspace median is denoted by T . The *support* of a probability measure $P \in \mathcal{P}(\mathbb{R}^d)$, usually denoted by $\text{supp}(P)$, is defined as the complement of the union of all open sets which are P -null sets, i.e. the smallest closed set A such that $P(\mathbb{R}^d \setminus A) = 0$. Using this notion, we say that P has *contiguous support* if $\text{supp}(P)$ cannot be separated by a zero P -mass slab of a non-empty interior. Here, by slab, we mean any set of the form $\{\mathbf{x} \in \mathbb{R}^d \mid a < \mathbf{x} \cdot \mathbf{u} < b\}$ for some $\mathbf{u} \in S^{d-1}$ and $a, b \in \mathbb{R}$. If we assume that a probability measure P has contiguous support and is smooth at some $\mathbf{x} \in \text{Med}(P)$, then $hD(\cdot; P)$ is strictly monotone. That is, for all $\alpha \in (0, \Pi(P))$, we have

$$\text{int}(D(\alpha; P)) = \{\mathbf{x} \in \mathbb{R}^d \mid hD(\mathbf{x}; P) > \alpha\}.$$

This has been proven by Laketa and Nagy [2022, Section 4.3]. Also, recall that P is said to be halfspace symmetric about \mathbf{x} if $P(H) = 1/2$ for every halfspace H with \mathbf{x} on its boundary (see Definition 2). Zuo and Serfling [2000b, Theorem 2.1] have shown that for every probability measure P halfspace symmetric about \mathbf{x} , we have $\text{Med}(P) = \{\mathbf{x}\}$ unless P is concentrated on a line and its distribution has more than one univariate median. That is, P is supported on some one-dimensional affine subspace of \mathbb{R}^d , and $\text{Med}(P)$ is a line segment in this subspace, symmetric around \mathbf{x} .

In order to determine the influence function of T at P , we need to calculate the halfspace median for the contaminated distribution $P_{(\varepsilon, \mathbf{x})}$. This is not an easy task in general, but as shown in the following theorem, it can be done after imposing an additional assumption on P .

Definition 22 Consider $P \in \mathcal{P}(\mathbb{R}^d)$, $\alpha \in (0, \Pi(P)]$ and $\mathbf{u} \in S^{d-1}$. Then, $r(\alpha, \mathbf{u})$ is defined as the radius of the α -depth region $D(\alpha; P)$ along the direction $\mathbf{u} \in S^{d-1}$, i.e. the length of

$$D(\alpha; P) \cap \{T(P) + s \mathbf{u} \mid s \geq 0\}.$$

In the subsequent theorem, we restrict ourselves to $d > 1$, which is the interesting case. For $d = 1$, the halfspace median coincides with the usual univariate median, and its robustness properties are well known. See, for example, [Huber and Ronchetti, 2009, Section 3.2].

Theorem 23 Let $d > 1$. Consider a probability measure $P \in \mathcal{P}(\mathbb{R}^d)$, which is halfspace symmetric about $\mathbf{x}_0 \in \mathbb{R}^d$, smooth at \mathbf{x}_0 and has contiguous support. Then for any $\varepsilon \in [0, 1/3)$ we have

$$T(P_{(\varepsilon, \mathbf{x})}) = \begin{cases} \mathbf{x} & \text{for } \mathbf{x} \in \text{int}\left(D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)\right), \\ \mathbf{x}_0 + \frac{1}{2} r\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}, \frac{\mathbf{x}-\mathbf{x}_0}{\|\mathbf{x}-\mathbf{x}_0\|}\right) \frac{\mathbf{x}-\mathbf{x}_0}{\|\mathbf{x}-\mathbf{x}_0\|} & \text{for } \mathbf{x} \notin \text{int}\left(D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)\right). \end{cases}$$

Note This theorem is due to Chen and Tyler [2002, Theorem 3.1]. However, the authors do not assume contiguous support of P . Without this assumption, the theorem does not hold, as demonstrated in Example 4. Additionally, Chen and Tyler [2002] assume that P is absolutely continuous w.r.t. the Lebesgue measure. This is not necessary. We replace this assumption with smoothness of P at its center. Finally, that paper does not explicitly state that $d > 1$ is required for the proof to be valid.

Proof. Using the affine equivariance of T , we can assume $\mathbf{x}_0 = \mathbf{0}$. If $\mathbf{x} = \mathbf{0}$, the claim is obvious. Let $\mathbf{x} = \mathbf{0}$.

- By the assumption of halfspace symmetry and the smoothness of P at $\mathbf{0}$, we have $P(H_{\mathbf{0},\mathbf{u}}) = 1/2$ for all $\mathbf{u} \in S^{d-1}$. Evidently, there exists $\mathbf{u} \in S^{d-1}$ such that $H_{\mathbf{0},\mathbf{u}}$ does not include \mathbf{x} . Therefore, $hD(\mathbf{0}; P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$.
- Consider any $\mathbf{y} \in \{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. Then, there is $\mathbf{u} \in S^{d-1}$ such that $H_{\mathbf{y},\mathbf{u}}$ does not include the set $\{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. Because P has contiguous support and $H_{\mathbf{y},\mathbf{u}}$ does not include $\mathbf{0}$, we have $P(H_{\mathbf{y},\mathbf{u}}) < 1/2$. Therefore, $P_{(\varepsilon,\mathbf{x})}(H_{\mathbf{y},\mathbf{u}}) = (1 - \varepsilon) P(H_{\mathbf{y},\mathbf{u}}) < \frac{1-\varepsilon}{2}$, hence $hD(\mathbf{y}; P_{(\varepsilon,\mathbf{x})}) < \frac{1-\varepsilon}{2}$.
- On the other hand, for any $\mathbf{y} \in \{\alpha\mathbf{x} / \alpha \in (0, 1)\}$, we can choose a sequence $\{\mathbf{u}_i\}_{i=1}^{\infty} \subset S^{d-1}$ such that for all $i \in \mathbb{N}$, $H_{\mathbf{y},\mathbf{u}_i}$ contains $\mathbf{0}$, does not contain \mathbf{x} and the distance between $\text{bd}(H_{\mathbf{y},\mathbf{u}_i})$ and $\mathbf{0}$ tends to 0 as $i \rightarrow \infty$. This is possible only because we assume that $d > 1$. As a result, $P_{(\varepsilon,\mathbf{x})}(H_{\mathbf{y},\mathbf{u}_i}) = (1 - \varepsilon)P(H_{\mathbf{y},\mathbf{u}_i}) \rightarrow \frac{1-\varepsilon}{2}$ as $i \rightarrow \infty$ and $hD(\mathbf{y}; P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$.

We have shown that $\Pi(P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$, $\text{Med}(P_{(\varepsilon,\mathbf{x})}) = \{\alpha\mathbf{x} / \alpha \in [0, 1]\}$ and $hD(\cdot; P_{(\varepsilon,\mathbf{x})})$ is bounded from above by $\frac{1-\varepsilon}{2}$ on $\{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. Now, we consider three different cases.

- If $\mathbf{x} \in \text{int}\left(D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)\right)$, then $hD(\mathbf{x}; P_{(\varepsilon,\mathbf{x})}) > (1 - \varepsilon)\frac{1-3\varepsilon}{2(1-\varepsilon)} + \varepsilon = \frac{1-\varepsilon}{2}$ since by the assumptions $hD(\cdot; P)$ is strictly monotone. Therefore, \mathbf{x} is the only point of the maximal depth, hence $T(P_{(\varepsilon,\mathbf{x})}) = \mathbf{x}$, as claimed.
- Let $\mathbf{x} \in \text{bd}\left(D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)\right)$. Using the strict monotonicity of $hD(\cdot; P)$, we have that $hD(\mathbf{x}; P) = \frac{1-3\varepsilon}{2(1-\varepsilon)}$. Therefore, $hD(\mathbf{x}; P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$, hence $\Pi(P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$. We already know that $hD(\mathbf{0}; P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$ and since the depth regions are convex, we have $\text{Med}(P_{(\varepsilon,\mathbf{x})}) = \{\alpha\mathbf{x} / \alpha \in [0, 1]\} = D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right) \cap \{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. Because the centroid of a line segment equals its center point, the claim follows readily.
- Finally, if $\mathbf{x} \in D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)$, then $hD(\mathbf{x}; P) < \frac{1-3\varepsilon}{2(1-\varepsilon)}$. As a consequence, $hD(\mathbf{x}; P_{(\varepsilon,\mathbf{x})}) < \frac{1-3\varepsilon}{2} + \varepsilon = \frac{1-\varepsilon}{2}$. Consider $\mathbf{y} \in \{\alpha\mathbf{x} / \alpha \in (0, 1)\}$. Much like before, if $\mathbf{y} \in D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)$, then $hD(\mathbf{y}; P_{(\varepsilon,\mathbf{x})}) < \frac{1-\varepsilon}{2}$. On the other hand, if $\mathbf{y} \in D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right)$, then $(1 - \varepsilon)P(H) = \frac{1-3\varepsilon}{2} = \frac{1-\varepsilon}{2} - \varepsilon$ for all halfspaces H with \mathbf{y} on boundary. Also, $(1 - \varepsilon)P(H) = \frac{1-\varepsilon}{2}$ for all such halfspaces that contain $\mathbf{0}$ and hence do not contain \mathbf{x} . As a consequence, $P_{(\varepsilon,\mathbf{x})}(H) = \frac{1-\varepsilon}{2}$ for all halfspaces H with \mathbf{y} on boundary, therefore $hD(\mathbf{y}; P_{(\varepsilon,\mathbf{x})}) = \frac{1-\varepsilon}{2}$. We have used the previous observation that $hD(\cdot; P_{(\varepsilon,\mathbf{x})})$ is bounded from above by $\frac{1-\varepsilon}{2}$ on $\{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. Summarizing this, we have $\text{Med}(P_{(\varepsilon,\mathbf{x})}) = D\left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P\right) \cap \{\alpha\mathbf{x} / \alpha \in [0, 1]\}$. The claim follows just like in the previous case.

Example 4 Let $P = P(\mathbb{R}^2)$ be the uniform distribution on the union of two rectangles with vertices

$$(-8, -6)^\top, (-5, -6)^\top, (-5, 6)^\top, (-8, 6)^\top \text{ and } (5, -6)^\top, (8, -6)^\top, (8, 6)^\top, (5, 6)^\top,$$

see Figure 3. Then, P is evidently halfspace symmetric about the origin $\mathbf{0}$ and is smooth at $\mathbf{0}$. Consider $\varepsilon = 1/4$ and $\mathbf{x} = (0, 4)^\top$. We have $1 - \varepsilon = 3/4$ and $\varepsilon/(1 - \varepsilon) = 1/3$. By considering all halfplanes with \mathbf{x} on their boundary, it is easy to see that $hD(\mathbf{x}; P) = 1/6$ and

$$\mathbf{x} \in \text{bd} \left(D \left(\frac{1 - 3\varepsilon}{2(1 - \varepsilon)}; P \right) \right) = \text{bd} \left(D \left(\frac{1}{6}; P \right) \right),$$

hence

$$r \left(\frac{1 - 3\varepsilon}{2(1 - \varepsilon)}, \frac{\mathbf{x}}{\|\mathbf{x}\|} \right) = \|\mathbf{x}\| = 4. \quad (1.5)$$

Further, let the points \mathbf{a} and \mathbf{b} be defined as in Figure 3 and denote by D the convex hull of the set $\{\mathbf{0}, \mathbf{a}, \mathbf{b}, \mathbf{x}\}$. Then, using the definition of $P_{(\varepsilon, \mathbf{x})}$, after considering all possible halfplanes, we realize that

- $hD(\mathbf{0}; P_{(\varepsilon, \mathbf{x})}) = (1 - \varepsilon)/2 = \frac{3}{8}$,
- $hD(\mathbf{x}; P_{(\varepsilon, \mathbf{x})}) = (1 - \varepsilon)hD(\mathbf{x}; P) + \varepsilon = \frac{3}{8}$,
- $hD(\mathbf{a}; P_{(\varepsilon, \mathbf{x})}) = hD(\mathbf{b}; P_{(\varepsilon, \mathbf{x})}) = (1 - \varepsilon)/2 = \frac{3}{8}$ and
- $hD(\mathbf{y}; P_{(\varepsilon, \mathbf{x})}) < \frac{1 - \varepsilon}{2}$ for all $\mathbf{y} \notin D$.

Therefore, using the quasiconcavity of the halfspace depth, we deduce that

$$D = D \left(\frac{3}{8}; P_{(\varepsilon, \mathbf{x})} \right).$$

It is also easy to see that for all points \mathbf{y} between $\mathbf{0}$ and \mathbf{x} , it holds that $hD(\mathbf{y}; P_{(\varepsilon, \mathbf{x})}) = \frac{1 - \varepsilon}{2} = \frac{3}{8}$. Therefore, using the quasiconcavity of $hD(\cdot; P_{(\varepsilon, \mathbf{x})})$ again, it holds that

$$\Pi(P_{(\varepsilon, \mathbf{x})}) = \frac{3}{8} \quad \text{and} \quad \text{Med}(P_{(\varepsilon, \mathbf{x})}) = D.$$

As a consequence, $T(P_{(\varepsilon, \mathbf{x})})$ is the centroid of D , marked as a red point in Figure 3. However, this point does not lie exactly in the middle between $\mathbf{0}$ and \mathbf{x} . This, together with (1.5), shows that the assumption of contiguous support is necessary in Theorem 23.

From the preceding Theorem 23, we observe that as we move \mathbf{x} farther away from \mathbf{x}_0 , the Tukey median of the contaminated distribution $P_{(\varepsilon, \mathbf{x})}$ (as a function of the contaminating point \mathbf{x}), initially increases linearly along a ray and subsequently decreases to a fixed value. This behavior will be illustrated in Example 5, where P is a Gaussian distribution. Theorem 23 also allows us, under certain assumptions, to determine the influence function and gross error sensitivity for the Tukey median.

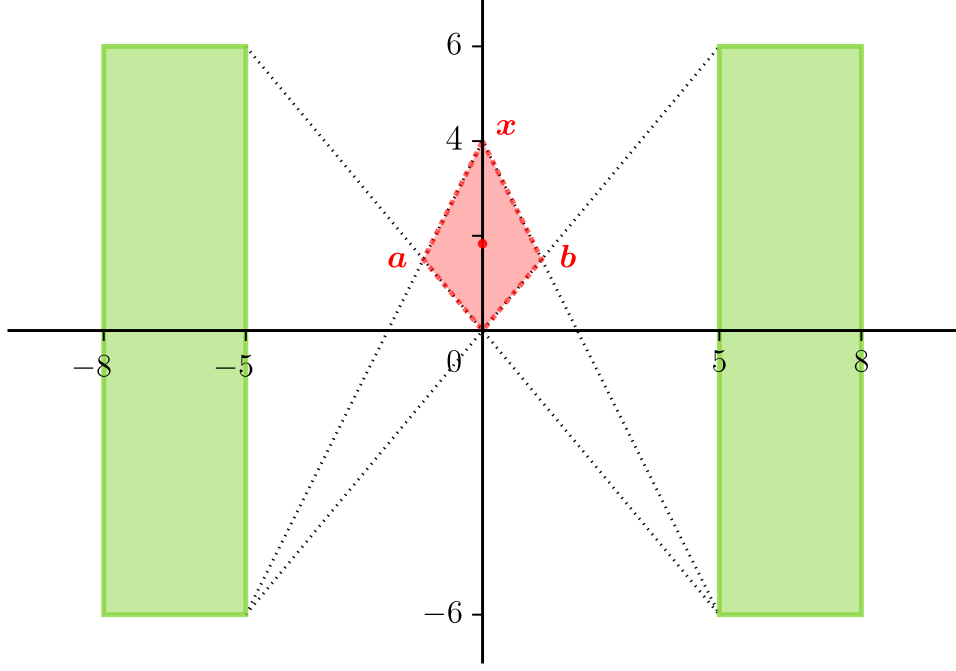


Figure 3: The situation from Example 4. The Tukey median of a uniform distribution on the union of two green rectangles, contaminated by \mathbf{x} , is marked by a red dot. This point does not lie exactly in the middle between $\mathbf{0}$ and \mathbf{x} , which contradicts [Chen and Tyler, 2002, Theorem 3.1].

Theorem 24 *Under the conditions of Theorem 23 it holds that*

$$IF(\mathbf{x}; T, P) = \begin{cases} \mathbf{0} & \text{for } \mathbf{x} = \mathbf{x}_0, \\ \lim_{\alpha \rightarrow 0^+} \left(\frac{1}{2\alpha} r \left(\frac{1}{2} - \alpha, \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0} \right) \right) \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0} & \text{for } \mathbf{x} \neq \mathbf{x}_0, \end{cases} \quad (1.6)$$

provided that the limit exists.

This theorem states that, under the specified assumptions, the influence function remains constant along rays originating from the center of symmetry. This implies boundedness of the influence function, indicating good robustness properties of Tukey's median.

Proof. For $\mathbf{x} = \mathbf{x}_0$, the claim is obvious. Consider $\mathbf{x} \neq \mathbf{x}_0$. Then, for ε small enough, we have $\mathbf{x} \in \text{int} \left(D \left(\frac{1-3\varepsilon}{2(1-\varepsilon)}; P \right) \right)$. Therefore, by Theorem 23, we have

$$\begin{aligned} IF(\mathbf{x}; T, P) &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} r \left(\frac{1-3\varepsilon}{2(1-\varepsilon)}, \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0} \right) \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{1}{2\varepsilon} r \left(\frac{1}{2} - \frac{\varepsilon}{1-\varepsilon}, \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0} \right) \frac{\mathbf{x} - \mathbf{x}_0}{\mathbf{x} - \mathbf{x}_0}, \end{aligned}$$

which is evidently equivalent to (1.6).

The formula for the gross error sensitivity follows immediately.

Theorem 25 *Let the conditions of Theorem 23 be satisfied and further assume that the limit*

$$\lim_{\alpha \rightarrow 0^+} \left(\frac{1}{\alpha} r \left(\frac{1}{2} - \alpha, \mathbf{u} \right) \right)$$

exists for all $\mathbf{u} \in S^{d-1}$. Then,

$$\gamma(T, P) = \sup_{\mathbf{u} \in S^{d-1}} \lim_{\alpha \rightarrow 0^+} \left(\frac{1}{2\alpha} r \left(\frac{1}{2} - \alpha, \mathbf{u} \right) \right).$$

Proof. The claim follows from Theorem 24.

Now we will demonstrate the application of the preceding theorems for the Tukey median of a contaminated Gaussian distribution.

Example 5 *Consider $P = N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. The assumptions of Theorem 23 are obviously satisfied for $\mathbf{x}_0 = \boldsymbol{\mu}$. Recall that the distribution function of the univariate standard Gaussian distribution and its density are denoted by Φ and φ , respectively. From Example 3 we know that*

$$D \left(\frac{1}{2} - \alpha; N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \right) = \left\{ \mathbf{y} \in \mathbb{R}^d / (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \leq \left(\Phi^{-1} \left(\frac{1}{2} + \alpha \right) \right)^2 \right\}.$$

For $\mathbf{x} = \boldsymbol{\mu}$, the width of this ellipsoid along the direction $(\mathbf{x} - \boldsymbol{\mu}) / \|\mathbf{x} - \boldsymbol{\mu}\|$ is

$$2r \left(\frac{1}{2} - \alpha, \frac{\mathbf{x} - \boldsymbol{\mu}}{\|\mathbf{x} - \boldsymbol{\mu}\|} \right) = 2 \Phi^{-1} \left(\frac{1}{2} + \alpha \right) \frac{\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\|}{\|\mathbf{x} - \boldsymbol{\mu}\|}.$$

Therefore, by Theorem 23 we have

$$T(P_{(\varepsilon, \mathbf{x})}) = \boldsymbol{\mu} + \frac{1}{2} \Phi^{-1} \left(\frac{1}{2} + \alpha \right) \frac{\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\|}{\|\mathbf{x} - \boldsymbol{\mu}\|^2} (\mathbf{x} - \boldsymbol{\mu})$$

and by Theorem 24

$$\begin{aligned} IF(\mathbf{x}; T, N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})) &= \lim_{\alpha \rightarrow 0^+} \left(\frac{1}{2\alpha} \Phi^{-1} \left(\frac{1}{2} + \alpha \right) \right) \frac{\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\|}{\|\mathbf{x} - \boldsymbol{\mu}\|^2} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \frac{\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\|}{2\varphi(\Phi^{-1}(1/2)) \|\mathbf{x} - \boldsymbol{\mu}\|^2} (\mathbf{x} - \boldsymbol{\mu}) = \frac{\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\|}{2\sqrt{2\pi} \|\mathbf{x} - \boldsymbol{\mu}\|^2} (\mathbf{x} - \boldsymbol{\mu}), \end{aligned}$$

where the limit was evaluated using the L'Hôpital's rule. For $\mathbf{x} = \boldsymbol{\mu}$, we have $IF(\boldsymbol{\mu}; T, N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = \mathbf{0}$. Furthermore, because the spectral norm of $\boldsymbol{\Sigma}^{1/2}$ equals $\sqrt{\sigma_{\max}(\boldsymbol{\Sigma})}$, we have $\|\boldsymbol{\Sigma}^{1/2}(\mathbf{x} - \boldsymbol{\mu})\| / \|\mathbf{x} - \boldsymbol{\mu}\| \leq \sqrt{\sigma_{\max}(\boldsymbol{\Sigma})}$. For \mathbf{x} such that $\mathbf{x} - \boldsymbol{\mu}$ is the corresponding eigenvector of $\boldsymbol{\Sigma}$ (or $\boldsymbol{\Sigma}^{1/2}$), we obtain equality. Therefore

$$\gamma(T, N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})) = \sqrt{\frac{\sigma_{\max}(\boldsymbol{\Sigma})}{8\pi}}.$$

This means that a larger maximal eigenvalue of the covariance matrix $\boldsymbol{\Sigma}$ suggests greater variability in the data, potentially making the estimator more sensitive to outliers or extreme values.

Now, we would like to determine the effect of ε -contamination by an arbitrary probability distribution Q . In this case, it is not possible to find any general expression for the halfspace median of the contaminated distribution $P_{(\varepsilon, Q)}$. However, as shown in the following theorem, the maximum possible bias for halfspace symmetric distributions is produced by point-mass contamination. For a general distribution, we are at least able to bound the maximum bias. This allows us to determine the breakdown point of the halfspace median. The subsequent theorem and its proof is taken from [Chen and Tyler, 2002, Theorem 4.1].

Theorem 26 For any $P \in \mathcal{P}(\mathbb{R}^d)$, $d > 1$ and $\varepsilon \in \left(0, \frac{(P)}{1+(P)}\right)$ we have

$$B(\varepsilon; T, P) = \sup_{\mathbf{u} \in S^{d-1}} r \left(\Pi(P) - \frac{\varepsilon}{(1-\varepsilon)}, \mathbf{u} \right) =: \left\| D \left(\Pi(P) - \frac{\varepsilon}{(1-\varepsilon)}; P \right) \right\|_r, \quad (1.7)$$

hence $\frac{(P)}{1+(P)} \geq \varepsilon(T, P)$. If P is smooth, then we also have $\varepsilon(T, P) \leq 1/3$. Finally, if P satisfies conditions of Theorem 23 and is smooth, then for all $\varepsilon \in (0, 1/3)$ it holds that

$$B(\varepsilon; T, P) = \left\| D \left(\frac{1}{2} - \frac{\varepsilon}{(1-\varepsilon)}; P \right) \right\|_r \quad (1.8)$$

and $\varepsilon(T, P) = 1/3$.

Proof.

- Fix $Q \in \mathcal{P}(\mathbb{R}^d)$. For all halfspaces H with $T(P)$ on its boundary, we have $P(H) \geq \Pi(P)$. Therefore, $hD(T(P); P_{(\varepsilon, Q)}) \geq (1-\varepsilon)\Pi(P)$. This implies $\Pi(P_{(\varepsilon, Q)}) \geq (1-\varepsilon)\Pi(P)$. Consider $\mathbf{y} \in D \left(\Pi(P) - \frac{\varepsilon}{(1-\varepsilon)}; P \right)$, then

$$\begin{aligned} hD(\mathbf{y}; P_{(\varepsilon, Q)}) &\geq (1-\varepsilon)hD(\mathbf{y}; P) + \varepsilon \\ &< (1-\varepsilon) \left(\Pi(P) - \frac{\varepsilon}{1-\varepsilon} \right) + \varepsilon \\ &= (1-\varepsilon)\Pi(P). \end{aligned}$$

This means that $\text{Med} \left(P_{(\varepsilon, Q)} \right) \in D \left(\Pi(P) - \frac{\varepsilon}{(1-\varepsilon)}; P \right)$. Since this depth region is bounded by Theorem 10, this proves the upper bound (1.7) for an arbitrary distribution P . Consequently we have $\frac{(P)}{1+(P)} \geq \varepsilon(T, P)$.

- For P which is smooth at all points, we can consider $Q = \delta_{\mathbf{x}}$ for $\mathbf{x} \in \mathbb{R}^d$ and $\varepsilon > 1/3$. Then, we have $hD(\mathbf{x}; P_{(\varepsilon, \mathbf{x})}) \geq \varepsilon > 1/3$. On the other hand, consider arbitrary $\mathbf{y} \in \mathbb{R}^d$, $\mathbf{y} \neq \mathbf{x}$, and find $\mathbf{u} \in S^{d-1}$ such that $\langle \mathbf{u}, \mathbf{x} - \mathbf{y} \rangle = 0$. Then, using the smoothness of P , we have that at least one of $H_{\mathbf{y}, \mathbf{u}}$, $H_{\mathbf{y}, -\mathbf{u}}$ has P -mass less than or equal to $1/2$. Without loss of generality, let $P(H_{\mathbf{y}, \mathbf{u}}) \leq 1/2$. Now, consider a sequence $\{\mathbf{u}_i\}_{i=1}^{\infty} \subset S^{d-1}$ such that $\mathbf{x} \in H_{\mathbf{y}, \mathbf{u}_i}$ and $\mathbf{u}_i \rightarrow \mathbf{u}$. Using the smoothness of P , we have $P_{(\varepsilon, \mathbf{x})}(H_{\mathbf{y}, \mathbf{u}_i}) \rightarrow (1-\varepsilon)P(H_{\mathbf{y}, \mathbf{u}}) = (1-\varepsilon)/2 < 1/3$. This implies that $hD(\mathbf{y}; P_{(\varepsilon, \mathbf{x})}) < 1/3 < hD(\mathbf{x}; P_{(\varepsilon, \mathbf{x})})$. Therefore, \mathbf{x} is the halfspace median of $P_{(\varepsilon, \mathbf{x})}$ and by letting $\mathbf{x} \rightarrow \mathbf{y}$, we can produce an arbitrarily large bias. Therefore, $\varepsilon(T, P) = 1/3$.

- Recall that under conditions of Theorem 23 we have $\Pi(P) = 1/2$, hence $\frac{\varepsilon(P)}{1 + \varepsilon(P)} = 1/3$. Therefore, for any $\varepsilon \in (0, 1/3)$ and $\mathbf{x} \in \text{int}\left(D\left(\frac{1}{2} - \frac{\varepsilon}{(1-\varepsilon)}; P\right)\right)$ we can use Theorem 23 to show that $T(P_{(\varepsilon, \mathbf{x})}) = \mathbf{x}$. As a consequence, considering $Q = \delta_{\mathbf{x}}$ proves equality (1.8). The equality $\varepsilon(T, P) = 1/3$ now follows from the first and the second part.

The previous theorem implies that the greater the maximal depth of a point w.r.t. a given distribution, the greater the breakdown point of Tukey's median. By Theorem 17 we know that for any $P \in \mathcal{P}(\mathbb{R}^d)$, the lower bound $\Pi(P) \geq 1/(d+1)$ holds. Together with the previous theorem, this implies that $\varepsilon(T, P) \geq 1/(d+2)$. However, the argument of Theorem 26 can be refined so that we obtain $\varepsilon(T, P) \geq 1/(d+1)$, see [Chen and Tyler, 2002, Theorem 4.2]. Note that this implies that for any smooth $P \in \mathcal{P}(\mathbb{R}^2)$ we have $\varepsilon(T, P) = 1/3$.

All robustness measures introduced above are based on the population distribution P . However, we can also investigate the robustness of the halfspace median based on finite samples. For example, a finite-sample version of the influence function is the sensitivity curve, see [Huber and Ronchetti, 2009, Section 1.5]. For a dataset $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, Donoho and Gasko [1992] define the finite-sample breakdown point as

$$\varepsilon(T, X) = \min_{m \in \mathbb{N}} \left\{ \frac{m}{n+m} / \sup \left\{ |T(X \cup Y^m) - T(X)| \mid Y^m = \{\mathbf{y}_1, \dots, \mathbf{y}_m\} \subset \mathbb{R}^d \right\} = \right\},$$

where $T(X)$ is the halfspace median calculated from X and $T(X \cup Y^m)$ is the halfspace median calculated from $X \cup Y^m$. That is, we contaminate out dataset X by m points Y^m . If, by some choice of Y^m , this shifts the halfspace median arbitrarily far from the original halfspace median, we say that T breaks down. Then, $\varepsilon(T, X)$ is the smallest contamination fraction $m/(n+m)$ such that T can break down. Donoho and Gasko [1992, Proposition 3.3] show that if X corresponds to a random sample from any absolutely continuous, centrally symmetric distribution in \mathbb{R}^d , $d > 2$, then $\varepsilon(T, X)$ converges to $1/3$ almost surely as $n \rightarrow \infty$. Furthermore, Donoho and Gasko [1992, Proposition 3.4] prove that $\varepsilon(T, X) \geq 1/(d+1)$ for any X in a general position (this means that no more than k points lie in any $(k-1)$ -dimensional affine subspace, $k = 1, \dots, d$). This corresponds to Theorem 26 above.

2. Scatter halfspace depth

In the first chapter of this thesis, we discussed the location halfspace depth. We showed that it allows us to order points with respect to some fixed distribution, and the deepest point can serve as a location estimator. This estimator possesses outstanding properties, especially those regarding robustness. However, in mathematical statistics, one is often interested in higher-order characteristics of probability distributions. A natural question is whether it is possible to use concepts similar to the ones discussed in the first chapter also in these situations. Therefore, in the second part of this thesis, we deal with scatter halfspace depth for matrices. We begin by defining the scatter parameter of a distribution.

Definition 27 (Scatter parameter) *Let $P \in \mathcal{P}(\mathbb{R}^d)$. Consider a functional $\Psi: \mathcal{P} \rightarrow \text{PD}_d$ satisfying $\Psi(P_{\mathbf{A}\mathbf{X}+\mathbf{b}}) = \mathbf{A}\Psi(P_{\mathbf{X}})\mathbf{A}^T$ for any $\mathbf{X} \in \mathcal{P}_{\mathbf{X}} \in \mathcal{P}$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{b} \in \mathbb{R}^d$ such that $P_{\mathbf{A}\mathbf{X}+\mathbf{b}} \in \mathcal{P}$. Then Ψ is called a scatter functional and $\Psi(P)$ is called a scatter parameter of $P \in \mathcal{P}$. The uniquely determined matrix $\Psi(P)^{1/2} \in \text{PD}_d$ is called a scale parameter of P .*

The functional assigning the covariance matrix to the appropriately integrable random vector is an example of a scatter functional. Another example (for $d = 1$) is the median squared deviation (2.1) defined below.

2.1 Axiomatic approach to scatter statistical depth

At the beginning of this work, we defined general location statistical depth. One of the desired properties stress that location depth should behave nicely with (centrally/angularly/halfspace) symmetric distributions. Analogously, scatter statistical depth should reflect the “shape” of the underlying distribution. Therefore, it should work well with so-called elliptically symmetric distributions, as defined below. Recall that by MSD we mean the *median squared deviation*, i.e.

$$\text{MSD}(Z) = \text{median}\left((Z - \text{median}(Z))^2\right) \quad (2.1)$$

for any $Z \in \mathcal{P} \in \mathcal{P}(\mathbb{R})$. In this definition, if the “median interval” is not a singleton, we always take its midpoint as the median (similarly as for the Tukey median).

Definition 28 (Elliptically symmetric distribution) *Consider a random vector $\mathbf{X} \in \mathcal{P} \in \mathcal{P}(\mathbb{R}^d)$. Then, we say that \mathbf{X} (or its distribution P) is spherically symmetric about the origin if and only if $\mathbf{O}\mathbf{X} \stackrel{D}{=} \mathbf{X}$ for all $d \times d$ orthogonal matrices \mathbf{O} . More generally, the random vector \mathbf{X} is said to have an elliptically symmetric distribution with location $\boldsymbol{\mu} \in \mathbb{R}^d$ and scatter $\boldsymbol{\Sigma} \in \text{PD}_d$ if and only if $\mathbf{X} \stackrel{D}{=} \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\mathbf{Z}$, where*

- $\mathbf{Z} \in \mathcal{P}(\mathbb{R}^d)$ is spherically symmetric about the origin,
- $\text{MSD}(Z_1) = 1$, where Z_1 is the first component of \mathbf{Z} .

Note that the definition of a spherically symmetric distribution implies that $\text{median}(Z_k^2) = 1$ for any $k \in \{1, \dots, d\}$. As a consequence of this condition, it can be shown that the location $\boldsymbol{\mu}$ and the scatter matrix $\boldsymbol{\Sigma}$ of an elliptically symmetric distribution are uniquely determined. It is also easy to see that every elliptically symmetric distribution with location $\boldsymbol{\mu}$ is centrally symmetric about its center $\boldsymbol{\mu}$.

An example of a spherically symmetric (about the origin) distribution is the d -variate Gaussian distribution $N_d(\mathbf{0}, \alpha \mathbf{I}_d)$ for any $\alpha > 0$. Indeed, if a random variable Z has univariate standard Gaussian distribution, then we see that

$$\begin{aligned} \text{MSD}(Z) &= \text{median}(Z^2) = \arg \max_{z>0} \left(\min \left\{ \mathbb{P}(Z^2 \leq z), \mathbb{P}(Z^2 \geq z) \right\} \right) \\ &= \arg \max_{z>0} \left(\min \left\{ \Phi(\sqrt{z}) - \Phi(-\sqrt{z}), 1 - \Phi(\sqrt{z}) + \Phi(-\sqrt{z}) \right\} \right) \\ &= \arg \max_{z>0} \left(\min \left\{ 2\Phi(\sqrt{z}) - 1, 2 - 2\Phi(\sqrt{z}) \right\} \right) \\ &= \arg \max_{z>0} \left(\min \left\{ \Phi(\sqrt{z}) - \frac{1}{2}, 1 - \Phi(\sqrt{z}) \right\} \right) = \left(\Phi^{-1}\left(\frac{3}{4}\right) \right)^2 =: \beta. \end{aligned}$$

Therefore, $\text{MSD}(Z/\sqrt{\beta}) = 1$. As a consequence, for all $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} \in \text{PD}_d$, the Gaussian distribution $N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}/\beta)$ is elliptically symmetric with location $\boldsymbol{\mu}$ and scatter $\boldsymbol{\Sigma}$. Note that β is defined such that it satisfies

$$\Phi(\sqrt{\beta}) = \frac{3}{4} \quad (2.2)$$

and this constant will be used in the following text.

The subsequent text deals with the space PD_d of all symmetric positive definite $d \times d$ matrices. Therefore, it is necessary to introduce a topology in this space. We will consider the Frobenius topology induced by the Frobenius distance, which is itself inherited from the Frobenius norm, i.e.

$$d_F(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) := \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_F = \sqrt{\text{tr}((\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2)^\top)}.$$

A real function $f: \text{PD}_d \rightarrow [0, \infty)$ will be called F -(semi)continuous if it is (semi)continuous w.r.t. the Frobenius topology on PD_d . Similarly, a set $A \subset \text{PD}_d$ is called F -(bounded/closed) if it is bounded/closed w.r.t. the Frobenius topology on PD_d . We also use the operator matrix norm in Section 3.5 about the minimax optimality of the scatter halfspace median. This norm (also called the spectral norm) is defined for any $d \times d$ symmetric matrix \mathbf{A} as

$$\mathbf{A}_{\text{op}} = \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \|\mathbf{A}\mathbf{u}\|,$$

which is equal to the square root of the largest singular value of \mathbf{A} . Because \mathbf{A} is symmetric, \mathbf{A}_{op} is equal to

$$\sqrt{\max\{|\lambda_1|, \dots, |\lambda_d|\}}$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of \mathbf{A} .

As Zuo and Serfling [2000a] define general location depth, Paidaveine and Van Bever [2018] similarly suggest the following scatter statistical depth definition.

Definition 29 (Scatter statistical depth) *A mapping $D^{\text{sc}}: \text{PD}_d \times P(\mathbb{R}^d) \rightarrow [0, 1]$ is called a scatter statistical depth (or simply scatter depth) if it is bounded and satisfies the following properties:*

- (B1) *For any random vector $\mathbf{X} \sim P_{\mathbf{X}} \in P(\mathbb{R}^d)$, $\Sigma \in \text{PD}_d$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{b} \in \mathbb{R}^d$ it holds that $D^{\text{sc}}(\mathbf{A}\Sigma\mathbf{A}^{\top}; P_{\mathbf{A}\mathbf{X}+\mathbf{b}}) = D^{\text{sc}}(\Sigma; P_{\mathbf{X}})$.*
- (B2) *$D^{\text{sc}}(\Sigma_0; P) = \max_{\Sigma \in \text{PD}_d} D^{\text{sc}}(\Sigma; P)$ holds true for any elliptically symmetric distribution $P \in P(\mathbb{R}^d)$ with scatter Σ_0 .*
- (B3) *If for any fixed $P \in P(\mathbb{R}^d)$ a matrix $\Sigma_0 \in \text{PD}_d$ maximizes $D^{\text{sc}}(\cdot; P)$, then for an arbitrary $\Sigma \in \text{PD}_d$ the mapping $\alpha \mapsto D^{\text{sc}}(\Sigma_0 + \alpha(\Sigma - \Sigma_0); P)$ is non-increasing on $[0, 1]$,*
- (B4) *Consider any fixed $P \in P(\mathbb{R}^d)$ and $\{\Sigma_n\}_{n=1}^{\infty} \subset \text{PD}_d$ such that either $d_F(\Sigma_n, \Sigma) \xrightarrow[n]{\text{---}} 0$ for some singular $d \times d$ matrix Σ , or $d_F(\Sigma_n, \mathbf{I}_d) \xrightarrow[n]{\text{---}} 0$. Then $D^{\text{sc}}(\Sigma_n; P) \xrightarrow[n]{\text{---}} 0$.*

Properties (B1)–(B4) are somewhat analogous to those for location statistical depth. Namely, (B1) states that depth should transform accordingly after the affine transformation of the underlying distribution. Property (B2) expresses that the scatter estimator based on maximal depth should be Fisher-consistent for elliptically symmetric distributions. Properties (B3) and (B4) claim that scatter depth should possess some monotonicity properties, and the depth of matrices near the “boundary” of the space PD_d should be close to zero.

Again, we can define the sample version of the scatter depth D^{sc} based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P \in P(\mathbb{R}^d)$ as $D^{\text{sc}}(\Sigma; \hat{P}_n)$ and a matrix maximizing this function can be used as an estimator of the scatter parameter of P .

2.2 Definition of scatter halfspace depth

The scatter halfspace depth was first introduced in [Chen et al., 2018]. There, it was studied only for centred Gaussian distributions with fixed location $\mathbf{0}$. The following general definition is provided in [Paindaveine and Van Bever, 2018].

Definition 30 (Scatter halfspace depth) *Let $\mathbf{X} \sim P \in P(\mathbb{R}^d)$ and $\Sigma \in \text{PD}_d$. Further, consider some location functional Ξ . The scatter halfspace depth of Σ w.r.t. P is defined as*

$$hD_{\Xi}^{\text{sc}}(\Sigma; P) := \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \min \left\{ \mathbb{P} \left(\left| \mathbf{X} - \Xi(P), \mathbf{u} \right| \leq \sqrt{\mathbf{u}^{\top} \Sigma \mathbf{u}} \right), \right. \\ \left. \mathbb{P} \left(\left| \mathbf{X} - \Xi(P), \mathbf{u} \right| \geq \sqrt{\mathbf{u}^{\top} \Sigma \mathbf{u}} \right) \right\}.$$

To interpret the scatter halfspace depth geometrically, we introduce the subsequent definition.

Definition 31 For $\mathbf{u} \in S^{d-1}$, $t \in [0, 1]$ and $\Sigma \in \text{PD}_d$, we define

$$\begin{aligned} H_{\mathbf{u},t}^{\text{in}} &= \left\{ \mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{u} \rangle \leq t \right\}, \\ H_{\mathbf{u},t}^{\text{out}} &= \left\{ \mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{u} \rangle \geq t \right\}, \\ H_{\mathbf{u},\Sigma}^{\text{in}} &= \left\{ \mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{u} \rangle \leq \overline{\mathbf{u}^\top \Sigma \mathbf{u}} \right\}, \\ H_{\mathbf{u},\Sigma}^{\text{out}} &= \left\{ \mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{u} \rangle \geq \overline{\mathbf{u}^\top \Sigma \mathbf{u}} \right\}. \end{aligned} \quad (2.3)$$

Using this notion, we can rewrite

$$hD_{\Xi}^{\text{sc}}(\Sigma; P) = \inf_{\mathbf{u} \in S^{d-1}} \min \left\{ P \left(H_{\mathbf{u},\Sigma}^{\text{in}} + \Xi(P) \right), P \left(H_{\mathbf{u},\Sigma}^{\text{out}} + \Xi(P) \right) \right\}, \quad (2.4)$$

where $A + \mathbf{z} = \{ \mathbf{x} + \mathbf{z} \mid \mathbf{x} \in A \}$ for any $A \subseteq \mathbb{R}^d$ and $\mathbf{z} \in \mathbb{R}^d$. Denote by E_{Σ} the surface

$$\Sigma^{1/2} S^{d-1} = \left\{ \Sigma^{1/2} \mathbf{u} \in \mathbb{R}^d \mid \mathbf{u} \in S^{d-1} \right\}$$

of an ellipsoid $\Sigma^{1/2} B_d(\mathbf{0}, 1)$. Then, it can be shown (see [Brabenec, 2021, Section 1.3] for reference) that the set $H_{\mathbf{u},\Sigma}^{\text{in}}$ consists exactly of all points in between two supporting hyperplanes of $\Sigma^{1/2} B_d(\mathbf{0}, 1)$ with normal vector \mathbf{u} . Similarly, $H_{\mathbf{u},\Sigma}^{\text{out}}$ consists of all points which are not in between these two hyperplanes. Our situation is shown in Figure 4. From (2.4) we can see that if the matrix Σ has a small scatter halfspace depth w.r.t. P , then in some direction $\mathbf{u} \in S^{d-1}$, at least one of the sets $H_{\mathbf{u},\Sigma}^{\text{in}} + \Xi(P)$, $H_{\mathbf{u},\Sigma}^{\text{out}} + \Xi(P)$ has small P -probability. Therefore, when looking for a matrix with a high scatter halfspace depth w.r.t. P , one has to find Σ such that in every direction $\mathbf{u} \in S^{d-1}$, both $H_{\mathbf{u},\Sigma}^{\text{in}} + \Xi(P)$ and $H_{\mathbf{u},\Sigma}^{\text{out}} + \Xi(P)$ have high probability P . This is the reason why the scatter halfspace depth captures the geometric characteristics of the elliptically symmetric distribution well. For a smooth P , $hD_{\Xi}^{\text{sc}}(\cdot; P)$ is obviously bounded from above by $1/2$.

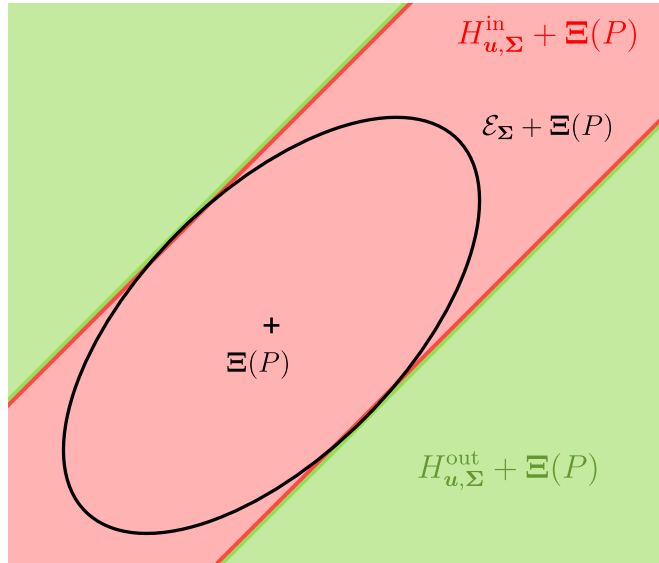


Figure 4: Interpretation of the definition of the scatter halfspace depth of Σ w.r.t. the distribution P : When determining the depth of Σ , we are looking for a direction $\mathbf{u} \in S^{d-1}$ in which at least one of the sets $H_{\mathbf{u},\Sigma}^{\text{in}} + \Xi(P)$ and $H_{\mathbf{u},\Sigma}^{\text{out}} + \Xi(P)$ has the smallest probability P .

For a better understanding, consider the univariate case $d = 1$. Then, for a fixed probability distribution $P \in \mathcal{P}(\mathbb{R})$, $X \sim P$ and $\sigma^2 > 0$, we have

$$hD^{\text{sc}}(\sigma^2; P) = \min \left\{ \mathbb{P} \left((X - \Xi(P))^2 \leq \sigma^2 \right), \mathbb{P} \left((X - \Xi(P))^2 \geq \sigma^2 \right) \right\}. \quad (2.5)$$

Note that if Ξ corresponds to the univariate median, then the point σ^2 , which maximizes (2.5), is precisely the median squared deviation (2.1). Further, we derive what the scatter halfspace depth of a matrix looks like with respect to a Gaussian distribution.

Example 6 Let $\mathbf{X} \sim P_{\mathbf{X}} \sim N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for an arbitrary $\boldsymbol{\mu} \in \mathbb{R}^d$ and $\boldsymbol{\Sigma} \in \text{PD}_d$. First, note that if Ξ is a location functional, then it must hold that $\Xi(P_{\mathbf{X}}) = \boldsymbol{\mu}$. Indeed, by the definition of the location functional, we have for all $\boldsymbol{\mu} \in \mathbb{R}^d$ that

$$\Xi(P_{\mathbf{X}}) = -\Xi(P_{-\mathbf{X}}) = -\Xi(P_{\mathbf{X}-2\boldsymbol{\mu}}) = -\Xi(P_{\mathbf{X}}) + 2\boldsymbol{\mu},$$

which implies $\Xi(P_{\mathbf{X}}) = \boldsymbol{\mu}$. Therefore, $\mathbf{X} - \Xi(P_{\mathbf{X}}) \sim N_d(\mathbf{0}, \boldsymbol{\Sigma})$ and

$$\mathbf{X} - \Xi(P_{\mathbf{X}}), \mathbf{u} \sim N_d(\mathbf{0}, \mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}).$$

Consider some $\boldsymbol{\Gamma} \in \text{PD}_d$. Then, we can write

$$\begin{aligned} \mathbb{P} \left(\left| \mathbf{X} - \Xi(P_{\mathbf{X}}), \mathbf{u} \right| \leq \sqrt{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}} \right) &= \mathbb{P} \left(\left| \frac{\mathbf{X} - \Xi(P_{\mathbf{X}}), \mathbf{u}}{\sqrt{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right| \leq \sqrt{\frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) \\ &= 2\Phi \left(\sqrt{\frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) - 1. \end{aligned}$$

Similarly,

$$\mathbb{P} \left(\left| \mathbf{X} - \Xi(P_{\mathbf{X}}), \mathbf{u} \right| \geq \sqrt{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}} \right) = 2 \left(1 - \Phi \left(\sqrt{\frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) \right).$$

As a consequence, we have

$$\begin{aligned} hD_{\Xi}^{\text{sc}}(\boldsymbol{\Gamma}; P_{\mathbf{X}}) &= 2 \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \min \left\{ \Phi \left(\sqrt{\frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) - \frac{1}{2}, 1 - \Phi \left(\sqrt{\frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) \right\} \\ &= 2 \min \left\{ \Phi \left(\sqrt{\inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) - \frac{1}{2}, 1 - \Phi \left(\sqrt{\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}}} \right) \right\}. \end{aligned}$$

Now, note that

$$\inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}} = \inf_{\mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}} \stackrel{\mathbf{v} = \boldsymbol{\Sigma}^{1/2} \mathbf{u}}{=} \inf_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^{\top} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Gamma} \boldsymbol{\Sigma}^{-1/2} \mathbf{v}}{\mathbf{v}^{\top} \mathbf{v}} = \inf_{\mathbf{v} \in \mathbb{S}^{d-1}} \mathbf{v}^{\top} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Gamma} \boldsymbol{\Sigma}^{-1/2} \mathbf{v}.$$

By the Cauchy-Schwarz inequality, this is equal to the smallest eigenvalue of $\boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Gamma} \boldsymbol{\Sigma}^{-1/2}$ denoted by $\sigma_{\min}(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Gamma} \boldsymbol{\Sigma}^{-1/2})$. Analogously

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{\mathbf{u}^{\top} \boldsymbol{\Gamma} \mathbf{u}}{\mathbf{u}^{\top} \boldsymbol{\Sigma} \mathbf{u}} = \sigma_{\max}(\boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Gamma} \boldsymbol{\Sigma}^{-1/2}).$$

Therefore, we can write

$$hD_{\Xi}^{\text{sc}}(\Gamma; P_{\mathbf{X}}) = 2 \min \left\{ \Phi \left(\sqrt{\sigma_{\min}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}})} \right) - \frac{1}{2}, 1 - \Phi \left(\sqrt{\sigma_{\max}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}})} \right) \right\}.$$

This is obviously maximized for Γ satisfying

$$\Phi \left(\sqrt{\sigma_{\min}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}})} \right) = \Phi \left(\sqrt{\sigma_{\max}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}})} \right) = \frac{3}{4}$$

which implies

$$\sigma_{\min}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}}) = \sigma_{\max}(\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}}) = \beta, \quad (2.6)$$

where β is defined in (2.2). Note that $\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}} \in \text{PD}_d$, therefore it is diagonalizable. From (2.6) we have

$$\Sigma^{-\frac{1}{2}} \Gamma \Sigma^{-\frac{1}{2}} = \mathbf{U}^{\top} \mathbf{\Lambda} \mathbf{U} = \beta \mathbf{I}_d,$$

where \mathbf{U} is an orthogonal matrix and $\mathbf{\Lambda}$ is a diagonal matrix with all diagonal entries equal to β . This implies that

$$\Gamma = \beta \Sigma,$$

which is the scatter parameter of $P_{\mathbf{X}}$. To sum up, the scatter halfspace depth w.r.t. the Gaussian distribution with a covariance matrix Σ is uniquely maximized for its scatter parameter $\beta \Sigma$. This will hold for any elliptically symmetric distribution, as stated in Theorem 36.

2.3 Basic properties of scatter halfspace depth

In this section, we introduce basic properties of the scatter halfspace depth. Similarly as in Chapter 1.3, we do not provide proofs, only references. This is because the purpose of this work is not to prove basic properties; rather, the work focuses mainly on minimax optimality of the scatter halfspace depth.

One of the most important properties of the scatter halfspace depth is that it is affine invariant, as stated in the following theorem.

Theorem 32 *The scatter halfspace depth is a ne invariant for any location functional Ξ . That is, for any $\mathbf{X} \sim P_{\mathbf{X}} \in \mathcal{P}(\mathbb{R}^d)$, $\Sigma \in \text{PD}_d$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{b} \in \mathbb{R}^d$, it holds that*

$$hD_{\Xi}^{\text{sc}}(\mathbf{A} \Sigma \mathbf{A}^{\top}; P_{\mathbf{A} \mathbf{X} + \mathbf{b}}) = hD_{\Xi}^{\text{sc}}(\Sigma; P_{\mathbf{X}}).$$

Recall that by $P_{\mathbf{A} \mathbf{X} + \mathbf{b}}$ we denote the distribution of $\mathbf{A} \mathbf{X} + \mathbf{b}$.

Proof. Refer to [Paindaveine and Van Bever, 2018, Theorem 2.1].

This property is useful because it implies that the scatter halfspace depth does not depend on our choice of the coordinate system. Next, similarly as for the location halfspace depth, we define a scatter α -depth region.

Definition 33 (Scatter α -depth region) *For any $P \subset P(\mathbb{R}^d)$ and $\alpha \in [0, 1]$, the scatter α -depth region of P is defined as*

$$D^{\text{sc}}(\alpha; P) := \{\Sigma \in \text{PD}_d \mid hD_{\Xi}^{\text{sc}}(\Sigma; P) \geq \alpha\} \subset \text{PD}_d.$$

An easy consequence of Theorem 32 is that if $\mathbf{X} \in P_{\mathbf{X}} \subset P(\mathbb{R}^d)$, $\mathbf{A} \in \text{GL}_d$ and $\mathbf{b} \in \mathbb{R}^d$, then

$$D^{\text{sc}}(\alpha; P_{\mathbf{A}\mathbf{X}+\mathbf{b}}) = \mathbf{A}D^{\text{sc}}(\alpha; P_{\mathbf{X}})\mathbf{A}^{\top}$$

holds for any $\alpha \in [0, 1]$. The subsequent theorem shows that the scatter halfspace depth is also quasiconcave. This is an analogous statement to Theorem 7, which establishes quasiconcavity of the location halfspace depth.

Theorem 34 *For any $P \subset P(\mathbb{R}^d)$, the function $hD^{\text{sc}}(\cdot; P)$ is quasiconcave. That is, for any $\Sigma_1, \Sigma_2 \in \text{PD}_d$ and $\alpha \in [0, 1]$ we have*

$$hD_{\Xi}^{\text{sc}}(\alpha\Sigma_1 + (1 - \alpha)\Sigma_2; P) \geq \min\{hD_{\Xi}^{\text{sc}}(\Sigma_1; P), hD_{\Xi}^{\text{sc}}(\Sigma_2; P)\}.$$

Proof. Refer to [Paindaveine and Van Bever, 2018, Theorem 3.3].

This theorem implies that $D^{\text{sc}}(\alpha; P)$ is always a convex subset of PD_d , which can be shown just like in (1.2) for location halfspace depth regions $D(\alpha; P)$. Further, we are interested in whether the scatter halfspace depth is upper F -semicontinuous or even F -continuous.

Theorem 35 *For any $P \subset P(\mathbb{R}^d)$, the function $hD^{\text{sc}}(\cdot; P)$ is always upper F -semicontinuous. If P is smooth at $\Xi(P)$, then $hD^{\text{sc}}(\cdot; P)$ is also F -continuous.*

Proof. Refer to [Paindaveine and Van Bever, 2018, Theorem 3.1].

The previous theorem implies that the sets $D^{\text{sc}}(\alpha; P), \alpha \in [0, 1]$, are always F -closed. Paindaveine and Van Bever [2018, Theorem 3.2] show that for $\alpha > 0$, these sets are also F -bounded. The authors of that paper claim that all bounded sets are also totally bounded in (PD_d, d_F) . Therefore, we have that for $\alpha > 0$, the set $D^{\text{sc}}(\alpha; P)$ is F -totally bounded and F -closed. But the space (PD_d, d_F) is not complete, hence we can not conclude that $D^{\text{sc}}(\alpha; P)$ is F -compact, in contrast with the case of the location halfspace depth. Therefore, Paindaveine and Van Bever [2018] also analyse properties of the scatter halfspace depth w.r.t. the so-called geodesic topology, see [Paindaveine and Van Bever, 2018, Section 4]. This is used, for example, to prove Theorem 40 stated below.

Similarly as for the location halfspace depth, arguably the most important property is that under certain conditions, the scatter halfspace depth is a scatter statistical depth in the sense of Definition 29.

Theorem 36 *For the scatter halfspace depth, properties (B1), (B2) and (B3) of Definition 29 are always met. Property (B4) is satisfied for P smooth at $\Xi(P)$.*

Proof. Property (B1) is implied by Theorem 32. For the proof of property (B2), refer to [Paindaveine and Van Bever, 2018, Theorem 5.1]. Property (B3) is implied by Theorem 34, see [Paindaveine and Van Bever, 2018, Section 5] for details. Finally, the property (B4) is proven in [Brabenec, 2021, Theorem 16 and Theorem 19].

2.4 Sample scatter halfspace depth and scatter halfspace median

Analogously as for the location halfspace depth, we estimate the scatter halfspace depth using the sample scatter halfspace depth, which is defined as follows.

Definition 37 (Sample scatter halfspace depth) *Consider a random sample $\{\mathbf{X}_i\}_{i=1}^n \stackrel{P}{\sim} P(\mathbb{R}^d)$, $\Sigma \in \text{PD}_d$ and some location functional Ξ . The sample scatter halfspace depth of Σ w.r.t. $\{\mathbf{X}_i\}_{i=1}^n$ is defined as*

$$hD_{\Xi}^{\text{sc}}(\Sigma; \{\mathbf{X}_i\}_{i=1}^n) := hD_{\Xi}^{\text{sc}}(\Sigma; \hat{P}_n) \\ = \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \min \left\{ \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{|\langle \mathbf{x}_i - \hat{\mu}_n, \mathbf{u} \rangle| \leq \overline{u^\top \Sigma u}\}, \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{|\langle \mathbf{x}_i - \hat{\mu}_n, \mathbf{u} \rangle| \leq \overline{u^\top \Sigma u}\} \right\},$$

where $\hat{\mu}_n = \Xi(\hat{P}_n)$.

It turns out that $hD_{\Xi}^{\text{sc}}(\Sigma; \{\mathbf{X}_i\}_{i=1}^n)$ is (under a certain condition for Ξ) a strongly uniformly consistent estimator of $hD_{\Xi}^{\text{sc}}(\Sigma; P)$. This claim is stated in the subsequent theorem.

Theorem 38 *Suppose we have a random sample $\{\mathbf{X}_i\}_{i=1}^n \stackrel{P}{\sim} P(\mathbb{R}^d)$. If $\hat{\mu}_n \xrightarrow[n]{a.s.} \Xi(P)$, then*

$$\sup_{\Sigma \in \text{PD}_d} |hD_{\Xi}^{\text{sc}}(\Sigma; \{\mathbf{X}_i\}_{i=1}^n) - hD_{\Xi}^{\text{sc}}(\Sigma; P)| \xrightarrow[n]{a.s.} 0. \quad (2.7)$$

Proof. Refer to [Paindaveine and Van Bever, 2018, Theorem 2.2].

A natural choice of Ξ is the Tukey median functional T . Then, by Theorem 18 it holds that $T(\hat{P}_n) \xrightarrow[n]{a.s.} T(P)$. As a consequence, $hD_T^{\text{sc}}(\cdot; \{\mathbf{X}_i\}_{i=1}^n)$ is always a strongly uniformly consistent estimator of $hD_T^{\text{sc}}(\cdot; P)$. In Chapter 3, we consider a model consisting of centered Gaussian distributions with an unknown covariance matrix. In that situation, we can use a constant location functional $\mathbf{0}$, which is obviously consistent. Consequently, in this case, the convergence (2.7) holds true.

Further, we define the maximal scatter halfspace depth and the set of matrices of maximal scatter depth as follows.

Definition 39 (Maximal scatter depth) *For $P \in \mathcal{P}(\mathbb{R}^d)$, we define the maximal scatter halfspace depth as*

$$\Pi^{\text{sc}}(P) := \sup_{\Sigma \in \text{PD}_d} hD_{\Xi}^{\text{sc}}(\Sigma; P)$$

and the set of matrices of maximal scatter halfspace depth as

$$\text{Med}^{\text{sc}}(P) := D^{\text{sc}}(\Pi^{\text{sc}}(P); P).$$

For the location halfspace depth, there is always a point with maximal location halfspace depth. This is not the case with the scatter halfspace depth, as noted in [Brabenec, 2021, Section 2.2]. Yet, under the assumption that P is smooth at its location $\Xi(P)$, the deepest matrix always exists. This is stated in the following theorem.

Theorem 40 *Let $P \subset P(\mathbb{R}^d)$. If Ξ is a location functional such that P is smooth at $\Xi(P)$, then there exists $\Sigma \in \text{PD}_d$ such that*

$$hD_{\Xi}^{\text{sc}}(\Sigma; P) = \Pi^{\text{sc}}(P),$$

therefore $\text{Med}^{\text{sc}}(P) = \Sigma$.

Proof. Refer to [Paindaveine and Van Bever, 2018, Theorem 4.3].

Similarly as for the deepest points of the location halfspace depth, we can have several deepest matrices of the scatter halfspace depth. Paindaveine and Van Bever [2018, Theorem 5.1] show that under some mild assumptions, the scatter halfspace depth w.r.t. an elliptically symmetric distribution is uniquely maximized for its scatter parameter.

Definition 41 (Scatter halfspace median) *Let $P \subset P(\mathbb{R}^d)$. Any matrix $\Sigma \in \text{PD}_d$ that maximizes $hD_{\Xi}^{\text{sc}}(\cdot; P)$ is called a scatter halfspace median¹ of P , denoted by $\Sigma^{\text{sc}}(P)$.*

To conclude this section, let us note that a *sample scatter halfspace median*

$$\widehat{\Sigma}_n^{\text{sc}} := \arg \max_{\Sigma \in \text{PD}_d} hD_{\Xi}^{\text{sc}}(\Sigma; \{\mathbf{X}_i\}_{i=1}^n)$$

always exists, because the sample scatter halfspace depth attains only finitely many values.

2.5 Robustness of scatter halfspace median

Since the concept of the scatter halfspace depth and the scatter halfspace median matrix is relatively new compared to the location halfspace depth, not much is known about its robustness properties so far. In this section, we briefly present one of the results of Louvet and Van Bever [2024]. That article primarily discusses the influence function of the scatter halfspace depth. Only in Section 4 of that article, the influence function of the scatter halfspace median is derived under certain conditions, as presented further. The following theorem can be found in [Louvet and Van Bever, 2024, Theorem 5].

¹If we only need one representative of this set of matrices, we can, for example, choose its centroid, similarly as for the location halfspace depth.

Theorem 42 Let $P \sim P(\mathbb{R}^d)$ correspond to the elliptically symmetric distribution with location $\boldsymbol{\mu} = \mathbf{0}$ and scatter $\boldsymbol{\Sigma}_0 \in \text{PD}_d$. That is, for $\mathbf{X} \sim P$, the canonical representation is $\mathbf{X} = \boldsymbol{\Sigma}_0^{1/2} \mathbf{Z}$, where $\mathbf{Z} \sim P(\mathbb{R}^d)$ is spherically symmetric about the origin. Denote the cumulative distribution function of Z_1 by F and assume that this distribution is absolutely continuous with density $f = F'$. Consider a functional $\Xi \in \mathbb{R}$ and let $f(1) > 0$. Then, $\boldsymbol{\Sigma}_0$ is the unique scatter halfspace median of P , i.e. the unique maximizer of $hD_{\Xi}^{\text{sc}}(\cdot; P)$. Furthermore, let $\mathbf{x} \in \mathbb{R}^d$ and $\varepsilon \in (0, 1)$.

- If $|\mathbf{u}, \mathbf{x}| < \overline{\mathbf{u}^\top \boldsymbol{\Sigma}_0 \mathbf{u}}$ for all $\mathbf{u} \in S^{d-1}$, then

$$\arg \max_{\boldsymbol{\Sigma} \in \text{PD}_d} hD_{\Xi}^{\text{sc}}(\boldsymbol{\Sigma}; P_{(\varepsilon, \mathbf{x})}) = \left(F^{-1} \left(\frac{3 - 4\varepsilon}{4 - 4\varepsilon} \right) \right)^2 \boldsymbol{\Sigma}_0.$$

As a consequence,

$$IF(\mathbf{x}; \boldsymbol{\Sigma}^{\text{sc}}, P) := \lim_{\varepsilon \rightarrow 0^+} \frac{\boldsymbol{\Sigma}^{\text{sc}}(P_{(\varepsilon, \mathbf{x})}) - \boldsymbol{\Sigma}^{\text{sc}}(P)}{\varepsilon} = \frac{-\boldsymbol{\Sigma}_0}{2f(1)}.$$

- If $|\mathbf{u}, \mathbf{x}| \geq \overline{\mathbf{u}^\top \boldsymbol{\Sigma}_0 \mathbf{u}}$ for some $\mathbf{u} \in S^{d-1}$, then

$$\sup_{\boldsymbol{\Sigma} \in \text{PD}_d} hD_{\Xi}^{\text{sc}}(\boldsymbol{\Sigma}; P_{(\varepsilon, \mathbf{x})}) = hD_{\Xi}^{\text{sc}}(\boldsymbol{\Sigma}_0; P_{(\varepsilon, \mathbf{x})}) = \frac{1 - \varepsilon}{2}.$$

Proof. Refer to [Louvet and Van Bever, 2024, Theorem 5].

This theorem assumes for simplicity that the distribution P is centered at the origin. Thus, we do not need to estimate the location of the distribution. Furthermore, the assumption that $f(1) > 0$ ensures that there is exactly one deepest matrix (scatter halfspace median) w.r.t. P , which is essential for us to be able to say anything about the influence function. If there were multiple deepest matrices, the situation would become complicated. From Theorem 42, we can see that if the contaminating point is inside the ellipsoid $\boldsymbol{\Sigma}_0^{1/2} B_d(\mathbf{0}, 1)$, then the only deepest matrix will still be $\boldsymbol{\Sigma}_0$, just rescaled by an appropriate constant depending on the level of contamination ε . However, if the point is outside that ellipsoid, then the deepest matrix is still $\boldsymbol{\Sigma}_0$, but it may not be the only one. This again leads to a complicated situation because to determine the influence function in this case, we would need to be able to select exactly one representative of those deepest matrices.

To summarize, from Theorem 42, we can see that the influence function is bounded (and constant), at least for contamination by points close to the center of symmetry. This, however, indicates good robust properties only partially.

3. Minimax optimality

3.1 Convergence rate

Consider a general parametric statistical model $\mathcal{M} = \{P_\theta \mid \theta \in \Theta\} \subset \mathcal{P}(\mathbb{R}^d)$, where Θ is some parametric space. Our task is to estimate the parameter θ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_\theta$. To do so, we construct an estimator $\hat{\theta}_n = \Xi(\hat{P}_n)$ using some statistical functional Ξ . Besides properties of $\hat{\theta}_n$ such as consistency or unbiasedness, we are often interested in how fast the estimator $\hat{\theta}_n$ converges to θ .

Example 7 Let $\mathcal{M} = \{N_d(\boldsymbol{\mu}, \mathbf{I}_d) \mid \boldsymbol{\mu} \in \mathbb{R}^d\}$ be the statistical model of interest. We consider the sample mean $\hat{\boldsymbol{\mu}}_n = \bar{\mathbf{X}}_n$. Then $\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu})$ is distributed as $N_d(\mathbf{0}, \mathbf{I}_d)$ and consequently $n \|\bar{\mathbf{X}}_n - \boldsymbol{\mu}\|^2 \sim \chi_d^2$. Note that $\mathbb{E} \left(n \|\bar{\mathbf{X}}_n - \boldsymbol{\mu}\|^2 \right) = d$. Utilizing Markov's inequality, we have that for any $a > 0$, it holds that

$$\mathbb{P} \left(n \|\bar{\mathbf{X}}_n - \boldsymbol{\mu}\|^2 \geq a \right) \leq \frac{d}{a} \quad \forall n \in \mathbb{N}.$$

Let $0 < c < 1$ and consider $a = d/c$. Rewriting the inequality above, we get

$$\mathbb{P} \left(n \|\bar{\mathbf{X}}_n - \boldsymbol{\mu}\|^2 \leq \frac{d}{c} \right) \geq c.$$

As a consequence, for all $n, d \in \mathbb{N}$ it holds that

$$\mathbb{P} \left(\|\bar{\mathbf{X}}_n - \boldsymbol{\mu}\|^2 < \frac{1}{c} \frac{d}{n} \right) \geq 1 - c. \quad (3.1)$$

Inequality (3.1) indicates that we need twice as many observations to obtain twice as accurate (in terms of the squared Euclidean norm) estimate with the same probability. On the other hand, the higher the dimension, the higher the number of observations required to obtain equally precise estimates.

The quantity d/n (considered as a function of n) in the example above is the so-called *convergence rate*. This depends on our choice of measuring the accuracy of estimates. Generally, this is done using a *loss function*.

Definition 43 (Loss function) A loss function is any measurable map

$$L: \Theta \times \Theta \rightarrow [0, \infty).$$

We consider a symmetric loss function such that $L(\theta_1, \theta_2) = 0$ if and only if $\theta_1 = \theta_2$, and we suppose that it satisfies the τ -triangle inequality (sometimes also called weak triangle inequality), i.e.

$$L(\theta_1, \theta_2) \leq \tau (L(\theta_1, \theta_3) + L(\theta_3, \theta_2)) \quad (3.2)$$

for some $\tau \geq 1$ and all $\theta_1, \theta_2, \theta_3 \in \Theta$.

A loss function is usually chosen as some metric or its square. If $L(\cdot, \cdot) = d(\cdot, \cdot)^2$ for some metric d on Θ , then (3.2) is satisfied for $\tau = 2$. To see this, note that the following inequality holds

$$d(\theta_1, \theta_2)^2 \leq (d(\theta_1, \theta_3) + d(\theta_3, \theta_2))^2 \leq 2(d(\theta_1, \theta_3)^2 + d(\theta_3, \theta_2)^2). \quad (3.3)$$

In this chapter, we drop the indices in sequences $\{a_n\}_{n=1}^{\infty}$, i.e., we write only $\{a_n\}$. By the term *rate*, we mean a positive real sequence. All rates will be indexed by n , the random sample size. To define the convergence rate, we consider the following notation.

Definition 44 (Order of rates) *Consider two rates $\{a_n\}, \{b_n\}$. We say that $\{a_n\}$ is of lower order than $\{b_n\}$ (writing $\{a_n\} \ll \{b_n\}$) if there exists an absolute constant $C > 0$ such that $a_n \leq C b_n$ for all $n \in \mathbb{N}$. We say that $\{a_n\}$ is of the same order as $\{b_n\}$ if $\{a_n\} \ll \{b_n\}$ and $\{b_n\} \ll \{a_n\}$. This is denoted by $\{a_n\} \asymp \{b_n\}$.*

Now, we are ready to define the convergence rate of an estimator.

Definition 45 (Convergence rate) *Consider an estimator $\hat{\theta}_n$, a loss function L and a rate $\{R_n\}$. We say that $\{R_n\}$ is the convergence rate of the estimator $\hat{\theta}_n$ in terms of the loss L if:*

1. For any $\delta > 0$ there exists an absolute constant $C > 0$ such that

$$\sup_{n \in \mathbb{N}} \sup_{\theta \in \Theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \geq C R_n \right) \leq \delta, \quad (3.4)$$

where the estimator $\hat{\theta}_n$ is based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{\theta}$.

2. If the property (3.4) above is satisfied for some other rate $\{R'_n\}$, then $\{R'_n\} \asymp \{R_n\}$.

Note We use C, D, C', \dots to represent universal constants. Note that these constants may vary in different sections and theorems.

This work considers convergence rates dependent on the sample size n , the dimension d , and the proportion of contamination ε (this will be explained in Section 3.3). Therefore, in Definition 45, by C we mean a constant not depending on n, d and ε . To summarize the definition above, a rate $\{R_n\}$ is said to be the convergence rate of an estimator $\hat{\theta}_n$ in terms of some loss L if $\{R_n\}$ is a rate of the lowest order such that $\{L(\hat{\theta}_n, \theta)/R_n\}$ is bounded in probability (or tight) uniformly over all $\theta \in \Theta$ and $n \in \mathbb{N}$. From the definition, we can see that a convergence rate is not uniquely defined. This is summarized in the following theorem.

Theorem 46 *Let $\hat{\theta}_n$ be an estimator and L be a loss function.*

1. If both $\{R_n\}$ and $\{Q_n\}$ are convergence rates of $\hat{\theta}_n$ in terms of the loss L , then $\{R_n\} \asymp \{Q_n\}$.
2. If $\{R_n\}$ is a convergence rate of $\hat{\theta}_n$ in terms of the loss L and $\{Q_n\} \ll \{R_n\}$, then $\{Q_n\}$ is also a convergence rate of $\hat{\theta}_n$ in terms of the loss L .

Proof.

1. By Definition 45 we have that $\{R_n\} \ll \{Q_n\}$ and $\{Q_n\} \ll \{R_n\}$, hence $\{R_n\} \ll \{Q_n\}$.
2. Let $C_1, C_2 > 0$ be such that for all $n \in \mathbb{N}$ it holds that $R_n \leq C_1 Q_n$ and $Q_n \leq C_2 R_n$. Let $\delta > 0$. Then, by the definition of the convergence rate, there exists $C > 0$ such that (3.4) is satisfied. Therefore, it also holds that

$$\sup_{n \in \mathbb{N}} \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \leq C C_1 Q_n \right) \leq \sup_{n \in \mathbb{N}} \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \leq C R_n \right) \leq \delta.$$

It is also easy to see that if the property (3.4) in Definition 45 is satisfied for some other rate $\{R_n\}$, then for all $n \in \mathbb{N}$ it holds that $Q_n \leq C_2 R_n \leq C_2 C_3 R_n$ for some $C_3 > 0$. Hence $\{Q_n\} \ll \{R_n\}$ and therefore $\{Q_n\}$ is a convergence rate of $\hat{\theta}$ in terms of L .

From the definition, it is also easy to see that estimators with lower order convergence rates (in terms of the loss L) converge faster to the true value of θ . Therefore, one can choose among different estimators based on their convergence rates. This task is called the *minimax estimation*.

Definition 47 (Minimax optimal estimator) *Let L be some loss function. Consider an estimator $\hat{\theta}_n$ with a convergence rate $\{R_n\}$ in terms of the loss L . This estimator (and its convergence rate) is called minimax optimal if, for all other estimators $\tilde{\theta}_n$ with the corresponding convergence rate $\{Q_n\}$, it holds that $\{R_n\} \ll \{Q_n\}$.*

In this chapter, we will derive convergence rates for the location halfspace median and scatter halfspace median in the Huber contamination model and show that these estimators are indeed minimax optimal in such settings.

The general strategy to find estimators that attain a minimax optimal convergence rate is based on two steps. First, we find some lower bound on the minimax optimal rate. That is, we search for some rate $\{Q_n\}$ such that for all estimators with convergence rate $\{R_n\}$ it holds that $\{Q_n\} \ll \{R_n\}$. This can be solved, for example, by converting this problem into a hypothesis testing task; refer to [Tsybakov, 2009, Chapter 2] for details. In the second step, we have to find some estimator with convergence rate $\{R_n\}$ such that $\{Q_n\} \ll \{R_n\}$, to show that $\{Q_n\}$ is attainable. The following theorem will be useful to obtain the lower bound of the minimax optimal rate.

Theorem 48 *Let $\{Q_n\}$ be a rate and L be some loss function. Assume that there exist $C, \delta > 0$ such that*

$$\inf_{n \in \mathbb{N}} \inf_{\hat{\theta}_n} \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \leq C Q_n \right) \leq \delta, \quad (3.5)$$

where the second infimum is taken over all estimators $\hat{\theta}_n$ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{\theta}$. Then, for every estimator $\hat{\theta}_n$ with a convergence rate $\{R_n\}$ in terms of L , it holds that $\{Q_n\} \ll \{R_n\}$.

Proof. For a contradiction, assume that there exists a convergence rate $\{R_n\}$ of some estimator $\hat{\theta}_n$ such that $\{Q_n\} \neq \{R_n\}$ does not hold. By the assumption (3.5), let $C, \delta > 0$ be such that

$$\inf_n \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \leq C Q_n \right) \geq \delta. \quad (3.6)$$

Because $\{R_n\}$ is a convergence rate of the estimator $\hat{\theta}_n$ in terms of L , by definition (3.4) there exists $C_2 > 0$ such that

$$\sup_n \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_n, \theta) \leq C_2 R_n \right) \geq \delta/2. \quad (3.7)$$

Because we assume that $\{Q_n\} \neq \{R_n\}$ is not true, there exists some $m \in \mathbb{N}$ such that $Q_m > C_2/C R_m$. However, this implies that $C Q_m > C_2 R_m$ and together with (3.6), we can bound

$$\sup_{\theta} P_{\theta} \left(L(\hat{\theta}_m, \theta) \leq C_2 R_m \right) \leq \sup_{\theta} P_{\theta} \left(L(\hat{\theta}_m, \theta) \leq C Q_m \right) \leq \delta.$$

This contradicts (3.7); thus, the theorem is proved.

Also, note that if we define the risk function as

$$R(\hat{\theta}_n, \theta) = \mathbb{E}_{\theta} \left(L(\hat{\theta}_n, \theta) \right),$$

then the other approach to minimax estimation is to minimize the maximal risk

$$\sup_{\theta} R(\hat{\theta}_n, \theta)$$

over all estimators $\hat{\theta}_n$. If the condition (3.5) is satisfied for some $C, \delta > 0$, then by the Markov inequality we have

$$\inf_n \inf_{\hat{\theta}_n} \sup_{\theta} \mathbb{E}_{\theta} \left(L(\hat{\theta}_n, \theta) \right) \leq \delta C Q_n.$$

Therefore, to bound the maximal risk from below, it is sufficient to find $C, \delta > 0$ such that (3.5) holds. This is the most common method in the literature; see [Tsybakov, 2009, Section 2.2]. However, to show that the location and scatter halfspace median are minimax optimal, it will be convenient to work directly with probabilities instead of risks.

3.2 Distance of probability measures

In this chapter, the following notions regarding the distance of probability measures will be useful.

Definition 49 (Total variation distance, Kullback-Leibler divergence) *Consider two probability measures $P, Q \in \mathcal{P}(\mathbb{R}^d)$. Then, the total variation distance of P and Q is defined as*

$$\text{TV}(P, Q) = \sup_F |P(F) - Q(F)|,$$

where the supremum is taken over all measurable sets $F \subseteq \mathbb{R}^d$. If $P \ll Q$, then the Kullback-Leibler divergence from Q to P is defined as

$$D(P||Q) = \int_{\mathbb{R}^d} \log \left(\frac{dP}{dQ} \right) dP.$$

It is easy to see that $\text{TV}(\cdot, \cdot)$ is a metric on $\mathcal{P}(\mathbb{R}^d)$. Also, note that for any two probability measures $P, Q \in \mathcal{P}(\mathbb{R}^d)$ there exists another probability measure $\mu \in \mathcal{P}(\mathbb{R}^d)$ which dominates both P and Q in the sense that $P \ll \mu$ and $Q \ll \mu$. For example, we can take $\mu = \frac{1}{2}P + \frac{1}{2}Q$ or some other convex combination. This allows us to formulate Scheffé's Theorem, which associates the total variation distance with L^1 distance of densities.

Theorem 50 (Scheffé's Theorem) *Let $P, Q \in \mathcal{P}(\mathbb{R}^d)$ and $\mu \in \mathcal{P}(\mathbb{R}^d)$ be such that $P, Q \ll \mu$. Denote $p = dP/d\mu$ and $q = dQ/d\mu$. Then, it holds that*

$$\text{TV}(P, Q) = \frac{1}{2} \int_{\mathbb{R}^d} |p - q| d\mu = 1 - \int_{\mathbb{R}^d} (p \wedge q) d\mu.$$

Proof. At first, note that the following equalities hold

$$\begin{aligned} (p - q) \mathbf{1}_{\{p > q\}} + (q - p) \mathbf{1}_{\{p < q\}} &= |p - q|, \\ (p - q) \mathbf{1}_{\{p > q\}} &= p - p \wedge q, \\ (q - p) \mathbf{1}_{\{p < q\}} &= q - p \wedge q. \end{aligned}$$

By integrating these over \mathbb{R}^d w.r.t. μ we get

$$\begin{aligned} \int_{p > q} (p - q) d\mu + \int_{p < q} (q - p) d\mu &= \int_{\mathbb{R}^d} |p - q| d\mu, \\ \int_{p > q} (p - q) d\mu &= 1 - \int_{\mathbb{R}^d} p \wedge q d\mu, \\ \int_{p < q} (q - p) d\mu &= 1 - \int_{\mathbb{R}^d} p \wedge q d\mu. \end{aligned} \tag{3.8}$$

The latter two integrals are equal. This implies

$$\int_{p > q} (p - q) d\mu = \int_{p < q} (q - p) d\mu = \frac{1}{2} \int_{\mathbb{R}^d} |p - q| d\mu. \tag{3.9}$$

Denoting $A = \{\mathbf{x} \in \mathbb{R}^d \mid p(\mathbf{x}) > q(\mathbf{x})\}$, for any measurable $F \subseteq \mathbb{R}^d$ it holds that

$$\begin{aligned} |P(F) - Q(F)| &= \left| \int_F (p - q) d\mu \right| = \int_{F \cap A} (p - q) d\mu + \int_{F \cap A^c} (q - p) d\mu \\ &= \int_{F \cap A} (p - q) d\mu - \int_{F \cap A^c} (p - q) d\mu = |P(A \cap F) - Q(A \cap F)|. \end{aligned}$$

Therefore, the supremum in the definition of the total variation distance is attained for the set A (or A^c by symmetry). Hence

$$\begin{aligned} \text{TV}(P, Q) &= |P(A) - Q(A)| = \int_{p > q} (p - q) d\mu \\ &\stackrel{(3.9)}{=} \frac{1}{2} \int_{\mathbb{R}^d} |p - q| d\mu \stackrel{(3.8)}{=} 1 - \int_{\mathbb{R}^d} (p \wedge q) d\mu. \end{aligned}$$

The total variation distance and Kullback-Leibler divergence are closely linked through Pinsker's inequality.

Theorem 51 (Pinsker's inequality) For any $P, Q \in \mathcal{P}(\mathbb{R}^d)$ it holds that

$$\text{TV}(P, Q)^2 \leq D(P\|Q)/2.$$

Proof. Refer to [Tsybakov, 2009, Lemma 2.5].

The following theorem determines the Kullback-Leibler divergence of two multivariate Gaussian distributions.

Theorem 52 For any $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \text{PD}_d$ it holds that

$$\begin{aligned} D(N_d(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \| N_d(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) \\ = \frac{1}{2} \left(\log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \text{tr}(\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - d \right). \end{aligned}$$

In particular, for $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \mathbf{I}_d$ we obtain

$$D(N_d(\boldsymbol{\mu}_1, \mathbf{I}_d) \| N_d(\boldsymbol{\mu}_2, \mathbf{I}_d)) = \frac{1}{2} \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2.$$

Proof. For $i = 1, 2$, denote by P_i the distribution $N_d(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ and by p_i its probability density function w.r.t. λ_d . It is easy to see that $dP_1/dP_2 = p_1/p_2$. Note that

$$\begin{aligned} \log \left(\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \right) \\ = \frac{1}{2} \left(\log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) - (\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + (\mathbf{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \right) \end{aligned}$$

for any $\mathbf{x} \in \mathbb{R}^d$. Let $\mathbf{X} \sim P_1$, then $D(P_1\|P_2)$ can be obviously calculated as

$$\begin{aligned} D(P_1\|P_2) &= \mathbb{E} \left[\log \left(\frac{p_1(\mathbf{X})}{p_2(\mathbf{X})} \right) \right] \\ &= \frac{1}{2} \left(\log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) - \mathbb{E}(\mathbf{X} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}_1^{-1}(\mathbf{X} - \boldsymbol{\mu}_1) + \mathbb{E}(\mathbf{X} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\mathbf{X} - \boldsymbol{\mu}_2) \right). \end{aligned}$$

Note that for a random vector $\mathbf{Y} \sim P(\mathbb{R}^d)$ with mean $\boldsymbol{\mu}$, variance matrix $\boldsymbol{\Sigma}$, and a matrix $\mathbf{A} \in \text{PD}_d$, it holds that

$$\mathbb{E}[\mathbf{Y}^\top \mathbf{A} \mathbf{Y}] = \boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu} + \text{tr}(\mathbf{A} \boldsymbol{\Sigma}).$$

Using this, we obtain

$$\begin{aligned} D(P_1\|P_2) \\ &= \frac{1}{2} \left(\log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) - \text{tr}(\mathbf{I}_d) + \text{tr}(\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right) \\ &= \frac{1}{2} \left(\log \left(\frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} \right) + \text{tr}(\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - d \right). \end{aligned}$$

3.3 Huber's contamination model

Recall that we are dealing with a statistical model $\mathcal{M} = \{P_\theta \mid \theta \in \Theta\} \subset P(\mathbb{R}^d)$. However, in this chapter, our task is to estimate θ from Huber's contamination model.

Definition 53 (Huber's contamination model) *Let $\varepsilon \in [0, 1]$ and $Q \in P(\mathbb{R}^d)$. We denote by*

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

the distribution P_θ that is ε -contaminated by Q . For a fixed $\varepsilon \in [0, 1)$, the statistical model

$$H_\varepsilon = \{P_{(\varepsilon, \theta, Q)} \mid \theta \in \Theta, Q \in P(\mathbb{R}^d)\}$$

is called Huber's contamination model.

Huber's contamination model assumes that the data may contain both "clean" observations from the assumed distribution P_θ and "contaminating" observations from some other distribution Q (outliers, faulty observations, ...).

Note *For $\varepsilon > 0$, the parameter θ is usually not identifiable in the model H_ε . This means that for $\theta_1, \theta_2 \in \Theta, \theta_1 \neq \theta_2$, one is able to find $Q_1, Q_2 \in P(\mathbb{R}^d)$ such that $P_{(\varepsilon, \theta_1, Q_1)} = P_{(\varepsilon, \theta_2, Q_2)}$. This will be used to prove Theorem 55.*

As indicated in Section 3.1, we must first find a lower bound to determine a minimax optimal convergence rate in Huber's contamination model H_ε . However, it turns out that if we can find a lower bound $\{R_n\}$ for a minimax optimal rate in the non-contaminated model \mathcal{M} , then $\{R_n \omega(\varepsilon, \Theta)\}$ is a lower bound for a minimax optimal rate in model H_ε , where $\omega(\varepsilon, \Theta)$ is the modulus of continuity of model \mathcal{M} .

Definition 54 (Modulus of continuity) *Given a statistical model \mathcal{M} , a loss function L and $\varepsilon \in [0, 1)$, the quantity*

$$\omega(\varepsilon, \Theta) = \sup \{L(\theta_1, \theta_2) \mid \text{TV}(P_{\theta_1}, P_{\theta_2}) \leq \varepsilon/(1 - \varepsilon); \theta_1, \theta_2 \in \Theta\}$$

is called the modulus of continuity of model \mathcal{M} .

The concept of the modulus of continuity dates back to the pioneering work of Donoho and Liu [1991]. The interpretation is that the quantity $\omega(\varepsilon, \Theta)$ assesses the ability of the loss function L to distinguish between two distributions $P_{\theta_1}, P_{\theta_2}$ that are close in terms of the total variation. In every model \mathcal{M} with identifiable parameter θ , $\omega(0, \Theta)$ equals 0.

The following theorem is one of the most important results of this work. It asserts that if we have a lower bound for the minimax optimal convergence rate in model \mathcal{M} , we can determine a lower bound for the minimax optimal convergence rate in model H_ε using the modulus of continuity $\omega(\varepsilon, \Theta)$. Recall that $\hat{\theta}_n = \Xi(\hat{P}_n)$ represents an estimator based on a statistical functional Ξ and an empirical measure \hat{P}_n . Further on, this empirical measure will correspond to a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ from a distribution P_θ (in model \mathcal{M}) or $P_{(\varepsilon, \theta, Q)}$ (in model H_ε). A lower index always indicates the underlying generating distribution, compare (3.10) and (3.11).

Theorem 55 Let $\delta > 0$ be an absolute constant, L be a loss function satisfying the τ -triangle inequality, and $\{\mathcal{R}_n^0\}$ be a rate such that in model \mathcal{M} it holds that

$$\inf_n \inf_{\hat{\theta}_n} \sup_{\theta} P_{\theta} \left(L \left(\hat{\theta}_n, \theta \right) \leq \mathcal{R}_n^0 \right) \geq \delta. \quad (3.10)$$

Then, for all $\varepsilon \in [0, 1)$, in model H_{ε} it holds that

$$\inf_n \inf_{\hat{\theta}_n} \sup_{\theta, Q} P_{(\varepsilon, \theta, Q)} \left(L \left(\hat{\theta}_n, \theta \right) \leq \mathcal{R}_n^{\varepsilon} \right) \geq \left(\delta - \frac{1}{2} \right), \quad (3.11)$$

where $\{\mathcal{R}_n^{\varepsilon}\} = \{\mathcal{R}_n^0 - \omega(\varepsilon, \Theta)\}$.

Proof. The proof is based on the proof of [Chen et al., 2018, Theorem 5.1]. However, the proof is elaborated with more detailed steps, and some ambiguities are clarified; especially the last part of the proof, where the τ -triangle inequality is used to obtain the desired lower bound.

The assertion is trivial if $\varepsilon = 0$. Consider fixed $n \in \mathbb{N}, \varepsilon \in (0, 1)$ and assume that (3.10) holds. We consider two cases.

Case 1: At first, assume $\mathcal{R}_n^0 \geq \omega(\varepsilon, \Theta)$, hence $\mathcal{R}_n^{\varepsilon} = \omega(\varepsilon, \Theta) = \mathcal{R}_n^0$. By setting $Q = P_{\theta}$, we have $P_{(\varepsilon, \theta, Q)} = P_{\theta}$. Therefore

$$\inf_{\hat{\theta}_n} \sup_{\theta, Q} P_{(\varepsilon, \theta, Q)} \left(L \left(\hat{\theta}_n, \theta \right) \leq \mathcal{R}_n^0 \right) = \inf_{\hat{\theta}_n} \sup_{\theta} P_{\theta} \left(L \left(\hat{\theta}_n, \theta \right) \leq \mathcal{R}_n^0 \right) \geq \delta - \left(\delta - \frac{1}{2} \right).$$

Hence, (3.11) holds with $\mathcal{R}_n^{\varepsilon} = \mathcal{R}_n^0$. The interesting case is the second one.

Case 2: Now, assume $\mathcal{R}_n^0 < \omega(\varepsilon, \Theta)$, hence $\mathcal{R}_n^{\varepsilon} = \omega(\varepsilon, \Theta) = \omega(\varepsilon, \Theta)$. The idea is to find $\theta_1, \theta_2 \in \Theta$ and $Q_1, Q_2 \in \mathcal{P}(\mathbb{R}^d)$ such that $P_{(\varepsilon, \theta_1, Q_1)} = P_{(\varepsilon, \theta_2, Q_2)}$ but $L(\theta_1, \theta_2)$ is close to $\omega(\varepsilon, \Theta)$.

By Definition 54 of the modulus of continuity, find θ_1, θ_2 such that

$$0 < \text{TV}(P_{\theta_1}, P_{\theta_2}) \leq \varepsilon / (1 - \varepsilon)$$

and $L(\theta_1, \theta_2)$ is arbitrarily close to $\omega(\varepsilon, \Theta)$. Then, there exists $\varepsilon \in (0, \varepsilon]$ such that

$$\text{TV}(P_{\theta_1}, P_{\theta_2}) = \frac{\varepsilon}{1 - \varepsilon}.$$

Let $\mu = \frac{1}{2}P_{\theta_1} + \frac{1}{2}P_{\theta_2}$ so that $P_{\theta_1}, P_{\theta_2} \ll \mu$ and define densities

$$p_{\theta_j} = \frac{dP_{\theta_j}}{d\mu}, \quad j \in \{1, 2\}.$$

Using these, we can define measures Q_1, Q_2 by their densities

$$\begin{aligned} \frac{dQ_1}{d\mu} &= \frac{(p_{\theta_2} - p_{\theta_1}) \mathbf{1}_{\{p_{\theta_2} \geq p_{\theta_1}\}}}{\text{TV}(P_{\theta_1}, P_{\theta_2})}, \\ \frac{dQ_2}{d\mu} &= \frac{(p_{\theta_1} - p_{\theta_2}) \mathbf{1}_{\{p_{\theta_1} \geq p_{\theta_2}\}}}{\text{TV}(P_{\theta_1}, P_{\theta_2})}. \end{aligned}$$

Now, we verify that Q_1, Q_2 are probability measures. This is done by the same analysis as was carried out in the proof of Scheffé's Theorem 50, see equalities (3.8) and (3.9). Namely, we have

$$\begin{aligned} \int_{\mathbb{R}^d} (p_{\theta_2} - p_{\theta_1}) \mathbf{1}_{\{p_{\theta_2} > p_{\theta_1}\}} d\mu &= \int_{\mathbb{R}^d} (p_{\theta_1} - p_{\theta_2}) \mathbf{1}_{\{p_{\theta_1} > p_{\theta_2}\}} d\mu \\ &= \frac{1}{2} \int_{\mathbb{R}^d} |p_{\theta_1} - p_{\theta_2}| d\mu = \text{TV}(P_{\theta_1}, P_{\theta_2}), \end{aligned}$$

where the last equality also follows from Scheffé's Theorem 50. Therefore, Q_1 and Q_2 are indeed probability measures. Now, consider contaminated probability distributions

$$P_j = (1 - \varepsilon)P_{\theta_j} + \varepsilon Q_j, \quad j \in \{1, 2\}.$$

Recall that $\text{TV}(P_{\theta_1}, P_{\theta_2}) = \varepsilon / (1 - \varepsilon)$. Hence, the density of P_1 w.r.t. μ is

$$\begin{aligned} \frac{dP_1}{d\mu} &= (1 - \varepsilon)p_{\theta_1} + \frac{\varepsilon}{\text{TV}(P_{\theta_1}, P_{\theta_2})} (p_{\theta_2} - p_{\theta_1}) \mathbf{1}_{\{p_{\theta_2} > p_{\theta_1}\}} \\ &= (1 - \varepsilon)(p_{\theta_1} + (p_{\theta_2} - p_{\theta_1}) \mathbf{1}_{\{p_{\theta_2} > p_{\theta_1}\}}) \\ &= (1 - \varepsilon)(p_{\theta_1} \vee p_{\theta_2}). \end{aligned}$$

Analogously

$$\frac{dP_2}{d\mu} = (1 - \varepsilon)(p_{\theta_1} \wedge p_{\theta_2}),$$

hence $P_1 = P_2$. To sum up, we have found parameters θ_1, θ_2 and distributions Q_1, Q_2 such that $L(\theta_1, \theta_2)$ is arbitrarily close to $\omega(\varepsilon, \Theta)$ but the contaminated distributions P_1 and P_2 generated by these parameters are identical, denoted by $P_1 = P_2 = P$.

Now, we can prove the desired inequality (3.11). First, note that for all estimators $\hat{\theta}_n$ we have, using the τ -triangle inequality of L (3.3), that

$$L(\theta_1, \theta_2) \leq \tau \left(L(\hat{\theta}_n, \theta_1) + L(\hat{\theta}_n, \theta_2) \right) \leq 2\tau \left(L(\hat{\theta}_n, \theta_1) \vee L(\hat{\theta}_n, \theta_2) \right),$$

which implies

$$\frac{1}{2\tau} L(\theta_1, \theta_2) \leq L(\hat{\theta}_n, \theta_1) \vee L(\hat{\theta}_n, \theta_2). \quad (3.12)$$

As a consequence, we can bound

$$\begin{aligned} \inf_{\hat{\theta}_n} \sup_{\theta, Q} P_{(\varepsilon, \theta, Q)} \left(L(\hat{\theta}_n, \theta) \leq \frac{1}{2\tau} L(\theta_1, \theta_2) \right) \\ &\leq \inf_{\hat{\theta}_n} \max_{j \in \{1, 2\}} P_j \left(L(\hat{\theta}_n, \theta_j) \leq \frac{1}{2\tau} L(\theta_1, \theta_2) \right) \\ &= \inf_{\hat{\theta}_n} \max_{j \in \{1, 2\}} P \left(L(\hat{\theta}_n, \theta_j) \leq \frac{1}{2\tau} L(\theta_1, \theta_2) \right) \\ &\leq \inf_{\hat{\theta}_n} \max_{j \in \{1, 2\}} P \left(L(\hat{\theta}_n, \theta_j) \leq \left(L(\hat{\theta}_n, \theta_1) \vee L(\hat{\theta}_n, \theta_2) \right) \right). \end{aligned} \quad (3.13)$$

Because $P_j \leq H_\varepsilon$ for $j = 1, 2$, the first inequality follows by taking the supremum over a smaller set. For the equality, we used that $P_1 = P_2 = P$, and in the second inequality, we used (3.12). However, for any estimator $\hat{\theta}_n$, at least one of the events

$$\left[L(\hat{\theta}_n, \theta_1) \leq L(\hat{\theta}_n, \theta_2) \right], \left[L(\hat{\theta}_n, \theta_2) \leq L(\hat{\theta}_n, \theta_1) \right]$$

has P -probability at least $1/2$. Therefore, (3.13) can be further bounded from below by $1/2$. In conclusion, because n was chosen arbitrarily, we have shown that

$$\inf_n \inf_{\hat{\theta}_n} \sup_{\theta, Q} P_{(\varepsilon, \theta, Q)} \left(L(\hat{\theta}_n, \theta) \geq \frac{1}{2\tau} L(\theta_1, \theta_2) \right) \geq \frac{1}{2}.$$

Because $\tau \geq 1$ is some fixed constant, it holds that

$$\frac{1}{2\tau} L(\theta_1, \theta_2) \leq L(\theta_1, \theta_2)$$

and θ_1, θ_2 can be chosen such that $L(\theta_1, \theta_2)$ is arbitrarily close to $\omega(\varepsilon, \Theta)$. This implies that (3.11) holds with

$$R_n^\varepsilon \leq \omega(\varepsilon, \Theta),$$

which concludes the proof.

Note *Chen et al. [2018] use Le Cam's two point testing method (see [Yu, 1997, Lemma 1]) to bound (3.13) from below. As we have shown, this is not necessary.*

To sum up, we have shown that to bound a minimax optimal rate in model H_ε from below, it suffices to find the lower bound for minimax optimal rate in model \mathcal{M} and determine the modulus of continuity of model \mathcal{M} .

3.4 Optimality of location halfspace median

Motivated by the previous theory for general parametric models, in this section, we consider a specific parametric model

$$\mathcal{M}^{\text{loc}} = \left\{ N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mid \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \in \text{PD}_d, \sigma_{\max}(\boldsymbol{\Sigma}) \leq M \right\},$$

where $M > 0$ is fixed. We consider all Gaussian distributions with covariance matrices having a bounded spectrum. The parameter of interest is the mean $\theta = \boldsymbol{\mu} \in \mathbb{R}^d = \Theta$. This is a location parameter, hence the superscript ‘‘loc’’ in \mathcal{M}^{loc} . Estimators of $\boldsymbol{\mu}$ are denoted by $\hat{\boldsymbol{\mu}}_n$ and the location halfspace median is denoted by $\hat{\boldsymbol{\mu}}_n^{\text{hs}}$. We consider the squared Euclidean loss, i.e.

$$L(\boldsymbol{\mu}_1, \boldsymbol{\mu}_2) = \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2.$$

Recall that this loss function satisfies the τ -triangle inequality for $\tau = 2$ (see inequality (3.3)). The corresponding Huber's contamination model for $\varepsilon \in [0, 1)$ is

$$H_\varepsilon^{\text{loc}} = \left\{ (1 - \varepsilon) N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \varepsilon Q \mid \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \in \text{PD}_d, \sigma_{\max}(\boldsymbol{\Sigma}) \leq M, Q \in P(\mathbb{R}^d) \right\}.$$

In this section, by $P_{(\boldsymbol{\mu}, \boldsymbol{\Sigma})}$ we denote the probability distribution $N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Further, by $P_{(\boldsymbol{\mu}, \boldsymbol{\Sigma}, Q)}$ we denote the contaminated probability distribution

$$(1 - \varepsilon) N_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + \varepsilon Q.$$

The amount of contamination ε is omitted in this notation and will be evident from the context.

First, we focus on establishing a lower bound for a minimax optimal convergence rate in model \mathcal{M}^{loc} . The following theorem is one of the essential tools for determining lower bounds in the minimax estimation; see [Tsybakov, 2009, Section 2.7] for reference.

Theorem 56 (Fano's inequality) *Let $M \leq N$, $M \geq 2$ and $\{P_1, \dots, P_M\} \subset \mathcal{P}(\mathbb{R}^d)$ with the corresponding parameters $\theta_1, \dots, \theta_M \in \mathbb{R}^p$. If $a, b > 0$ are such that $\|\theta_i - \theta_j\| \geq a$ and $D(P_i||P_j) \leq b$ for all $i \neq j$, then for all $n \in \mathbb{N}$ it holds that*

$$\inf_{\hat{\mu}_n} \frac{1}{M} \sum_{i=1}^M P_i \left(\|\theta_i - \hat{\mu}_n\| \geq a/2 \right) \leq 1 - \frac{nb + \log(2)}{\log(M)}, \quad (3.14)$$

where the infimum is taken over all estimators $\hat{\mu}_n$ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$. In (3.14), the probability

$$P_i \left(\|\theta_i - \hat{\mu}_n\| \geq a/2 \right)$$

considers that the estimator $\hat{\mu}_n$ is calculated using a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ from the distribution P_i .

Proof. The claim follows from the proof of [Yu, 1997, Lemma 3]. The proof utilizes advanced methods from information theory.

Using Fano's inequality, we can deduce the following result.

Theorem 57 *In model \mathcal{M}^{loc} , there exists an absolute constant $\delta > 0$ such that*

$$\inf_{n \in \mathbb{N}} \inf_{\hat{\mu}_n} \sup_{P_{(\mu, \Sigma)} \in \mathcal{M}^{\text{loc}}} P_{(\mu, \Sigma)} \left(\|\hat{\mu}_n - \mu\|^2 \geq R_n^0 \right) \geq \delta \quad (3.15)$$

holds for some $\{R_n^0\} \subset \mathbb{R}^+$ with $R_n^0 \geq d/n$. The infimum is taken over all estimators $\hat{\mu}_n$ of μ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{(\mu, \Sigma)}$.

Note *Chen et al. [2018] assure that the claim of Theorem 57 is a classic result of minimax estimation theory. The authors provide [Ma and Wu, 2015] as a reference for this assertion. However, I did not find the stated claim in this paper or in other literature. Therefore, the proof of this theorem is the original result of this work. The guide [Scarlett and Cevher, 2021] was helpful for correctly applying Fano's inequality.*

Proof. Let us consider

$$R_n^0 = \frac{\kappa}{64} \frac{d}{n} \quad \text{where} \quad \kappa = \frac{\log(2)}{4} > 0. \quad (3.16)$$

We will prove that (3.15) holds with $\delta = 1/4$ using Fano's inequality. Obviously, for all $n \in \mathbb{N}$, we can bound

$$\inf_{\hat{\mu}_n} \sup_{P_{(\mu, \Sigma)} \in \mathcal{M}^{\text{loc}}} P_{(\mu, \Sigma)} \left(\|\hat{\mu}_n - \mu\|^2 \geq R_n^0 \right) \leq \inf_{\hat{\mu}_n} \sup_{\mu \in \mathcal{I}_d} P_{(\mu, \mathbf{I}_d)} \left(\|\hat{\mu}_n - \mu\|^2 \geq R_n^0 \right),$$

where

$$\Theta_0 = \left\{ \boldsymbol{\mu} \in \mathbb{R}^d \mid \|\boldsymbol{\mu}\| \leq \sqrt{\frac{\kappa d}{n}} \right\}. \quad (3.17)$$

Therefore, to prove this theorem, it is sufficient to show that for all $n \in \mathbb{N}$, the inequality

$$\inf_{\hat{\boldsymbol{\mu}}_n} \sup_{\boldsymbol{\mu}_0} P_{(\boldsymbol{\mu}, \mathbf{I}_d)} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|^2 \leq R_n^0 \right) \geq \frac{1}{4} \quad (3.18)$$

holds.

Fix $n \in \mathbb{N}$ and let $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_M\} \subset \Theta_0$ be a set of the maximum cardinality such that for all $i \neq j$ we have

$$\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\| \geq \frac{1}{4} \sqrt{\frac{\kappa d}{n}}. \quad (3.19)$$

The set $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_M\}$ is a so-called $\frac{1}{4} \sqrt{\frac{\kappa d}{n}}$ -packing of Θ_0 . Then:

1. By Theorem 52 and our choice of Θ_0 in (3.17), we can bound

$$D(P_{(\boldsymbol{\mu}_i, \mathbf{I}_d)} \| P_{(\boldsymbol{\mu}_j, \mathbf{I}_d)}) = \frac{1}{2} \|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2 \leq \frac{1}{2} (\|\boldsymbol{\mu}_i\| + \|\boldsymbol{\mu}_j\|)^2 \leq 2 \frac{\kappa d}{n}. \quad (3.20)$$

2. By Theorem A5 in the Appendix, we can bound the cardinality of any $\frac{1}{4} \sqrt{\frac{\kappa d}{n}}$ -packing of Θ_0 from below by 4^d , i.e.

$$M \geq 4^d. \quad (3.21)$$

Using this, we bound the left side of (3.18) from below as follows

$$\begin{aligned} & \inf_{\hat{\boldsymbol{\mu}}_n} \sup_{\boldsymbol{\mu}_0} P_{(\boldsymbol{\mu}, \mathbf{I}_d)} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|^2 \leq R_n^0 \right) \\ & \quad \inf_{\hat{\boldsymbol{\mu}}_n} \max_{\{1, \dots, M\}} P_{(\boldsymbol{\mu}_i, \mathbf{I}_d)} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}_i\|^2 \leq R_n^0 \right) \\ & \quad \inf_{\hat{\boldsymbol{\mu}}_n} \frac{1}{M} \sum_{i=1}^M P_{(\boldsymbol{\mu}_i, \mathbf{I}_d)} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}_i\|^2 \leq R_n^0 \right) \\ & \stackrel{(3.16)}{=} \inf_{\hat{\boldsymbol{\mu}}_n} \frac{1}{M} \sum_{i=1}^M P_{(\boldsymbol{\mu}_i, \mathbf{I}_d)} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}_i\| \leq \frac{1}{2} \left(\frac{1}{4} \sqrt{\frac{\kappa d}{n}} \right) \right) \\ & \stackrel{\text{Fano}}{\geq} 1 - \frac{2n \frac{\kappa d}{n} + \log(2)}{\log(M)} \\ & \stackrel{(3.21)}{\geq} 1 - \frac{2n \frac{\kappa d}{n} + \log(2)}{\log(4^d)} = 1 - \frac{\kappa}{\log(2)} - \frac{1}{2d} \\ & \stackrel{d \geq 1}{\geq} 1 - \frac{\kappa}{\log(2)} - \frac{1}{2} \stackrel{(3.16)}{=} \frac{1}{2} - \frac{1}{4} = \frac{1}{4}. \end{aligned}$$

In the first inequality, we take the supremum over a smaller set. In the second inequality, we use that a maximum is always bounded from below by an average. In the third inequality, we use Fano's inequality (stated in Theorem 56) for $P_i = P_{(\boldsymbol{\mu}_i, \mathbf{I}_d)}$ and $\boldsymbol{\mu}_i = \boldsymbol{\mu}_i$, $i = 1, \dots, M$. The assumptions of Fano's inequality

are satisfied by (3.19), (3.20) and (3.21). This proves the theorem.

Now we can use Theorem 55 to determine a lower bound of a minimax optimal convergence rate in model $H_\varepsilon^{\text{loc}}$.

Theorem 58 *There exists an absolute constant $\delta > 0$ such that for any $\varepsilon \in [0, 1)$*

$$\inf_n \inf_{\hat{\boldsymbol{\mu}}_n} \sup_{P_{(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathcal{Q})} \in H_\varepsilon^{\text{loc}}} P_{(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathcal{Q})} \left(\|\hat{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\|^2 \geq R_n^\varepsilon \right) \geq \delta \quad (3.22)$$

holds for some $\{R_n^\varepsilon\} \asymp (d/n) \varepsilon^2$. The infimum is taken over all estimators $\hat{\boldsymbol{\mu}}_n$ of $\boldsymbol{\mu}$ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \mathcal{Q})}$.

Proof. Without loss of generality, let $M = 1$. To use Theorem 55, we must first bound the modulus of continuity $\omega(\varepsilon, \Theta)$. Note that by Pinsker's inequality given in Theorem 51 and using the Kullback-Leibler divergence of two Gaussian distributions (see Theorem 52), we have for any $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^d$ that

$$\text{TV} \left(P_{(\boldsymbol{\mu}_1, \mathbf{I}_d)}, P_{(\boldsymbol{\mu}_2, \mathbf{I}_d)} \right)^2 = \frac{1}{2} D(P_{(\boldsymbol{\mu}_1, \mathbf{I}_d)} \| P_{(\boldsymbol{\mu}_2, \mathbf{I}_d)}) = \frac{1}{4} \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2.$$

Therefore, if $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|/2 \leq \varepsilon$, then

$$\text{TV} \left(P_{(\boldsymbol{\mu}_1, \mathbf{I}_d)}, P_{(\boldsymbol{\mu}_2, \mathbf{I}_d)} \right) \leq \frac{\varepsilon}{1 - \varepsilon}.$$

As a consequence

$$\omega(\varepsilon, \Theta) = \sup \left\{ \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2 / \text{TV} \left(P_{(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}, P_{(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} \right) \mid \frac{\varepsilon}{1 - \varepsilon} \right\} \\ \sup \left\{ \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|^2 / \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|/2 \leq \varepsilon \right\} = 4\varepsilon^2.$$

The inequality follows from the fact that the first supremum is taken over all $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \text{PD}_d$ such that $\sigma_{\max}(\boldsymbol{\Sigma}_1), \sigma_{\max}(\boldsymbol{\Sigma}_2) \leq M$ and the second supremum is taken only over $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \mathbf{I}_d$ and $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2$ such that $\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|/2 \leq \varepsilon$. We have shown that $\omega(\varepsilon, \Theta) \leq 4\varepsilon^2$, hence $(d/n) \omega(\varepsilon, \Theta) \leq (d/n) (4\varepsilon^2)$. By Theorem 55 and Theorem 57, (3.22) holds for some $\{R_n^\varepsilon\} \asymp (d/n) \omega(\varepsilon, \Theta)$, hence also for some $\{R_n^\varepsilon\} \asymp (d/n) \varepsilon^2$.

We demonstrated that a minimax optimal convergence rate in model $H_\varepsilon^{\text{loc}}$ is bounded from below by $(d/n) \varepsilon^2$. Next, we will show that this rate is indeed attained for the location halfspace median $\hat{\boldsymbol{\mu}}_n^{\text{HS}}$. This indicates that the location halfspace median is indeed minimax optimal in model $H_\varepsilon^{\text{loc}}$.

Theorem 59 *Assume that $\varepsilon \in (0, 1/5)$. Then, for any $\delta \in (0, 1/2)$ there exist absolute constants $C, D > 0$ such that for all $n \in \mathbb{N}$ satisfying*

$$\left(\frac{d}{n} + \frac{\log(1/\delta)}{n} \right) < D \quad (3.23)$$

we have

$$\inf_{P_{(\mu, \Sigma, Q)}} \sup_{H_\varepsilon^{\text{loc}}} P_{(\mu, \Sigma, Q)} \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} - \boldsymbol{\mu} \right\|^2 < C \left(\left(\frac{d}{n} \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right) \geq 1 - 2\delta,$$

hence

$$\sup_{P_{(\mu, \Sigma, Q)}} \sup_{H_\varepsilon^{\text{loc}}} P_{(\mu, \Sigma, Q)} \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} - \boldsymbol{\mu} \right\|^2 \geq C \left(\left(\frac{d}{n} \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right) < 2\delta.$$

Proof. The proof is based on the proof of [Chen et al., 2018, Theorem 2.1]. However, the theorem is proven down to finer details, and some ambiguities are clarified. Especially the last part of the proof, where the Lipschitz continuity of the quantile function of the standard Gaussian distribution is utilized, is rigorously argued. The proof is divided into four parts.

Part 1: Auxiliary observations

First, we prepare the following observations that will be useful in deriving the intended bound.

- (L1) By the affine equivariance of the location halfspace median, we can assume that $\boldsymbol{\mu} = \mathbf{0}$. Also, without loss of generality, we can take $\boldsymbol{\Sigma} = \mathbf{I}_d$. Otherwise, we would transform the random sample by the linear mapping $\boldsymbol{x} \mapsto \boldsymbol{\Sigma}^{-1/2} \boldsymbol{x}$ and note that

$$\left\| \boldsymbol{\Sigma}^{-1/2} \hat{\boldsymbol{\mu}}_n^{\text{hs}} \right\|^2 = \left(\hat{\boldsymbol{\mu}}_n^{\text{hs}} \right)^\top \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_n^{\text{hs}} = \left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} \right\|^2 / M.$$

The inequality holds since all eigenvalues of $\boldsymbol{\Sigma}$ are bounded from above by M , hence all eigenvalues of $\boldsymbol{\Sigma}^{-1}$ are bounded from below by M^{-1} . Therefore, it suffices to find C under the assumption $\boldsymbol{\Sigma} = \mathbf{I}_d$ and then multiply this constant by M .

- (L2) Recall that $hD(\boldsymbol{x}; N_d(\mathbf{0}, \mathbf{I}_d)) = 1 - \Phi(\|\boldsymbol{x}\|)$ for all $\boldsymbol{x} \in \mathbb{R}^d$. See Example 3 for derivation.
- (L3) Considering a random sample $\{\boldsymbol{X}_i\}_{i=1}^n$ from a contaminated distribution $P_{(\mathbf{0}, \mathbf{I}_d, Q)}$, we can decompose $\{\boldsymbol{X}_i\}_{i=1}^n = \{\boldsymbol{Y}_i\}_{i=1}^{n_1} \cup \{\boldsymbol{Z}_i\}_{i=1}^{n_2}$ where, marginally, $n_2 \sim \text{Binomial}(n, \varepsilon)$, $n_1 = n - n_2$, and conditionally on n_1, n_2 we have that $\{\boldsymbol{Y}_i\}_{i=1}^{n_1}$ is a random sample from $N_d(\mathbf{0}, \mathbf{I}_d)$ and $\{\boldsymbol{Z}_i\}_{i=1}^{n_2}$ is a random sample from Q .

- (L4) By Theorem A2 in the Appendix, we have with probability at least $1 - \delta$ that

$$\sup_{\boldsymbol{x} \in \mathbb{R}^d} |hD(\boldsymbol{x}; N_d(\mathbf{0}, \mathbf{I}_d)) - hD(\boldsymbol{x}; \{\boldsymbol{Y}_i\}_{i=1}^{n_1})| \leq \sup_{H \in \mathcal{H}_d} \left| P_{(\mathbf{0}, \mathbf{I}_d)}(H) - \hat{P}_{n_1}(H) \right| \leq \sqrt{\frac{1440\pi e}{1 - e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}},$$

where \hat{P}_{n_1} is the empirical distribution of $\{\boldsymbol{Y}_i\}_{i=1}^{n_1}$ and \mathcal{H}_d is the system of all closed halfspaces in \mathbb{R}^d , i.e.

$$\mathcal{H}_d := \left\{ H_{\boldsymbol{x}, \boldsymbol{u}} \mid \boldsymbol{x} \in \mathbb{R}^d, \boldsymbol{u} \in S^{d-1} \right\}. \quad (3.24)$$

(L5) By the definition of the sample halfspace depth, it follows that

$$n_1 hD(\mathbf{x}; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - n hD(\mathbf{x}; \{\mathbf{X}_i\}_{i=1}^n) - n_2 \geq n_1 hD(\mathbf{x}; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - n_2 \quad (3.25)$$

holds for all $\mathbf{x} \in \mathbb{R}^d$. For example, to see the first inequality, note that

$$\inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \sum_{i=1}^{n_1} \mathbf{1}_{\{\mathbf{Y}_i, \mathbf{u} \cdot \mathbf{x}, \mathbf{u}\}} \geq \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \sum_{i=1}^n \mathbf{1}_{\{\mathbf{X}_i, \mathbf{u} \cdot \mathbf{x}, \mathbf{u}\}} - n_2. \quad (3.26)$$

This inequality holds because, for a fixed $\mathbf{u} \in \mathbb{S}^{d-1}$, the left side of (3.26) represents the number of observations \mathbf{Y}_i in $H_{\mathbf{x}, \mathbf{u}}$, which is always greater than or equal to the number of observations \mathbf{X}_i in $H_{\mathbf{x}, \mathbf{u}}$ without n_2 . This is because some of the observations \mathbf{Z}_i can also lie in $H_{\mathbf{x}, \mathbf{u}}$. The second inequality is proven analogously.

(L6) By Theorem A4 in the Appendix, if

$$\sqrt{\frac{\log(1/\delta)}{2n}} < 1/5 \quad (3.27)$$

holds for $\delta > 0$, then

$$\frac{n_2}{n_1} \leq \frac{\varepsilon}{1 - \varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} < \frac{2}{3}$$

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

Part 2: Choice of constant D

Now, we define a constant D so that the proof of the statement works. In particular, we need to ensure that (3.27) holds and also that

$$80 \sqrt{\frac{3\pi e}{1 - e^{-1}}} \sqrt{\frac{d}{n}} + \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}} < \frac{1}{5} \quad (3.28)$$

is true, which will be useful later. For the entire subsequent proof, fix $\delta \in (0, 1/2)$ and consider $n \in \mathbb{N}$ such that the assumption (3.23) is satisfied for

$$D = \left(\frac{\frac{1}{5}}{80 \sqrt{\frac{3\pi e}{1 - e^{-1}}} + \frac{7}{2}} \right)^2 > 0.$$

Because $a + b < D$ implies $a - b < D$ for all $a, b > 0$, it is easy to see that also

$$\sqrt{\frac{d}{n}} - \sqrt{\frac{\log(1/\delta)}{n}} < \bar{D}. \quad (3.29)$$

Therefore, by our choice of D and (3.29), we can bound

$$80 \sqrt{\frac{3\pi e}{1 - e^{-1}}} \sqrt{\frac{d}{n}} + \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}} < \left(80 \sqrt{\frac{3\pi e}{1 - e^{-1}}} + \frac{7}{2} \right) \bar{D} = \frac{1}{5}.$$

As a consequence, inequality (3.28) holds. Also, inequality (3.28) implies that

$$\sqrt{\frac{\log(1/\delta)}{2n}} < \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}} < \frac{1}{5},$$

hence, using the assumption $\varepsilon < 1/5$, (L6) gives that

$$\mathbb{P} \left[\frac{n_2}{n_1} \leq \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} < \frac{2}{3} \right] \geq 1 - \delta. \quad (3.30)$$

Further, note that

$$\frac{n_2}{n_1} < \frac{2}{3} \quad 3n_2 < 2n_1 \quad 3(n - n_1) < 2n_1 \quad n_1 > 3n/5. \quad (3.31)$$

To sum up, we find out that for our fixed δ and a large enough n , at least $3/5$ of all observations are non-contaminating with probability at least $1 - \delta$. Also, the inequality (3.28) holds. Both of these facts will be used later in the proof.

Part 3: Bound for the depth of $\hat{\boldsymbol{\mu}}_n^{\text{hs}}$

We derive the following series of inequalities using the above observations and our choice of D . These hold with probability at least $1 - \delta$ conditionally on the marginal counts n_1, n_2 . We have

$$\begin{aligned} & 1 - \Phi \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} \right\| \right) \stackrel{\text{(L2)}}{=} hD(\hat{\boldsymbol{\mu}}_n^{\text{hs}}; N_d(\mathbf{0}, \mathbf{I}_d)) \\ & \stackrel{\text{(L4)}}{\leq} hD(\hat{\boldsymbol{\mu}}_n^{\text{hs}}; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \stackrel{\text{(L5)}}{\leq} \frac{n}{n_1} hD(\hat{\boldsymbol{\mu}}_n^{\text{hs}}; \{\mathbf{X}_i\}_{i=1}^n) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \leq \frac{n}{n_1} hD(\mathbf{0}; \{\mathbf{X}_i\}_{i=1}^n) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \stackrel{\text{(L5)}}{\leq} hD(\mathbf{0}; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \stackrel{\text{(L4)}}{\leq} hD(\mathbf{0}; N_d(\mathbf{0}, \mathbf{I}_d)) - \frac{n_2}{n_1} - 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{2\log(1/\delta)}{n_1}} \\ & \stackrel{\text{(L2)}}{=} \frac{1}{2} - \frac{n_2}{n_1} - 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} - \sqrt{\frac{2\log(1/\delta)}{n_1}}, \end{aligned} \quad (3.32)$$

where the third inequality follows from the fact that $\hat{\boldsymbol{\mu}}_n^{\text{hs}}$ is the maximizer of $hD(\cdot; \{\mathbf{X}_i\}_{i=1}^n)$. Rewriting (3.32), we have that

$$\mathbb{P} \left[\Phi \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} \right\| \right) \leq \frac{1}{2} + \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \mid n_1, n_2 \right] \geq 1 - \delta.$$

However, by taking the expectation on both sides (and considering its monotonicity), we obtain

$$\mathbb{P} \left[\Phi \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} \right\| \right) \leq \frac{1}{2} + \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \right] \geq 1 - \delta.$$

Now, we combine this result with inequality (3.30). Note that for any two random events A, B we have that $\mathbb{P}(A \cap B) \geq \mathbb{P}(A) + \mathbb{P}(B) - 1$. Therefore, it holds that

$$\mathbb{P} \left[\Phi \left(\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| \right) \leq \frac{1}{2} + \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}}, \right. \\ \left. \frac{n_2}{n_1} \leq \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} \leq \frac{2}{3} \right] \geq 1 - 2\delta. \quad (3.33)$$

Now, under the condition of the second random event in (3.33), we can further upper bound

$$\begin{aligned} \Phi \left(\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| \right) &\stackrel{(3.32)}{\leq} \frac{1}{2} + \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \\ &\stackrel{(0)=\frac{1}{2}}{=} \Phi(0) + \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \\ &\stackrel{(3.30)}{\leq} \Phi(0) + \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \\ &\stackrel{(3.31)}{\leq} \Phi(0) + \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{5(d+1)}{3n}} + \sqrt{\frac{10\log(1/\delta)}{3n}} \\ &= \Phi(0) + \frac{\varepsilon}{1-\varepsilon} + 40\sqrt{\frac{6\pi e}{1-e^{-1}}} \sqrt{\frac{d+1}{n}} + \sqrt{\frac{1585\log(1/\delta)}{288n}} \\ &\stackrel{d+1 \leq 2d}{\leq} \Phi(0) + \frac{\varepsilon}{1-\varepsilon} + 40\sqrt{\frac{6\pi e}{1-e^{-1}}} \sqrt{\frac{2d}{n}} + \sqrt{\frac{1585}{288}} \sqrt{\frac{\log(1/\delta)}{n}} \\ &\stackrel{\varepsilon < 1/5}{<} \Phi(0) + \frac{5}{4}\varepsilon + 80\sqrt{\frac{3\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}}. \end{aligned}$$

In the last inequality, we also used the fact that $\sqrt{1585/288} < 7/2$. Ultimately, we have

$$\mathbb{P} \left[\Phi \left(\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| \right) - \Phi(0) \leq \frac{5}{4}\varepsilon + 80\sqrt{\frac{3\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}}, \right. \\ \left. \frac{n_2}{n_1} \leq \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} \leq \frac{2}{3} \right] \geq 1 - 2\delta.$$

For any random events A, B , we have $\mathbb{P}(A \cap B) \geq \mathbb{P}(A)$. Therefore, the preceding inequality implies that

$$\Phi \left(\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| \right) - \Phi(0) < \frac{5}{4}\varepsilon + 80\sqrt{\frac{3\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + \frac{7}{2} \sqrt{\frac{\log(1/\delta)}{n}} \quad (3.34)$$

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

Part 4: Conclusion

We assumed that n is chosen such that (3.28) holds. Therefore, using this and the assumption $\varepsilon < 1/5$, from (3.34) we can deduce

$$\Phi \left(\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| \right) - \Phi(0) < \frac{5}{4}\varepsilon + \frac{1}{5} < \frac{1}{4} + \frac{1}{5} = \frac{9}{20}.$$

Now we can use the fact that $\Phi(0) = 1/2$ and the quantile function of the standard Gaussian distribution is Lipschitz continuous on the interval $[1/2, 19/20]$ with the Lipschitz constant

$$C = [\varphi(\Phi^{-1}(19/20))]^{-1},$$

where φ is the probability density function of $\mathcal{N}_1(0, 1)$. Therefore, from (3.34) we can further deduce

$$\begin{aligned} \|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\| &< C \left(\frac{5}{4}\varepsilon + 80\sqrt{\frac{3\pi e}{1-e^{-1}}}\sqrt{\frac{d}{n}} + \frac{7}{2}\sqrt{\frac{\log(1/\delta)}{n}} \right) \\ &C \left(\left(\frac{5}{4} + 80\sqrt{\frac{3\pi e}{1-e^{-1}}} \right) \left(\sqrt{\frac{d}{n}} \quad \varepsilon \right) + \frac{7}{2}\sqrt{\frac{\log(1/\delta)}{n}} \right), \end{aligned}$$

hence, using 2-triangle inequality (see (3.3)), we have

$$\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\|^2 \leq C \left(\left(\frac{d}{n} \quad \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right),$$

where C is defined as

$$2(C)^2 \left(\left(\frac{5}{4} + 80\sqrt{\frac{3\pi e}{1-e^{-1}}} \right)^2 + \frac{49}{4} \right).$$

All in all, we have proven that

$$\mathbb{P} \left[\|\hat{\boldsymbol{\mu}}_n^{\text{hs}}\|^2 \leq C \left(\left(\frac{d}{n} \quad \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right] \geq 1 - 2\delta, \quad (3.35)$$

which concludes the proof.

Finally, the minimax optimality of Tukey's median follows.

Theorem 60 *For $\varepsilon < 1/5$, the location halfspace median $\hat{\boldsymbol{\mu}}_n^{\text{hs}}$ is the minimax optimal estimator of $\boldsymbol{\mu}$ in terms of the squared Euclidean loss in Huber's contamination model $H_\varepsilon^{\text{loc}}$. The minimax optimal convergence rate is $d/n \quad \varepsilon^2$.*

Proof. From Theorem 58, we know that the minimax optimal convergence rate is bounded from below by $(d/n) \quad \varepsilon^2$. The statement follows after we prove that this rate is indeed attained for the location halfspace median. However, this is an immediate consequence of Theorem 59 after we recognize the subsequent two observations:

1. The following inequality holds

$$\begin{aligned} \left(\left(\frac{d}{n} \quad \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) &\leq \left(\left(\frac{d}{n} \quad \varepsilon^2 \right) + \frac{d \log(1/\delta)}{n} \right) \\ &= (1 + \log(1/\delta)) \left(\frac{d}{n} \quad \varepsilon^2 \right) = \tilde{C} \left(\frac{d}{n} \quad \varepsilon^2 \right) \end{aligned}$$

for an absolute constant $\tilde{C} = (1 + \log(1/\delta)) > 0$ independent of n, d and ε . Therefore, by Theorem 59

$$\inf_{P_{(\mu, \Sigma, Q)} \in \mathcal{H}_\varepsilon^{\text{loc}}} P_{(\mu, \Sigma, Q)} \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} - \boldsymbol{\mu} \right\|^2 < C \tilde{C} \left(\frac{d}{n} \varepsilon^2 \right) \right) \\ \inf_{P_{(\mu, \Sigma, Q)} \in \mathcal{H}_\varepsilon^{\text{loc}}} P_{(\mu, \Sigma, Q)} \left(\left\| \hat{\boldsymbol{\mu}}_n^{\text{hs}} - \boldsymbol{\mu} \right\|^2 < C \left(\left(\frac{d}{n} \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right) \\ 1 - 2\delta$$

holds.

- By (3.23), Theorem 59 holds only for $n > (d - \log(\delta))/D$. However, this is not a problem. There are only finitely many $n \leq (d - \log(\delta))/D$. Therefore, when we look closely at the proof of Theorem 59, specifically at equation (3.35), we can always choose a larger constant $C > 0$ such that this inequality holds even for $n \leq (d - \log(\delta))/D$.

Note As noted in [Chen et al., 2018, Remark 2.1], Theorem 59 remains valid for any $\varepsilon < 1/3 - c_0$, where c_0 can be an arbitrarily small constant. The threshold $1/3$ represents the maximum breakdown point for Tukey’s median; see Theorem 26.

To conclude, we have demonstrated that Tukey’s median is a minimax optimal estimator in Huber’s contamination model \mathcal{M}^{loc} . This highlights another of its robustness properties, in addition to those introduced in Chapter 1. Specifically, the minimax optimality underlines Tukey’s median’s effectiveness in handling outliers and contaminated data, making it a highly reliable choice for statistical estimation in various real-world scenarios.

3.5 Optimality of scatter halfspace median

In this section, we consider the following parametric model

$$\mathcal{M}^{\text{cov}} = \{N_d(\mathbf{0}, \boldsymbol{\Sigma}) \mid \boldsymbol{\Sigma} \in \text{PD}_d, \sigma_{\max}(\boldsymbol{\Sigma}) \leq M\},$$

which consists of all centered Gaussian distributions with a covariance matrix whose spectrum is bounded by some fixed constant $M > 0$. The parameter of interest is the covariance matrix $\theta = \boldsymbol{\Sigma}$, hence the superscript “cov” in \mathcal{M}^{cov} . We have

$$\Theta = \{\boldsymbol{\Sigma} \in \text{PD}_d \mid \sigma_{\max}(\boldsymbol{\Sigma}) \leq M\}.$$

The error of an estimator is measured by the squared spectral norm

$$L(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) = \left\| \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2 \right\|_{\text{op}}^2.$$

Same as the squared Euclidean norm, this loss function also satisfies the τ -triangle inequality for $\tau = 2$ (see inequality (3.3)). The following theorem presents various methods to determine L .

Theorem 61 For any $\Gamma_1, \Gamma_2 \in \text{PD}_d$ it holds that

$$\|\Gamma_1 - \Gamma_2\|_{\text{op}} = \max\{|\lambda_1|, \dots, |\lambda_d|\} = \max_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \mathbf{u}^\top (\Gamma_1 - \Gamma_2) \mathbf{u} \right|,$$

where $\lambda_1, \dots, \lambda_d$ are the eigenvalues of $\Gamma_1 - \Gamma_2$.

Proof. We know that the spectral norm of a matrix is equal to its largest singular value. The matrix $\Gamma_1 - \Gamma_2$ is the difference between two symmetric matrices. Hence, it is symmetric. For a symmetric matrix, the singular values are the absolute values of its eigenvalues, which gives us the first equality.

To prove the second equality, we decompose

$$\Gamma_1 - \Gamma_2 = \mathbf{Q}^\top \Lambda \mathbf{Q},$$

where \mathbf{Q} is an orthogonal matrix and Λ is a diagonal matrix with the eigenvalues of $\Gamma_1 - \Gamma_2$ on the diagonal. Using this decomposition, we can calculate

$$\begin{aligned} \max_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \mathbf{u}^\top (\Gamma_1 - \Gamma_2) \mathbf{u} \right| &= \max_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \mathbf{u}^\top \mathbf{Q}^\top \Lambda \mathbf{Q} \mathbf{u} \right| \\ &= \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \mathbf{v}^\top \Lambda \mathbf{v} \right| \\ &= \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \sum_{i=1}^d \lambda_i v_i^2 \right| \\ &= \max\{|\lambda_1|, \dots, |\lambda_d|\} \max_{\mathbf{v} \in \mathbb{S}^{d-1}} \left| \sum_{i=1}^d v_i^2 \right| \\ &= \max\{|\lambda_1|, \dots, |\lambda_d|\}. \end{aligned}$$

The second equality above holds because the mapping $\mathbf{u} \mapsto \mathbf{Q} \mathbf{u}$ maps \mathbb{S}^{d-1} onto \mathbb{S}^{d-1} . Note that equality is achieved when \mathbf{v} is the canonical vector corresponding to the largest eigenvalue in absolute value. This proves the theorem.

The corresponding Huber's contamination model for \mathcal{M}^{cov} is

$$H_\varepsilon^{\text{cov}} = \left\{ (1 - \varepsilon) N_d(\mathbf{0}, \Sigma) + \varepsilon Q \mid \Sigma \in \text{PD}_d, \sigma_{\max}(\Sigma) \leq M, Q \in P(\mathbb{R}^d) \right\}.$$

In contrast to the previous section, here we use P_Σ to represent the probability distribution $N_d(\mathbf{0}, \Sigma)$, and $P_{(\Sigma, Q)}$ to denote the contaminated probability distribution $(1 - \varepsilon)N_d(\mathbf{0}, \Sigma) + \varepsilon Q$. The contamination level ε is not explicitly stated, as it will be clear from the context, just like before. Recall that we denote by $hD^{\text{sc}}(\cdot; \cdot)$ the scatter halfspace depth with the fixed location $\mathbf{0}$, that is

$$hD^{\text{sc}}(\Gamma; P) = \inf_{\mathbf{u} \in \mathbb{S}^{d-1}} \min \left\{ P(H_{\mathbf{u}, \Gamma}^{\text{in}}), P(H_{\mathbf{u}, \Gamma}^{\text{out}}) \right\}$$

for any $\Gamma \in \text{PD}_d$ and $P \in P(\mathbb{R}^d)$, where $H_{\mathbf{u}, \Gamma}^{\text{in}}$ and $H_{\mathbf{u}, \Gamma}^{\text{out}}$ are defined in (2.3). Estimators of Σ are denoted by $\hat{\Sigma}_n$. The deepest matrix is denoted by $\hat{\Sigma}_n^{\text{hs}}$, i.e.

$$\hat{\Sigma}_n^{\text{hs}} = \arg \max_{\Sigma \in \text{PD}_d} hD^{\text{sc}}(\Gamma; \hat{P}_n)$$

and by $\widehat{\Sigma}_n^{\text{cov}}$ we denote the estimator of covariance matrix

$$\widehat{\Sigma}_n^{\text{cov}} = \frac{1}{\beta} \widehat{\Sigma}_n^{\text{sc}},$$

where $\beta > 0$ solves $\Phi(\bar{\beta}) = 3/4$, see Example 6.

First, we focus on establishing a lower bound for the minimax optimal convergence rate in the model \mathcal{M}^{cov} .

Theorem 62 *In model \mathcal{M}^{cov} , there exists an absolute constant $\delta > 0$ such that*

$$\inf_n \inf_{\widehat{\Sigma}_n} \sup_{P_{\Sigma} \in \mathcal{M}^{\text{cov}}} P_{\Sigma} \left(\left\| \widehat{\Sigma}_n - \Sigma \right\|_{\text{op}}^2 \geq R_n^0 \right) \geq \delta$$

holds for some $\{R_n^0\} \asymp d/n$. The infimum is taken over all estimators $\widehat{\Sigma}_n$ of Σ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{\Sigma}$.

Proof. This claim follows from the proof of Theorem 6 and Remark 2 in [Ma and Wu, 2015]. Note that the spectral norm $\|\cdot\|_{\text{op}}$ is a unitarily invariant norm in terms of this paper, which is assumed by that theorem.

As in the section on the optimality of the location halfspace median, we can use this to determine a lower bound for the minimax optimal convergence rate in model $\mathcal{H}_{\varepsilon}^{\text{cov}}$.

Theorem 63 *There exists an absolute constant $\delta > 0$ such that for any $\varepsilon \in [0, 1)$*

$$\inf_n \inf_{\widehat{\Sigma}_n} \sup_{P_{(\Sigma, Q)} \in \mathcal{H}_{\varepsilon}^{\text{cov}}} P_{(\Sigma, Q)} \left(\left\| \widehat{\Sigma}_n - \Sigma \right\|_{\text{op}}^2 \geq R_n^{\varepsilon} \right) \geq \delta$$

holds for some $\{R_n^{\varepsilon}\} \asymp (d/n) \varepsilon^2$. The infimum is taken over all estimators $\widehat{\Sigma}_n$ of Σ based on a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n \sim P_{(\Sigma, Q)}$.

Proof. Same as in the proof of Theorem 58, it suffices to bound the modulus of continuity $\omega(\varepsilon, \Theta)$ from below by ε^2 in order to use Theorem 55 and Theorem 62.

Without loss of generality, assume $M \geq 1 + \varepsilon$. As in the proof of [Chen et al., 2018, Theorem 3.2], consider two matrices $\Sigma_a = \mathbf{I}_d$ and $\Sigma_b = \mathbf{I}_d + \varepsilon \mathbf{M}$, where \mathbf{M} is a matrix with the only non-zero entry $M_{(1,1)} = 1$. Then, obviously

$$\Sigma_a - \Sigma_b \Big|_{\text{op}}^2 = \varepsilon \mathbf{M} \Big|_{\text{op}}^2 = \varepsilon^2.$$

Using Pinsker's inequality given in Theorem 51 and the Kullback-Leibler di-

vergence of two Gaussian distributions (see Theorem 52), we can bound

$$\begin{aligned}
\text{TV}(P_{\Sigma_a}, P_{\Sigma_b})^2 &= \frac{1}{2} D(P_{\Sigma_a} // P_{\Sigma_b}) \\
&= \frac{1}{4} \left(\log \left(\frac{|\Sigma_b|}{|\Sigma_a|} \right) + \text{tr}(\Sigma_b^{-1} \Sigma_a) - d \right) \\
&= \frac{1}{4} \left(\log(1 + \varepsilon) + \frac{1}{1 + \varepsilon} + (d - 1) - d \right) \\
&= \frac{1}{4} \left(\log(1 + \varepsilon) - \frac{\varepsilon}{1 + \varepsilon} \right) \\
&= \frac{1}{4} \left(\varepsilon - \frac{\varepsilon}{1 + \varepsilon} \right) \\
&= \frac{\varepsilon^2}{4(1 + \varepsilon)} \left(\frac{\varepsilon}{1 - \varepsilon} \right)^2.
\end{aligned}$$

In the second inequality, we used the fact that $\log(1 + x) \geq x$ for $x > 0$, and in the last inequality, we used the assumption $\varepsilon \in [0, 1)$. Therefore

$$\text{TV}(P_{\Sigma_a}, P_{\Sigma_b}) \leq \varepsilon / (1 - \varepsilon).$$

As a consequence, the modulus of continuity can be bounded from below as follows

$$\omega(\varepsilon, \Theta) = \sup \left\{ \|\Sigma_1 - \Sigma_2\|_{\text{op}}^2 / \text{TV}(P_{\Sigma_1}, P_{\Sigma_2}) \leq \frac{\varepsilon}{1 - \varepsilon}; \Sigma_1, \Sigma_2 \in \Theta \right\}$$

$$\|\Sigma_a - \Sigma_b\|_{\text{op}} = \varepsilon^2.$$

This concludes the proof.

Now, we show that the convergence rate $(d/n) \leq \varepsilon^2$ in model $H_\varepsilon^{\text{cov}}$ is indeed achieved for the scatter halfspace median.

Theorem 64 *Assume that $\varepsilon \in (0, 1/5)$. Then, for any $\delta \in (0, 1/2)$ there exist absolute constants $C, D > 0$ such that for all $n \in \mathbb{N}$ satisfying*

$$\left(\frac{d}{n} + \frac{\log(1/\delta)}{n} \right) < D \tag{3.36}$$

we have

$$\inf_{P_{(\Sigma, Q)} \in H_\varepsilon^{\text{loc}}} P_{(\Sigma, Q)} \left(\left\| \widehat{\Sigma}_n^{\text{cov}} - \Sigma \right\|_{\text{op}}^2 < C \left(\left(\frac{d}{n} + \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right) \geq 1 - 2\delta,$$

hence

$$\sup_{P_{(\Sigma, Q)} \in H_\varepsilon^{\text{loc}}} P_{(\Sigma, Q)} \left(\left\| \widehat{\Sigma}_n^{\text{cov}} - \Sigma \right\|_{\text{op}}^2 \geq C \left(\left(\frac{d}{n} + \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right) \right) < 2\delta.$$

Proof. This proof is largely inspired by [Chen et al., 2018, Theorem 3.1]. It closely resembles the proof of Theorem 59, offering more precise details and clarifying some ambiguities present in [Chen et al., 2018]. As before, the proof is structured into four parts.

Part 1: Auxiliary observations

First, we make the following observations, which will be useful in deriving the intended bound.

(S1) By Example 6, for any $\Sigma, \Gamma \in \text{PD}_d$ and $\mathbf{u} \in \mathbb{S}^{d-1}$ we have

$$P_{\Sigma}(H_{\mathbf{u},\Gamma}^{\text{in}}) = 1 - P_{\Sigma}(H_{\mathbf{u},\Gamma}^{\text{out}}) = 2 \left(\Phi \left(\sqrt{\frac{\mathbf{u}^{\top} \Gamma \mathbf{u}}{\mathbf{u}^{\top} \Sigma \mathbf{u}}} \right) - \frac{1}{2} \right)$$

and

$$\max_{\Gamma \in \text{PD}_d} hD^{\text{sc}}(\Gamma; P_{\Sigma}) = hD^{\text{sc}}(\beta \Sigma; P_{\Sigma}) = \frac{1}{2}.$$

(S2) Given a random sample $\{\mathbf{X}_i\}_{i=1}^n$ from a contaminated distribution $P_{(\Sigma, Q)}$, we can decompose $\{\mathbf{X}_i\}_{i=1}^n = \{\mathbf{Y}_i\}_{i=1}^{n_1} \cup \{\mathbf{Z}_i\}_{i=1}^{n_2}$. Here, marginally, it holds that $n_2 \sim \text{Binomial}(n, \varepsilon)$ and $n_1 = n - n_2$. Conditionally on n_1 and n_2 , $\{\mathbf{Y}_i\}_{i=1}^{n_1}$ is a random sample from P_{Σ} , while $\{\mathbf{Z}_i\}_{i=1}^{n_2}$ is a random sample from Q .

(S3) By Theorem A3 in the Appendix, we have with probability at least $1 - \delta$ that

$$\begin{aligned} & \sup_{\Gamma \in \text{PD}_d} |hD^{\text{sc}}(\Gamma; P_{\Sigma}) - hD^{\text{sc}}(\Gamma; \{\mathbf{Y}_i\}_{i=1}^{n_1})| \\ & \sup_{\Gamma \in \text{PD}_d, \mathbf{u} \in \mathbb{S}^{d-1}} \max \left\{ \left| P_{\Sigma}(H_{\mathbf{u},\Gamma}^{\text{in}}) - \widehat{P}_{n_1}(H_{\mathbf{u},\Gamma}^{\text{in}}) \right|, \left| P_{\Sigma}(H_{\mathbf{u},\Gamma}^{\text{out}}) - \widehat{P}_{n_1}(H_{\mathbf{u},\Gamma}^{\text{out}}) \right| \right\} \\ & \sup_{\mathbf{u} \in \mathbb{S}^{d-1}, t \geq 0} \max \left\{ \left| P_{\Sigma}(H_{\mathbf{u},t}^{\text{in}}) - \widehat{P}_{n_1}(H_{\mathbf{u},t}^{\text{in}}) \right|, \left| P_{\Sigma}(H_{\mathbf{u},t}^{\text{out}}) - \widehat{P}_{n_1}(H_{\mathbf{u},t}^{\text{out}}) \right| \right\} \\ & \sqrt{\frac{1440\pi e}{1 - e^{-1}}} \sqrt{\frac{2d + 3}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}}, \end{aligned}$$

where \widehat{P}_{n_1} is the empirical distribution of $\{\mathbf{Y}_i\}_{i=1}^{n_1}$.

(S4) By the definition of the sample scatter halfspace depth, it follows that for any $\Gamma \in \text{PD}_d$, we have

$$n_1 hD^{\text{sc}}(\Gamma; \{\mathbf{Y}_i\}_{i=1}^{n_1}) \leq n hD^{\text{sc}}(\Gamma; \{\mathbf{X}_i\}_{i=1}^n) - n_2 \leq n_1 hD^{\text{sc}}(\Gamma; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - n_2.$$

This can be proven exactly as (3.25) in the proof of Theorem 59.

(S5) By Theorem A4 in the Appendix, if

$$\sqrt{\frac{\log(1/\delta)}{2n}} < 1/5 \tag{3.37}$$

holds for $\delta > 0$, then

$$\frac{n_2}{n_1} \leq \frac{\varepsilon}{1 - \varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} < \frac{2}{3}$$

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

Part 2: Choice of constant D

To ensure the proof of the statement is valid, we need to define a constant D . Specifically, it needs to satisfy both (3.37) and the following condition

$$20\sqrt{\frac{30\pi e}{1-e^{-1}}}\sqrt{\frac{d}{n}} + 2\sqrt{\frac{\log(1/\delta)}{n}} < \frac{1}{16}, \quad (3.38)$$

which will be useful later. For the remainder of the proof, fix $\delta \in (0, 1/2)$ and consider $n \in \mathbb{N}$ such that the assumption (3.36) holds with

$$D = \left(\frac{\frac{1}{16}}{20\sqrt{\frac{30\pi e}{1-e^{-1}}} + 2} \right)^2 > 0.$$

Since $a + b < D$ implies $a - b < D$ for all $a, b > 0$, it follows that

$$\sqrt{\frac{d}{n}} - \sqrt{\frac{\log(1/\delta)}{n}} < \bar{D}. \quad (3.39)$$

Therefore, by our choice of D and using (3.39), we can bound

$$20\sqrt{\frac{30\pi e}{1-e^{-1}}}\sqrt{\frac{d}{n}} + 2\sqrt{\frac{\log(1/\delta)}{n}} < \left(20\sqrt{\frac{30\pi e}{1-e^{-1}}} + 2 \right) \bar{D} = \frac{1}{16}.$$

Consequently, inequality (3.38) holds. Additionally, (3.38) implies that

$$\sqrt{\frac{\log(1/\delta)}{2n}} < 2\sqrt{\frac{\log(1/\delta)}{n}} < \frac{1}{16} < \frac{1}{5},$$

thus, under the assumption $\varepsilon < 1/5$, (S5) gives that

$$\frac{n_2}{n_1} - \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12}\sqrt{\frac{\log(1/\delta)}{2n}} < \frac{2}{3} \quad (3.40)$$

holds with probability at least $1 - \delta$. Furthermore, note that

$$\frac{n_2}{n_1} < \frac{2}{3} \quad 3n_2 < 2n_1 \quad 3(n - n_1) < 2n_1 \quad n_1 > 3n/5 \quad (3.41)$$

holds, much like before in the proof of Theorem 59. To summarize, for our fixed $\delta > 0$ and a sufficiently large $n \in \mathbb{N}$, at least $3/5$ of all observations are non-contaminating with probability at least $1 - \delta$. Moreover, inequality (3.38) holds. Both of these facts will be used later in the proof.

Part 3: Bound for the depth of $\hat{\Sigma}_n^{\text{hs}}$

Based on the earlier observations and our selection of D , we establish the following sequence of inequalities. These are valid with a probability of at least $1 - \delta$ given

the marginal counts n_1, n_2 for all $\Sigma \in \Theta$. We have

$$\begin{aligned}
hD^{\text{sc}}(\widehat{\Sigma}_n^{\text{hs}}; P_{\Sigma}) &\stackrel{\text{(S3)}}{=} hD^{\text{sc}}(\widehat{\Sigma}_n^{\text{hs}}; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\
&\stackrel{\text{(S4)}}{\leq} \frac{n}{n_1} hD^{\text{sc}}(\widehat{\Sigma}_n^{\text{hs}}; \{\mathbf{X}_i\}_{i=1}^n) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\
&\quad - \frac{n}{n_1} hD^{\text{sc}}(\beta\Sigma; \{\mathbf{X}_i\}_{i=1}^n) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\
&\stackrel{\text{(S4)}}{\leq} hD^{\text{sc}}(\beta\Sigma; \{\mathbf{Y}_i\}_{i=1}^{n_1}) - \frac{n_2}{n_1} - \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{\log(1/\delta)}{2n_1}} \\
&\stackrel{\text{(S3)}}{\leq} hD^{\text{sc}}(\beta\Sigma; P_{\Sigma}) - \frac{n_2}{n_1} - 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{2\log(1/\delta)}{n_1}} \\
&\stackrel{\text{(S1)}}{=} \frac{1}{2} - \frac{n_2}{n_1} - 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} - \sqrt{\frac{2\log(1/\delta)}{n_1}}. \tag{3.42}
\end{aligned}$$

The third inequality follows because $\widehat{\Sigma}_n^{\text{hs}}$ is the maximizer of $hD^{\text{sc}}(\cdot; \{\mathbf{X}_i\}_{i=1}^n)$. Rewriting (3.42), we have that

$$\mathbb{P} \left[\frac{1}{2} - hD^{\text{sc}}(\widehat{\Sigma}_n^{\text{hs}}; P_{\Sigma}) \leq \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \mid n_1, n_2 \right] \geq 1 - \delta.$$

However, by taking the expectation on both sides (and considering its monotonicity), we obtain

$$\mathbb{P} \left[\frac{1}{2} - hD^{\text{sc}}(\widehat{\Sigma}_n^{\text{hs}}; P_{\Sigma}) \leq \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \right] \geq 1 - \delta.$$

Using (S1), this means that uniformly over all $\mathbf{u} \in \mathbb{S}^{d-1}$ the inequality

$$\begin{aligned}
\frac{1}{2} - \min \left\{ P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right), 1 - P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right) \right\} \\
\leq \frac{n_2}{n_1} + 2\sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{2\log(1/\delta)}{n_1}} \tag{3.43}
\end{aligned}$$

holds with probability at least $1 - \delta$. Note that for all $x \in \mathbb{R}$ the equality $1/2 - \min \{x, 1 - x\} = |1/2 - x|$ holds. Therefore, for all $\mathbf{u} \in \mathbb{S}^{d-1}$ we have

$$\begin{aligned}
\frac{1}{2} - \min \left\{ P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right), 1 - P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right) \right\} &= \left| \frac{1}{2} - P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right) \right| \\
&\stackrel{\text{(S1)}}{=} \left| P_{\Sigma}(H_{\mathbf{u}, \beta\Sigma}^{\text{in}}) - P_{\Sigma} \left(H_{\mathbf{u}, \widehat{\Sigma}_n^{\text{hs}}}^{\text{in}} \right) \right| \\
&\stackrel{\text{(S1)}}{=} 2 \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^{\text{T}} \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^{\text{T}} \Sigma \mathbf{u}}} \right) \right|. \tag{3.44}
\end{aligned}$$

Combining (3.43) with (3.44), we have that

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^{\text{T}} \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^{\text{T}} \Sigma \mathbf{u}}} \right) \right| \leq \frac{n_2}{2n_1} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}}$$

holds with probability at least $1 - \delta$. Now, we combine this result with inequality (3.40). Same as in the proof of Theorem 59, we obtain

$$\mathbb{P} \left[\sup_{\mathbf{u} \in S^{d-1}} \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}} \right) \right| \leq \frac{n_2}{2n_1} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}}, \right. \\ \left. \frac{n_2}{n_1} \leq \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} \leq \frac{2}{3} \right] \geq 1 - 2\delta. \quad (3.45)$$

Under the condition of the second random event in (3.45), we can further upper bound

$$\begin{aligned} & \sup_{\mathbf{u} \in S^{d-1}} \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}} \right) \right| \stackrel{(3.45)}{\leq} \frac{n_2}{2n_1} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \stackrel{(3.40)}{\leq} \frac{\varepsilon}{2(1-\varepsilon)} + \frac{25}{24} \sqrt{\frac{\log(1/\delta)}{2n}} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{2d+3}{n_1}} + \sqrt{\frac{\log(1/\delta)}{2n_1}} \\ & \stackrel{(3.41)}{\leq} \frac{\varepsilon}{2(1-\varepsilon)} + \frac{25}{24} \sqrt{\frac{\log(1/\delta)}{2n}} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{5(2d+3)}{3n}} + \sqrt{\frac{5 \log(1/\delta)}{6n}} \\ & \stackrel{d \geq 1}{\leq} \frac{\varepsilon}{2(1-\varepsilon)} + \frac{25}{24} \sqrt{\frac{\log(1/\delta)}{2n}} + \sqrt{\frac{1440\pi e}{1-e^{-1}}} \sqrt{\frac{25d}{3n}} + \sqrt{\frac{5 \log(1/\delta)}{6n}} \\ & = \frac{\varepsilon}{2(1-\varepsilon)} + 20 \sqrt{\frac{30\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + \left(\frac{25}{24} \sqrt{\frac{1}{2}} + \sqrt{\frac{5}{6}} \right) \sqrt{\frac{\log(1/\delta)}{n}} \\ & < \frac{\varepsilon}{2(1-\varepsilon)} + 20 \sqrt{\frac{30\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + 2 \sqrt{\frac{\log(1/\delta)}{n}} \\ & \stackrel{\varepsilon < 1/5}{<} \frac{5}{8} \varepsilon + 20 \sqrt{\frac{30\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + 2 \sqrt{\frac{\log(1/\delta)}{n}}. \end{aligned}$$

In the first strict inequality, we used the fact that

$$\left(\frac{25}{24} \sqrt{\frac{1}{2}} + \sqrt{\frac{5}{6}} \right) < 2.$$

Ultimately, we have

$$\mathbb{P} \left[\sup_{\mathbf{u} \in S^{d-1}} \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}} \right) \right| \leq \frac{5}{8} \varepsilon + 20 \sqrt{\frac{30\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + 2 \sqrt{\frac{\log(1/\delta)}{n}}, \right. \\ \left. \frac{n_2}{n_1} \leq \frac{\varepsilon}{1-\varepsilon} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} \leq \frac{2}{3} \right] \geq 1 - 2\delta.$$

This further implies (by the same argument as in the proof of Theorem 59) that

$$\sup_{\mathbf{u} \in S^{d-1}} \left| \Phi \left(\sqrt{\beta} \right) - \Phi \left(\sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}} \right) \right| \leq \frac{5}{8} \varepsilon + 20 \sqrt{\frac{30\pi e}{1-e^{-1}}} \sqrt{\frac{d}{n}} + 2 \sqrt{\frac{\log(1/\delta)}{n}} \quad (3.46)$$

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

Part 4: Conclusion

We assumed that n is chosen such that (3.38) holds. Therefore, using this and the assumption $\varepsilon < 1/5$, from (3.46) we can deduce that for all $\Sigma \in \Theta$ it holds that

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \Phi\left(\sqrt{\beta}\right) - \Phi\left(\sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}}\right) \right| < \frac{5}{8}\varepsilon + \frac{1}{16} < \frac{1}{8} + \frac{1}{16} = \frac{3}{16}. \quad (3.47)$$

Now we can use the fact that $\Phi(\bar{\beta}) = 3/4$ and that the quantile function of the standard Gaussian distribution is Lipschitz continuous on the interval

$$\left[\frac{3}{4} - \frac{3}{16}, \frac{3}{4} + \frac{3}{16}\right] = \left[\frac{9}{16}, \frac{15}{16}\right]$$

with the Lipschitz constant

$$C = \left[\varphi\left(\Phi^{-1}(15/16)\right)\right]^{-1} > 0.$$

Therefore, from (3.46) we can further deduce

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \sqrt{\beta} - \sqrt{\frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}}} \right| < C \left(\frac{5}{8}\varepsilon + 20\sqrt{\frac{30\pi e}{1-e^{-1}}}\sqrt{\frac{d}{n}} + 2\sqrt{\frac{\log(1/\delta)}{n}} \right). \quad (3.48)$$

In contrast to the proof of Theorem 59, we need to apply the trick involving Lipschitz continuity once again. Specifically, we observe that the upper bound of the right side of (3.48) can be constrained similarly to (3.47) by $3C/16 > 0$. The function $x \mapsto x^2$ is Lipschitz continuous on the interval

$$\left[\sqrt{\beta} - 3C/16, \sqrt{\beta} + 3C/16\right]$$

with the Lipschitz constant

$$C = 2\left(\sqrt{\beta} + 3C/16\right) > 0.$$

As a consequence, with probability at least $1 - 2\delta$, we have that

$$\sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \beta - \frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}} \right| < C C \left(\frac{5}{8}\varepsilon + 20\sqrt{\frac{30\pi e}{1-e^{-1}}}\sqrt{\frac{d}{n}} + 2\sqrt{\frac{\log(1/\delta)}{n}} \right). \quad (3.49)$$

We can finally establish the required bound by summarizing all the arguments

above. With probability at least $1 - \delta$, for all $\Sigma \in \Theta$ it holds that

$$\begin{aligned}
& \left\| \widehat{\Sigma}_n^{\text{cov}} - \Sigma \right\|_{\text{op}}^2 \stackrel{\text{Thm. 61}}{=} \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \mathbf{u}^\top (\Sigma - \widehat{\Sigma}_n^{\text{cov}}) \mathbf{u} \right|^2 \\
&= \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \left| \mathbf{u}^\top \Sigma \mathbf{u} - \mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} / \beta \mathbf{u} \right|^2 = \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{M^2}{\beta^2} \left| \frac{\mathbf{u}^\top \beta \Sigma \mathbf{u} - \mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{M} \right|^2 \\
&\stackrel{\Sigma}{=} \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{M^2}{\beta^2} \left| \frac{\mathbf{u}^\top \beta \Sigma \mathbf{u} - \mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}} \right|^2 = \sup_{\mathbf{u} \in \mathbb{S}^{d-1}} \frac{M^2}{\beta^2} \left| \beta - \frac{\mathbf{u}^\top \widehat{\Sigma}_n^{\text{hs}} \mathbf{u}}{\mathbf{u}^\top \Sigma \mathbf{u}} \right|^2 \\
(3.49) \quad & \left(\frac{M C C}{\beta} \right)^2 \left(\frac{5}{8} \varepsilon + 20 \sqrt{\frac{30\pi e}{1 - e^{-1}}} \sqrt{\frac{d}{n}} + 2 \sqrt{\frac{\log(1/\delta)}{n}} \right)^2 \\
& \left(\frac{M C C}{\beta} \right)^2 \left(\left(\frac{5}{8} + 20 \sqrt{\frac{30\pi e}{1 - e^{-1}}} \right) \left(\sqrt{\frac{d}{n}} + \varepsilon \right) + 2 \sqrt{\frac{\log(1/\delta)}{n}} \right)^2 \\
(3.3) \quad & 2 \left(\frac{M C C}{\beta} \right)^2 \left(\left(\frac{5}{8} + 20 \sqrt{\frac{30\pi e}{1 - e^{-1}}} \right)^2 \left(\frac{d}{n} + \varepsilon^2 \right) + \frac{4 \log(1/\delta)}{n} \right) \\
& C \left(\left(\frac{d}{n} + \varepsilon^2 \right) + \frac{\log(1/\delta)}{n} \right),
\end{aligned}$$

where

$$C = 2 \left(\frac{M C C}{\beta} \right)^2 \left(\left(\frac{5}{8} + 20 \sqrt{\frac{30\pi e}{1 - e^{-1}}} \right)^2 + 4 \right) > 0.$$

In the first inequality, we used the Cauchy-Schwarz inequality and the assumption that $\sigma_{\max}(\Sigma) \leq M$. This concludes the proof.

Minimax optimality of the scatter halfspace median follows immediately.

Theorem 65 *For $\varepsilon < 1/5$, the scatter halfspace median $\widehat{\Sigma}_n^{\text{cov}}$ is the minimax optimal estimator of Σ in terms of the squared operator loss in Huber's contamination model $H_\varepsilon^{\text{cov}}$. The minimax optimal convergence rate is $d/n + \varepsilon^2$.*

Proof. This claim follows from Theorem 64 utilizing precisely the same arguments used in the proof of Theorem 60.

We have established that the scatter halfspace median is the minimax optimal estimator of the covariance matrix when the observations $\mathbf{X}_1, \dots, \mathbf{X}_n$ are generated from a centered Gaussian distribution. Chen et al. [2018] provide two extensions. First, Theorem 65 can be extended to elliptically symmetric distributions, as shown in [Chen et al., 2018, Theorem 4.1]. Secondly, the assumption of a centered generating distribution can also be relaxed. The key idea is that if $\mathbf{X}_1, \dots, \mathbf{X}_n \sim N_d(\boldsymbol{\mu}, \Sigma)$, then for all $i, j \in \{1, \dots, n\}, i \neq j$,

$$Y_{ij} = \frac{1}{2} (\mathbf{X}_i - \mathbf{X}_j) \sim N_d(\mathbf{0}, \Sigma),$$

as follows from the basic properties of the Gaussian distribution. This motivates the use of the U-version of scatter depth. That is, we compute

$$\hat{\Sigma}_n^{\text{hs,U}} = \arg \max_{\mathbf{\Gamma} \text{ PD}_d} hD^{\text{sc}}(\mathbf{\Gamma}; \{Y_{ij}\}_{i < j})$$

and define the estimator of Σ as

$$\hat{\Sigma}_n^{\text{cov,U}} = \frac{1}{\beta} \hat{\Sigma}_n^{\text{hs,U}}.$$

It turns out that the optimality results also hold for these U-estimators. For details, refer to Section A of the supplementary material in [Chen et al., 2018].

Conclusion

In this thesis, we examined the fundamental properties of location halfspace depth (Chapter 1) and scatter halfspace depth (Chapter 2), with an emphasis on the robustness of the corresponding halfspace medians. In Chapter 3, we introduced the minimax optimality of location and scatter halfspace medians. This allowed us to meet the objectives of the thesis.

We utilized tools from mathematical statistics, probability theory, measure theory, empirical processes, geometry, and combinatorics. The primary contribution of this thesis is the establishment of a unified framework for analyzing the convergence rates of estimators and the minimax optimality of estimators. We also identified and corrected several deficiencies in the literature, such as Theorem 23 (see Example 4), and Theorem 57. The proofs presented in this thesis include numerous additional details and clarify various ambiguities in the proofs from the original papers.

Future work could build on this thesis in several directions. For example, it would be beneficial to determine other measures of robustness for the scatter halfspace median, such as the breakdown point. Additionally, it remains an open question whether location and scatter halfspace medians are minimax optimal in Huber's contamination model when the contaminated distribution is not Gaussian or elliptically symmetric. Another interesting question is whether the scatter halfspace median would still be minimax optimal if we also had to non-trivially estimate the location parameter.

Bibliography

- Patrick Billingsley. *Probability and measure*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons, Inc., New York, third edition, 1995. A Wiley-Interscience Publication.
- Tomáš Brabenec. Hloubka variančních matic. Master's thesis, Charles University, Faculty of Mathematics and Physics, Prague, 2021.
- Mengjie Chen, Chao Gao, and Zhao Ren. Robust covariance and scatter matrix estimation under Huber's contamination model. *Ann. Statist.*, 46(5):1932–1960, 2018.
- Zhiqiang Chen and David E. Tyler. The influence function and maximum bias of Tukey's median. *Ann. Statist.*, 30(6):1737–1759, 2002.
- Zhiqiang Chen and David E. Tyler. On the behavior of Tukey's depth and median under symmetric stable distributions. *J. Statist. Plann. Inference*, 122(1-2):111–124, 2004.
- David L. Donoho and Miriam Gasko. Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *Ann. Statist.*, 20(4):1803–1827, 1992.
- David L. Donoho and Richard C. Liu. Geometrizing rates of convergence. II, III. *Ann. Statist.*, 19(2):633–667, 668–701, 1991.
- Thomas P. Hettmansperger and Hannu Oja. Affine invariant multivariate multi-sample sign tests. *J. Roy. Statist. Soc. Ser. B*, 56(1):235–249, 1994.
- Peter J. Huber and Elvezio M. Ronchetti. *Robust statistics*. Wiley Series in Probability and Statistics. John Wiley & Sons, Inc., Hoboken, NJ, second edition, 2009.
- Rebecka Jörnsten. Clustering and classification based on the L_1 data depth. *J. Multivariate Anal.*, 90(1):67–89, 2004.
- Petra Laketa and Stanislav Nagy. Halfspace depth for general measures: the ray basis theorem and its consequences. *Statist. Papers*, 63(3):849–883, 2022.
- Gaëtan Louvet and Germain Van Bever. The influence function of scatter halfspace depth. In *Recent Advances in Econometrics and Statistics*. Springer Nature Switzerland, 2024.
- Zongming Ma and Yihong Wu. Volume ratio, sparsity, and minimaxity under unitarily invariant norms. *IEEE Trans. Inform. Theory*, 61(12):6939–6956, 2015.
- Karl Mosler and Pavlo Mozharovskyi. Choosing among notions of multivariate depth statistics. *Statist. Sci.*, 37(3):348–368, 2022.
- Stanislav Nagy, Carsten Schütt, and Elisabeth M. Werner. Halfspace depth and floating body. *Stat. Surv.*, 13:52–118, 2019.

- Davy Paindaveine and Germain Van Bever. Halfspace depths for scatter, concentration and shape matrices. *Ann. Statist.*, 46(6B):3276–3307, 2018.
- Peter J. Rousseeuw and Ida Ruts. The depth function of a population distribution. *Metrika*, 49(3):213–244, 1999.
- Jonathan Scarlett and Volkan Cevher. *An Introductory Guide to Fano's Inequality with Applications in Statistical Estimation*, pages 487–528. Cambridge University Press, Cambridge, 2021.
- Robert Serfling. Depth functions in nonparametric multivariate inference. In *Data depth: robust multivariate analysis, computational geometry and applications*, volume 72 of *DIMACS Ser. Discrete Math. Theoret. Comput. Sci.*, pages 1–16. Amer. Math. Soc., Providence, RI, 2006.
- Alexandre B. Tsybakov. *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York, 2009.
- John W. Tukey. Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians (Vancouver, B.C., 1974)*, Vol. 2, pages 523–531. Canad. Math. Congr., Montreal, QC, 1975.
- Roman Vershynin. *High-dimensional probability*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2018.
- Bin Yu. Assouad, Fano, and Le Cam. In *Festschrift for Lucien Le Cam*, pages 423–435. Springer, New York, 1997.
- Yijun Zuo and Robert Serfling. General notions of statistical depth function. *Ann. Statist.*, 28(2):461–482, 2000a.
- Yijun Zuo and Robert Serfling. On the performance of some robust nonparametric location measures relative to a general notion of multivariate symmetry. *J. Statist. Plann. Inference*, 84(1-2):55–79, 2000b.

Auxiliary theorems

Theorem A1 (Lebesgue's dominated convergence theorem) *Let (X, \mathcal{A}, μ) be a space with measure. Consider measurable functions $\{f_n\}_{n=1}, f, g: X \rightarrow \mathbb{R}$. Assume that for all $n \in \mathbb{N}$ we have $|f_n| \leq g$ μ -a.e., where g is integrable. If $f_n \rightarrow f$ μ -a.e. as $n \rightarrow \infty$, then $\{f_n\}_{n=1}, f$ are integrable and*

$$\int_X f_n(x) d\mu(x) \rightarrow \int_X f(x) d\mu(x) \quad \text{as } n \rightarrow \infty.$$

Proof. Refer to [Billingsley, 1995, Theorem 16.4].

Theorem A2 (Rate of convergence for halfspaces) *Let $P = P(\mathbb{R}^d)$ and consider the empirical distribution \hat{P}_n based on a random sample of size n from P . Then for all $\delta \in (0, 1/2)$ the inequality*

$$\sup_{H \in H_d} |P(H) - \hat{P}_n(H)| \leq \sqrt{\frac{1440\pi e}{1 - e^{-1}}} \sqrt{\frac{d+1}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

holds with probability at least $1 - \delta$, where by H_d we denote the system of all closed halfspaces in \mathbb{R}^d , see (3.24).

Proof. This can be proven analogously to [Chen et al., 2018, Lemma 7.3] using that the VC dimension of H_d is $d + 1$.

Theorem A3 (Rate of convergence for slabs) *Let $P = P(\mathbb{R}^d)$ and consider the empirical distribution \hat{P}_n based on a random sample of size n from P . Then for all $\delta \in (0, 1/2)$ the inequality*

$$\sup_{\mathbf{u} \in S^{d-1}, t \geq 0} \max \left\{ |P(H_{\mathbf{u},t}^{\text{in}}) - \hat{P}_n(H_{\mathbf{u},t}^{\text{in}})|, |P(H_{\mathbf{u},t}^{\text{out}}) - \hat{P}_n(H_{\mathbf{u},t}^{\text{out}})| \right\} \leq \sqrt{\frac{1440\pi e}{1 - e^{-1}}} \sqrt{\frac{2d+3}{n}} + \sqrt{\frac{\log(1/\delta)}{2n}}$$

holds with probability at least $1 - \delta$. The sets $H_{\mathbf{u},t}^{\text{in}}$ and $H_{\mathbf{u},t}^{\text{out}}$ are defined in (2.3).

Proof. Refer to [Chen et al., 2018, Lemma 7.3].

Theorem A4 *Let $N \sim \text{Binomial}(n, p)$ and assume $p < 1/5$. Then, for every $\delta > 0$ satisfying $\sqrt{\frac{\log(1/\delta)}{2n}} < 1/5$, we have*

$$\frac{N}{n - N} \leq \frac{p}{1 - p} + \frac{25}{12} \sqrt{\frac{\log(1/\delta)}{2n}} < \frac{2}{3}$$

with probability at least $1 - \delta$.

Proof. Refer to [Chen et al., 2018, Lemma 7.1].

Theorem A5 Let $d \in \mathbb{N}$, $r > 0$ and $s > 1$. Consider $\{\mathbf{x}_1, \dots, \mathbf{x}_M\} \subset B_d(\mathbf{0}, r)$ a set of the maximum cardinality M such that

$$\|\mathbf{x}_i - \mathbf{x}_j\| \geq r/s$$

holds for all $i \neq j$. Then, $M \leq s^d$.

Proof. Refer to [Vershynin, 2018, Proposition 4.2.12].