

Robust estimation of a mean in a multivariate Gaussian model: Part 1



Frejus, December 17, 2018

Arnak S. Dalalyan
ENSAE ParisTech / CREST

1. Various models of contamination

General notation

We first introduce the notation that are common to all the models of contamination considered in this talk.

- Number of observations : n .
- Dimension of the unknown parameter μ^* : p .
- Observations $(\mathbf{X}_1, \dots, \mathbf{X}_n) \sim P_n$.
- Number of outliers (possibly random): $s \in \{1, \dots, n\}$.
- Set of outliers: $S \subset \{1, \dots, n\}$.
- Proportion of outliers: $\varepsilon = \mathbf{E}[s/n] = \mathbf{E}[|S|/n]$.

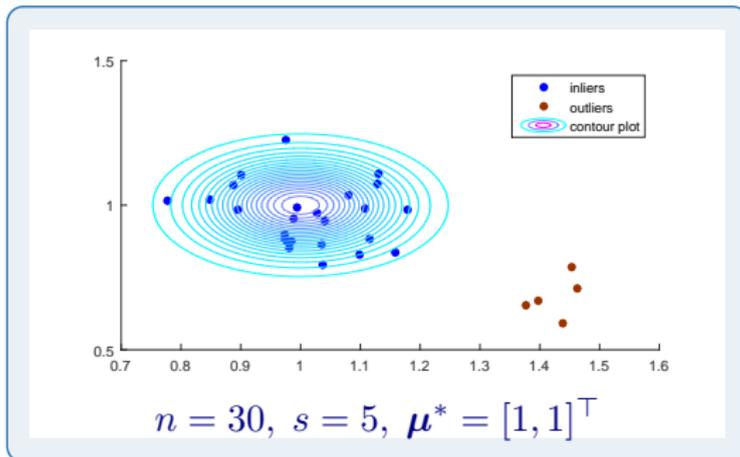
Setting (informal)

Among the n observations $\mathbf{X}_1, \dots, \mathbf{X}_n$, there is a small number s of outliers. If we remove the outliers, all the other \mathbf{X}_i 's are iid drawn from a reference distribution P_{μ^*} .

Gaussian model with unknown mean

Assumption (model for inliers)

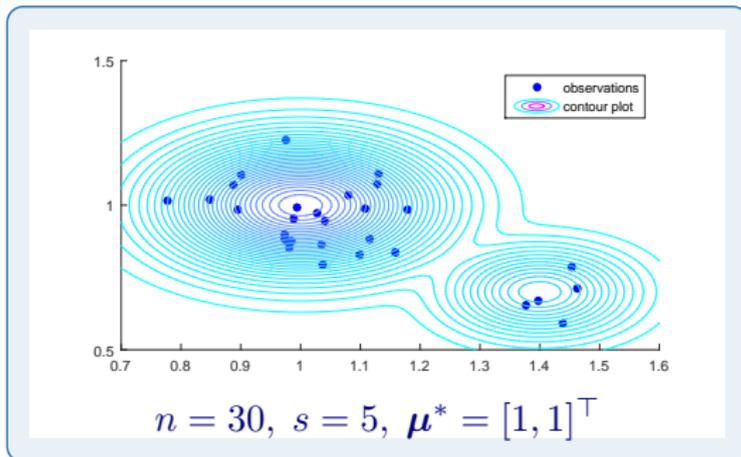
Throughout this presentation, we assume that the reference distribution P_{μ^*} is p -variate Gaussian $\mathcal{N}_p(\mu^*, \mathbf{I}_p)$. The goal is to estimate the parameter $\mu^* \in \mathbb{R}^p$.



Gaussian model with unknown mean

Assumption (model for inliers)

Throughout this presentation, we assume that the reference distribution P_{μ^*} is p -variate Gaussian $\mathcal{N}_p(\mu^*, \mathbf{I}_p)$. The goal is to estimate the parameter $\mu^* \in \mathbb{R}^p$.



Huber's contamination

Assumption (HC model for outliers)

There are unobserved iid random variables $Z_1, \dots, Z_n \sim \mathcal{B}(\varepsilon)$ and a distribution \mathbf{Q} , such that

$$\mathcal{L}(\mathbf{X}_i | Z_i = 0) = \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p), \quad \mathcal{L}(\mathbf{X}_i | Z_i = 1) = \mathbf{Q},$$

the observations \mathbf{X}_i corresponding to different i 's are independent.
This is equivalent to

$$\mathbf{P}_n = \{(1 - \varepsilon)\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p) + \varepsilon\mathbf{Q}\}^{\otimes n}.$$

In this model,

$$\underbrace{S = \{i : Z_i = 1\}}_{\text{set of outliers}} \text{ and } \underbrace{s \sim \mathcal{B}(n, \varepsilon)}_{\text{nb of outliers}}$$

are both random.

Huber's contamination

Assumption (HC model for outliers)

There are unobserved iid random variables $Z_1, \dots, Z_n \sim \mathcal{B}(\varepsilon)$ and a distribution \mathbf{Q} , such that

$$\mathcal{L}(\mathbf{X}_i | Z_i = 0) = \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p), \quad \mathcal{L}(\mathbf{X}_i | Z_i = 1) = \mathbf{Q},$$

the observations \mathbf{X}_i corresponding to different i 's are independent. This is equivalent to

$$\mathbf{P}_n = \{(1 - \varepsilon)\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p) + \varepsilon\mathbf{Q}\}^{\otimes n}.$$

In this model,

$$\underbrace{S = \{i : Z_i = 1\}}_{\text{set of outliers}} \text{ and } \underbrace{s \sim \mathcal{B}(n, \varepsilon)}_{\text{nb of outliers}}$$

are both random.

i	Z_i	$\mathbf{X}_i \sim$
1	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
2	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
3	1	\mathbf{Q}
4	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
5	1	\mathbf{Q}
6	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
7	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
\vdots	\vdots	\vdots
30	0	$\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$
$s =$	5	$\sim \mathcal{B}(30, 0.2)$

Huber's contamination

Assumption (HC model for outliers)

There are unobserved iid random variables $Z_1, \dots, Z_n \sim \mathcal{B}(\varepsilon)$ and a distribution Q , such that

$$\mathcal{L}(X_i | Z_i = 0) = \mathcal{N}_p(\mu^*, \mathbf{I}_p), \quad \mathcal{L}(X_i | Z_i = 1) = Q,$$

the observations X_i corresponding to different i 's are independent. This is equivalent to

$$P_n = \{(1 - \varepsilon)\mathcal{N}_p(\mu^*, \mathbf{I}_p) + \varepsilon Q\}^{\otimes n}.$$

In this model,

$$\underbrace{S = \{i : Z_i = 1\}}_{\text{set of outliers}} \text{ and } \underbrace{s \sim \mathcal{B}(n, \varepsilon)}_{\text{nb of outliers}}$$

are both random.

i	Z_i	$X_i \sim$
1	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
2	1	Q
3	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
4	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
5	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
6	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
7	1	Q
\vdots	\vdots	\vdots
30	0	$\mathcal{N}_p(\mu^*, \mathbf{I}_p)$
$s =$	6	$\sim \mathcal{B}(30, 0.2)$

Huber's contamination

Assumption (HC model for outliers)

There are unobserved iid random variables $Z_1, \dots, Z_n \sim \mathcal{B}(\varepsilon)$ and a distribution Q , such that

$$\mathcal{L}(X_i | Z_i = 0) = \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p), \quad \mathcal{L}(X_i | Z_i = 1) = Q,$$

the observations X_i corresponding to different i 's are independent. This is equivalent to

$$P_n = \{(1 - \varepsilon)\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p) + \varepsilon Q\}^{\otimes n}.$$

In this model,

$$\underbrace{S = \{i : Z_i = 1\}}_{\text{set of outliers}} \text{ and } \underbrace{s \sim \mathcal{B}(n, \varepsilon)}_{\text{nb of outliers}}$$

are both random.

We write

$$P_n \in \mathcal{M}_n^{\text{HC}}(p, \varepsilon, \boldsymbol{\mu}^*).$$

for the model of Huber's contamination.

Huber's deterministic contamination

Assumption (HDC model for outliers)

There is a set $S \subset \{1, \dots, n\}$ of cardinality $s = \lceil n\varepsilon \rceil$ and a distribution Q , such that

$$\{\mathbf{X}_i : i \in S^c\} \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p) \quad \perp\!\!\!\perp \quad \{\mathbf{X}_i : i \in S\} \stackrel{\text{iid}}{\sim} Q.$$

- Similar to HC: the outliers are iid.
- Different from HC: the set of outliers is deterministic.

Remark *The number of outliers s should be smaller than $n/2$, otherwise Q would be the reference distribution and $\mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$ the contamination.*

We write $P_n \in \mathcal{M}_n^{\text{HDC}}(p, \varepsilon, \boldsymbol{\mu}^*)$.

Assumption (PC model for outliers)

There is a set $S \subset \{1, \dots, n\}$ of cardinality $s = \lfloor n\varepsilon \rfloor$ and a collection of vectors $\{\boldsymbol{\mu}_i : i \in S\}$, such that

$$\{\mathbf{X}_i : i \in S^c\} \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p) \quad \perp\!\!\!\perp \quad \{\mathbf{X}_i : i \in S\} \sim \bigotimes_{i \in S} \mathcal{N}_p(\boldsymbol{\mu}_i, \mathbf{I}_p).$$

- Similar to HC & HDC: the outliers are independent.
- Different from HC & HDC: the outliers might have different distributions.

We write $P_n \in \mathcal{M}_n^{\text{PC}}(p, \varepsilon, \boldsymbol{\mu}^*)$.

Assumption (AC model for outliers)

For a sequence $Y_i \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\boldsymbol{\mu}^*, \mathbf{I}_p)$, $i = 1, \dots, n$, and a random set $S \subset \{1, \dots, n\}$ of cardinality $s = \lceil n\varepsilon \rceil$ we have

$$\mathbf{X}_i = \mathbf{Y}_i, \quad \forall i \in S^c.$$

- The set S **is not** independent of $\{\mathbf{Y}_i : i = 1, \dots, n\}$.
- The observations $\{\mathbf{X}_i : i \in S\}$ may have arbitrary dependence structure.

We write $P_n \in \mathcal{M}_n^{\text{AC}}(p, \varepsilon, \boldsymbol{\mu}^*)$.

Relation between the models

$$\mathcal{M}_n^{\text{HDC}}(p, 2\varepsilon, \boldsymbol{\mu}^*)$$

$$\mathcal{M}_n^{\text{HC}}(p, \varepsilon, \boldsymbol{\mu}^*)$$

$$\mathcal{M}_n^{\text{PC}}(p, 2\varepsilon, \boldsymbol{\mu}^*)$$

$$\mathcal{M}_n^{\text{AC}}(p, 2\varepsilon, \boldsymbol{\mu}^*)$$

2. Problem formulation and overview of results

Historical approach

Breakdown point

- Assume the unknown parameter μ^* is in \mathbb{R}^p .
- Let $\hat{\mu}$ be an estimator of μ^* . Thus,

$$\hat{\mu} : \bigcup_{n=1}^{\infty} \mathcal{X}^n \rightarrow \mathbb{R}^p.$$

- The breakdown point ε_n^* of $\hat{\mu}$ is defined by

$$\varepsilon_n^* = \frac{1}{n} \min \left\{ s \in \{1, \dots, n\} : \sup_{y_1, \dots, y_s} \|\hat{\mu}(\mathbf{x}_{1:(n-s)}, \mathbf{y}_{1:s})\| = +\infty \right\}.$$

- Drawbacks:
 - does not take into account the impact of “mild” outliers,
 - meaningless if the parameter space is bounded,
 - does not depend on the norm under consideration,
 - ...

Minimax approach

In expectation

- A more informative way of quantifying the robustness is the evaluation of the worst-case risk and its comparison to the minimax risk.
- Worst-case risk of an estimator $\hat{\mu}_n$:

$$R_{n,p,\varepsilon}^*(\hat{\mu}_n) = \sup_{\mu^*} \sup_{P_n \in \mathcal{M}_n^*(p,\varepsilon,\mu^*)} \mathbf{E}_{\mathbf{X} \sim P_n} [\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2].$$

- Here, $\mathcal{M}_n^*(p,\varepsilon,\mu^*)$ is one of the 4 models of contamination considered in previous slides.
- For instance, $R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n)$ is the minimax risk for Huber's contamination model.
- Minimax risk:

$$R_{n,p,\varepsilon}^* = \inf_{\hat{\mu}_n} R_{n,p,\varepsilon}^*(\hat{\mu}_n).$$

Minimax approach

In expectation

- A more informative way of quantifying the robustness is the evaluation of the worst-case risk and its comparison to the minimax risk.
- Worst-case risk of an estimator $\hat{\mu}_n$:

$$R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n) = \sup_{\mu^*} \sup_{P_n \in \mathcal{M}_n^{\text{HC}}(p,\varepsilon,\mu^*)} \mathbf{E}_{\mathbf{X} \sim P_n} [\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2].$$

- Here, $\mathcal{M}_n^*(p, \varepsilon, \mu^*)$ is one of the 4 models of contamination considered in previous slides.
- For instance, $R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n)$ is the minimax risk for Huber's contamination model.
- Minimax risk:

$$R_{n,p,\varepsilon}^{\text{HC}} = \inf_{\hat{\mu}_n} R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n).$$

Minimax approach

In expectation

- A more informative way of quantifying the robustness is the evaluation of the worst-case risk and its comparison to the minimax risk.
- Worst-case risk of an estimator $\hat{\mu}_n$:

$$R_{n,p,\varepsilon}^{\text{HDC}}(\hat{\mu}_n) = \sup_{\mu^*} \sup_{P_n \in \mathcal{M}_n^{\text{HDC}}(p,\varepsilon,\mu^*)} \mathbf{E}_{\mathbf{X} \sim P_n} [\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2].$$

- Here, $\mathcal{M}_n^*(p, \varepsilon, \mu^*)$ is one of the 4 models of contamination considered in previous slides.
- For instance, $R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n)$ is the minimax risk for Huber's contamination model.
- Minimax risk:

$$R_{n,p,\varepsilon}^{\text{HDC}} = \inf_{\hat{\mu}_n} R_{n,p,\varepsilon}^{\text{HDC}}(\hat{\mu}_n).$$

Minimax approach

In expectation

- A more informative way of quantifying the robustness is the evaluation of the worst-case risk and its comparison to the minimax risk.
- Worst-case risk of an estimator $\hat{\mu}_n$:

$$R_{n,p,\varepsilon}^{\text{PC}}(\hat{\mu}_n) = \sup_{\mu^*} \sup_{P_n \in \mathcal{M}_n^{\text{PC}}(p,\varepsilon,\mu^*)} \mathbf{E}_{\mathbf{X} \sim P_n} [\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2].$$

- Here, $\mathcal{M}_n^*(p, \varepsilon, \mu^*)$ is one of the 4 models of contamination considered in previous slides.
- For instance, $R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n)$ is the minimax risk for Huber's contamination model.
- Minimax risk:

$$R_{n,p,\varepsilon}^{\text{PC}} = \inf_{\hat{\mu}_n} R_{n,p,\varepsilon}^{\text{PC}}(\hat{\mu}_n).$$

Minimax approach

In expectation

- A more informative way of quantifying the robustness is the evaluation of the worst-case risk and its comparison to the minimax risk.
- Worst-case risk of an estimator $\hat{\mu}_n$:

$$R_{n,p,\varepsilon}^{\text{AC}}(\hat{\mu}_n) = \sup_{\mu^*} \sup_{P_n \in \mathcal{M}_n^{\text{AC}}(p,\varepsilon,\mu^*)} \mathbf{E}_{\mathbf{X} \sim P_n} [\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2].$$

- Here, $\mathcal{M}_n^*(p, \varepsilon, \mu^*)$ is one of the 4 models of contamination considered in previous slides.
- For instance, $R_{n,p,\varepsilon}^{\text{HC}}(\hat{\mu}_n)$ is the minimax risk for Huber's contamination model.
- Minimax risk:

$$R_{n,p,\varepsilon}^{\text{AC}} = \inf_{\hat{\mu}_n} R_{n,p,\varepsilon}^{\text{AC}}(\hat{\mu}_n).$$

Minimax approach

In deviation

- Most results in the literature provide bounds on the deviation, not for the expectation.
- Fix a confidence level $\delta \in (0, 1)$.
- Worst-case deviation of an estimator $\hat{\mu}_n$: $r_{n,p,\varepsilon}^*(\hat{\mu}_n)$ is solution to

$$\begin{aligned} & \text{minimize} && r \\ & \text{subject to} && \mathbf{P}_{\mathbf{X} \sim \mathcal{P}_n} (\|\hat{\mu}_n(\mathbf{X}) - \mu^*\|_2^2 > r) \leq \delta \\ & && \forall \mu^* \in \mathbb{R}^p, \forall \mathcal{P}_n \in \mathcal{M}_n^*(p, \varepsilon, \mu^*). \end{aligned}$$

Clearly, $r_{n,p,\varepsilon}^*(\hat{\mu}_n)$ depends on δ , but we will not be interested in this dependence.

- Minimax risk:

$$r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}_n} r_{n,p,\varepsilon}^*(\hat{\mu}_n).$$

- Tchebychev's inequality yields $\delta r_{n,p,\varepsilon}^*(\hat{\mu}_n) \leq R_{n,p,\varepsilon}^*(\hat{\mu}_n)$.

Common robust estimators of the mean

- The most common robust estimators of the mean are perhaps the coordinatewise median, the geometric median and the Huber's estimator.
- All these estimators can be defined as an M -estimator:

$$\hat{\boldsymbol{\mu}}_n \in \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^p} \sum_{i=1}^n \Psi(\mathbf{X}_i - \boldsymbol{\mu})$$

with

$$\Psi(\mathbf{x}) = \begin{cases} \|\mathbf{x}\|_1, & \text{coordinatewise median,} \\ \|\mathbf{x}\|_2, & \text{geometric median,} \\ \frac{\|\mathbf{x}\|_2^2}{2} \wedge \lambda(\|\mathbf{x}\|_2 - 0.5\lambda), & \text{Huber's estimator.} \end{cases}$$

- In all the three cases, the function Ψ is convex and the estimator is computable in polynomial time.

Overview of the results

Minimax rates in deviation

- Lower bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \geq c\left(\frac{p}{n} + \varepsilon^2\right)$$

where c is a constant depending only on δ .

- Upper bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \leq C\left(\frac{p}{n} + \varepsilon^2\right)$$

where C is a constant depending only on δ .

- It is attained by Tukey's median, which is not computationally tractable.
- The coordinatewise median, the geometric median and the Huber estimator are sub-optimal:

$$r_{n,p,\varepsilon}^{\text{HC}}(\hat{\boldsymbol{\mu}}) \geq \bar{c}\left(\frac{p}{n} + p\varepsilon^2\right).$$

Overview of the results

Minimax rates in deviation

- Lower bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \geq c\left(\frac{p}{n} + \varepsilon^2\right)$$

where c is a constant depending only on δ .

- Upper bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \leq C\left(\frac{p}{n} + \varepsilon^2\right)$$

where C is a constant depending only on δ .

- It is attained by Tukey's median, which is not computationally tractable.
- The coordinatewise median, the geometric median and the Huber estimator are sub-optimal:

$$r_{n,p,\varepsilon}^{\text{HC}}(\hat{\boldsymbol{\mu}}) \geq \bar{c}\left(\frac{p}{n} + p\varepsilon^2\right).$$

Overview of the results

Minimax rates in deviation

- Lower bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \geq c\left(\frac{p}{n} + \varepsilon^2\right)$$

where c is a constant depending only on δ .

- Upper bound on the minimax risk (Chen et al., 2015):

$$r_{n,p,\varepsilon}^{\text{HC}} \leq C\left(\frac{p}{n} + \varepsilon^2\right)$$

where C is a constant depending only on δ .

- It is attained by Tukey's median, which is not tractable.
- The coordinatewise median, the geometric median estimator are sub-optimal:

There is an extra factor p .

$$r_{n,p,\varepsilon}^{\text{HC}}(\hat{\boldsymbol{\mu}}) \geq \bar{c}\left(\frac{p}{n} + p\varepsilon^2\right).$$

Overview of the results

Tractable estimators

We will present three tractable estimators that improve on the coordinatewise median.

- 1 The ellipsoid method ([Diakonikolas et al., 2016](#)).
- 2 The spectral method ([Lai et al., 2016](#)).
- 3 The iterative soft thresholding ([Collier and Dalalyan, 2017](#)).

3. The minimax rate

Minimax lower bound

Theorem 1 (Chen et al., 2015)

There is a constant $c > 0$ such that for every $\varepsilon \in [0, 1]$ and every $\delta \in (0, 1/2)$, it holds that

$$r_{n,p,\varepsilon}^{\text{HC}} \geq c \left(\frac{p}{n} + \varepsilon^2 \right).$$

Some remarks

- By Tchebychev's inequality, $\frac{p}{n} + \varepsilon^2$ is also a lower bound for the minimax risk in expectation.
- By inclusion, $\frac{p}{n} + \varepsilon^2$ is also a lower bound for the minimax risk in models **HDC** and **AC**.
- The same lower bound $\frac{p}{n} + \varepsilon^2$ holds true for the model **PC**.

Proof of the lower bound 1

- ① From the classic parametric minimax theory:

$$r_{n,p,\varepsilon}^{\text{HC}} \gtrsim \frac{p}{n}.$$

- ② Thus, we need only to show that

$$r_{n,p,\varepsilon}^{\text{HC}} \gtrsim \varepsilon^2.$$

- ③ Main steps of the proof:

- Reduction to dimension 1: $r_{n,p,\varepsilon}^{\text{HC}} \geq r_{n,1,\varepsilon}^{\text{HC}}$.
- Construct a probability density function f_ε such that

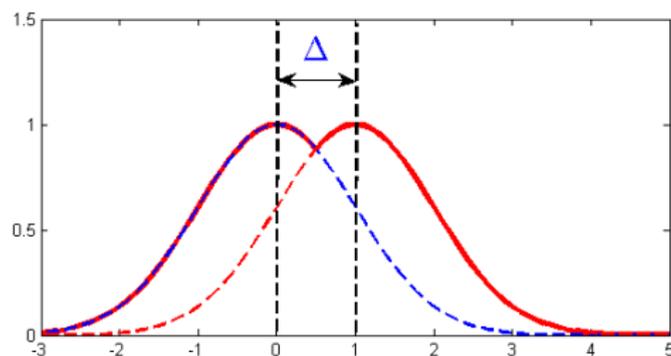
$$\begin{aligned} f_\varepsilon^{\otimes n} &\in \mathcal{M}_n^{\text{HC}}(1, \varepsilon, 0) \\ f_\varepsilon^{\otimes n} &\in \mathcal{M}_n^{\text{HC}}(1, \varepsilon, \Delta_\varepsilon) \end{aligned} \quad \text{with} \quad \Delta_\varepsilon \asymp \varepsilon.$$

- Parameter values $\mu^* = 0$ and $\mu^* = \Delta_\varepsilon$ are indistinguishable from the observations $\mathbf{X}_1, \dots, \mathbf{X}_n \sim f_\varepsilon^{\otimes n}$.

- Therefore $r_{n,p,\varepsilon}^{\text{HC}} \gtrsim \|\Delta_\varepsilon - 0\|_2^2 \asymp \varepsilon^2$.

Proof of the lower bound 2

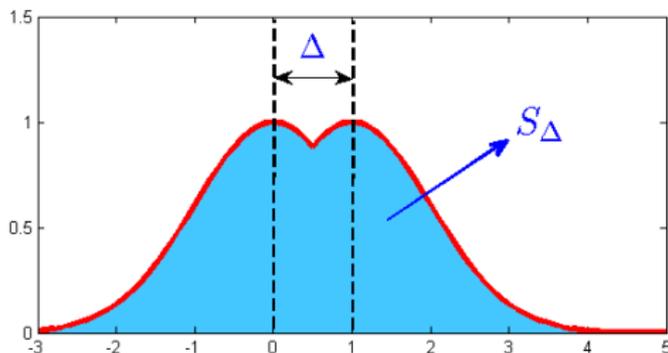
For a $\Delta > 0$, define $f_{\Delta}^{\circ} = \varphi_0 \vee \varphi_{\Delta}$.



Proof of the lower bound 2

For a $\Delta > 0$, define $f_{\Delta}^{\circ} = \varphi_0 \vee \varphi_{\Delta}$.

We have $S_{\Delta} = \int f_{\Delta}^{\circ}(x) dx = 1 + a\Delta + O(\Delta^2)$.

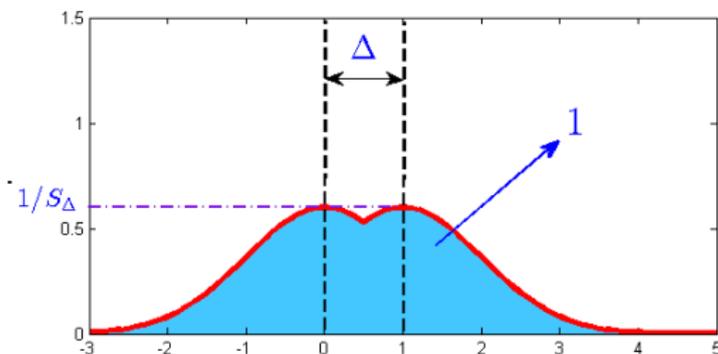


Proof of the lower bound 2

For a $\Delta > 0$, define $f_{\Delta}^{\circ} = \varphi_0 \vee \varphi_{\Delta}$.

We have $S_{\Delta} = \int f_{\Delta}^{\circ}(x) dx = 1 + a\Delta + O(\Delta^2)$.

Then, $f_{\Delta} = f_{\Delta}^{\circ}/S_{\Delta}$ is a pdf.



Proof of the lower bound 2

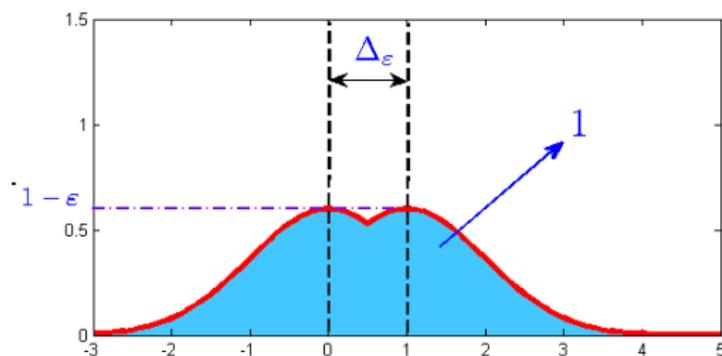
For a $\Delta > 0$, define $f_{\Delta}^{\circ} = \varphi_0 \vee \varphi_{\Delta}$.

We have $S_{\Delta} = \int f_{\Delta}^{\circ}(x) dx = 1 + a\Delta + O(\Delta^2)$.

Then, $f_{\Delta} = f_{\Delta}^{\circ}/S_{\Delta}$ is a pdf.

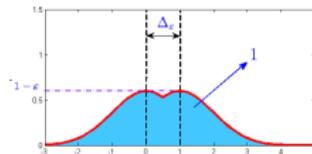
We choose Δ_{ε} so that $1/S_{\Delta_{\varepsilon}} = 1 - \varepsilon$

and set $f_{\varepsilon} = f_{\Delta_{\varepsilon}}$.



Proof of the lower bound 3

$$f_\varepsilon = (1 - \varepsilon)(\varphi_0 \vee \varphi_{\Delta_\varepsilon})$$



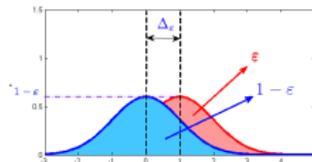
Proof of the lower bound 3

$$f_\varepsilon = (1 - \varepsilon)(\varphi_0 \vee \varphi_{\Delta_\varepsilon})$$

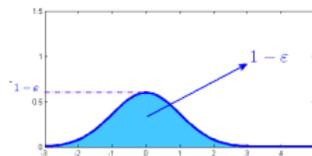
||

$$(1 - \varepsilon)\varphi_0$$

+



||



+

Proof of the lower bound 3

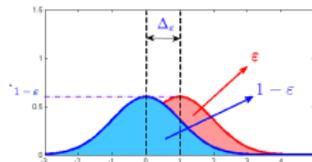
$$f_\varepsilon = (1 - \varepsilon)(\varphi_0 \vee \varphi_{\Delta_\varepsilon})$$

||

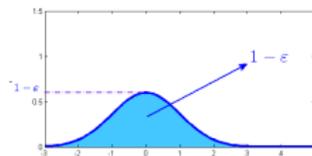
$$(1 - \varepsilon)\varphi_0$$

+

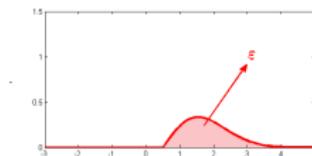
$$\varepsilon q$$



||



+



Proof of the lower bound 3

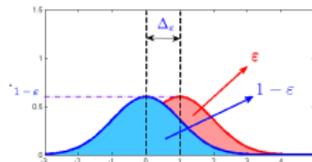
$$f_\varepsilon = (1 - \varepsilon)(\varphi_0 \vee \varphi_{\Delta_\varepsilon})$$

||

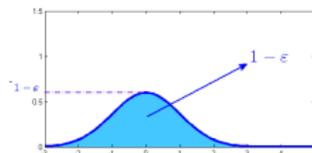
$$(1 - \varepsilon)\varphi_0$$

+

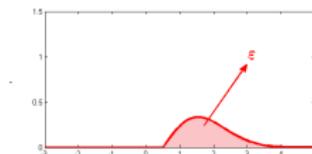
$$\varepsilon q$$



||



+



$$f_\varepsilon = \begin{matrix} (1 - \varepsilon)\varphi_0 & + & \varepsilon q \\ (1 - \varepsilon)\varphi_{\Delta_\varepsilon} & + & \varepsilon q' \end{matrix}$$

Minimax upper bound

Theorem 2 (Chen et al., 2015)

There are two constants $C_1, C_2 > 0$ such that

- for every $\varepsilon \leq 1/5$
- for every $p \leq C_1 n$
- for every $\delta \geq e^{-C_1 n}$,

it holds that

$$r_{n,p,\varepsilon}^{\text{HC}} \leq C_2 \left(\frac{p}{n} + \varepsilon^2 + \frac{\log 1/\delta}{n} \right).$$

Some remarks:

- The upper bound is attained by Tukey's median.
- The condition $\varepsilon \leq 1/5$ can be replaced by $\varepsilon \leq 1/3 - c'$, with an arbitrarily small $c' > 0$.
- The estimator does not rely on the knowledge of ε .

Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$

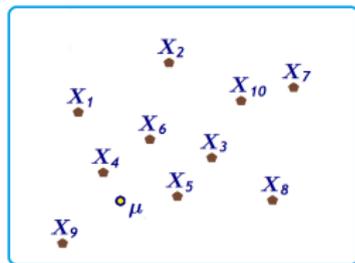
Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\mu}_n^{\text{TM}} \in \arg \max_{\mu \in \mathbb{R}^p} \mathcal{D}(\mu, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\mu, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \mu).$$



depth of μ in $\mathbf{X}_{1:10}$?

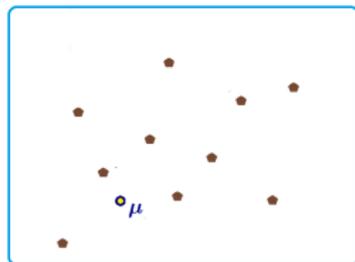
Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$



depth of $\boldsymbol{\mu}$ in $\mathbf{X}_{1:10}$?

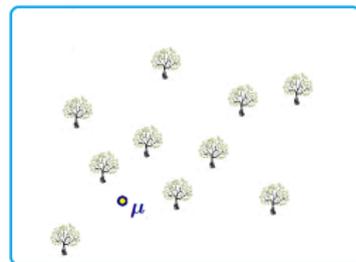
Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$



depth of $\boldsymbol{\mu}$ in the forest ?

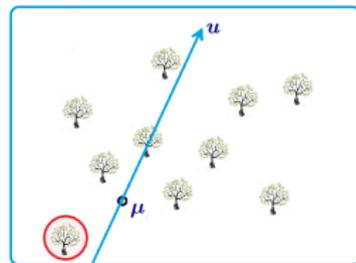
Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$



$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = 1.$$

Tukey's median

- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$



depth of $\boldsymbol{\mu}$ in the forest ?

Tukey's median

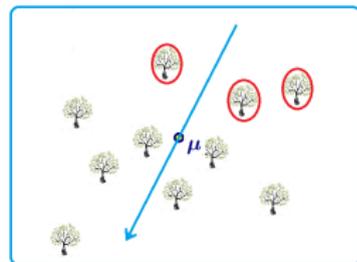
- The upper bound is attained by Tukey's median.
- Tukey's median is any maximizer of Tukey's depth:

$$\hat{\boldsymbol{\mu}}_n^{\text{TM}} \in \arg \max_{\boldsymbol{\mu} \in \mathbb{R}^p} \mathcal{D}(\boldsymbol{\mu}, \{\mathbf{X}_{1:n}\}).$$

- Tukey's (halfspace) depth is

$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = \min_{\mathbf{u} \in \mathbb{S}_1} \sum_{i=1}^n \mathbb{1}(\mathbf{u}^\top \mathbf{X}_i \leq \mathbf{u}^\top \boldsymbol{\mu}).$$

- $\hat{\boldsymbol{\mu}}_n^{\text{TM}}$ is computationally intractable for large p .



$$\mathcal{D}(\boldsymbol{\mu}, \mathbf{X}_{1:n}) = 3.$$

Summary

- We introduced four models of contamination by outliers:
 - Huber's contamination $\mathcal{M}_n^{\text{HC}}(p, \varepsilon, \mu^*)$.
 - Huber's deterministic contamination $\mathcal{M}_n^{\text{HDC}}(p, \varepsilon, \mu^*)$.
 - Parameter contamination $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$.
 - Adversarial contamination $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$.
- We have defined the worst case risks in expectation and in deviation, $R_{n,p,\varepsilon}^*(\hat{\mu})$ and $r_{n,p,\varepsilon}^*(\hat{\mu})$.
- We have defined the minimax risks $R_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} R_{n,p,\varepsilon}^*(\hat{\mu})$.

Summary

- We introduced four models of contamination by outliers:
 - Huber's contamination $\mathcal{M}_n^{\text{HC}}(p, \varepsilon, \mu^*)$.
 - Huber's deterministic contamination $\mathcal{M}_n^{\text{HDC}}(p, \varepsilon, \mu^*)$.
 - Parameter contamination $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$.
 - Adversarial contamination $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$.
- We have defined the worst case risks in expectation and in deviation, $R_{n,p,\varepsilon}^*(\hat{\mu})$ and $r_{n,p,\varepsilon}^*(\hat{\mu})$.
- We have defined the minimax risks $r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} r_{n,p,\varepsilon}^*(\hat{\mu})$.

Summary

- We introduced four models of contamination by outliers:
 - Huber's contamination $\mathcal{M}_n^{\text{HC}}(p, \varepsilon, \mu^*)$.
 - Huber's deterministic contamination $\mathcal{M}_n^{\text{HDC}}(p, \varepsilon, \mu^*)$.
 - Parameter contamination $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$.
 - Adversarial contamination $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$.
- We have defined the worst case risks in expectation and in deviation, $R_{n,p,\varepsilon}^*(\hat{\mu})$ and $r_{n,p,\varepsilon}^*(\hat{\mu})$.
- We have defined the minimax risks $r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} r_{n,p,\varepsilon}^*(\hat{\mu})$.
- For every $\varepsilon < 1/3 - \square$, we have $r_{n,p,\varepsilon}^* \asymp \frac{p}{n} + \varepsilon^2$.
- This minimax rate is obtained by Tukey's median, which is hard to compute for large p .

Summary

- We introduced four models of contamination by outliers:
 - Huber's contamination $\mathcal{M}_n^{\text{HC}}(p, \varepsilon, \mu^*)$.
 - Huber's deterministic contamination $\mathcal{M}_n^{\text{HDC}}(p, \varepsilon, \mu^*)$.
 - Parameter contamination $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$.
 - Adversarial contamination $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$.
- We have defined the worst case risks in expectation and in deviation, $R_{n,p,\varepsilon}^*(\hat{\mu})$ and $r_{n,p,\varepsilon}^*(\hat{\mu})$.
- We have defined the minimax risks $r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} r_{n,p,\varepsilon}^*(\hat{\mu})$.
- For every $\varepsilon < 1/3 - \square$, we have $r_{n,p,\varepsilon}^* \asymp \frac{p}{n} + \varepsilon^2$.
- This minimax rate is obtained by Tukey's median, which is hard to compute for large p .

Question

What is the smallest rate of the worst-case risk that can be obtained by an estimator computable in $\text{poly}(n, p, 1/\varepsilon)$ time?

References I

- M. Chen, C. Gao, and Z. Ren. Robust Covariance and Scatter Matrix Estimation under Huber's Contamination Model. ArXiv e-prints, to appear in the Annals of Statistics, 2015.
- Olivier Collier and Arnak S. Dalalyan. Minimax estimation of a multidimensional linear functional in sparse gaussian models and robust estimation of the mean. submitted 1712.05495, arXiv, December 2017. URL <https://arxiv.org/abs/1712.05495>.
- Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high dimensions without the computational intractability. In IEEE 57th Annual Symposium on Foundations of Computer Science, FOCS 2016, USA, pages 655–664, 2016.
- K. A. Lai, A. B. Rao, and S. Vempala. Agnostic estimation of mean and covariance. In 2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS), pages 665–674, Oct 2016.