

# Robust estimation of a mean in a multivariate Gaussian model: Part 2



Frejus, December 18, 2018

Arnak S. Dalalyan  
ENSAE ParisTech / CREST

## Quick Recap

# 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  $\mu^*$ :  $p$ .
- Observations  $(X_1, \dots, X_n) \sim P_n$ .
- Number of outliers (possibly random):  $s \in \{1, \dots, n\}$ .
- Set of outliers:  $S \subset \{1, \dots, n\}$ .
- Proportion of outliers:  $\varepsilon = \mathbf{E}[s/n] = \mathbf{E}[|S|/n]$ .

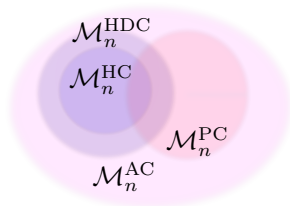
## Setting (informal)

Among the  $n$  observations  $X_1, \dots, X_n$ , there is a small number  $s$  of outliers. If we remove the outliers, all the other  $X_i$ 's are iid drawn from a reference distribution  $\mathcal{N}_p(\mu^*, \mathbf{I}_p)$ .

# Summary of the first lecture

- We have introduced four models of contamination:  $\mathcal{M}_n^\square(p, \varepsilon, \mu^*)$ .

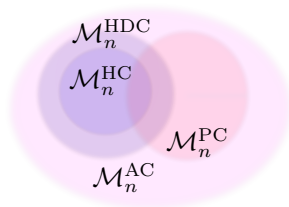
- Huber's Cont.:  $\square = \text{HC}$ .
- Huber's deterministic:  $\square = \text{HDC}$ .
- Parameter Cont.:  $\square = \text{PC}$ .
- Adversarial Cont.:  $\square = \text{AC}$ .



# Summary of the first lecture

- We have introduced four models of contamination:  $\mathcal{M}_n^\square(p, \varepsilon, \boldsymbol{\mu}^*)$ .

- Huber's Cont.:  $\square = \text{HC}$ .
- Huber's deterministic:  $\square = \text{HDC}$ .
- Parameter Cont.:  $\square = \text{PC}$ .
- Adversarial Cont.:  $\square = \text{AC}$ .

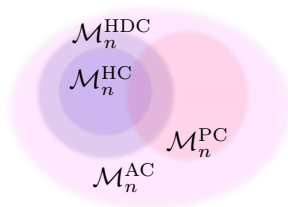


- We have defined the worst case risk  $r_{n,p,\varepsilon}^*(\hat{\boldsymbol{\mu}})$  and the minimax risk  $r_{n,p,\varepsilon}^* = \inf_{\hat{\boldsymbol{\mu}}} r_{n,p,\varepsilon}^*(\hat{\boldsymbol{\mu}})$ .

# Summary of the first lecture

- We have introduced four models of contamination:  $\mathcal{M}_n^\square(p, \varepsilon, \mu^*)$ .

- Huber's Cont.:  $\square = \text{HC}$ .
- Huber's deterministic:  $\square = \text{HDC}$ .
- Parameter Cont.:  $\square = \text{PC}$ .
- Adversarial Cont.:  $\square = \text{AC}$ .

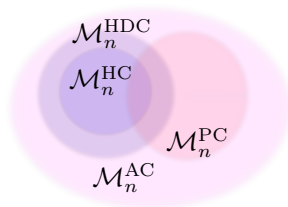


- We have defined the worst case risk  $r_{n,p,\varepsilon}^*(\hat{\mu})$  and the minimax risk  $r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} r_{n,p,\varepsilon}^*(\hat{\mu})$ .
- We have seen that  $\forall \varepsilon < 1/3 - \square$ , we have  $r_{n,p,\varepsilon}^* \asymp \frac{p}{n} + \varepsilon^2$ .
- This minimax rate is obtained by Tukey's median.

# Summary of the first lecture

- We have introduced four models of contamination:  $\mathcal{M}_n^\square(p, \varepsilon, \mu^*)$ .

- Huber's Cont.:  $\square = \text{HC}$ .
- Huber's deterministic:  $\square = \text{HDC}$ .
- Parameter Cont.:  $\square = \text{PC}$ .
- Adversarial Cont.:  $\square = \text{AC}$ .



- We have defined the worst case risk  $r_{n,p,\varepsilon}^*(\hat{\mu})$  and the minimax risk  $r_{n,p,\varepsilon}^* = \inf_{\hat{\mu}} r_{n,p,\varepsilon}^*(\hat{\mu})$ .
- We have seen that  $\forall \varepsilon < 1/3 - \square$ , we have  $r_{n,p,\varepsilon}^* \asymp \frac{p}{n} + \varepsilon^2$ .
- This minimax rate is obtained by Tukey's median.

## 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?

## **4. Robust estimation by the ellipsoid method**



# Worst-case risk bound

## Ellipsoid method for robust estimation

### Theorem 3 (Diakonikolas et al., 2016)

Let  $\delta \in (0, 1/2)$ . There are constants  $c, C > 0$  such that, for every  $\varepsilon \leq c$ , on a set of probability  $\geq 1 - \delta$ , the *ellipsoid method for robust estimation* terminates in  $\text{poly}(n, p, 1/\varepsilon)$  steps and outputs a weight vector  $\hat{\mathbf{w}} \in [0, 1]^n$  such that the mean

$$\hat{\boldsymbol{\mu}}_n^{Ell} = \sum_{i=1}^n \hat{w}_i \mathbf{X}_i$$

satisfies  $\|\hat{\boldsymbol{\mu}}_n^{Ell} - \boldsymbol{\mu}^*\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 \log(1/\varepsilon) + \frac{\log(1/\delta)}{n} \right)$ .

- Valid for  $\mathcal{M}_n^{AC}(p, \varepsilon, \boldsymbol{\mu}^*)$ .
- Complexity of 1 step:  $O(np^2)$ .

# Worst-case risk bound

## Ellipsoid method for robust estimation

### Theorem 3 (Diakonikolas et al., 2016)

Let  $\delta \in (0, 1/2)$ . There are constants  $c, C > 0$  such that, for every  $\varepsilon \leq c$ , on a set of probability  $\geq 1 - \delta$ , the *ellipsoid method for robust estimation* terminates in  $\text{poly}(n, p, 1/\varepsilon)$  steps and outputs a weight vector  $\hat{\mathbf{w}} \in [0, 1]^n$  such that the mean

$$\hat{\boldsymbol{\mu}}_n^{Ell} = \sum_{i=1}^n \hat{w}_i \mathbf{X}_i$$

satisfies  $\|\hat{\boldsymbol{\mu}}_n^{Ell} - \boldsymbol{\mu}^*\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 \frac{\log(1/\varepsilon)}{n} + \frac{\log(1/\delta)}{n} \right).$

- Valid for  $\mathcal{M}_n^{AC}(p, \varepsilon, \boldsymbol{\mu}^*)$ .
- Complexity of 1 step:  $O(np^2)$ .
- Extra factor  $\log(1/\varepsilon)$ .

# Ellipsoid Algorithm

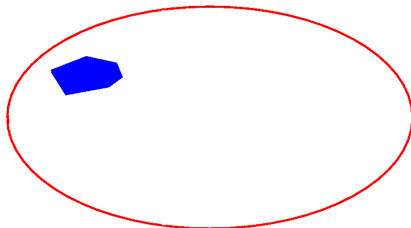
The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .

# Ellipsoid Algorithm

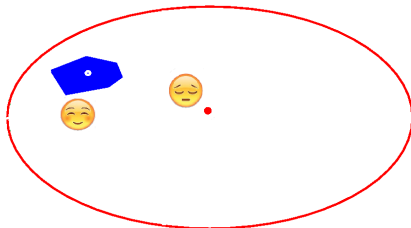
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- One needs
  - an ellipsoid containing  $\mathcal{P}$ ,

# Ellipsoid Algorithm

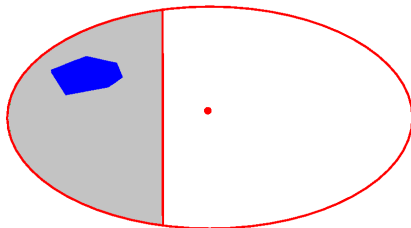
The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- One needs
  - an ellipsoid containing  $\mathcal{P}$ ,
  - a membership oracle,

# Ellipsoid Algorithm

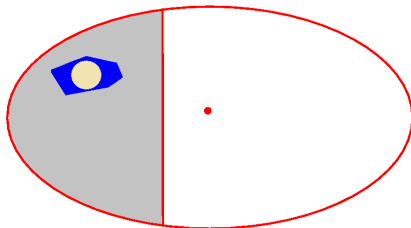
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- One needs
  - an ellipsoid containing  $\mathcal{P}$ ,
  - a membership oracle,
  - a separation oracle,

# Ellipsoid Algorithm

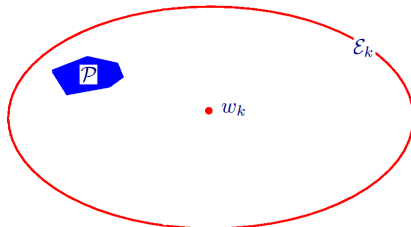
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- One needs
  - an ellipsoid containing  $\mathcal{P}$ ,
  - a membership oracle,
  - a separation oracle,
  - a lower bound on the volume of  $\mathcal{P}$ ,

# Ellipsoid Algorithm

The classic one

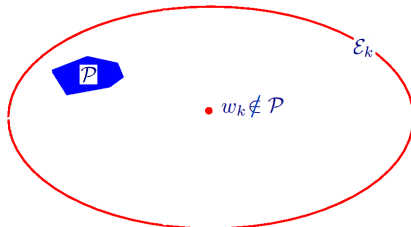


- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,



# Ellipsoid Algorithm

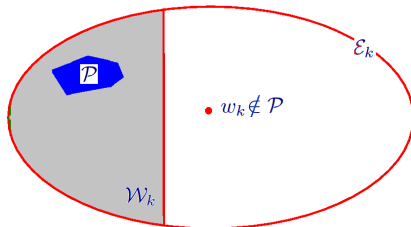
The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,

# Ellipsoid Algorithm

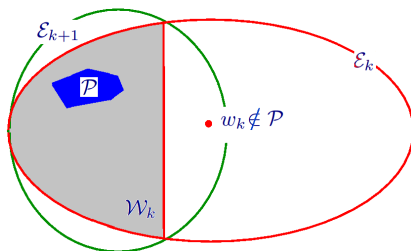
The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .

# Ellipsoid Algorithm

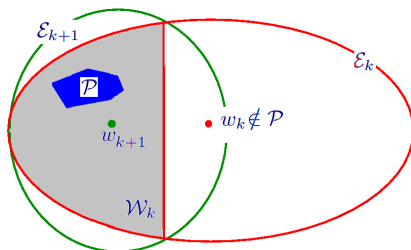
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .

# Ellipsoid Algorithm

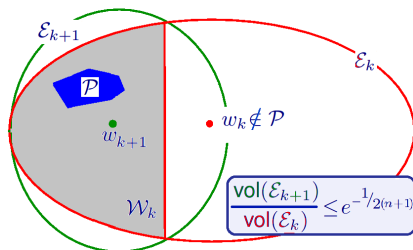
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .
- Since  $\text{vol}(\mathcal{E}_{k+1}) \leq e^{-1/2(n+1)} \text{vol}(\mathcal{E}_k)$ , the algo will stop after  $O(n^2)$  steps,

# Ellipsoid Algorithm

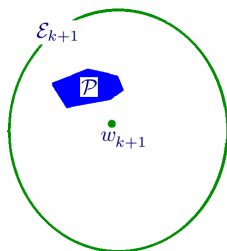
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .
- Since  $\text{vol}(\mathcal{E}_{k+1}) \leq e^{-1/2(n+1)} \text{vol}(\mathcal{E}_k)$ , the algo will stop after  $O(n^2)$  steps, if  $\log \text{vol}(\mathcal{P}) \geq -cn$ .

# Ellipsoid Algorithm

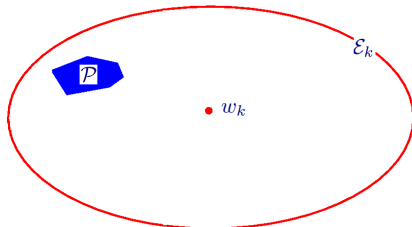
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .
- Since  $\text{vol}(\mathcal{E}_{k+1}) \leq e^{-1/2(n+1)} \text{vol}(\mathcal{E}_k)$ , the algo will stop after  $O(n^2)$  steps, if  $\log \text{vol}(\mathcal{P}) \geq -cn$ .

# Ellipsoid Algorithm

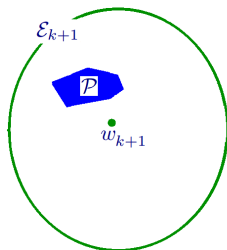
## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .
- Since  $\text{vol}(\mathcal{E}_{k+1}) \leq e^{-1/2(n+1)} \text{vol}(\mathcal{E}_k)$ , the algo will stop after  $O(n^2)$  steps, if  $\log \text{vol}(\mathcal{P}) \geq -cn$ .

# Ellipsoid Algorithm

## The classic one



- Nemirovski-Yudin (1976), Chor (1976-77), Khachiyan (1979).
- The goal is to find a point  $w$  belonging to the polytope  $\mathcal{P}$ .
- Repeat  $K$  times:
  - if the center  $w_k \in \mathcal{P}$ , stop and output  $w_k$ ,
  - else compute the minimal volume ellipsoid  $\mathcal{E}_{k+1}$  containing  $\mathcal{W}_k$ .
- Since  $\text{vol}(\mathcal{E}_{k+1}) \leq e^{-1/2(n+1)} \text{vol}(\mathcal{E}_k)$ , the algo will stop after  $O(n^2)$  steps, if  $\log \text{vol}(\mathcal{P}) \geq -cn$ .



# Ellipsoid method for robust estimation

**Goal:** use the ellipsoid algorithm for approximating the ideal weight vector  $\mathbf{w}^*$  defined by  $w_i^* = \mathbb{1}(i \in S^c)/(n-s)$ ,  $i = 1, \dots, n$ .

- Candidate weights  $\Omega = \{\mathbf{w} \in [0, \frac{1}{n-2s}]^n : \mathbf{w}^\top \mathbf{1}_n = 1\}$ .
- Good weights [note that  $\mathbf{w}^* \in \Omega^* \subset \Omega$ ]

$$\Omega^* = \left\{ \mathbf{w} \in \Omega : \left\| \sum_{i=1}^n w_i (\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})(\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})^\top - \mathbf{I}_p \right\|_{sp}^2 \leq \tau \right\}.$$

with  $\tau = C\left(\frac{p}{n} + \varepsilon^2 \log(1/\varepsilon) + \frac{\log(1/\delta)}{n}\right)$ .

# Ellipsoid method for robust estimation

**Goal:** use the ellipsoid algorithm for approximating the ideal weight vector  $\mathbf{w}^*$  defined by  $w_i^* = \mathbb{1}(i \in S^c)/(n-s)$ ,  $i = 1, \dots, n$ .

- Candidate weights  $\Omega = \{\mathbf{w} \in [0, \frac{1}{n-2s}]^n : \mathbf{w}^\top \mathbf{1}_n = 1\}$ .
- Good weights [note that  $\mathbf{w}^* \in \Omega^* \subset \Omega$ ]

$$\Omega^* = \left\{ \mathbf{w} \in \Omega : \left\| \sum_{i=1}^n w_i (\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})(\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})^\top - \mathbf{I}_p \right\|_{sp}^2 \leq \tau \right\}.$$

with  $\tau = C\left(\frac{p}{n} + \varepsilon^2 \log(1/\varepsilon) + \frac{\log(1/\delta)}{n}\right)$ .

- With probability  $\geq 1 - \delta$ ,
  - $\mathbf{w}^* \in \Omega^*$ ,
  - $\forall \mathbf{w} \notin \Omega^*$ , one can linearly separate  $\mathbf{w}$  from  $\mathbf{w}^*$ .

# Ellipsoid method for robust estimation

**Goal:** use the ellipsoid algorithm for approximating the ideal weight vector  $\mathbf{w}^*$  defined by  $w_i^* = \mathbb{1}(i \in S^c)/(n-s)$ ,  $i = 1, \dots, n$ .

- Candidate weights  $\Omega = \{\mathbf{w} \in [0, \frac{1}{n-2s}]^n : \mathbf{w}^\top \mathbf{1}_n = 1\}$ .
- Good weights [note that  $\mathbf{w}^* \in \Omega^* \subset \Omega$ ]

$$\Omega^* = \left\{ \mathbf{w} \in \Omega : \left\| \sum_{i=1}^n w_i (\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})(\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})^\top - \mathbf{I}_p \right\|_{sp}^2 \leq \tau \right\}.$$

with  $\tau = C\left(\frac{p}{n} + \varepsilon^2 \log(1/\varepsilon) + \frac{\log(1/\delta)}{n}\right)$ .

- With probability  $\geq 1 - \delta$ ,
  - $\mathbf{w}^* \in \Omega^*$ ,
  - $\forall \mathbf{w} \notin \Omega^*$ , one can linearly separate  $\mathbf{w}$  from  $\mathbf{w}^*$ .
  - no lower bound on  $\text{vol}(\Omega^*)$ ,
  - no separation between  $\mathbf{w}$  and  $\Omega^*$ .

## Ellipsoid method for robust estimation

**Goal:** use the ellipsoid algorithm for approximating the ideal weight vector  $\mathbf{w}^*$  defined by  $w_i^* = \mathbb{1}(i \in S^c)/(n-s)$ ,  $i = 1, \dots, n$ .

- Candidate weights  $\Omega = \{\mathbf{w} \in [0, \frac{1}{n-2s}]^n : \mathbf{w}^\top \mathbf{1}_n = 1\}$ .
- Good weights [note that  $\mathbf{w}^* \in \Omega^* \subset \Omega$ ]

$$\Omega^* = \left\{ \mathbf{w} \in \Omega : \left\| \sum_{i=1}^n w_i (\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})(\mathbf{X}_i - \bar{\mathbf{X}}_{\mathbf{w}})^\top - \mathbf{I}_p \right\|_{sp}^2 \leq \tau \right\}.$$

with  $\tau = C\left(\frac{p}{n} + \varepsilon^2 \log(1/\varepsilon) + \frac{\log(1/\delta)}{n}\right)$ .

- With probability  $\geq 1 - \delta$ ,

**Does it really terminate  
in polynomial time ?**

separate  $\mathbf{w}$  from  $\mathbf{w}^*$ .

- no lower bound on  $\text{vol}(\Omega^*)$ ,
- no separation between  $\mathbf{w}$  and  $\Omega^*$ .

## **5. Robust estimation by the spectral method**

# Spectral method for robust estimation

## Finite sample guarantees

### Theorem 4 (Lai et al., 2016)

There are constants  $\alpha, c, C > 0$  such that, for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - 1/p^\alpha$ , the *spectral method for robust estimation* runs in  $\text{poly}(n, p, 1/\varepsilon)$ -time and outputs a vector  $\hat{\mu}_n^{\text{Sp}}$  such that

$$\|\hat{\mu}_n^{\text{Sp}} - \mu^*\|_2^2 \leq C \left( \frac{p(\log p)^2 \log n}{n} + \varepsilon^2 \log p \right).$$

- Valid for  $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity:  $O(np^2)$ .

# Spectral method for robust estimation

## Finite sample guarantees

### Theorem 4 (Lai et al., 2016)

There are constants  $\alpha, c, C > 0$  such that, for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - 1/p^\alpha$ , the *spectral method for robust estimation* runs in  $\text{poly}(n, p, 1/\varepsilon)$ -time and outputs a vector  $\hat{\mu}_n^{\text{Sp}}$  such that

$$\|\hat{\mu}_n^{\text{Sp}} - \mu^*\|_2^2 \leq C \left( \frac{p(\log p)^2 \log n}{n} + \varepsilon^2 \log p \right).$$

- Valid for  $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity:  $O(np^2)$ .
- Extra factors  $(\log p)^2 \log n$  and  $\log p$ .

# Spectral method for robust estimation

## Finite sample guarantees

### Theorem 4 (Lai et al., 2016)

There are constants  $\alpha, c, C > 0$  such that, for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - 1/p^\alpha$ , the *spectral method for robust estimation* runs in  $\text{poly}(n, p, 1/\varepsilon)$ -time and outputs a vector  $\hat{\mu}_n^{\text{Sp}}$  such that

$$\|\hat{\mu}_n^{\text{Sp}} - \mu^*\|_2^2 \leq C \left( \frac{p(\log p)^2 \log n}{n} + \varepsilon^2 \log p \right).$$

- Valid for  $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity:  $O(np^2)$ .
- Extra factors  $(\log p)^2 \log n$  and  $\log p$ .
- The decay of the probability is polynomial in  $p$  (versus exponential for other methods).



# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $X_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $X_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
- 3 ... so on...

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
- 3 ... so on...
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
- 3 ... so on...
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
  - $\mathcal{D}'_1 := \mathcal{D}_1 \setminus \{\mathbf{X} : \|\mathbf{X} - \text{Med}(\mathcal{D}_1)\|_2 > C\sqrt{p \log n}\}$ .
  
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
  
- 3 ... so on...
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
  - $\mathcal{D}'_1 := \mathcal{D}_1 \setminus \{\mathbf{X} : \|\mathbf{X} - \text{Med}(\mathcal{D}_1)\|_2 > C\sqrt{p \log n}\}$ .
  - Compute the SVD of the  $\text{Cov}(\mathcal{D}'_1)$ .
  
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
  
- 3 ... so on...
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$



# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
  - $\mathcal{D}'_1 := \mathcal{D}_1 \setminus \{\mathbf{X} : \|\mathbf{X} - \text{Med}(\mathcal{D}_1)\|_2 > C\sqrt{p \log n}\}$ .
  - Compute the SVD of the  $\text{Cov}(\mathcal{D}'_1)$ .
  - Let  $V_1$  be the subspace of  $p/2$  smallest eigenvectors.
  
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
  
  
  
  
  
  
  
  
  
- 3 ... so on...
  
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\boldsymbol{\mu}}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
  - $\mathcal{D}'_1 := \mathcal{D}_1 \setminus \{\mathbf{X} : \|\mathbf{X} - \text{Med}(\mathcal{D}_1)\|_2 > C\sqrt{p \log n}\}$ .
  - Compute the SVD of the  $\text{Cov}(\mathcal{D}'_1)$ .
  - Let  $V_1$  be the subspace of  $p/2$  smallest eigenvectors.
  - Set  $\Pi_{V_1}(\hat{\boldsymbol{\mu}}^{\text{Sp}}) := \text{Mean}(\Pi_{V_1} \mathcal{D}'_1)$  and  $\mathcal{D}'_2 := \Pi_{V_1^\perp} \mathcal{D}_2$ .
- 2 Define  $\hat{\boldsymbol{\mu}}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
- 3 ... so on...
- 4 Define  $\hat{\boldsymbol{\mu}}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$

# Spectral method for robust estimation

## The algorithm

Start by data-splitting  $\mathbf{X}_{1:n} = \mathcal{D}_1, \dots, \mathcal{D}_k$  with  $k = \lceil \log_2 p \rceil$ .

- 1 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_1$  of dimension  $p/2$  using  $\mathcal{D}_1$ .
  - $\mathcal{D}'_1 := \mathcal{D}_1 \setminus \{\mathbf{X} : \|\mathbf{X} - \text{Med}(\mathcal{D}_1)\|_2 > C\sqrt{p \log n}\}$ .
  - Compute the SVD of the  $\text{Cov}(\mathcal{D}'_1)$ .
  - Let  $V_1$  be the subspace of  $p/2$  smallest eigenvectors.
  - Set  $\Pi_{V_1}(\hat{\mu}^{\text{Sp}}) := \text{Mean}(\Pi_{V_1} \mathcal{D}'_1)$  and  $\mathcal{D}'_2 := \Pi_{V_1^\perp} \mathcal{D}_2$ .
- 2 Define  $\hat{\mu}_n^{\text{Sp}}$  on a subspace  $V_2 \subset V_1^\perp$  of dimension  $p/4$  using  $\mathcal{D}_2$ .
  - $\mathcal{D}''_2 := \mathcal{D}'_2 \setminus \{\mathbf{X} \in V_1^\perp : \|\mathbf{X} - \text{Med}(\mathcal{D}'_2)\|_2 > C\sqrt{p \log n}\}$ .
  - Compute the SVD of the  $\text{Cov}(\mathcal{D}''_2)$ .
  - Let  $V_2$  be the subspace of  $p/4$  smallest eigenvectors.
  - Set  $\Pi_{V_2}(\hat{\mu}^{\text{Sp}}) := \text{Mean}(\Pi_{V_2} \mathcal{D}''_2)$  and  $\mathcal{D}'_3 := \Pi_{(V_1 \oplus V_2)^\perp} \mathcal{D}_3$ .
- 3 ... so on...
- 4 Define  $\hat{\mu}_n^{\text{Sp}}$  on  $V_k = (V_1 \oplus \dots \oplus V_{k-1})^\perp$  using  $\mathcal{D}_k$ .
  - Compute the median of  $\{\Pi_{V_k} \mathbf{X} : \mathbf{X} \in \mathcal{D}_k\}$

## **5. Iterative group soft thresholding**

# Parameter contamination

## Some notation

- We observe  $X_1, \dots, X_n$  in  $\mathbb{R}^p$  such that

$$X_i = \mu^* + \theta_i^* + \xi_i, \quad \xi_i \stackrel{\text{iid}}{\sim} \mathcal{N}_p(\mathbf{0}, \mathbf{I}_p).$$

- **Goal:** estimate the vector  $\mu^*$ .
- **Sparsity assumption:** most vectors  $\theta_i^*$  are equal to zero.
- $S = \{i : \|\theta_i^*\|_2 > 0\}$  is considered as the set of outliers.
- Vectors  $\theta_i$  are unknown nuisance parameters.
- **Matrix notation:**  $\mathbf{X} = \mu^* \mathbf{1}_n^\top + \Theta^* + \xi.$
- Auxiliary problem: estimate  $L_n(\Theta^*) = \frac{1}{n} \sum_{i=1}^n \theta_i^*$ .

## Naive idea: Group-lasso estimator

- Group lasso (Chesneau and Hebiri, 2008; Lin and Zhang, 2006; Lounici, Pontil, van de Geer, and Tsybakov, 2011; Meier, van de Geer, and Bühlmann, 2009; Yuan and Lin, 2006):

$$(\hat{\mu}, \hat{\Theta}) \in \arg \min_{\mu, \Theta} \left\{ \sum_{i=1}^n \|X_i - \mu - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}.$$

- The above optimization problem is convex and can be solved efficiently even when  $p$  and  $n$  are large.  $\hat{\mu}$  is exactly the Huber M-estimator (Donoho and Montanari, 2016).

### Theorem

If  $s \leq n/32$  and  $\lambda_i = 6\sqrt{p}$ , then, with prob.  $\geq 1 - \delta$ ,

$$\|L_n(\hat{\Theta}) - L_n(\Theta^*)\|_2^2 \lesssim \varepsilon^2 p, \quad \|\hat{\mu} - \mu^*\|_2^2 \lesssim \frac{p}{n} + \varepsilon^2 p + \frac{\log(2/\delta)}{n}.$$

# Idea behind iterative group soft thresholding

- Group lasso:

$$(\hat{\mu}, \hat{\Theta}) \in \arg \min_{\mu, \Theta} \left\{ \sum_{i=1}^n \|X_i - \mu - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}.$$

# Idea behind iterative group soft thresholding

- Group lasso:

$$(\hat{\mu}, \hat{\Theta}) \in \arg \min_{\mu, \Theta} \left\{ \sum_{i=1}^n \|X_i - \mu - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}.$$

- We have

$$\hat{\mu} = \arg \min_{\mu} \left\{ \sum_{i=1}^n \|X_i - \mu - \hat{\Theta}_i\|_2^2 \right\} = L_n(\mathbf{X}) - L_n(\hat{\Theta}).$$



# Idea behind iterative group soft thresholding

- Group lasso:

$$(\hat{\mu}, \hat{\Theta}) \in \arg \min_{\mu, \Theta} \left\{ \sum_{i=1}^n \|X_i - \mu - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}.$$

- We have

$$\hat{\mu} = \arg \min_{\mu} \left\{ \sum_{i=1}^n \|X_i - \mu - \hat{\Theta}_i\|_2^2 \right\} = L_n(\mathbf{X}) - L_n(\hat{\Theta}).$$

- If we set  $Z_i = X_i - \{L_n(\mathbf{X}) - L_n(\hat{\Theta})\}$ , we get

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \sum_{i=1}^n \|Z_i - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}$$

which is the group-soft-thresholding (GST) estimator applied to  $\mathbf{Z}$ .

# Idea behind iterative group soft thresholding

- Group lasso:

$$(\hat{\mu}, \hat{\Theta}) \in \arg \min_{\mu, \Theta} \left\{ \sum_{i=1}^n \|X_i - \mu - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}.$$

- We have

$$\hat{\mu} = \arg \min_{\mu} \left\{ \sum_{i=1}^n \|X_i - \mu - \hat{\Theta}_i\|_2^2 \right\} = L_n(\mathbf{X}) - L_n(\hat{\Theta}).$$

- If we set  $Z_i = X_i - \{L_n(\mathbf{X}) - L_n(\hat{\Theta})\}$ , we get

$$\hat{\Theta} \in \arg \min_{\Theta} \left\{ \sum_{i=1}^n \|Z_i - \theta_i\|_2^2 + \sum_{i=1}^n \lambda_i \|\theta_i\|_2 \right\}$$

which is the group-soft-thresholding (GST) estimator applied to  $Z$ .

- Some prior work on linear functional estimation suggests to choose

$$\lambda_i = \frac{2p^{1/4} \|Z_i\|_2}{(\|Z_i\|_2^2 - p)_+^{1/2}}.$$

Unfortunately, we can not do that since  $Z$  depends on  $\hat{\Theta}$ .

# Iterative group soft thresholding

## Algorithm

### Algorithm of IGST

- Start with an estimator  $\hat{\Theta}^0$ , for instance, group lasso.
- for  $k = 1, \dots, K$ , do

$$1) \mathbf{Z}_i = \mathbf{X}_i - \{L_n(\mathbf{X}) - L_n(\hat{\Theta}^{k-1})\},$$

$$2) \lambda_i = \frac{2p^{1/4} \|\mathbf{Z}_i\|_2}{(\|\mathbf{Z}_i\|_2^2 - p)_+^{1/2}},$$

$$3) \hat{\Theta}_i^k = \text{GST}(\mathbf{Z}_i, \lambda_i).$$

- Final estimator:

$$\hat{\mu}^{\text{IGST}} = L_n(\mathbf{X}) - L_n(\hat{\Theta}^K).$$

# Iterative group soft thresholding

## Risk bound

### Theorem 5 (Collier and Dalalyan, 2017)

Let  $\nu > 0$  and assume the IGST estimator is run for  $K = \log_2(1/\nu) + \log \log p$  iterations.

There are constants  $c, C > 0$  such that for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - e^{-p/8}$ , the IGST estimator satisfies

$$\|\hat{\mu}^{\text{IGST}} - \mu\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 + (p\varepsilon^4)^{1-\nu} \right).$$

# Iterative group soft thresholding

## Risk bound

### Theorem 5 (Collier and Dalalyan, 2017)

Let  $\nu > 0$  and assume the IGST estimator is run for  $K = \log_2(1/\nu) + \log \log p$  iterations.

There are constants  $c, C > 0$  such that for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - e^{-p/8}$ , the IGST estimator satisfies

$$\|\hat{\mu}^{\text{IGST}} - \mu\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 + (p\varepsilon^4)^{1-\nu} \right).$$

- Valid for  $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$ .

# Iterative group soft thresholding

## Risk bound

### Theorem 5 (Collier and Dalalyan, 2017)

Let  $\nu > 0$  and assume the IGST estimator is run for  $K = \log_2(1/\nu) + \log \log p$  iterations.

There are constants  $c, C > 0$  such that for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - e^{-p/8}$ , the IGST estimator satisfies

$$\|\hat{\mu}^{\text{IGST}} - \mu\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 + (p\varepsilon^4)^{1-\nu} \right).$$

- Valid for  $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity  $O(np \log \log p)$

# Iterative group soft thresholding

## Risk bound

### Theorem 5 (Collier and Dalalyan, 2017)

Let  $\nu > 0$  and assume the IGST estimator is run for  $K = \log_2(1/\nu) + \log \log p$  iterations.

There are constants  $c, C > 0$  such that for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - e^{-p/8}$ , the IGST estimator satisfies

$$\|\hat{\mu}^{\text{IGST}} - \mu\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 + (p\varepsilon^4)^{1-\nu} \right).$$

- Valid for  $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity  $O(np \log \log p)$
- Exponential in  $p$  decay of probability.

# Iterative group soft thresholding

## Risk bound

### Theorem 5 (Collier and Dalalyan, 2017)

Let  $\nu > 0$  and assume the IGST estimator is run for  $K = \log_2(1/\nu) + \log \log p$  iterations.

There are constants  $c, C > 0$  such that for every  $\varepsilon \leq c$ , in an event of probability  $\geq 1 - e^{-p/8}$ , the IGST estimator satisfies

$$\|\hat{\mu}^{\text{IGST}} - \mu\|_2^2 \leq C \left( \frac{p}{n} + \varepsilon^2 + (p\varepsilon^4)^{1-\nu} \right).$$

- Valid for  $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$ .
- Overall complexity  $O(np \log \log p)$
- Exponential in  $p$  decay of probability.
- The rate is optimal in the regime  $\varepsilon = O(p^{-1/2} \vee n^{-1/4})$ .



# Summary

- The ellipsoid method for robust estimation:
  - achieves the minimax rate on  $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$  up to an extra factor  $\log(1/\varepsilon)$ .
  - Complexity of 1 iteration  $O(np^2)$ .
  - Exponential in  $p$  decay of probability.
  - Poly number of iterations ?
- The spectral method for robust estimation:
  - achieves the minimax rate on  $\mathcal{M}_n^{\text{AC}}(p, \varepsilon, \mu^*)$  up to an extra factor  $(\log p)^2 \log n$ .
  - Overall complexity  $O(np^2)$ .
  - Polynomial in  $p$  decay of probability.
- The iterative group-soft-thresholding:
  - achieves the minimax rate on  $\mathcal{M}_n^{\text{PC}}(p, \varepsilon, \mu^*)$  without any extra factor when  $\varepsilon = O(p^{-1/2} \vee n^{-1/4})$ .
  - Overall complexity  $O(np)$ .
  - Exponential in  $p$  decay of probability.



# References I

- Christophe Chesneau and Mohamed Hebiri. Some theoretical results on the grouped variables Lasso. Math. Methods Statist., 17(4):317–326, 2008.
- 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.
- David Donoho and Andrea Montanari. High dimensional robust m-estimation: asymptotic variance via approximate message passing. Probability Theory and Related Fields, 166(3):935–969, Dec 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.

## References II

- Yi Lin and Hao Helen Zhang. Component selection and smoothing in multivariate nonparametric regression. Ann. Statist., 34(5):2272–2297, 2006.
- Karim Lounici, Massimiliano Pontil, Sara van de Geer, and Alexandre B. Tsybakov. Oracle inequalities and optimal inference under group sparsity. Ann. Statist., 39(4):2164–2204, 2011.
- Lukas Meier, Sara van de Geer, and Peter Bühlmann. High-dimensional additive modeling. Ann. Statist., 37(6B):3779–3821, 2009.
- Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables. J. R. Stat. Soc. Ser. B Stat. Methodol., 68(1):49–67, 2006.