

# Algorithmic High-Dimensional Robust Statistics

Ilias Diakonikolas (USC)

TTI Chicago  
August 2018

Can we develop learning algorithms that are *robust* to a *constant* fraction of *corruptions* in the data?

## MOTIVATION

- **Model Misspecification/Robust Statistics:** Any model only approximately valid. Need *stable* estimators [Fisher 1920, Huber 1960s, Tukey 1960s]
- **Outlier Removal:** Natural outliers in real datasets (e.g., biology). Hard to detect in several cases [Rosenberg *et al.*, Science'02; Li *et al.*, Science'08; Paschou *et al.*, Journal of Medical Genetics'10]
- **Reliable/Adversarial/Secure ML:** Data poisoning attacks (e.g., crowdsourcing) [Biggio *et al.* ICML'12, ...]

# DETECTING OUTLIERS IN REAL DATASETS

- High-dimensional datasets tend to be inherently noisy.

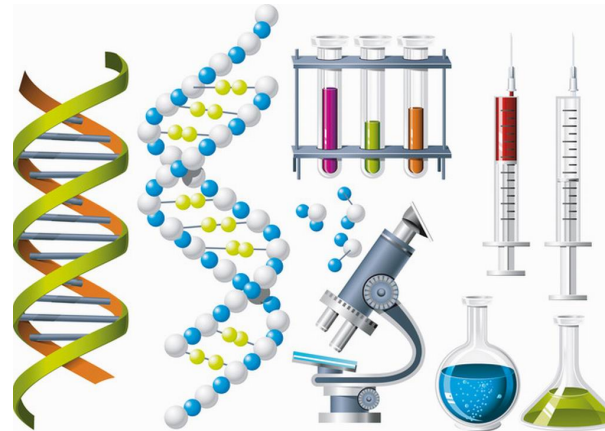
Biological Datasets: POPRES project,  
HGDP datasets

[November *et al.*, Nature'08];

[Rosenberg *et al.*, Science'02];

[Li *et al.*, Science'08];

[Paschou *et al.*, Medical Genetics'10]

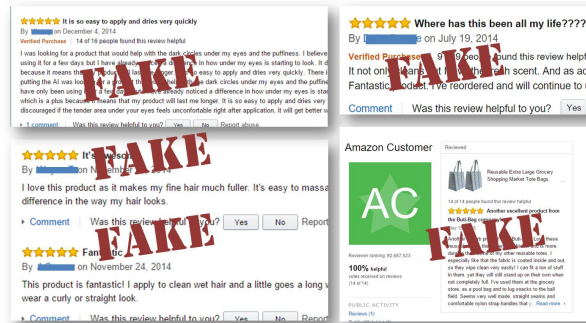


- Outliers: either interesting or can contaminate statistical analysis

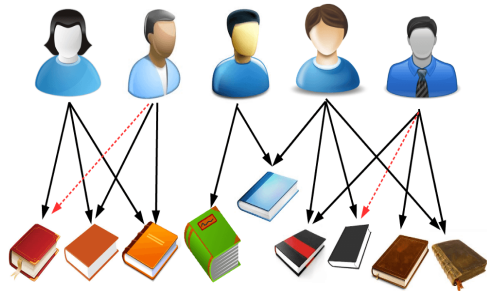
# DATA POISONING

Fake Reviews [Mayzlin et al. '14]

## So Many Misleading, "Fake" Reviews



## Recommender Systems:



[Li et al. '16]

## Crowdsourcing:



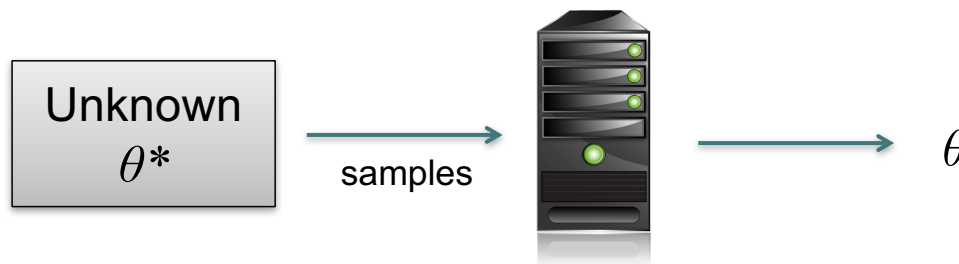
[Wang et al. '14]

## Malware/spam:



[Nelson et al. '08]

# THE STATISTICAL LEARNING PROBLEM



- *Input:* sample generated by a **probabilistic model** with unknown  $\theta^*$
- *Goal:* estimate parameters  $\theta$  so that  $\theta \approx \theta^*$

**Question 1: Is there an *efficient* learning algorithm?**

**Main performance criteria:**

- Sample size
- Running time
- **Robustness**

**Question 2: Are there *tradeoffs* between these criteria?**

## ROBUSTNESS IN A GENERATIVE MODEL

### Contamination Model:

Let  $\mathcal{F}$  be a family of probabilistic models.

We say that a set of  $N$  samples is  $\epsilon$ -corrupted from  $\mathcal{F}$  if it is generated as follows:

- $N$  samples are drawn from an unknown  $F \in \mathcal{F}$
- An omniscient adversary inspects these samples and changes arbitrarily an  $\epsilon$ -fraction of them.

cf. Huber's contamination model [1964]

## MODELS OF ROBUSTNESS

- Oblivious/Adaptive Adversary
- Adversary can: add corrupted samples, subtract uncorrupted samples or both.
- Six Distinct Models:

	Oblivious	Adaptive
Additive Errors	Huber's Contamination Model $P = (1 - \epsilon)G + \epsilon B$	Additive Contamination ("Data Poisoning")
Subtractive Errors	$P = (1 - \epsilon)G - \epsilon L$	Subtractive Contamination
Additive and Subtractive Errors	Hampel's Contamination $d_{TV}(P, G) \leq \epsilon$ $P = G - \epsilon L + \epsilon B$	Strong Contamination ("Nasty Learning Model")

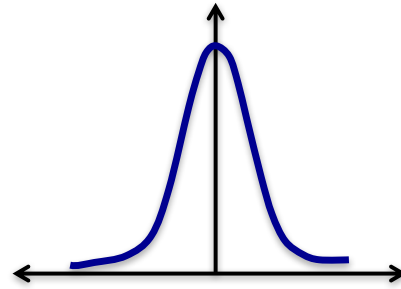


## EXAMPLE: PARAMETER ESTIMATION

Given samples from an unknown distribution:

e.g., a 1-D Gaussian

$$\mathcal{N}(\mu, \sigma^2)$$



how do we accurately estimate its parameters?

**empirical mean:**

$$\frac{1}{N} \sum_{i=1}^N X_i \rightarrow \mu$$

**empirical variance:**

$$\frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2 \rightarrow \sigma^2$$



R. A. Fisher

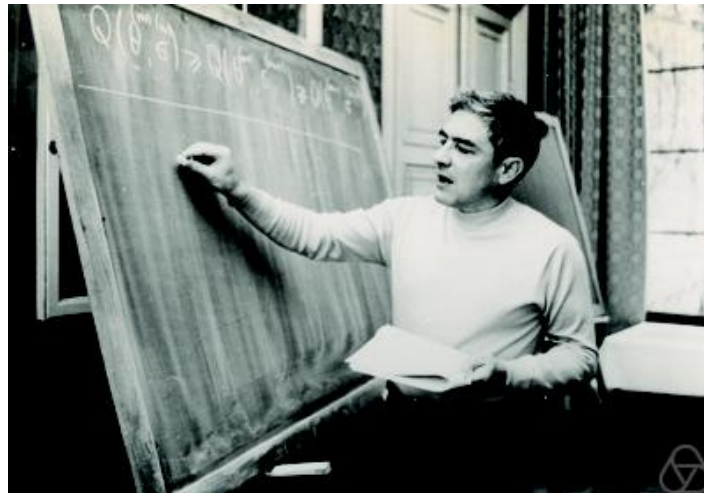
The **maximum likelihood estimator** is asymptotically efficient (1910-1920)



J. W. Tukey

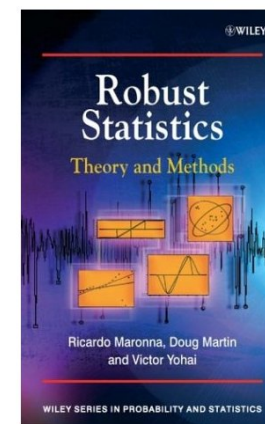
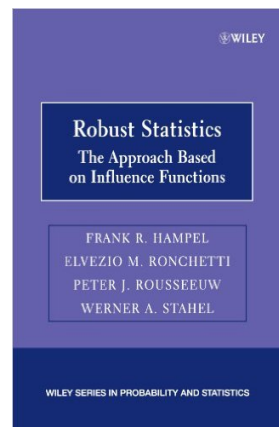
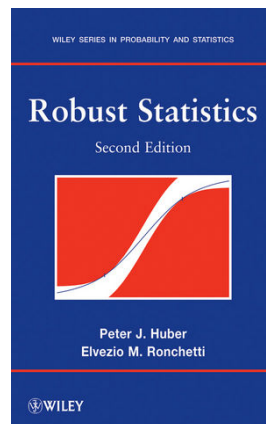
What about **errors** in the model itself? (1960)

Peter J. Huber



“Robust Estimation of a Location Parameter”  
Annals of Mathematical Statistics, 1964.

# ROBUST STATISTICS



What estimators behave well in a **neighborhood** around the model?

## ROBUST ESTIMATION: ONE DIMENSION

Given **corrupted** samples from a *one-dimensional* Gaussian, can we accurately estimate its parameters?

- A single corrupted sample can arbitrarily corrupt the empirical mean and variance.
- But the **median** and **interquartile range** work.

**Fact [Folklore]:** Given a set  $S$  of  $N$   $\epsilon$ -corrupted samples from a one-dimensional Gaussian

$$\mathcal{N}(\mu, \sigma^2)$$

with high constant probability we have that:

$$|\hat{\mu} - \mu| \leq O\left(\epsilon + \sqrt{1/N}\right) \cdot \sigma$$

where  $\hat{\mu} = \text{median}(S)$ .

---

What about robust estimation in high-dimensions?

# GAUSSIAN ROBUST MEAN ESTIMATION

**Robust Mean Estimation:** Given an  $\epsilon$ -corrupted set of samples from an **unknown mean**, identity covariance Gaussian  $\mathcal{N}(\mu, I)$  in  $d$  dimensions, recover  $\hat{\mu}$  with

$$\|\hat{\mu} - \mu\|_2 = O(\epsilon) .$$

**Remark:** Optimal rate of convergence with  $N$  samples is

$$O(\epsilon) + O\left(\sqrt{d/N}\right)$$

[Tukey'75, Donoho'82]

## PREVIOUS APPROACHES: ROBUST MEAN ESTIMATION

<b>Unknown Mean</b>	Error Guarantee	Running Time
Pruning	$\Theta(\epsilon\sqrt{d})$ ✗	$O(dN)$ ✓
Geometric Median	$\Theta(\epsilon\sqrt{d})$ ✗	$\text{poly}(d, N)$ ✓
Tukey Median	$\Theta(\epsilon)$ ✓	NP-Hard ✗
Tournament	$\Theta(\epsilon)$ ✓	$N^{O(d)}$ ✗



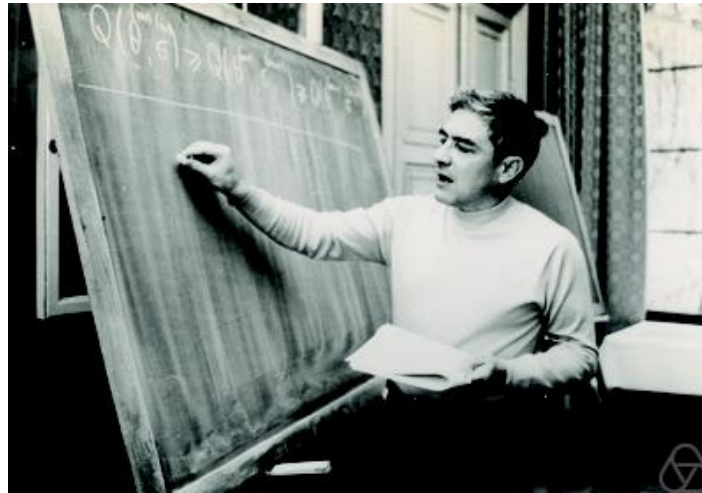
---

All known estimators are either **hard to compute** or  
can tolerate a **negligible fraction of corruptions**.

---

Is robust estimation algorithmically possible in high-dimensions?

Peter J. Huber, 1975



“[...] Only simple algorithms (i.e., with **a low degree of computational complexity**) will survive the onslaught of huge data sets. This runs counter to recent developments in computational robust statistics. **It appears to me that none of the above problems will be amenable to a treatment through theorems and proofs.** They will have to be attacked by heuristics and judgment, and by alternative “what if” analyses.[...]”

Robust Statistical Procedures, 1996, *Second Edition*.

# THIS TALK

Robust estimation in high-dimensions is algorithmically possible!

- First computationally efficient robust estimators that can tolerate a **constant** fraction of corruptions.
- General methodology to detect outliers in high dimensions.

**Meta-Theorem (Informal):** Can obtain *dimension-independent* error guarantees, as long as good data has nice concentration.

**[D-Kamath-Kane-Li-Moitra-Stewart, FOCS'16]**

Can tolerate a ***constant*** fraction of corruptions:

- Mean and Covariance Estimation
- Mixtures of Spherical Gaussians, Mixtures of Balanced Product Distributions

**[Lai-Rao-Vempala, FOCS'16]**

Can tolerate a ***mild sub-constant*** (*inverse logarithmic*) fraction of corruptions:

- Mean and Covariance Estimation
- Independent Component Analysis, SVD

## THIS TALK: ROBUST GAUSSIAN MEAN ESTIMATION

**Theorem:** There are polynomial time algorithms with the following behavior:  
Given  $\epsilon > 0$  and a set of  $N = \tilde{O}(d/\epsilon^2)$   $\epsilon$ -corrupted samples from a  $d$ -dimensional Gaussian  $\mathcal{N}(\mu, I)$ , the algorithms find  $\hat{\mu} \in \mathbb{R}^d$  that with high probability satisfies:

- **[LRV'16]:**

$$\|\mu - \hat{\mu}\|_2 = O(\epsilon\sqrt{\log d})$$

in *additive* contamination model.

- **[DKKLMS'16]:**

$$\|\mu - \hat{\mu}\|_2 = O(\epsilon\sqrt{\log(1/\epsilon)})$$

in *strong* contamination model.

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- Sample Complexity versus Robustness
- Certificate of Robustness
- Recursive Dimension Halving
- Iterative Filtering, Soft Outlier Removal
- Extensions

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- **Sample Complexity versus Robustness**
- Certificate of Robustness
- Recursive Dimension Halving
- Iterative Filtering, Soft Outlier Removal
- Extensions

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

# INFORMATION-THEORETIC LIMITS ON ROBUST ESTIMATION (I)

**Proposition:** Any robust mean estimator for  $\mathcal{N}(\mu, I)$  has error  $\Omega(\epsilon)$ , even in Huber's model.

We start with the following claim:

**Claim:** Let  $P_1, P_2$  be such that  $d_{\text{TV}}(P_1, P_2) = \epsilon/(1 - \epsilon)$ . There exist noise distributions  $B_1, B_2$  such that  $(1 - \epsilon)P_1 + \epsilon B_1 = (1 - \epsilon)P_2 + \epsilon B_2$ .

**Proof:**

Can write

$$P_i = \left(1 - \frac{\epsilon}{1 - \epsilon}\right) P + \frac{\epsilon}{1 - \epsilon} Q_i$$

Take  $B_1 = Q_2$  and  $B_2 = Q_1$ . In this case,

$$(1 - \epsilon)P_1 + \epsilon B_1 = (1 - \epsilon)P_2 + \epsilon B_2 = (1 - 2\epsilon)P + \epsilon(Q_1 + Q_2).$$



## INFORMATION-THEORETIC LIMITS ON ROBUST ESTIMATION (II)

**Proposition:** Any robust mean estimator for  $\mathcal{N}(\mu, I)$  has error  $\Omega(\epsilon)$ , even in Huber's model.

**Proof:**

Need similar construction where  $P_1, P_2$  are unit variance Gaussians.

Let  $P_i = \mathcal{N}(\mu_i, 1)$  such that  $d_{\text{TV}}(P_1, P_2) = \epsilon/(1 - \epsilon)$ .

Since  $d_{\text{TV}}(\mathcal{N}(\mu_1, 1), \mathcal{N}(\mu_2, 1)) \leq |\mu_1 - \mu_2|/2$ , this implies that

$$|\mu_1 - \mu_2| = \Omega(\epsilon).$$

More careful, calculation shows that constant in  $O(\cdot)$  is  $\sqrt{\pi/2} - o(1)$ .

**Remark:** Under different assumptions on good data, we obtain different functions of  $\epsilon$ .

## SAMPLE EFFICIENT ROBUST MEAN ESTIMATION (I)

**Proposition:** There is an algorithm that uses  $N = O(d/\epsilon^2)$   $\epsilon$ -corrupted samples from  $\mathcal{N}(\mu, I)$  and outputs  $\tilde{\mu} \in \mathbb{R}^d$  that with probability at least 9/10 satisfies  $\|\tilde{\mu} - \mu\|_2 = O(\epsilon)$ .

**Main Idea:** To robustly learn the mean of  $\mathcal{N}(\mu, I)$ , it suffices to learn the mean of *all* its 1-dimensional projections. (cf. Tukey median)

**Basic Fact:**  $\|\tilde{\mu} - \mu\|_2 = \max_{v: \|v\|_2=1} |v \cdot \tilde{\mu} - v \cdot \mu|$

Suppose that we can estimate  $v \cdot \mu$  for each  $v \in \mathbb{R}^d$ ,  $\|v\|_2 = 1$ , i.e., find  $\{\hat{\mu}_v\}_v$  such that for all  $v \in \mathbb{R}^d$  with  $\|v\|_2 = 1$  we have  $|\hat{\mu}_v - \mu \cdot v| \leq \delta$ . Then, we can learn  $\mu$  within error  $2\delta$ .

Consider infinite size LP: Find  $x \in \mathbb{R}^d$  such that *for all*  $v \in \mathbb{R}^d$  with  $\|v\|_2 = 1$ :  $|\hat{\mu}_v - v \cdot x| \leq \delta$ . Let  $x^*$  be any feasible solution. Then

$$\|x^* - \mu\|_2 = \max_{v: \|v\|_2=1} |v \cdot x^* - v \cdot \mu| \leq \max_{v: \|v\|_2=1} |v \cdot x^* - \hat{\mu}_v| + \max_{v: \|v\|_2=1} |v \cdot \mu - \hat{\mu}_v| \leq 2\delta.$$

## SAMPLE EFFICIENT ROBUST MEAN ESTIMATION (II)

**Main Idea:** To robustly learn the mean of  $\mathcal{N}(\mu, I)$ , it suffices to learn the mean of “all” its 1-dimensional projections.

Suffices to consider a  $\gamma$ -net  $C$  over all possible directions, where  $\gamma$  is a small positive constant.

This gives the following *finite* LP:

Find  $x \in \mathbb{R}^d$  such that for all  $v \in C$ , we have  $|\hat{\mu}_v - v \cdot x| \leq \delta$ .

Let  $x^*$  be any feasible solution. Let  $u \in C$  such that  $\|u - \frac{\mu - x^*}{\|\mu - x^*\|_2}\|_2 \leq \gamma$ .

Then

$$\|x^* - \mu\|_2 = \left| \left( \left( \frac{\mu - x^*}{\|\mu - x^*\|_2} - u \right) + u \right) \cdot (x^* - \mu) \right| \leq \gamma \|x^* - \mu\|_2 + 2\delta$$

or

$$\|x^* - \mu\|_2 \leq \frac{2\delta}{1 - \gamma}.$$

## SAMPLE EFFICIENT ROBUST MEAN ESTIMATION (III)

**Main Idea:** To robustly learn the mean of  $\mathcal{N}(\mu, I)$ , it suffices to learn the mean of “all” its 1-dimensional projections.

So, for  $\gamma = 1/2$ , any feasible solution to the LP has  $\|x^* - \mu\|_2 \leq 4\delta$ .

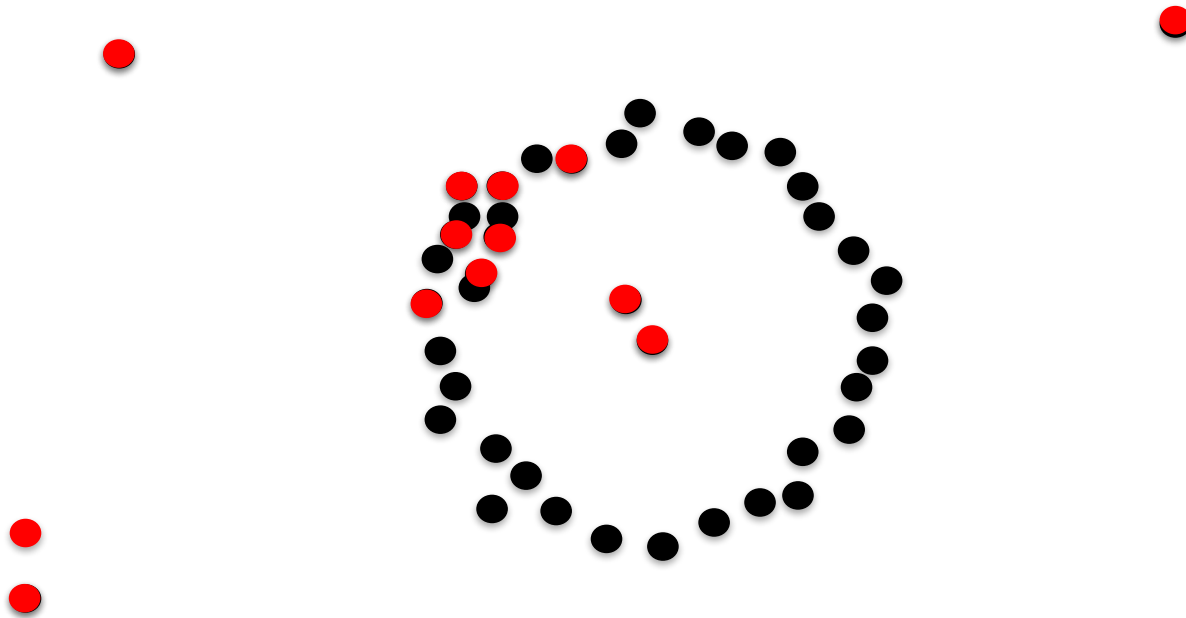
To bound the sample complexity, note that the empirical median satisfies  $\delta = O(\epsilon)$  with probability at least  $1 - \tau$  after  $O((1/\epsilon^2) \log(1/\tau))$  samples.

We need union bound over all  $v \in C$ . Since  $|C| = (1/\gamma)^{O(d)} = 2^{O(d)}$ , for  $\tau = 1/(10|C|)$  our algorithm works with probability at least 9/10.

Thus, sample complexity will be  $N = O(d/\epsilon^2)$ .

**Runtime:**  $\text{poly}(N, 2^d)$ .

# OUTLIER DETECTION ?



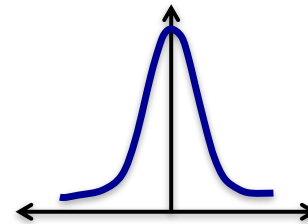
# ON THE EFFECT OF CORRUPTIONS

**Question:** What is the effect of additive and subtractive corruptions?

Let's study the simplest possible example of  $\mathcal{N}(\mu, 1)$ .

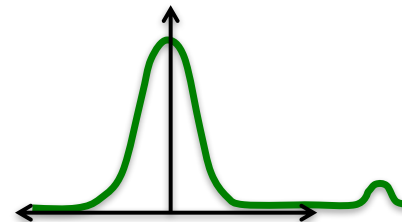
**Subtractive** errors at rate  $\epsilon$  can:

- Move the mean by at most  $O(\epsilon\sqrt{\log(1/\epsilon)})$
- Increase the variance by  $O(\epsilon)$  and decrease it by at most  $O(\epsilon \log(1/\epsilon))$



**Additive** errors at rate  $\epsilon$  can:

- Move the mean arbitrarily
- Increase the variance arbitrarily and decrease it by at most  $O(\epsilon)$



# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

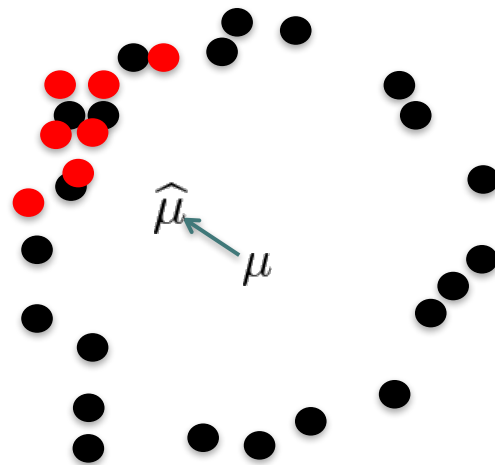
- Sample Complexity versus Robustness
- **Certificate of Robustness**
- Recursive Dimension Halving
- Iterative Filtering, Soft Outlier Removal
- Extensions

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

# CERTIFICATE OF ROBUSTNESS FOR EMPIRICAL ESTIMATOR

Detect when the empirical estimator *may* be compromised



$$\hat{\mu} \triangleq \frac{1}{N} \sum_{i=1}^N X_i$$

- = uncorrupted
- = corrupted

There is *no* direction of large ( $> 1$ ) variance



# COMPARISON OF THREE APPROACHES

## **Three Algorithmic Approaches:**

- Recursive Dimension-Halving [LRV'16]
- Iterative Filtering [DKKLMS'16]
- Soft Outlier Removal [DKKLMS'16]

## **Commonalities:**

- Rely on Spectrum of Empirical Covariance to Robustly Estimate the Mean
- Certificate of Robustness for the Empirical Estimator

## **Exploiting the Certificate:**

- Recursive Dimension-Halving: Find “good” large subspace.
- Iterative Filtering: Check condition on entire space. If violated, filter outliers.
- Soft Outlier Removal: Convex optimization via approximate separation oracle.

**Key Lemma:** Let  $X_1, X_2, \dots, X_N$  be an  $\epsilon$ -corrupted set of samples from  $\mathcal{N}(\mu, I)$  and  $N = \Omega(d/\epsilon^2)$ , then for

$$(1) \hat{\mu} \triangleq \frac{1}{N} \sum_{i=1}^N X_i \quad (2) \hat{\Sigma} \triangleq \frac{1}{N} \sum_{i=1}^N (X_i - \hat{\mu})(X_i - \hat{\mu})^T$$

with high probability we have:

- **[LRV'16]:**

$$\|\hat{\Sigma}\|_2 \leq 1 + O(\epsilon) \quad \rightarrow \quad \|\hat{\mu} - \mu\|_2 \leq O(\epsilon)$$

in **additive** contamination model

- **[DKKLMS'16]:**

$$\|\hat{\Sigma}\|_2 \leq 1 + O(\epsilon \log(1/\epsilon)) \quad \rightarrow \quad \|\hat{\mu} - \mu\|_2 \leq O(\epsilon \sqrt{\log(1/\epsilon)})$$

in **strong** contamination model

**Take-away:** An adversary needs to corrupt the second empirical moment in order to corrupt the first empirical moment

## PROOF OF KEY LEMMA: ADDITIVE CORRUPTIONS (I)

Let  $S = \{X_1, \dots, X_N\}$  be a multi-set of additively  $\epsilon$ -corrupted samples from  $\mathcal{N}(\mu, I)$ . Can assume wlog that  $\mu = \mathbf{0}$ .

Note that  $S = G \cup B$ , where  $G$  is the uncorrupted set of samples and  $B$  is the set of added corrupted samples.

Will express the empirical mean and covariance as sum of two terms, one depending on  $G$  and one on  $B$ .

Let  $\hat{\mu}_G = (1/|G|) \cdot \sum_{i \in I_G} X_i$ , similarly define  $\hat{\mu}_B$ . Have

$$\hat{\mu} = (1 - \epsilon)\hat{\mu}_G + \epsilon\hat{\mu}_B .$$

When  $N \rightarrow \infty$ , we have that  $\hat{\mu}_G = \mu = \mathbf{0}$  .

Therefore,

$$\hat{\mu} = \epsilon\hat{\mu}_B .$$

## PROOF OF KEY LEMMA: ADDITIVE CORRUPTIONS (II)

Recall that  $\mu = \mathbf{0}$  by assumption.

We argued that  $\widehat{\mu} = \epsilon \widehat{\mu}_B$ .

Let's express  $\widehat{\Sigma}$  in similar form.

By definition,  $\widehat{\Sigma} = (1/N) \sum_{i \in [N]} X_i X_i^T - \widehat{\mu} \widehat{\mu}^T$

Define  $\widehat{\Sigma}_G = (1/|G|) \sum_{i \in I_G} X_i X_i^T - \widehat{\mu}_G \widehat{\mu}_G^T$  and similarly  $\widehat{\Sigma}_B$ .

Recall that since  $N \rightarrow \infty$ , we have  $\widehat{\mu}_G = \mu = \mathbf{0}$ .

Similarly, we have that  $\widehat{\Sigma}_G = I$ .

Therefore,

$$(1/N) \sum_{i \in I_G} X_i X_i^T = (1 - \epsilon)I.$$

Also by the definition of  $\widehat{\Sigma}_B$  we get

$$(1/N) \sum_{i \in I_B} X_i X_i^T = \epsilon \widehat{\Sigma}_B + \epsilon \widehat{\mu}_B \widehat{\mu}_B^T.$$

## PROOF OF KEY LEMMA: ADDITIVE CORRUPTIONS (III)

Putting everything together,

$$\widehat{\Sigma} = (1 - \epsilon)I + \epsilon\widehat{\Sigma}_B + (\epsilon - \epsilon^2)\widehat{\mu}_B\widehat{\mu}_B^T.$$

Can now finish argument. Recall that  $\|\widehat{\Sigma}\|_2 = \max_{v:\|v\|_2=1} v^T\widehat{\Sigma}v$ .

Note that  $v^T\widehat{\Sigma}v = (1 - \epsilon) + \epsilon(v^T\widehat{\Sigma}_Bv) + (\epsilon - \epsilon^2)v^T(\widehat{\mu}_B\widehat{\mu}_B^T)v$ .

Choosing  $v = \widehat{\mu}_B/\|\widehat{\mu}_B\|_2$  gives

$$\|\widehat{\Sigma}\|_2 \geq (1 - \epsilon) + (\epsilon - \epsilon^2)\|\widehat{\mu}_B\|_2^2.$$

In conclusion, if  $\|\widehat{\Sigma}\|_2 \leq 1 + \delta$ , then  $\|\widehat{\mu}_B\|_2^2 \leq O(1 + \delta/\epsilon)$

Since  $\widehat{\mu} = \epsilon\widehat{\mu}_B$ , we have shown the following implication:

$$\|\widehat{\Sigma}\|_2 \leq 1 + \delta \quad \longrightarrow \quad \|\widehat{\mu} - \mu\|_2 \leq O(\epsilon + \sqrt{\epsilon\delta}).$$

Choosing  $\delta = O(\epsilon)$  gives the lemma.

## PROOF OF KEY LEMMA: ADDITIVE CORRUPTIONS (IV)

So far assumed we are in infinite sample regime.

Essentially same argument holds in finite sample setting.  
The following concentration inequalities suffice:

For  $N = \Omega(d/\epsilon^2)$ , with high probability we have that

$$\|\mu - \hat{\mu}_G\|_2 \ll \epsilon$$

$$\|\hat{\Sigma}_G - I\|_2 \ll \epsilon$$

## PROOF OF KEY LEMMA: STRONG CORRUPTIONS (I)

Let  $S = \{X_1, \dots, X_N\}$  be a multi-set of  $\epsilon$ -corrupted samples from  $\mathcal{N}(\mu, I)$ . Can assume wlog that  $\mu = \mathbf{0}$ .

Note that  $S = (G \setminus L) \cup B$ , where  $G$  is the uncorrupted set of samples,  $B$  is the added corrupted samples, and  $L \subset G$  is the subtracted set of samples.

Will express empirical mean and covariance as sum of three terms, depending on  $G$ ,  $B$  and  $L$ .

Let  $\hat{\mu}_G = (1/|G|) \cdot \sum_{i \in I_G} X_i$ . Similarly define  $\hat{\mu}_B$  and  $\hat{\mu}_L$ .

We have

$$\hat{\mu} = \hat{\mu}_G - \epsilon \hat{\mu}_L + \epsilon \hat{\mu}_B .$$

When  $N \rightarrow \infty$ , we have that  $\hat{\mu}_G = \mu = \mathbf{0}$ .

Therefore,

$$\hat{\mu} = \epsilon(\hat{\mu}_B - \hat{\mu}_L) .$$

## PROOF OF KEY LEMMA: STRONG CORRUPTIONS (II)

Recall that  $\mu = \mathbf{0}$  by assumption.

We argued that  $\widehat{\mu} = \epsilon(\widehat{\mu}_B - \widehat{\mu}_L)$ .

Will express  $\widehat{\Sigma}$  in similar form.

By definition,  $\widehat{\Sigma} = (1/N) \sum_{i \in [N]} X_i X_i^T - \widehat{\mu} \widehat{\mu}^T$

Define  $\widehat{\Sigma}_G = (1/|G|) \sum_{i \in I_G} X_i X_i^T - \widehat{\mu}_G \widehat{\mu}_G^T$ . Similarly define  $\widehat{\Sigma}_B$ .

Let  $\widehat{M}_L = (1/|L|) \sum_{i \in I_L} X_i X_i^T$ .

Recall that since  $N \rightarrow \infty$ , we have  $\widehat{\mu}_G = \mu = \mathbf{0}$ .

Similarly, we have that  $\widehat{\Sigma}_G = I$ .

Therefore,

$$(1/N) \sum_{i \in I_G} X_i X_i^T = I.$$

Also, by the definition of  $\widehat{\Sigma}_B$  and  $\widehat{M}_L$  we get

$$(1/N) \sum_{I \in I_B} X_i X_i^T = \epsilon \widehat{\Sigma}_B + \epsilon \widehat{\mu}_B \widehat{\mu}_B^T \quad (1/N) \sum_{I \in I_L} X_i X_i^T = \epsilon \widehat{M}_L$$



## PROOF OF KEY LEMMA: STRONG CORRUPTIONS (III)

Putting everything together,

$$\widehat{\Sigma} = I + \epsilon \widehat{\Sigma}_B + \epsilon \widehat{\mu}_B \widehat{\mu}_B^T - \epsilon \widehat{M}_L - \epsilon^2 (\widehat{\mu}_B - \widehat{\mu}_L)(\widehat{\mu}_B - \widehat{\mu}_L)^T .$$

To finish argument, need to bound  $\widehat{M}_L$  and  $\widehat{\mu}_L$  .

**Claim:** Have  $\|\widehat{M}_L\|_2 = O(\log(1/\epsilon))$  and  $\|\widehat{\mu}_L\|_2 = O(\sqrt{\log(1/\epsilon)})$  .

Assuming the claim holds, we get

$$\widehat{\Sigma} = I + \epsilon \widehat{\Sigma}_B + (\epsilon - \epsilon^2) \widehat{\mu}_B \widehat{\mu}_B^T + O(\epsilon \log(1/\epsilon)) .$$

This gives

$$\|\widehat{\Sigma}\|_2 \geq 1 + (\epsilon - \epsilon^2) \|\widehat{\mu}_B\|_2^2 - O(\epsilon \log(1/\epsilon)) .$$

## PROOF OF KEY LEMMA: STRONG CORRUPTIONS (IV)

We can now finish the argument.  
We have shown that

$$\|\widehat{\Sigma}\|_2 \geq 1 + (\epsilon - \epsilon^2)\|\widehat{\mu}_B\|_2^2 - O(\epsilon \log(1/\epsilon)) .$$

Suppose that  $\|\widehat{\Sigma}\|_2 \leq 1 + \delta$ . Then

$$\|\widehat{\mu}_B\|_2 \leq O\left(\sqrt{\delta/\epsilon} + \sqrt{\log(1/\epsilon)}\right)$$

Since  $\widehat{\mu} = \epsilon(\widehat{\mu}_B - \widehat{\mu}_L)$ , the final error is

$$\begin{aligned} \|\widehat{\mu}\|_2 &\leq \epsilon\|\widehat{\mu}_B\|_2 + \epsilon\|\widehat{\mu}_L\|_2 \\ &\leq O\left(\sqrt{\delta\epsilon} + \epsilon\sqrt{\log(1/\epsilon)}\right) . \end{aligned}$$

For  $\delta = \Theta(\epsilon \log(1/\epsilon))$ , lemma follows.

## PROOF OF KEY LEMMA: STRONG CORRUPTIONS (V)

Recall that  $\widehat{M}_L := (1/|L|) \sum_{i \in I_L} X_i X_i^T = \mathbf{E}_{X \sim_{UL}}[X X^T]$ . Remains to prove:

**Claim:** We have  $\|\widehat{M}_L\|_2 = O(\log(1/\epsilon))$  and  $\|\widehat{\mu}_L\|_2 = O(\sqrt{\log(1/\epsilon)})$ .

**Proof:** By definition have  $\|\widehat{M}_L\|_2 = \max_{v: \|v\|_2=1} |v^T \widehat{M}_L v| = \max_{v: \|v\|_2=1} \mathbf{E}_{X \sim_{UL}}[(v \cdot X)^2]$ .

Since  $L \subset G$ , for any event,  $|L| \cdot \Pr_{X \sim_{UL}}[X \in \mathcal{E}] \leq |S| \cdot \Pr_{X \sim_{UG}}[X \in \mathcal{E}]$ .

For any unit vector  $v$ :

$$\begin{aligned} \mathbf{E}_{X \sim_{UL}}[(v \cdot X)^2] &= 2 \int_0^{O(\sqrt{d})} \Pr_{X \sim_{UL}}[|v \cdot X| > T] T dT \\ &\leq 2 \int_0^{O(\sqrt{d})} \min\{1, (1/\epsilon) \cdot \Pr_{X \sim_{UG}}[|v \cdot X| > T]\} T dT \\ &\leq 2 \int_0^{O(\sqrt{\log(1/\epsilon)})} T dT + (1/\epsilon) \cdot \int_{O(\sqrt{\log(1/\epsilon)})}^{O(\sqrt{d})} e^{-T^2/2} T dT \\ &= O(\log(1/\epsilon)) + O(1). \end{aligned}$$

Finally, by definition we have that  $\|\widehat{\mu}_L\|_2^2 \leq \|\widehat{M}_L\|_2$ .

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- Sample Complexity versus Robustness
- Certificate of Robustness
- **Recursive Dimension Halving**
- Iterative Filtering, Soft Outlier Removal
- Extensions

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

## RECURSIVE DIMENSION-HALVING [LRV'16]

### Recursive Procedure:

**Step #1:** Find large subspace where “standard” estimator works.

**Step #2:** Recurse on complement.

(If dimension is small, use brute-force.)

Combine Results.

Can reduce dimension by factor of 2 in each recursive step.

## FINDING A GOOD SUBSPACE (I)

- Good subspace  $\mathbf{G}$ : one where the empirical mean works.
- By **Key Lemma**, sufficient condition is:

Projection of empirical covariance on  $\mathbf{G}$  has no large eigenvalues.

- Also want  $\mathbf{G}$  to be “high-dimensional”.
- **How do we find such a subspace?**

## FINDING A GOOD SUBSPACE (II)

**Good Subspace Lemma:** Let  $X_1, X_2, \dots, X_N$  be an *additively*  $\epsilon$ -corrupted set of  $N = \Omega(d \log d / \epsilon^2)$  samples from  $\mathcal{N}(\mu, I)$ . After *naïve pruning*, we have that

$$\lambda_{d/2}(\hat{\Sigma}) \leq 1 + O(\epsilon)$$

---

**Corollary:** Let  $W$  be the span of the bottom  $d/2$  eigenvalues of  $\hat{\Sigma}$ . Then  $W$  is a good subspace.

## PROOF OF GOOD SUBSPACE LEMMA (I)

Let  $S = \{X_1, \dots, X_N\}$  be a multi-set of additively  $\epsilon$ -corrupted samples from  $\mathcal{N}(\mu, I)$ . Can assume wlog that  $\mu = \mathbf{0}$ .

Note that  $S = G \cup B$ , where  $G$  is the uncorrupted set of samples and  $B$  is the added corrupted samples. Let  $S'$  be the subset of  $S$  obtained after naïve pruning. We know that  $S' = G \cup B'$ , where  $B' \subseteq B$ , and each  $x \in S'$  satisfies  $\|x\|_2 = O(\sqrt{d})$ .

Let  $\widehat{\Sigma}_{S'} = (1/|S'|) \sum_{i \in I_{S'}} X_i X_i^T - \widehat{\mu}_{S'} \widehat{\mu}_{S'}^T$  be the empirical covariance of  $S'$  and  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$  be its spectrum.

Want to show that  $\lambda_{d/2} \leq 1 + O(\epsilon)$ .

This follows from the following claims:

**Claim 1:**  $\lambda_1 \geq 1 - O(\epsilon)$ .

**Claim 2:**  $\text{Tr}(\widehat{\Sigma}_{S'}) \leq d(1 + O(\epsilon))$ .



## PROOF OF GOOD SUBSPACE LEMMA (II)

Let  $\widehat{\Sigma}_{S'} = (1/|S'|) \sum_{i \in I_{S'}} X_i X_i^T - \widehat{\mu}_{S'} \widehat{\mu}_{S'}^T$  be the empirical covariance of  $S'$  and  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$  be its spectrum.

**Claim 1:**  $\lambda_1 \geq 1 - O(\epsilon)$  .

**Claim 2:**  $\text{Tr}(\widehat{\Sigma}_{S'}) \leq d(1 + O(\epsilon))$  .

By Claim 1, 
$$A = \sum_{i=1}^{d/2} \lambda_i \geq (d/2)(1 - O(\epsilon))$$

Moreover, 
$$B = \sum_{i=d/2+1}^d \lambda_i \geq (d/2)\lambda_{d/2}$$

By Claim 2, 
$$A + B \leq d(1 + O(\epsilon))$$

Therefore, 
$$B \leq (d/2)(1 + O(\epsilon))$$

which gives 
$$\lambda_{d/2} \leq 1 + O(\epsilon)$$
 .

## PROOF OF GOOD SUBSPACE LEMMA (III)

Let  $\widehat{\Sigma}_{S'} = (1/|S'|) \sum_{i \in I_{S'}} X_i X_i^T - \widehat{\mu}_{S'} \widehat{\mu}_{S'}^T$  be the empirical covariance of  $S'$  and  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$  be its spectrum.

**Claim 1:**  $\lambda_1 \geq 1 - O(\epsilon)$ .

**Proof:** Recall that  $S' = G \cup B'$ , where  $G$  is the uncorrupted set of samples and  $B'$  is a subset of the added corrupted samples. Therefore,

$$\widehat{\Sigma}_{S'} = (1 - \epsilon)I + \epsilon \widehat{\Sigma}_{B'} + (\epsilon - \epsilon^2) \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T$$

Denoting  $M = \epsilon \widehat{\Sigma}_{B'} + (\epsilon - \epsilon^2) \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T$ , we have that

$$\lambda_{\min}(\widehat{\Sigma}_{S'}) \geq (1 - \epsilon) + \min_{v: \|v\|_2=1} v^T M v \geq 1 - \epsilon.$$

## PROOF OF GOOD SUBSPACE LEMMA (IV)

Let  $\widehat{\Sigma}_{S'} = (1/|S'|) \sum_{i \in I_{S'}} X_i X_i^T - \widehat{\mu}_{S'} \widehat{\mu}_{S'}^T$  be the empirical covariance of  $S'$  and  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$  be its spectrum.

**Claim 2:**  $\text{Tr}(\widehat{\Sigma}_{S'}) \leq d(1 + O(\epsilon))$ .

**Proof:** Recall that

$$\widehat{\Sigma}_{S'} = (1 - \epsilon)I + \epsilon \widehat{\Sigma}_{B'} + (\epsilon - \epsilon^2) \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T$$

Thus,

$$\text{Tr}(\widehat{\Sigma}_{S'}) \leq d(1 - \epsilon) + \epsilon \text{Tr}(\widehat{\Sigma}_{B'} + \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T)$$

Note that

$$\widehat{\Sigma}_{B'} + \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T = (1/|B'|) \sum_{i \in I_{B'}} X_i X_i^T$$

Moreover, for every  $x \in B' \subseteq S'$  we have  $\|x\|_2 = O(\sqrt{d})$ .

Thus,

$$\text{Tr}(\widehat{\Sigma}_{B'} + \widehat{\mu}_{B'} \widehat{\mu}_{B'}^T) = O(d).$$

## RECURSIVE DIMENSION-HALVING ALGORITHM [LRV'16]

Algorithm works as follows:

- Remove gross outliers (e.g., naïve pruning).
- Let  $W, V$  be the span of bottom  $d/2$  and upper  $d/2$  eigenvalues of  $\hat{\Sigma}$  respectively.
- Use empirical mean on  $W$ .
- Recurse on  $V$ . (If the dimension is one, use median.)

$O(\log d)$  levels of the recursion  $\rightarrow$  final error of  $O(\epsilon\sqrt{\log d})$

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- Sample Complexity versus Robustness
- Certificate of Robustness
- Recursive Dimension Halving
- **Iterative Filtering, Soft Outlier Removal**
- Extensions

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

## ITERATIVE FILTERING [DKKLMS'16]

### **Iterative Two-Step Procedure:**

**Step #1:** Find certificate of robustness of “standard” estimator

**Step #2:** If certificate is violated, detect and remove outliers

Iterate on “cleaner” dataset.

General recipe that works for fairly general settings.

Let's see how this works for robust mean estimation.

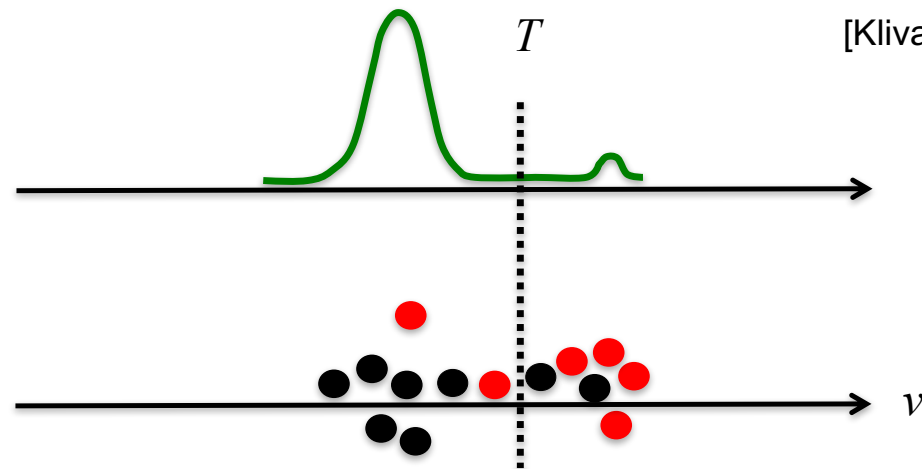
# FILTERING SUBROUTINE

Either output empirical mean, or remove many outliers.

**Filtering Approach:** Suppose that:

$$\|\hat{\Sigma}\|_2 \geq 1 + \Omega(\epsilon \log(1/\epsilon))$$

Let  $v^*$  be the direction of maximum variance.



[Klivans-Long-Servedio'09]

## FILTERING SUBROUTINE

Either output empirical mean, or remove many outliers.

**Filtering Approach:** Suppose that:

$$\|\widehat{\Sigma}\|_2 \geq 1 + \Omega(\epsilon \log(1/\epsilon))$$

Let  $v^*$  be the direction of maximum variance.

- Project all the points on the direction of  $v^*$
- Find a threshold  $T$  such that

$$\Pr_{X \sim \mathcal{U}S}[|v^* \cdot X - \text{median}(\{v^* \cdot x, x \in S\})| > T + 1] \geq 8 \cdot e^{-T^2/2}.$$

- Throw away all points  $x$  such that

$$|v^* \cdot x - \text{median}(\{v^* \cdot x, x \in S\})| > T + 1$$

- Iterate on new dataset.



## FILTERING SUBROUTINE: ANALYSIS SKETCH

Either output empirical mean, or remove many outliers.

**Filtering Approach:** Suppose that:

$$\|\hat{\Sigma}\|_2 \geq 1 + \Omega(\epsilon \log(1/\epsilon))$$

**Claim:** We remove more corrupted than uncorrupted points.

After a number of iterations, we have removed all corrupted points.

Eventually the empirical mean works

## FILTERING SUBROUTINE: PSEUDO-CODE

**Input:**  $\epsilon$ -corrupted set  $S$  from  $\mathcal{N}(\mu, I)$

**Output:** Set  $S' \subseteq S$  that is  $\epsilon'$ -corrupted, for some  $\epsilon' < \epsilon$   
OR robust estimate of the unknown mean  $\mu$

1. Let  $\hat{\mu}_S, \hat{\Sigma}_S$  be the empirical mean and covariance of the set  $S$ .
2. **If**  $\|\hat{\Sigma}_S\|_2 \leq 1 + C\epsilon \log(1/\epsilon)$ , for an appropriate constant  $C > 0$ :  
**Output**  $\hat{\mu}_S$
3. **Otherwise**, let  $(\lambda^*, v^*)$  be the top eigenvalue-eigenvector pair of  $\hat{\Sigma}_S$ .
4. Find  $T > 0$  such that

$$\Pr_{X \sim \mathcal{U}S}[|v^* \cdot X - \text{median}(\{v^* \cdot x, x \in S\})| > T + 1] \geq 8 \cdot e^{-T^2/2}.$$

5. **Return**

$$S' = \{x \in S : |v^* \cdot x - \text{median}(\{v^* \cdot x, x \in S\})| \leq T + 1\}.$$

## SKETCH OF CORRECTNESS (I)

**Claim:** Can always find a threshold satisfying the Condition of Step 4.

**Proof:**

By contradiction. Suppose that for all  $T > 0$  we have

$$\Pr_{X \sim U_S}[|v^* \cdot X - \text{median}(\{v^* \cdot x, x \in S\})| > T + 1] < 8 \cdot e^{-T^2/2}.$$

Will use this to show that  $\lambda^* = \|\hat{\Sigma}_S\|_2$  is smaller than it was assumed to be.

Since the median is a robust estimator of the mean, it follows that for all  $T > 0$

$$\Pr_{X \sim U_S}[|v^* \cdot X - \mu| > T + 2] < 8 \cdot e^{-T^2/2}.$$

Since  $B \subset S$ , for any event  $\mathcal{E}$ ,  $|B| \cdot \Pr_{X \sim U_B}[X \in \mathcal{E}] \leq |S| \cdot \Pr_{X \sim U_S}[X \in \mathcal{E}]$

Therefore,

$$\Pr_{X \sim U_B}[|v^* \cdot (X - \mu)| > T] \leq (1/\epsilon) \cdot \Pr_{X \sim U_S}[|v^* \cdot (X - \mu)| > T]$$

## SKETCH OF CORRECTNESS (II)

Assume wlog  $\mu = 0$ . Recall that

$$\widehat{\Sigma} = I + \epsilon \widehat{\Sigma}_B + (\epsilon - \epsilon^2) \widehat{\mu}_B \widehat{\mu}_B^T + O(\epsilon \log(1/\epsilon)) .$$

So, it suffices to show that  $\widehat{M}_B := \widehat{\Sigma}_B + \widehat{\mu}_B \widehat{\mu}_B^T = \mathbf{E}_{X \sim UB}[XX^T]$  has small  $v^*$ -variance, i.e., that  $\mathbf{E}_{X \sim UB}[(v^* \cdot X)^2]$  is small.

We have

$$\begin{aligned} \mathbf{E}_{X \sim UB}[(v^* \cdot X)^2] &= 2 \int_0^{O(\sqrt{d})} \mathbf{Pr}_{X \sim UB}[|v^* \cdot X| > T] T dT \\ &\leq O(1) + 2 \int_2^{O(\sqrt{d})} \mathbf{Pr}_{X \sim UB}[|v^* \cdot X| > T] T dT \\ &\leq O(1) + 2 \int_2^{O(\sqrt{d})} \min\{1, (1/\epsilon) \cdot \mathbf{Pr}_{X \sim US}[|v^* \cdot X| > T]\} T dT \\ &\leq O(1) + 2 \int_2^{O(\sqrt{\log(1/\epsilon)})} T dT + 16 \int_{O(\sqrt{\log(1/\epsilon)})}^{O(\sqrt{d})} T e^{-(T-2)^2/2} dT \\ &= O(\log(1/\epsilon)) + O(1) . \end{aligned}$$

# SUMMARY: ROBUST MEAN ESTIMATION VIA FILTERING

## **Certificate of Robustness:**

“Spectral norm of empirical covariance is what it should be.”

## **Exploiting the Certificate:**

- Check if certificate is satisfied.
- If violated, find “subspace” where behavior of outliers different than behavior of inliers.
- Use it to detect and remove outliers.
- Iterate on “cleaner” dataset.

## SOFT OUTLIER REMOVAL

Let

$$S_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : 0 \leq w_i \leq \frac{1}{(1-2\epsilon)N} \right\}$$

Let  $\delta = \Theta(\epsilon \log(1/\epsilon))$ . Consider the convex set

$$\mathcal{C}_\delta = \left\{ w \in S_{N,\epsilon} : \left\| \sum_{i=1}^N w_i (X_i - \mu)(X_i - \mu)^T - I \right\|_2 \leq \delta \right\}$$

**Algorithm:**

- Find  $w^* \in \mathcal{C}_\delta$
- Output  $\hat{\mu}_{w^*} = \sum_{i=1}^N w_i^* X_i$ .
- Adaptation of key lemma gives: For all  $w \in \mathcal{C}_\delta$ , we have:

$$\|\hat{\Sigma}_w\|_2 \leq 1 + \delta \quad \longrightarrow \quad \|\hat{\mu}_w - \mu\|_2 \leq O(\epsilon \sqrt{\log(1/\epsilon)})$$

## SOFT OUTLIER REMOVAL

Let

$$S_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : 0 \leq w_i \leq \frac{1}{(1-2\epsilon)N} \right\}$$

Let  $\delta = \Theta(\epsilon \log(1/\epsilon))$ . Consider the convex set

$$\mathcal{C}_\delta = \left\{ w \in S_{N,\epsilon} : \left\| \sum_{i=1}^N w_i (X_i - \mu)(X_i - \mu)^T - I \right\|_2 \leq \delta \right\}$$

**Algorithm:**

- Find  $w^* \in \mathcal{C}_\delta$
- Output  $\hat{\mu}_{w^*} = \sum_{i=1}^N w_i^* X_i$ .

**Main Issue:**  $\mu$  unknown.

## APPROXIMATE SEPARATION ORACLE

**Input:**  $\epsilon$ -corrupted set  $S$  and weight vector  $w$

**Output:** Separation oracle for  $\mathcal{C}_\delta$

- Let  $\delta = \Theta(\epsilon \log(1/\epsilon))$
- Let  $\hat{\mu}_w = \sum_{i=1}^N w_i X_i$  and  $\hat{\Sigma}_w = \sum_{i=1}^N w_i X_i X_i^T - \hat{\mu}_w \hat{\mu}_w^T$
- Let  $(\lambda^*, v^*)$  be the top eigenvalue-eigenvector pair of  $\hat{\Sigma}_w$ .
- If  $\lambda^* \leq 1 + \delta$ , return “YES”.
- Otherwise, return the hyperplane  $L : \mathbb{R}^d \rightarrow \mathbb{R}$  with

$$L(u) = \sum_{i=1}^N u_i ((X_i - \hat{\mu}_w) \cdot v^*)^2 - \lambda^* .$$



# DETERMINISTIC REGULARITY CONDITIONS

Convex program only requires the following conditions:

- For all  $w \in S_{N,\epsilon}$ , the following hold:

$$\left\| \sum_{i \in I_G} w_i (X_i - \mu)(X_i - \mu)^T - I \right\|_2 \leq \delta_1 := \Theta(\epsilon \log(1/\epsilon))$$

$$\left\| \sum_{i \in I_G} w_i (X_i - \mu) \right\|_2 \leq \delta_2 := \Theta(\epsilon \sqrt{\log(1/\epsilon)})$$

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- Sample Complexity versus Robustness
- Certificate of Robustness
- Recursive Dimension Halving
- Iterative Filtering, Soft Outlier Removal
- **Extensions**

## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

## ROBUST MEAN ESTIMATION: *SUB-GAUSSIAN CASE*

**Problem:** Given data  $x_1, \dots, x_N \in \mathbb{R}^d$ , of which  $(1 - \epsilon)N$  come from some distribution  $D$ , estimate mean  $\mu$  of  $D$ .

**Theorem [DKKLMS'17]:** Let  $\epsilon < 1/2$ . If  $N = \Omega(d/\epsilon^2)$  and  $D$  is sub-Gaussian with identity covariance, then can efficiently recover  $\hat{\mu}$  with,

$$\|\hat{\mu} - \mu\|_2 = O(\epsilon \sqrt{\log(1/\epsilon)}) .$$

Information-theoretically **optimal error**, even in one-dimension.

## OPTIMAL GAUSSIAN ROBUST MEAN ESTIMATION?

**Recall [DKKLS'16]:** There is a  $\text{poly}(d/\epsilon)$  time algorithm for robustly learning  $\mathcal{N}(\mu, I)$  within error

$$O(\epsilon\sqrt{\log(1/\epsilon)}) .$$

**Open Question:** Is there a  $\text{poly}(d/\epsilon)$  time algorithm for robustly learning  $\mathcal{N}(\mu, I)$  within error  $O(\epsilon)$ ?

How about

$$o(\epsilon\sqrt{\log(1/\epsilon)}) ?$$

## GAUSSIAN ROBUST MEAN ESTIMATION: ADDITIVE ERRORS

**Theorem [DKKLMS'18]** There is a polynomial time algorithm with the following behavior: Given  $\epsilon > 0$  and  $N = \text{poly}(d/\epsilon)$  corrupted samples from an unknown mean, identity covariance Gaussian distribution on  $\mathbb{R}^d$ , the algorithm finds a hypothesis mean  $\hat{\mu}$  that satisfies

$$\|\mu - \hat{\mu}\|_2 \leq \sqrt{\pi} \cdot \epsilon + o(\epsilon)$$

in *additive* contamination model.

- Robustness guarantee optimal up to  $\sqrt{2}$  factor!
- For any univariate projection, mean robustly estimated by median.

## GENERALIZED FILTERING: ADDITIVE CORRUPTIONS

- *Univariate* filtering based on tails not sufficient to remove the incurred  $\Omega(\epsilon\sqrt{\log(1/\epsilon)})$  error, even for additive errors.
- **Generalized Filter Idea:** Filter using top -  $k$  eigenvectors of empirical covariance.
- **Key Observation:** Suppose that  $\|\mu - \hat{\mu}\|_2 \geq \epsilon$ . Then either
  - (1)  $\hat{\Sigma}$  has  $k$  eigenvalues at least  $1 + \Omega(\epsilon)$ , or
  - (2) The error comes from a  $k$ -dimensional subspace.
- Choose  $k = \Theta(\log(1/\epsilon))$ .

# COMPUTATIONAL LIMITATIONS TO ROBUST MEAN ESTIMATION

**Theorem [DKS'17]** Suppose  $d \geq \text{polylog}(1/\epsilon)$ . Any *Statistical Query\** algorithm that learns an  $\epsilon$ -corrupted Gaussian  $\mathcal{N}(\mu, I)$  in the *strong* contamination model within distance

$$o(\epsilon \sqrt{\log(1/\epsilon)})$$

requires runtime

$$d^{\omega(1)} .$$

\*Instead of accessing samples from distribution  $D$ , a Statistical Query algorithm can adaptively query  $\mathbb{E}_{x \sim D}[f(x)]$ , for any  $f : \mathbb{R}^d \rightarrow [0, 1]$

**Take-away:** Any asymptotic improvement in error guarantee over [DKKLMS'16] algorithms may require super-polynomial time.

## ROBUST MEAN ESTIMATION: GENERAL CASE

**Problem:** Given data  $x_1, \dots, x_N \in \mathbb{R}^d$ , of which  $(1 - \epsilon)N$  come from some distribution  $D$ , estimate mean  $\mu$  of  $D$ .

**Theorem [DKKLMS'17, CSV'18]** Let  $\epsilon < 1/2$ . If  $N = \Omega(d/\epsilon)$ , and  $D$  has covariance  $\Sigma \preceq \sigma^2 \cdot I$ , then we can efficiently recover  $\hat{\mu}$  with,

$$\|\hat{\mu} - \mu\|_2 = O(\sigma \cdot \sqrt{\epsilon}).$$

- **Sample-optimal**, even without corruptions.
- Information-theoretically **optimal error**, even in one-dimension.
- Adaptation of Iterative Filtering.



## ROBUST COVARIANCE ESTIMATION

**Problem:** Given data  $x_1, \dots, x_N \in \mathbb{R}^d$ , of which  $(1 - \epsilon)N$  come from some distribution  $D$ , estimate covariance  $\Sigma$  of  $D$ .

**Theorem:** Let  $\epsilon < 1/2$ . If  $N = \Omega(d^2/\epsilon^2)$ , then can efficiently recover  $\hat{\Sigma}$  such that

$$\|\Sigma^{-1/2}(\hat{\Sigma} - \Sigma)\Sigma^{-1/2}\|_F = f(\epsilon),$$

where  $f$  depends on the concentration of  $D$ .

**Main Idea:** Use *fourth-order moment tensors* !

# OUTLINE

## **Part I: Introduction**

- Motivation
- Robust Statistics in Low and High Dimensions
- This Talk

## **Part II: High-Dimensional Robust Mean Estimation**

- Sample Complexity versus Robustness
- Certificate of Robustness
- Recursive Dimension Halving
- Iterative Filtering, Soft Outlier Removal
- Extensions

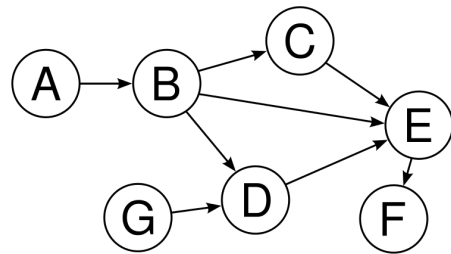
## **Part III: Summary and Conclusions**

- Beyond Robust Statistics: Unsupervised and Supervised Learning
- Future Directions

## SUMMARY AND CONCLUSIONS

- High-Dimensional Computationally Efficient Robust Estimation is Possible!
- First Computationally Efficient Robust Estimators with **Dimension-Independent** Error Guarantees.
- General Methodologies for High-Dimensional Estimation Problems.

# BEYOND ROBUST STATISTICS: ROBUST *UNSUPERVISED* LEARNING

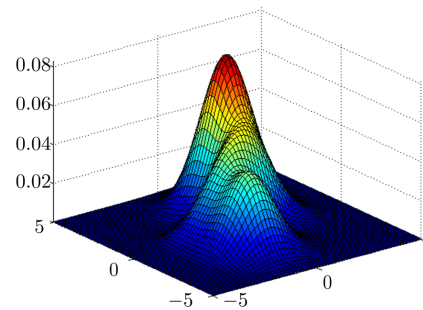


Robustly Learning Graphical Models  
[Cheng-D-Kane-Stewart'16,  
D-Kane-Stewart'18]

Clustering in Mixture Models  
[Charikar-Steinhardt-Valiant'17,  
D-Kane-Stewart'18,  
Hopkins-Li'18,  
Kothari-Steinhardt-Steurer'18]

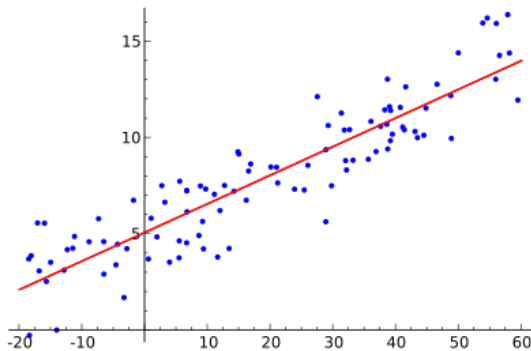
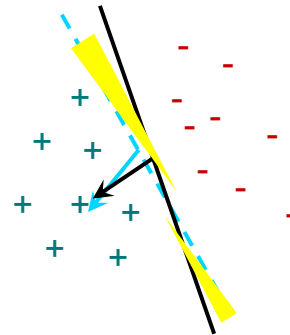


Computational/Statistical-Robustness Tradeoffs  
[D-Kane-Stewart'17, D-Kong-Stewart'18]

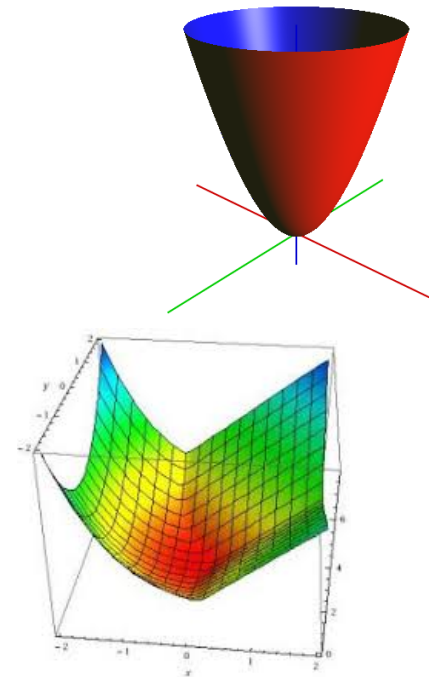


# BEYOND ROBUST STATISTICS: ROBUST *SUPERVISED* LEARNING

Malicious PAC Learning  
[Klivans-Long-Servedio'10,  
Awasthi-Balcan-Long'14,  
**D-Kane-Stewart'18**]



Robust Linear Regression  
[**D-Kong-Stewart'18**,  
Klivans-Kothari-Meka'18]



Stochastic (Convex) Optimization  
[Prasad-Suggala-Balakrishnan-Ravikumar'18,  
**D-Kamath-Kane-Li-Steinhardt-Stewart'18**]

## RELATED WORKS

- [Known Structure Bayes Nets](#) [Cheng-D-Kane-Stewart'16]
- [Sparse models \(e.g., sparse PCA, sparse regression\)](#) [Li'17, Du-Balakrishan-Singh'17, Liu-Shen-Li-Caramanis'18]
- [List-Decodable Learning](#) [Charikar-Steinhardt-Valiant '17, Meister-Valiant'18]
- [Robust PAC Learning](#) [Klivans-Long-Servedio'10, Awasthi-Balcan-Long'14, D-Kane-Stewart'18]
- [“Robust estimation via SoS” \(higher moments, learning mixture models\)](#) [Hopkins-Li'18, Kothari-Steinhardt-Steurer'18]
- [“SoS Free” learning of mixture models](#) [D-Kane-Stewart'18]
- [Robust Regression](#) [Klivans-Kothari-Meka'18, D-Kong-Stewart'18]
- [Robust Stochastic Optimization](#) [Prasad-Suggala-Balakrishnan-Ravikumar'18, D-Kamath-Kane-Li-Steinhardt-Stewart'18]
- ...

## FUTURE DIRECTIONS

### General Algorithmic Theory of Robustness

How can we robustly learn rich representations of data, based on natural hypotheses about the structure in data?

Can we robustly *test* our hypotheses about structure in data before learning?

#### **Concrete Challenges:**

- Richer Families of Problems and Models
- Connections to Non-convex Optimization
- Relation to other Notions of Algorithmic Stability  
(Differential Privacy, Adaptive Data Analysis)
- Further applications (ML Security, Computer Vision).

**Thank you!**  
**Questions?**