

Algorithmic Robust Statistics

Ilias Diakonikolas (UW Madison)
Mathematical Methods of Statistics, CIRM
December 2023

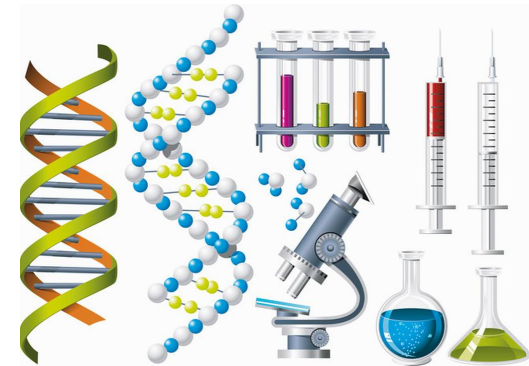
Can we develop learning algorithms that are *robust* to a *constant* fraction of *corruptions* in the data?

MOTIVATION

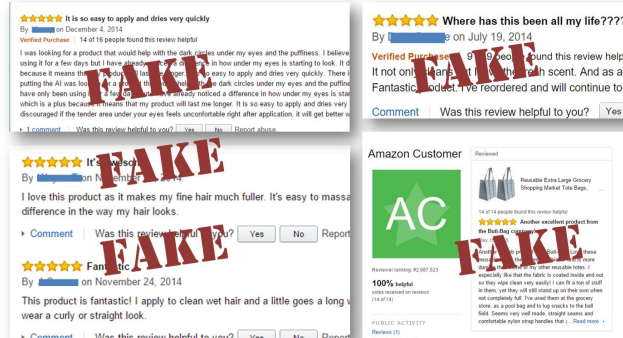
- **Model Misspecification/Robust Statistics**
[Fisher 1920s, Tukey 1960s, Huber 1960s]

- **Outlier Detection/Removal**

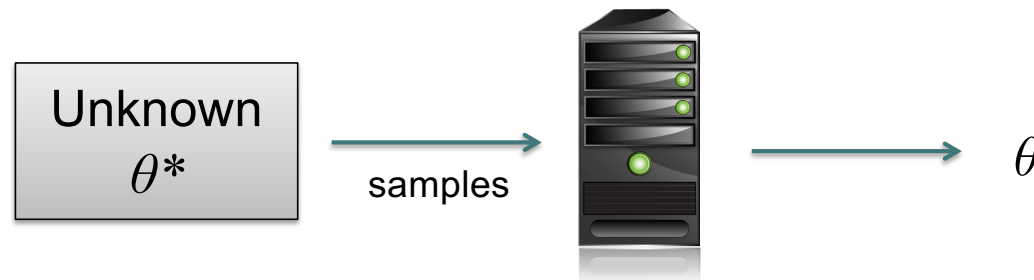
- **Adversarial/Secure ML**



So Many Misleading, “Fake” Reviews



THE STATISTICAL LEARNING PROBLEM



- *Input:* sample generated by a **statistical model** with unknown θ^*
- *Goal:* estimate parameters θ 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?

(OUTLIER-) ROBUSTNESS

Strong Contamination Model:

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

We say that a set of N samples is ϵ -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 ϵ -fraction of them.

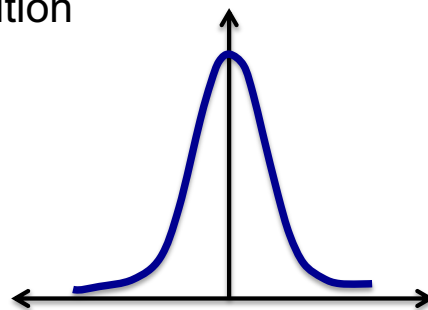
cf. Huber's contamination model [1964]

EXAMPLE: PARAMETER ESTIMATION

Given i.i.d. 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:

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

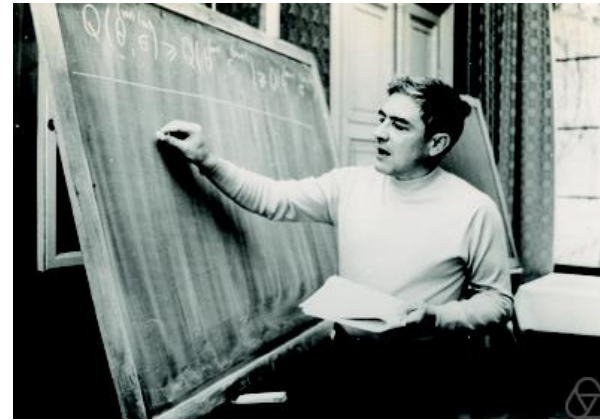
empirical variance:

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



John W. Tukey

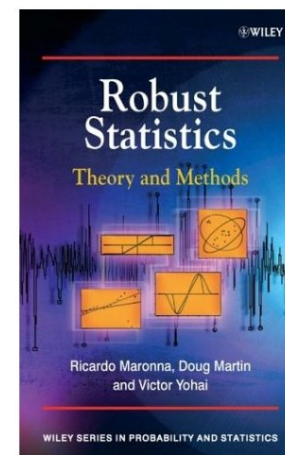
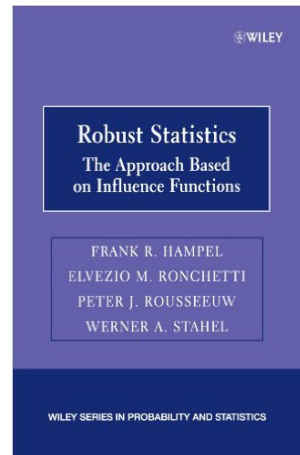
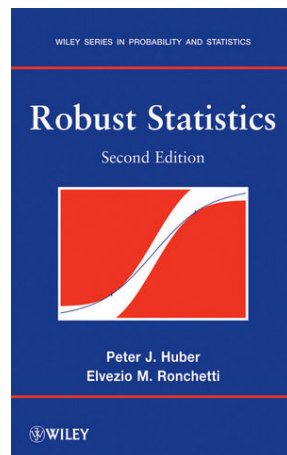
Model Misspecification
(1960s)



Peter J. Huber

Robust Estimation of Location
(1964)

ROBUST STATISTICS



What estimators behave well in the presence of outliers?

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 ϵ -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*?

HIGH-DIMENSIONAL ROBUST MEAN ESTIMATION

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

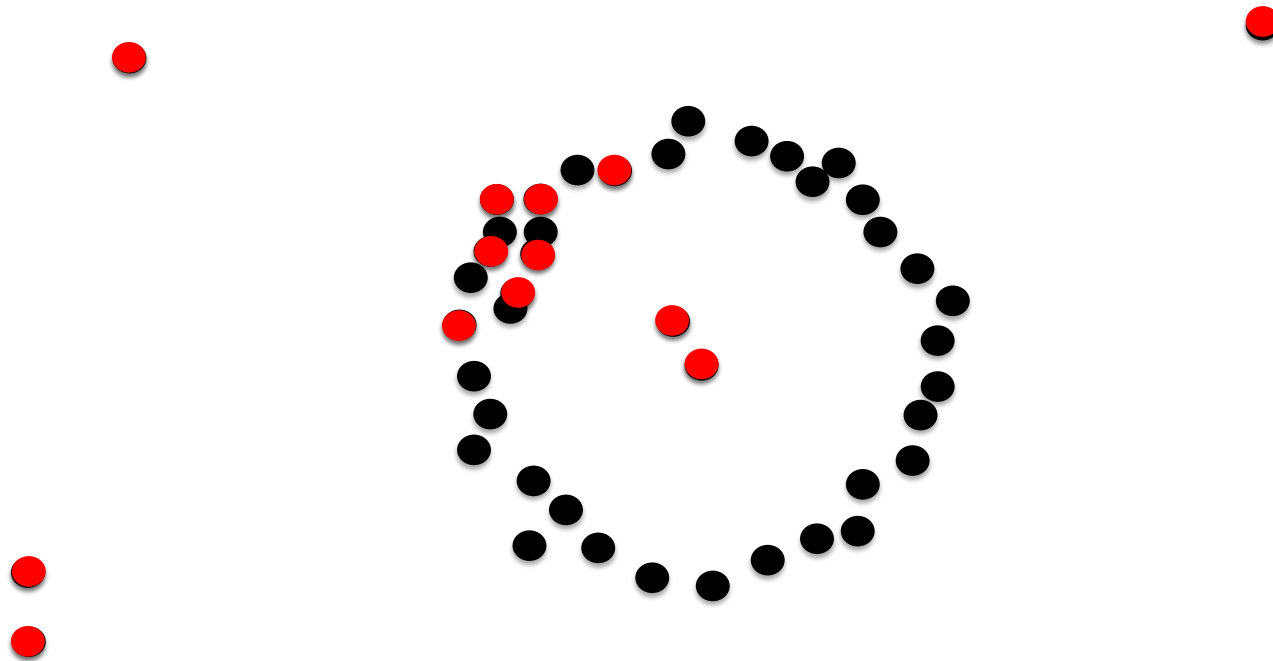
$$\|\hat{\mu} - \mu\|_2 = O(\epsilon) + O(\sqrt{d/N})$$

Remark: Above convergence rate is optimal [Tukey'75, Donoho'82]

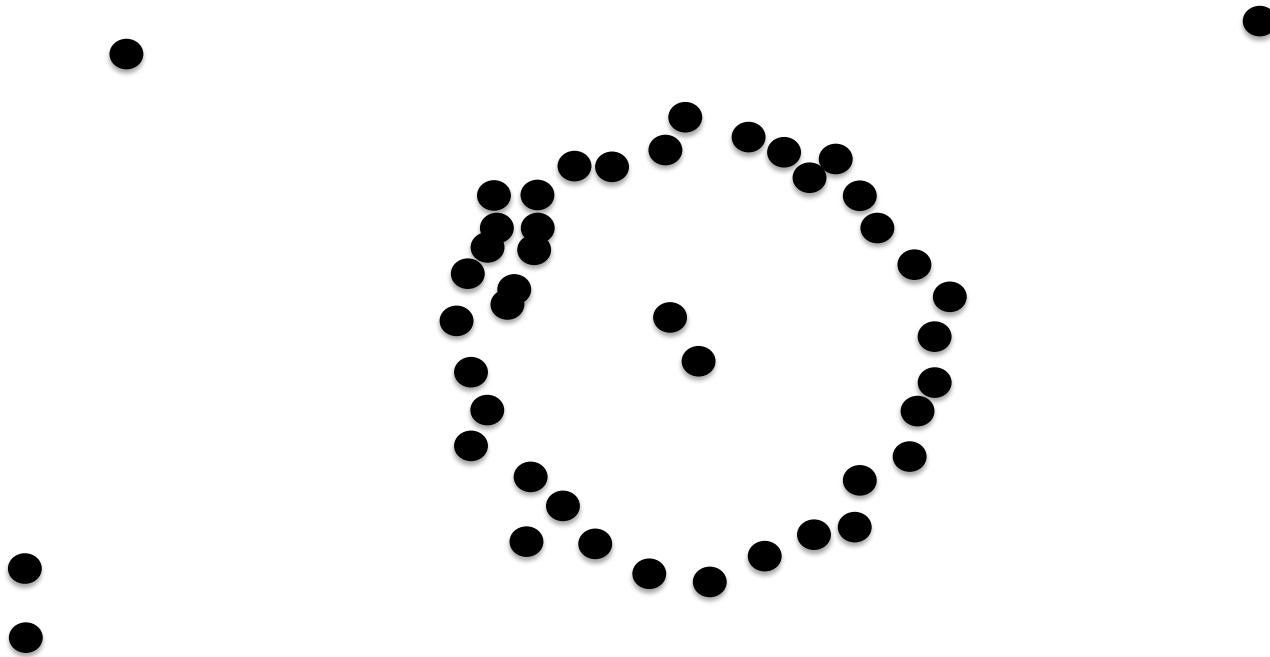
PREVIOUS APPROACHES: ROBUST MEAN ESTIMATION

Estimator	Error Rate	Running Time
Distance-Based Pruning	$\Theta(\epsilon\sqrt{d})$ ✗	$O(dN)$ ✓
Coordinate-wise Median	$\Theta(\epsilon\sqrt{d})$ ✗	$O(dN)$ ✓
Geometric Median	$\Theta(\epsilon\sqrt{d})$ ✗	$\text{poly}(d, N)$ ✓
Tukey Median	$\Theta(\epsilon)$ ✓	NP-Hard ✗
Tournament	$\Theta(\epsilon)$ ✓	$N^{O(d)}$ ✗

DISTANCE-BASED PRUNING



DISTANCE-BASED PRUNING = NAÏVE OUTLIER REMOVAL

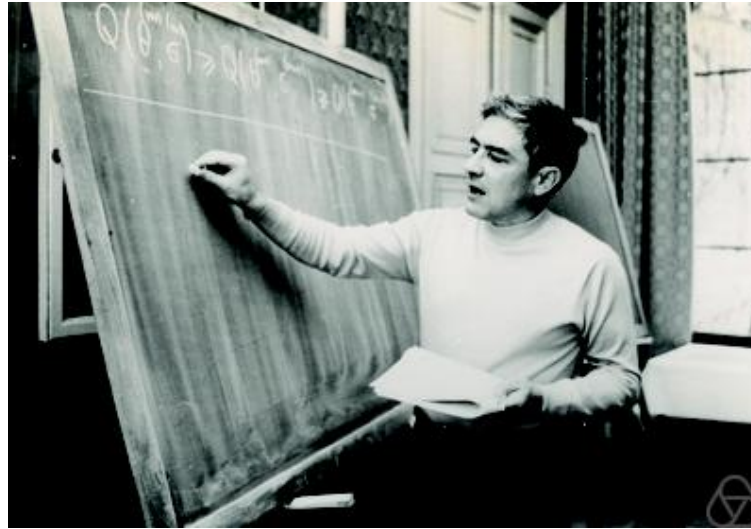


HIGH-DIMENSIONAL ROBUST STATISTICS: 1960-2016

All known estimators either **require exponential time to compute**
or can tolerate a **negligible fraction of outliers**.

Is robust estimation *algorithmically* possible in high-dimensions?

Peter J. Huber, 1975



“The bad news is that with **all currently known algorithms the effort of computing those estimates increases exponentially in d . We might say they break down by failing to give a timely answer!**

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.[...]”

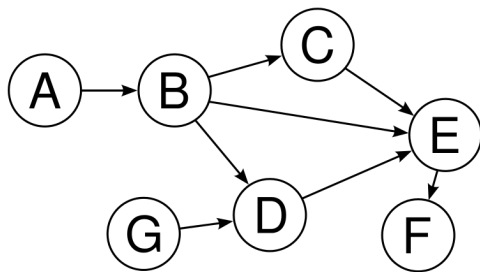
Robust Statistical Procedures, 1996, *Second Edition*.

Meta-Theorem [D-Kamath-Kane-Li-Moitra-Stewart'16]

Efficient robust estimators with *dimension-independent* error for robust mean and covariance estimation, if inlier distribution has bounded moments/nice concentration.

Related results by [Lai-Rao-Vempala'16]

ROBUST *UNSUPERVISED* LEARNING

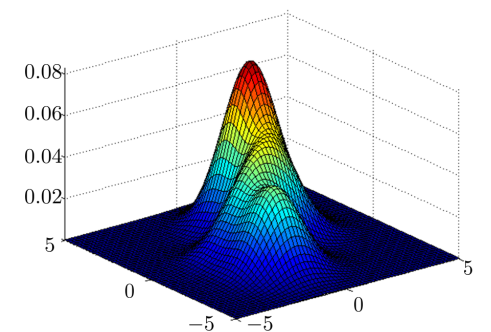


Robustly Learning Graphical Models



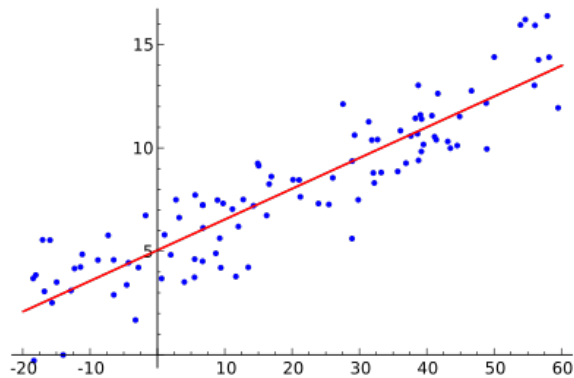
Computational/Statistical-Robustness Tradeoffs

List-decodable Learning and
Robustly Learning Mixture Models

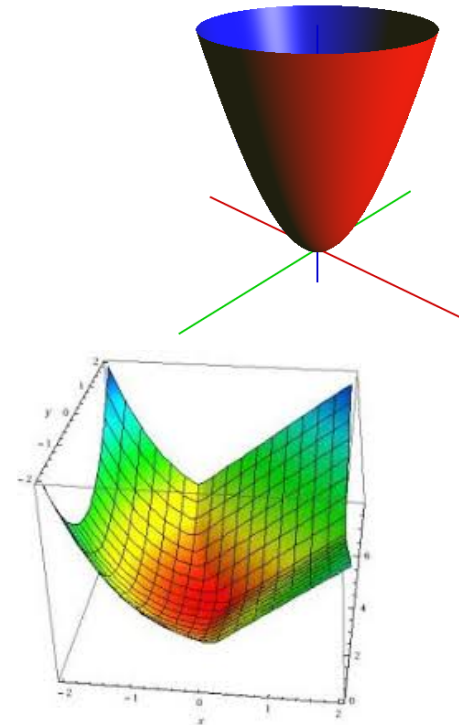
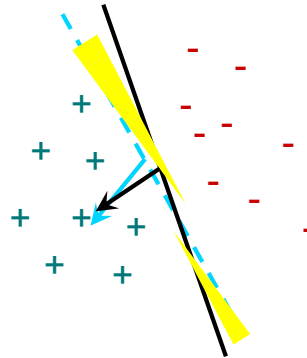


ROBUST *SUPERVISED* LEARNING

Robust Supervised
Learning



Robust Regression

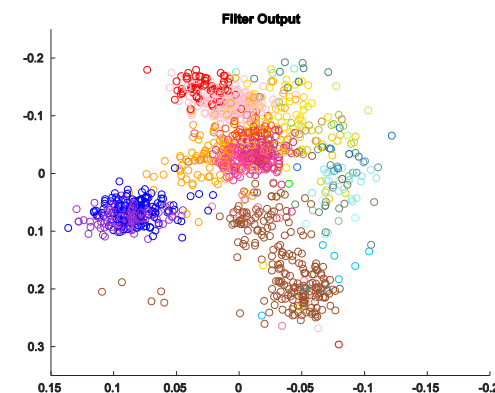
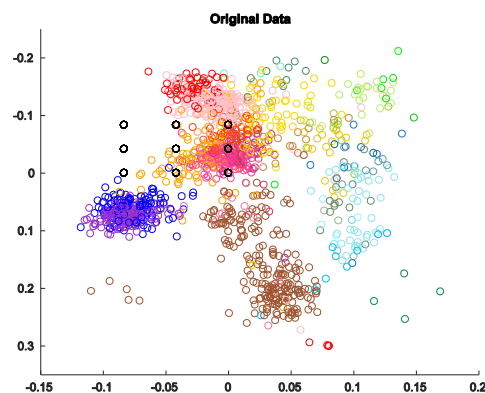


Stochastic Convex Optimization

APPLICATIONS

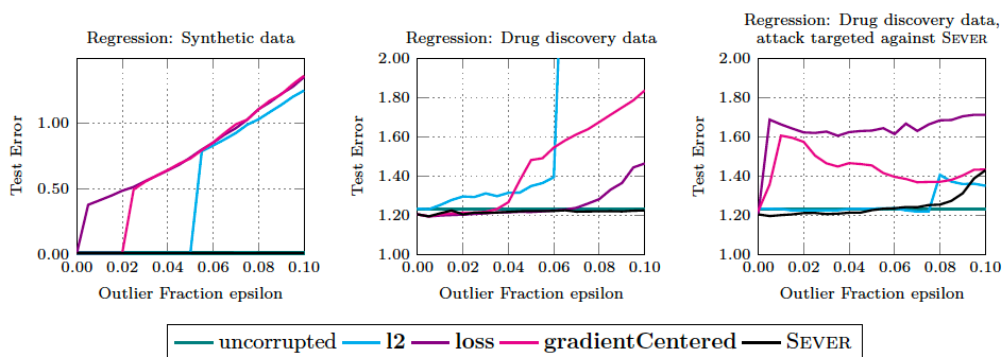
[D-Kamath-Kane-Li-Moitra-Stewart, ICML'17]

Detecting Patterns in Biological Data

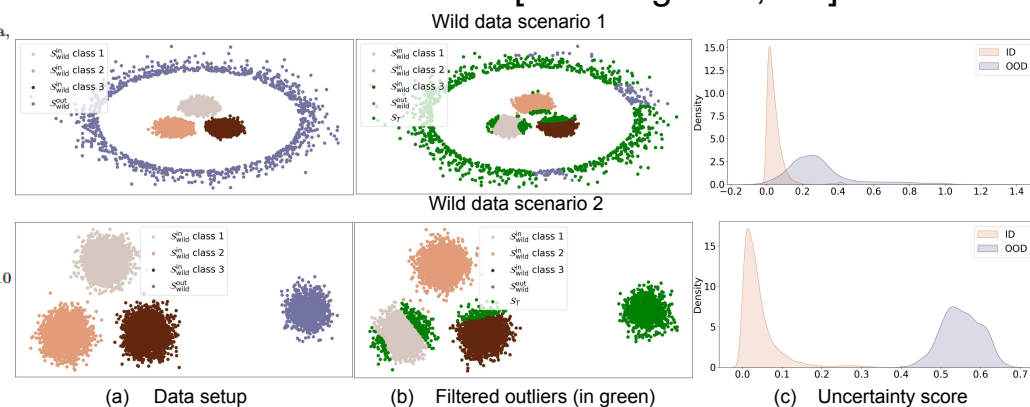


Provable Defenses against Data Poisoning

[D-Kamath-Kane-Li-Moitra-Steinhardt, ICML'19]



OOD Detection [Du-Fang-D-Li, '23]

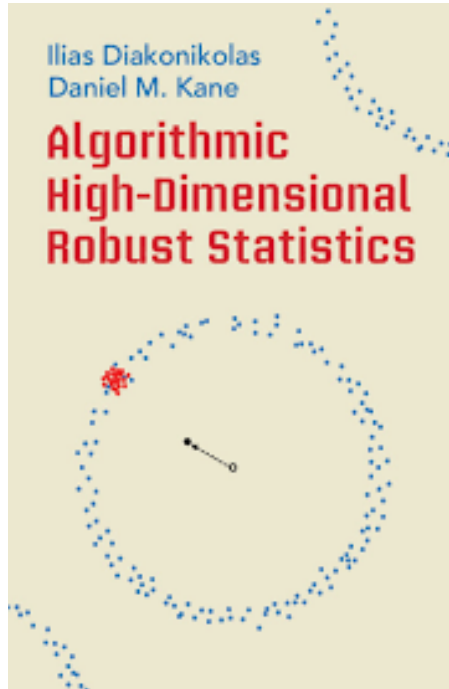


SUBSEQUENT WORKS

- **Sparse Models** [Balakrishnan-Du-Li-Singh'17, D-Karmalkar-Kane-Price-Stewart'19, D-Kane-Lee-Pensia'22,...]
- **Graphical Models** [Cheng-D-Kane-Stewart'18, D-Kane-Stewart-Sun'21, D-Kane-Sun'22]
- **Robust Regression/Classification** [D-Kane-Stewart'18, Klivans-Kothari-Meka'18, D-Kong-Stewart'19 Bakshi-Prasad'21, ...]
- **Robust Stochastic Optimization** [Prasad-Suggala-Balakrishnan-Ravikumar'19, D-Kamath-Kane-Li-Steinhardt-Stewart'19, ...]
- **Robust Estimation via SoS** [Hopkins-Li'18, Kothari-Steinhardt-Steurer'18, Karmalkar-Klivans-Kothari'19, Raghavendra-Yau'19, Bakshi-Kothari'20, D-Hopkins-Kane-Karmalkar'20, Liu-Moitra'21, Bakshi-D-Jia-Kane-Kothari-Vempala'21, Ivkov-Kothari'22, ...]
- **Near-Linear Time Algorithms** [Chen-D-Ge'18, Cheng-D-Ge-Woodruff'19, Depersin-Lecue'19, Dong-Hopkins-Li'19, Li-Ye'20, Cherapanamjeri-Mohanty-Yau'20, D-Kane-Koongsgard-Li-Tian'21, ...]
- **Computational-Statistical Tradeoffs** [D-Kane-Stewart'17, D-Kong-Stewart'19, Hopkins-Li'19, ...]
- **Connections to Non-Convex Optimization** [Chen-D-Ge-Soltanolkotabi'20, Zhu-Jiao-Steinhardt'20, ...]
- **List-Decodable Learning** [Charikar-Steinhardt-Valiant'17, D-Kane-Stewart'18, Meister-Valiant'18, Karmalkar-Klivans-Kothari'19, Raghavendra-Yau'19, D-Kane-Koongsgard'20, D-Kane-Koongsgard-Li-Tian'21, D-Kane-Karmalkar-Pensia-Pittas'22]
- **Applications in Data Analysis** [D-Kamath-Kane-Li-Moitra-Stewart'17, Tran-Li-Madry'18, D-Kamath-Kane-Li-Steinhardt-Stewart'19, Hayase-Kong-Somani-Oh'21, Du-Fang-D-Li'23, ...]

Ilias Diakonikolas
Daniel M. Kane

Algorithmic High-Dimensional Robust Statistics



HIGH-DIMENSIONAL ROBUST MEAN ESTIMATION

ROBUST MEAN ESTIMATION: GAUSSIAN CASE

Problem: Given an ϵ -corrupted set of points $x_1, \dots, x_N \in \mathbb{R}^d$ from an unknown distribution D in a known family \mathcal{F} , estimate the mean μ of D .

Theorem 1: Let $\epsilon < 1/2$. If D is a spherical Gaussian, there is an efficient algorithm that outputs an estimate $\hat{\mu}$ that with high probability satisfies

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

in the **additive contamination** model.

First-term of RHS Independent of d !

[D-Kamath-Kane-Li-Moitra-Stewart, SODA'18; D-Kane-Pensia-Pittas, NeurIPS'23]

ROBUST MEAN ESTIMATION: *SUB-GAUSSIAN* CASE

Problem: Given an ϵ -corrupted set of points $x_1, \dots, x_N \in \mathbb{R}^d$ from an unknown distribution D in a known family \mathcal{F} , estimate the mean μ of D .

Theorem 2: Let $\epsilon < 1/2$. If D is a spherical *sub-Gaussian*, there is an efficient algorithm that outputs an estimate $\hat{\mu}$ that with high probability satisfies

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

in the **strong contamination** model.

Information-theoretically **optimal error**.

[D-Kamath-Kane-Li-Moitra-Stewart, FOCS'16, ICML'17; D-Kane-Pensia-Pittas, ICML'22]

ROBUST MEAN ESTIMATION: BOUNDED COVARIANCE CASE

Problem: Given an ϵ -corrupted set of points $x_1, \dots, x_N \in \mathbb{R}^d$ from an unknown distribution D in a known family \mathcal{F} , estimate the mean μ of D .

Theorem 3: Let $\epsilon < 1/2$. If D has covariance $\Sigma \preceq I$, there is an efficient algorithm that outputs an estimate $\hat{\mu}$ that with high probability satisfies

$$\|\hat{\mu} - \mu\|_2 = O(\sqrt{\epsilon} + \sqrt{d/N}) .$$

in the **strong contamination** model.

Information-theoretically **optimal error**.

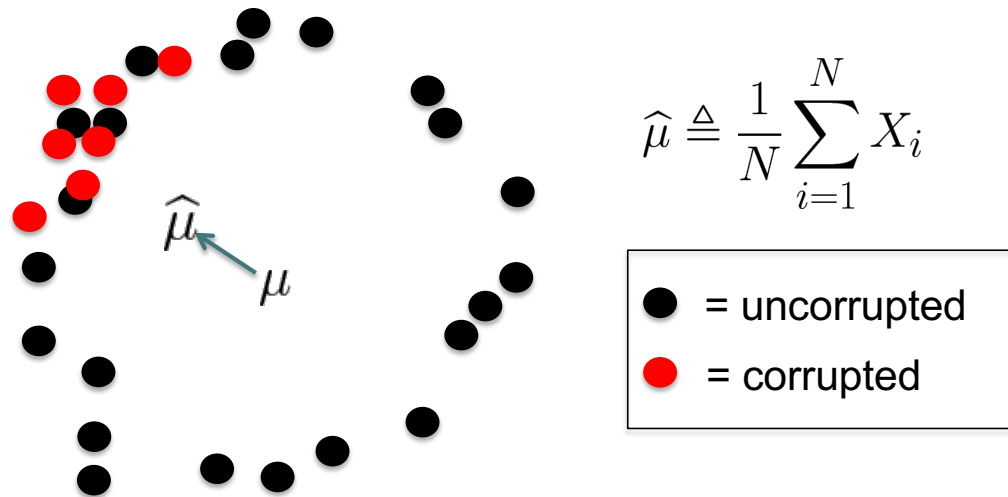
[D-Kamath-Kane-Li-Moitra-Stewart, ICML'17; Steinhardt, Charikar, Valiant, ITCS'18]

CERTIFICATE OF ROBUSTNESS FOR EMPIRICAL ESTIMATOR

Idea #1: If the empirical covariance is “close to what it should be”, then the empirical mean works.

CERTIFICATE FOR EMPIRICAL MEAN

Detect when the empirical estimator *may* be compromised



There is *no* direction of large empirical variance

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

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

with high probability we have:

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

in **strong** contamination model.

Idea #2: Removing *any* ϵ - fraction of inliers does not move the empirical mean and covariance by much.

Idea #3: Iteratively “remove outliers” to “fix” the empirical covariance.

ITERATIVE FILTERING

Iterative Two-Step Procedure:

Step #1: Test certificate of robustness of “standard” estimator

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

Iterate on “cleaner” dataset.

General recipe that works in general settings.

We'll see how this works for robust mean estimation.

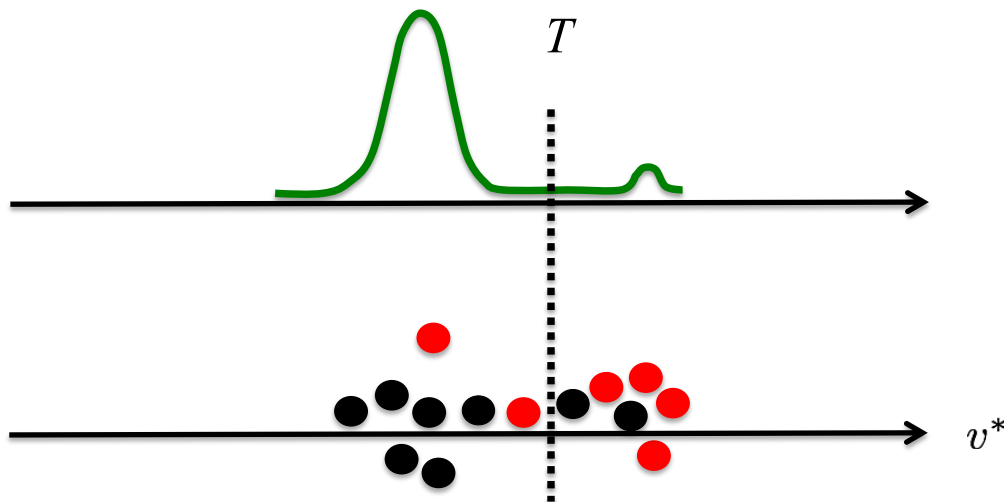
FILTERING SUBROUTINE

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.



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.

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

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

- 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: In each iteration, we remove more outliers than inliers.

After a bounded number of iterations, we stop removing points.

Eventually the empirical mean works

Runtime: $\tilde{O}(Nd^2)$

STABILITY CONDITION

Definition Fix $0 < \epsilon < 1/2$ and $\delta \geq \epsilon$. A set $S \subset \mathbb{R}^d$ is (ϵ, δ) -stable with respect to μ if for all $v \in \mathbb{S}^{d-1}$ and every $S' \subseteq S$ such that $|S'| \geq (1 - \epsilon)|S|$, we have:

- $\left| \frac{1}{|S'|} \sum_{x \in S'} v \cdot (x - \mu) \right| \leq \delta \iff \|\mu_{S'} - \mu\|_2 \leq \delta$
- $\left| \frac{1}{|S'|} \sum_{x \in S'} (v \cdot (x - \mu))^2 - 1 \right| \leq \delta^2/\epsilon \iff \|\bar{\Sigma}_{S'} - I\|_2 \leq \delta^2/\epsilon$

- Intended for inlier distributions with $\Sigma \preceq I$
- Similar definition for distributions as opposed to datasets.
- A sufficiently large clean sample from a well-behaved distribution is stable with high probability.

EFFICIENT ROBUST MEAN ESTIMATION UNDER STABILITY

General Theorem Let \mathcal{S} be (ϵ, δ) –stable with respect to a vector μ , and T an ϵ -corruption of \mathcal{S} . There is an efficient algorithm that given ϵ, δ, T it computes an estimate $\hat{\mu}$ such that

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

Fact A set of N i.i.d. samples from a well-behaved distribution is (ϵ, δ) –stable with high probability.

- For identity covariance sub-Gaussians, $\delta \sim \epsilon \sqrt{\log(1/\epsilon)}$ and $N \gg d/\delta^2$
- For identity covariance sub-exponentials, $\delta \sim \epsilon \log(1/\epsilon)$ and $N \gg d/\delta^2$
- For identity covariance with bounded k –th central moments ($k \geq 4$), $\delta \sim \epsilon^{1-1/k}$ and $N \gg d(\log d)/\delta^2$
- For *bounded* covariance distributions, $\delta \sim \sqrt{\epsilon}$ and $N \gg d(\log d)/\delta^2$
(after removing ϵ - fraction of inliers)

CERTIFICATE FOR EMPIRICAL MEAN

Lemma Let S be (ϵ, δ) -stable with respect to μ , and T be an ϵ -corruption of S .

If $\|\Sigma_T\|_2 \leq 1 + \lambda$, for $\lambda \geq 0$, then

$$\|\mu_T - \mu\|_2 \leq O(\delta + \sqrt{\epsilon\lambda})$$

Proof Let X, Y be uniform distribution over S, T respectively. Can write $Y = (1 - \epsilon)X' + \epsilon E$, where X' is ϵ -subtraction of X .

$$\Sigma_Y = (1 - \epsilon)\Sigma_{X'} + \epsilon\Sigma_E + \epsilon(1 - \epsilon)(\mu_{X'} - \mu_E)(\mu_{X'} - \mu_E)^\top$$

Let v be normalized version of $\mu_{X'} - \mu_E$.

$$\begin{aligned} 1 + \lambda &\geq v^\top \Sigma_Y v = (1 - \epsilon)v^\top \Sigma_{X'} v + \epsilon v^\top \Sigma_E v + \epsilon(1 - \epsilon)v^\top (\mu_{X'} - \mu_E)(\mu_{X'} - \mu_E)^\top v \\ &\geq (1 - \epsilon)(1 - \delta^2/\epsilon) + \epsilon(1 - \epsilon)\|\mu_{X'} - \mu_E\|_2^2 \\ &\geq 1 - O(\delta^2/\epsilon) + (\epsilon/2)\|\mu_{X'} - \mu_E\|_2^2 \end{aligned}$$

Rearranging

$$\|\mu_{X'} - \mu_E\|_2 = O(\delta/\epsilon + \sqrt{\lambda/\epsilon})$$

CERTIFICATE FOR EMPIRICAL MEAN

Lemma Let S be (ϵ, δ) -stable with respect to μ , and T be an ϵ -corruption of S .

If $\|\Sigma_T\|_2 \leq 1 + \lambda$, for $\lambda \geq 0$, then

$$\|\mu_T - \mu\|_2 \leq O(\delta + \sqrt{\epsilon\lambda})$$

Proof Let X, Y be uniform distribution over S, T respectively. Can write $Y = (1 - \epsilon)X' + \epsilon E$, where X' is ϵ -subtraction of X .

$$\|\mu_{X'} - \mu_E\|_2 = O(\delta/\epsilon + \sqrt{\lambda/\epsilon})$$

For the means, have that $\mu_T = \mu_Y = (1 - \epsilon)\mu_{X'} + \epsilon\mu_E$.

$$\begin{aligned} \|\mu_T - \mu\|_2 &= \|(1 - \epsilon)\mu_{X'} + \epsilon\mu_E - \mu\|_2 = \|\mu_{X'} - \mu + \epsilon(\mu_E - \mu_{X'})\|_2 \\ &\leq \|\mu_{X'} - \mu\|_2 + \epsilon\|\mu_{X'} - \mu_E\|_2 \\ &= O(\delta) + \epsilon \cdot O(\delta/\epsilon + \sqrt{\lambda/\epsilon}) \end{aligned}$$



RANDOMIZED FILTERING: IDEA

Main Idea: Suppose we can find $f : T \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\sum_{x \in T} f(x) \geq 2 \sum_{x \in S} f(x) .$$

Then we can randomly filter by removing each point $x \in T$ with probability $\propto f(x)$.

Need this property to hold across iterations, assuming certificate not satisfied.

Condition Given any $T' \subseteq T$ such that $|T' \cap S| \geq (1 - 4\epsilon)|S|$, if $\|\Sigma_{T'}\|_2 \geq 1 + \lambda$ there is an explicit $f : T' \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\sum_{x \in T'} f(x) \geq 2 \sum_{x \in T' \cap S} f(x)$$

RANDOMIZED FILTERING: PROPERTIES

Condition Given any $T' \subseteq T$ such that $|T' \cap S| \geq (1 - 4\epsilon)|S|$, if $\|\Sigma_{T'}\|_2 \geq 1 + \lambda$ there is an explicit $f : T' \rightarrow \mathbb{R}_{\geq 0}$ such that

$$\sum_{x \in T'} f(x) \geq 2 \sum_{x \in T' \cap S} f(x)$$

Theorem If condition holds, there is an efficient randomized algorithm that computes an estimate $\hat{\mu}$ such that with high probability

$$\|\hat{\mu} - \mu_X\|_2 = O(\delta + \sqrt{\epsilon\lambda})$$

RANDOMIZED FILTERING

Randomized Filtering Pseudocode

1. Compute $\nu = \|\Sigma_T\|_2$
2. If $\nu \leq 1 + \lambda$, return μ_T
3. Else
 - Compute the function f .
 - Remove each $x \in T$ with probability $f(x) / \max_{x \in T} f(x)$
 - Return to Step 1 with new set T .

RANDOMIZED FILTERING: ANALYSIS

At least one point is removed in each iteration, so algorithm runs in polynomial time.

Claim With probability at least $2/3$, throughout the algorithm have that $|S \cap T_i| \geq (1 - 4\epsilon)|S|$.

Proof Consider

Have
$$d(T_i) := |(S \cap T) \setminus T_i| + |T_i \setminus S| .$$

$$d(T_i) - d(T_{i-1}) = (\# \text{Inliers removed in iteration } i) - (\# \text{Outliers removed in iteration } i)$$

$$\mathbf{E}[d(T_i) - d(T_{i-1})] = \sum_{x \in S \cap T_i} f(x) - \sum_{x \in T_i \setminus S} f(x) = 2 \sum_{x \in S \cap T_i} f(x) - \sum_{x \in T_i} f(x) \leq 0 .$$

Since $d(T_i) \geq 0$ and $\mathbf{E}[d(T_i)] \leq \mathbf{E}[d(T_0)] \leq \epsilon|S|$, by Ville's inequality

$$\Pr[\max_i d(T_i) > 3\epsilon|S|] \leq 1/3 .$$

This implies that $|S \cap T_i| \geq (1 - 4\epsilon)|S|$ throughout. ■

FINDING f : UNIVERSAL FILTERING

Proposition Let S be $(2\epsilon, \delta)$ -stable and T be an ϵ -corruption of S . Suppose that $\|\Sigma_T\|_2 = 1 + \lambda > 1 + 8\delta^2/\epsilon$. There exists an efficient algorithm that given ϵ, δ, T it computes a function $f : T \rightarrow \mathbb{R}_{\geq 0}$ such that

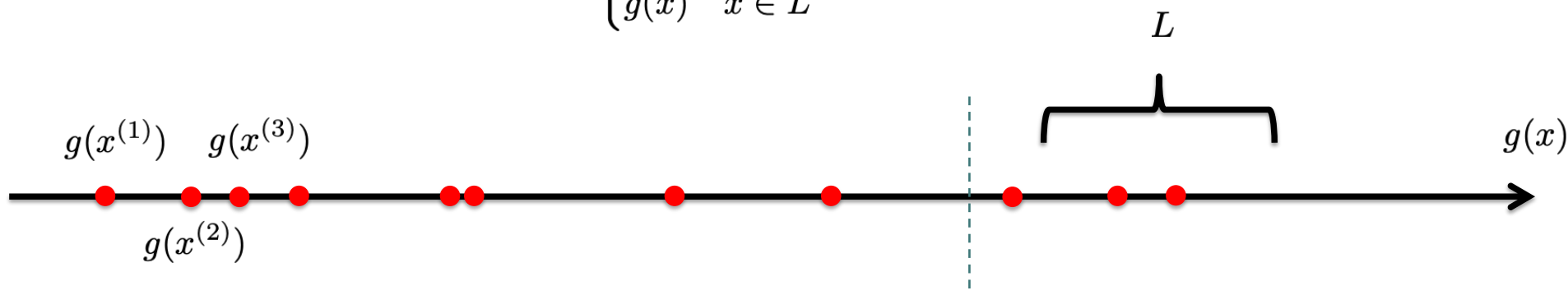
$$\sum_{x \in T} f(x) \geq 2 \sum_{x \in T \cap S} f(x).$$

Proof Define the function $g(x) = (v \cdot (x - \mu_T))^2$, where v is the top eigenvector.

Let L be the set of $\epsilon \cdot |T|$ points $x \in T$ for which $g(x)$ is largest.

Then

$$f(x) = \begin{cases} 0 & x \notin L \\ g(x) & x \in L \end{cases}$$



UNIVERSAL FILTERING: ANALYSIS

- By definition $\sum_{x \in T} g(x) = |T| \mathbf{Var}[v \cdot T] = |T|(1 + \lambda)$

and $\sum_{x \in S} g(x) = |S|(\mathbf{Var}[v \cdot S] + (v \cdot (\mu_T - \mu_S))^2)$

- By stability and our lemma $\sum_{x \in S} g(x)$ is small so that

$$\sum_{x \in T \setminus S} g(x) \geq \sum_{x \in T} g(x) - \sum_{x \in S} g(x) \geq (2/3)|S|\lambda.$$

- By the definition of L and λ

$$\sum_{x \in T} f(x) = \sum_{x \in L} g(x) \geq \sum_{x \in T \setminus S} g(x)$$

- Similarly

$$\begin{aligned} \sum_{x \in S \cap T} f(x) &= \sum_{x \in S \cap L} g(x) = \sum_{x \in S} g(x) - \sum_{x \in S \setminus L} g(x) \\ &\leq 2|S|\delta^2/\epsilon + |S|O(\delta^2 + \epsilon\lambda) \end{aligned}$$

WEIGHTED FILTERING

Assign *weights* to the samples so that weighted empirical mean works.

For $w : T \rightarrow \mathbb{R}_+$

$$\mu_w[T] := \frac{1}{\|w\|_1} \sum_{x \in T} w_x x \quad \Sigma_w[T] := \frac{1}{\|w\|_1} \sum_{x \in T} w_x (x - \mu_w)(x - \mu_w)^\top$$

Weighted Filtering Pseudocode

1. Set $t = 1$ and $w_x^{(1)} = 1/|T|$ for $x \in T$
2. While $\|\Sigma_{w^{(t)}}[T]\|_2 > 1 + \lambda$
 - Compute the function f .
 - Set $f_{\max} = \max\{f(x) \mid x \in T \text{ and } w_x^{(t)} \neq 0\}$
 - Set $w_x^{(t+1)} = w_x^{(t)}(1 - f(x)/f_{\max})$
 - Set t to $t + 1$
3. Return $\mu_{w^{(t)}}$

NON-CONVEX OPTIMIZATION FORMULATION (I)

- Consider the convex set:

$$\Delta_{T,\epsilon} = \left\{ w \in \mathbb{R}_{\geq 0}^T \text{ with } \|w\|_1 = 1 \text{ and } w_x \leq \frac{1}{|T|(1-\epsilon)} \right\}$$

Lemma: Let T be an ϵ -corruption of a $(3\epsilon, \delta)$ -stable set. For any $w \in \Delta_{T,\epsilon}$, if

$$\|\Sigma_w[T]\|_2 \leq 1 + \lambda \quad \rightarrow \quad \|\mu_w[T] - \mu\|_2 = O(\delta + \sqrt{\epsilon\lambda})$$

Non-Convex Optimization Formulation:

$$\min_{w \in \Delta_{T,\epsilon}} \|\Sigma_w[T]\|_2$$

NON-CONVEX OPTIMIZATION FORMULATION (II)

Problem Formulation:

Assign *weights* to the samples so that weighted empirical mean works.

Let $\Delta_{T,\epsilon} = \left\{ w \in \mathbb{R}_{\geq 0}^T \text{ with } \|w\|_1 = 1 \text{ and } w_x \leq \frac{1}{|T|(1-\epsilon)} \right\}$

Non-Convex Optimization Formulation:

$$\min_{w \in \Delta_{T,\epsilon}} \|\Sigma_w[T]\|_2$$

Algorithmic Approaches:

- This is what filtering does!
- Ellipsoid Method [DKKLMS'16]
- Bi-level optimization [Cheng-D-Ge'18]
- **Gradient Descent** [Cheng-D-Ge-Soltanolkotabi'20]

CONCRETE OPEN PROBLEMS

- **Design *near-linear time* algorithms for robust statistics tasks**

Robust Mean Estimation [Cheng-D-Ge, SODA'19; Dong-Hopkins-Li, NeurIPS'19; Depersin-Lecue'19]

Robust Covariance Estimation [Cheng-D-Ge-Woodruff, COLT'19]

Clustering mixture models [D-Kane-Koongsgard-Li-Tian, STOC'22]

Robust sparse estimation?

- ***Can we design robust estimators using first-order methods?***

Robust Mean Estimation [Cheng-D-Ge-Soltanolkotabi, ICML'20; Zhu et al. 2020]

More general tasks?

- **Obtain low-memory streaming robust learning algorithms**

[D-Kane-Pensia-Pittas, ICML'22] *Tradeoffs between memory and sample size?*

- **Robust *Online* Estimation?**