

Recent Advances in High-Dimensional Robust Statistics

Ilias Diakonikolas (UW Madison)

ICML 2020 Tutorial

July 2020

PART III: EXPERIMENTS AND EXTENSIONS

OUTLINE

Part III

- General Framework for Robust Mean Estimation
- Experiments
- Robust Stochastic Optimization
- Learning with Majority of Outliers

OUTLINE

Part III

- **General Framework for Robust Mean Estimation**
- Experiments
- Robust Stochastic Optimization
- Learning with Majority of Outliers

NON-CONVEX OPTIMIZATION FORMULATION (I)

Optimization Formulation:

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

Let

$$\Delta_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : \|w\|_1 = 1 \text{ and } 0 \leq w_i \leq \frac{1}{(1-\epsilon)N} \right\}$$
$$\hat{\mu}_w = \sum_{i=1}^N w_i X_i \quad \text{and} \quad \hat{\Sigma}_w = \sum_{i=1}^N w_i (X_i - \hat{\mu}_w)(X_i - \hat{\mu}_w)^T .$$

Generalization of Key Lemma: For any $w \in \Delta_{N,2\epsilon}$

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

NON-CONVEX OPTIMIZATION FORMULATION (II)

Notation: $\Delta_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : \|w\|_1 = 1 \text{ and } 0 \leq w_i \leq \frac{1}{(1-\epsilon)N} \right\}$

$$\hat{\mu}_w = \sum_{i=1}^N w_i X_i \quad \hat{\Sigma}_w = \sum_{i=1}^N w_i (X_i - \hat{\mu}_w)(X_i - \hat{\mu}_w)^T .$$

Generalization of Key Lemma

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

Non-Convex Formulation:

$$\min_w \|\hat{\Sigma}_w\|_2 \text{ subject to } w \in \Delta_{N,2\epsilon}$$

NON-CONVEX OPTIMIZATION FORMULATION (III)

Notation: $\Delta_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : \|w\|_1 = 1 \text{ and } 0 \leq w_i \leq \frac{1}{(1-\epsilon)N} \right\}$

$$\hat{\mu}_w = \sum_{i=1}^N w_i X_i \quad \hat{\Sigma}_w = \sum_{i=1}^N w_i (X_i - \hat{\mu}_w)(X_i - \hat{\mu}_w)^T .$$

Non-Convex Formulation:

$$\min_w \|\hat{\Sigma}_w\|_2 \text{ subject to } w \in \Delta_{N,2\epsilon}$$

Algorithmic Approaches:

- This is what filtering does!
- Ellipsoid Method [DKKLMS'16]
- Bi-level optimization [Cheng-D-Ge'18] (near-linear time!)
- **Gradient Descent** [Cheng-D-Ge-Soltanolkotabi, ICML'20]

ROBUST MEAN ESTIMATION VIA GRADIENT-DESCENT

Notation: $\Delta_{N,\epsilon} = \left\{ w \in \mathbb{R}^N : \|w\|_1 = 1 \text{ and } 0 \leq w_i \leq \frac{1}{(1-\epsilon)N} \right\}$

$$\hat{\mu}_w = \sum_{i=1}^N w_i X_i \quad \hat{\Sigma}_w = \sum_{i=1}^N w_i (X_i - \hat{\mu}_w)(X_i - \hat{\mu}_w)^T .$$

Non-Convex Formulation:

$$\min_w \|\hat{\Sigma}_w\|_2 \text{ subject to } w \in \Delta_{N,2\epsilon}$$

Theorem [Cheng-D-Ge-Soltanolkotabi, ICML'20, **Paper Id: #6611**]

Any approximate stationary point w defines $\hat{\mu}_w$ that is close to μ .

See also [Zhu et al., Arxiv, May 2020]

OUTLINE

Part III

- General Framework for Robust Mean Estimation
- **Experiments**
- Robust Stochastic Optimization
- Learning with Majority of Outliers

EXPERIMENTS

Being Robust (in High Dimensions) Can Be Practical

D., Kamath, Kane, Li, Moitra, Stewart, ICML'17

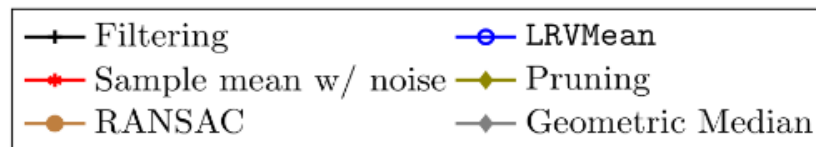
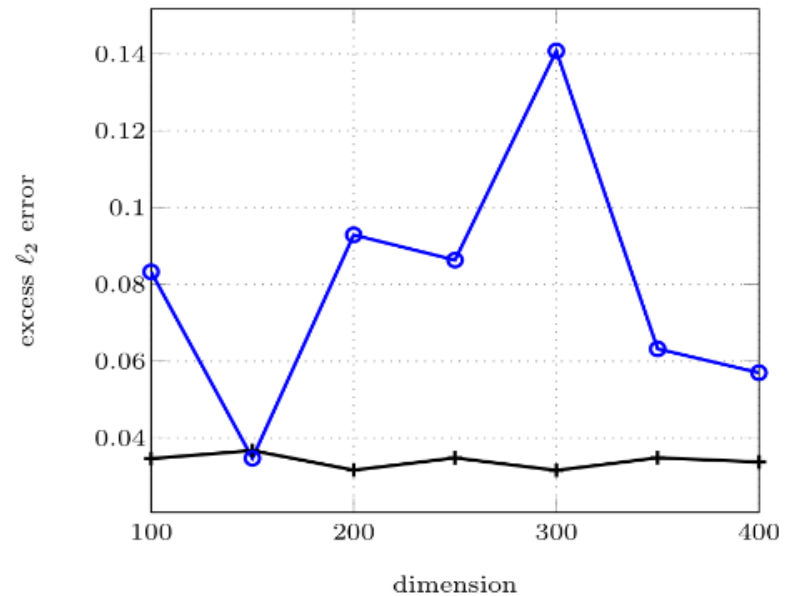
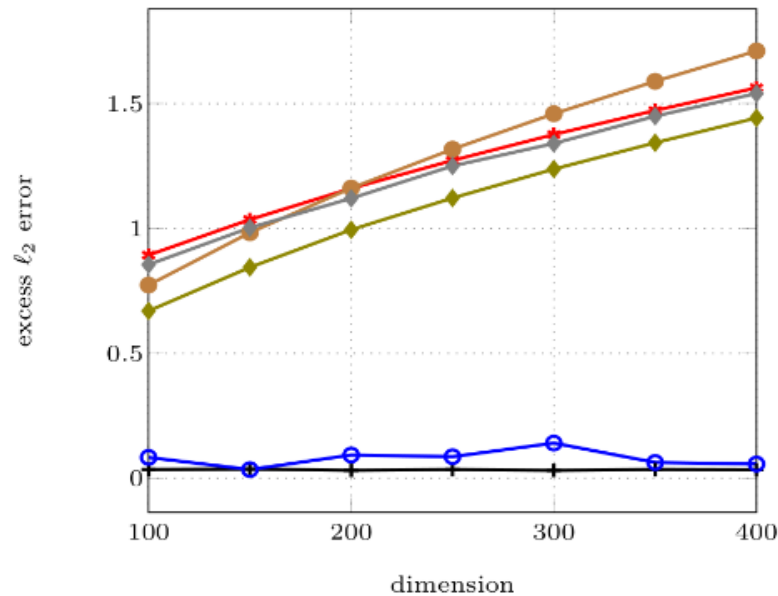
SYNTHETIC EXPERIMENTS: UNKNOWN MEAN

Error rates on synthetic data (**unknown mean**):

$$\mathcal{N}(\mu, I) + 10\% \text{ noise}$$

SYNTHETIC EXPERIMENTS: UNKNOWN MEAN

Error rates on synthetic data (**unknown mean**):



ROBUST COVARIANCE ESTIMATION

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 covariance of D .

Theorem: Let $\epsilon < 1/2$. We can efficiently recover $\hat{\Sigma}$ such that

$$\|\Sigma^{-1/2}(\hat{\Sigma} - \Sigma)\Sigma^{-1/2}\|_F \leq f(\epsilon) + \tilde{O}(d/\sqrt{N}),$$

where f depends on the concentration of D .

If D is a Gaussian, then $f(\epsilon) = O(\epsilon \log(1/\epsilon))$.

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

SYNTHETIC EXPERIMENTS: UNKNOWN COVARIANCE (I)

Error rates on synthetic data (**unknown covariance, isotropic**):

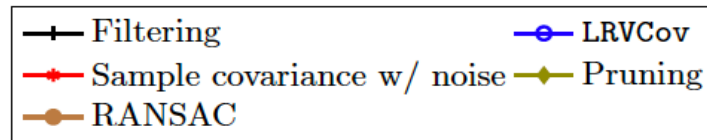
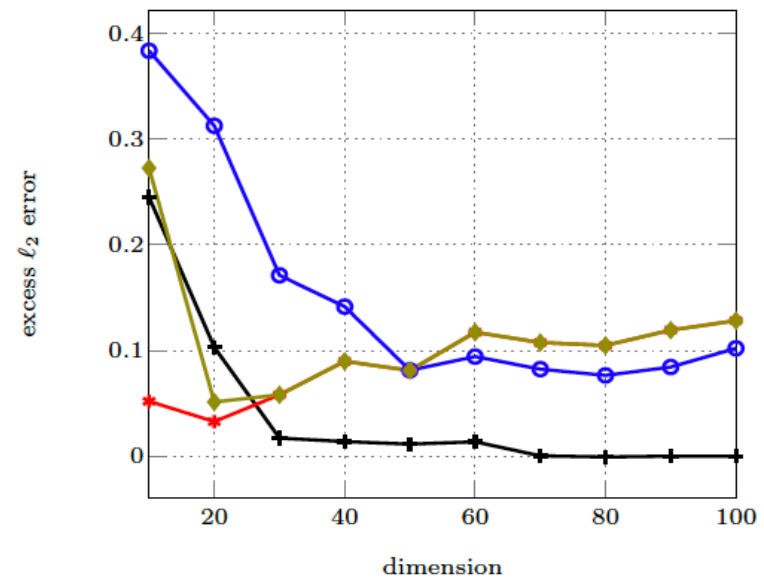
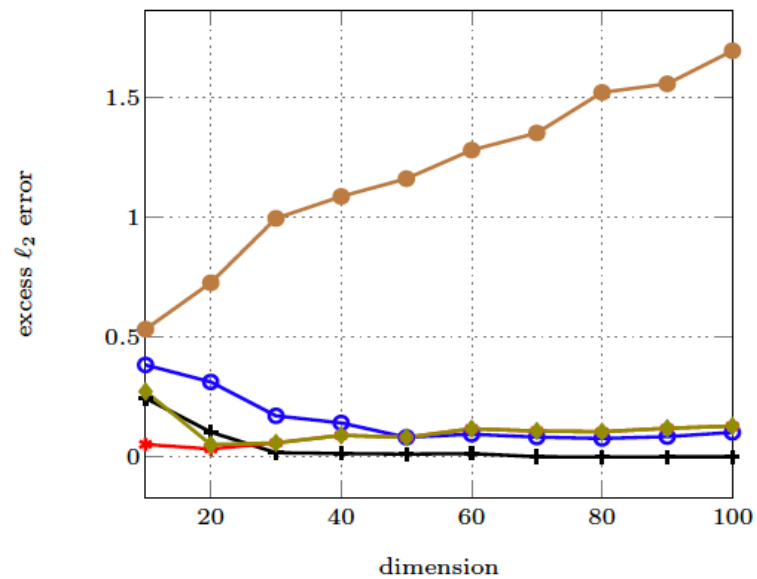
$\mathcal{N}(0, \Sigma) + 10\% \text{ noise}$



close to identity

SYNTHETIC EXPERIMENTS: UNKNOWN COVARIANCE (I)

Error rates on synthetic data (**unknown covariance, isotropic**):



SYNTHETIC EXPERIMENTS: UNKNOWN COVARIANCE (II)

Error rates on synthetic data (**unknown covariance**):

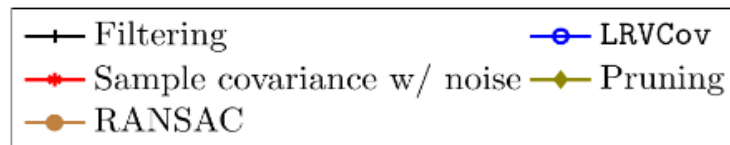
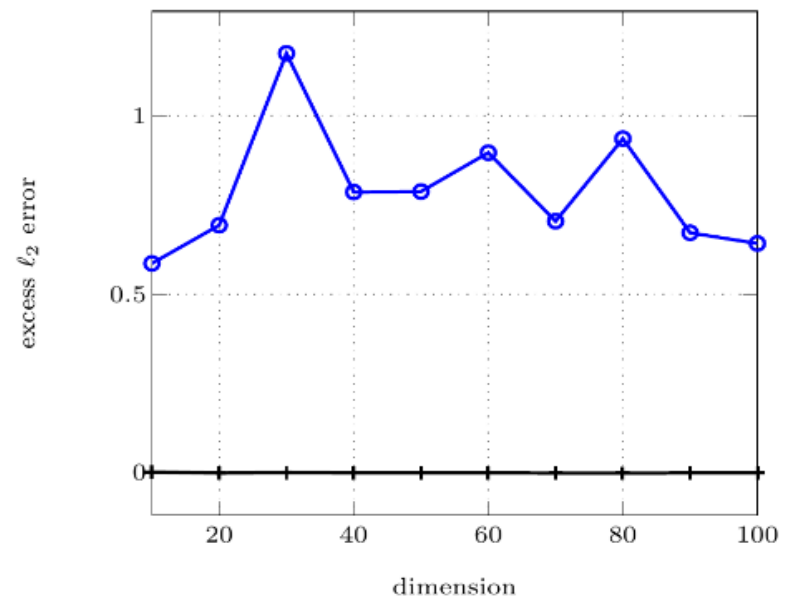
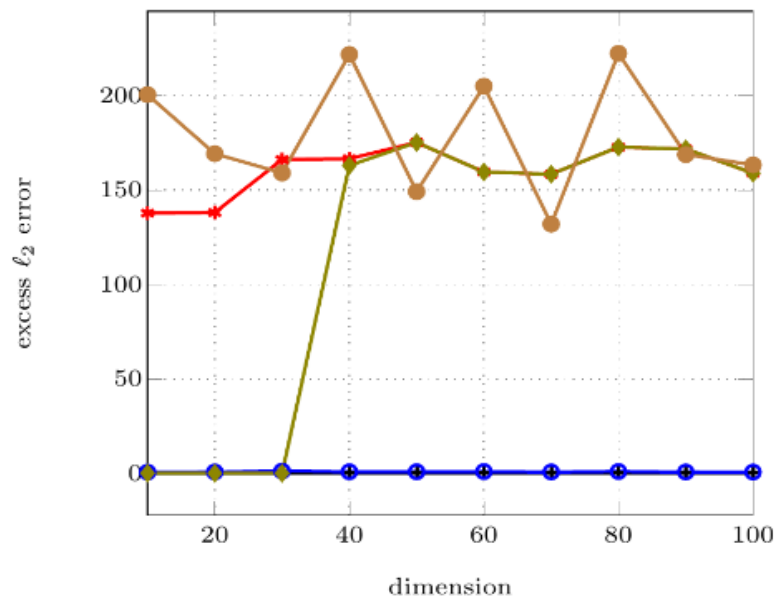
$\mathcal{N}(0, \Sigma) + 10\% \text{ noise}$



far from identity

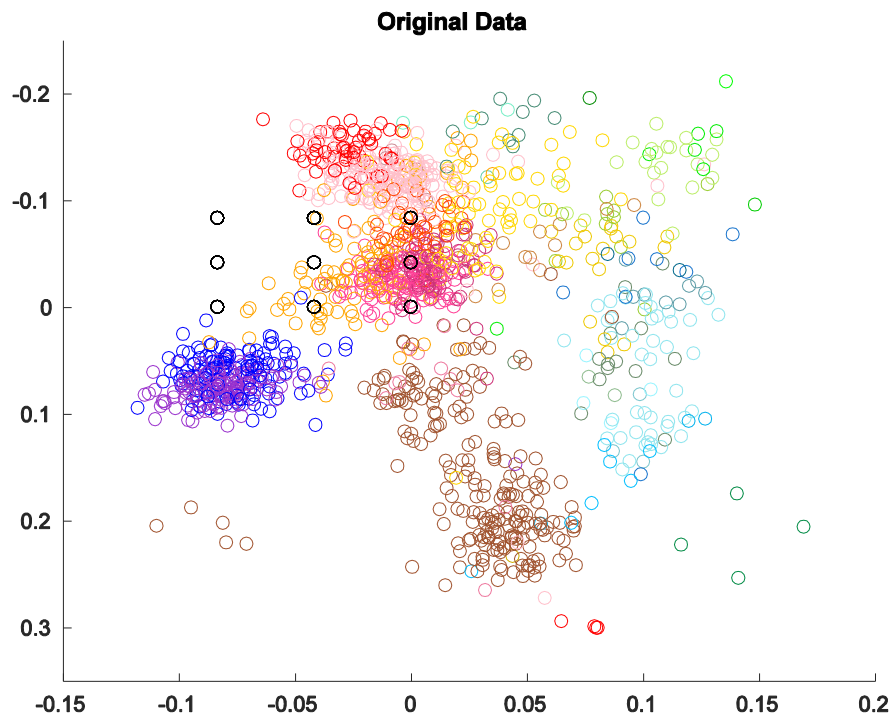
SYNTHETIC EXPERIMENTS: UNKNOWN COVARIANCE (II)

Error rates on synthetic data (**unknown covariance, anisotropic**):



REAL DATA EXPERIMENTS

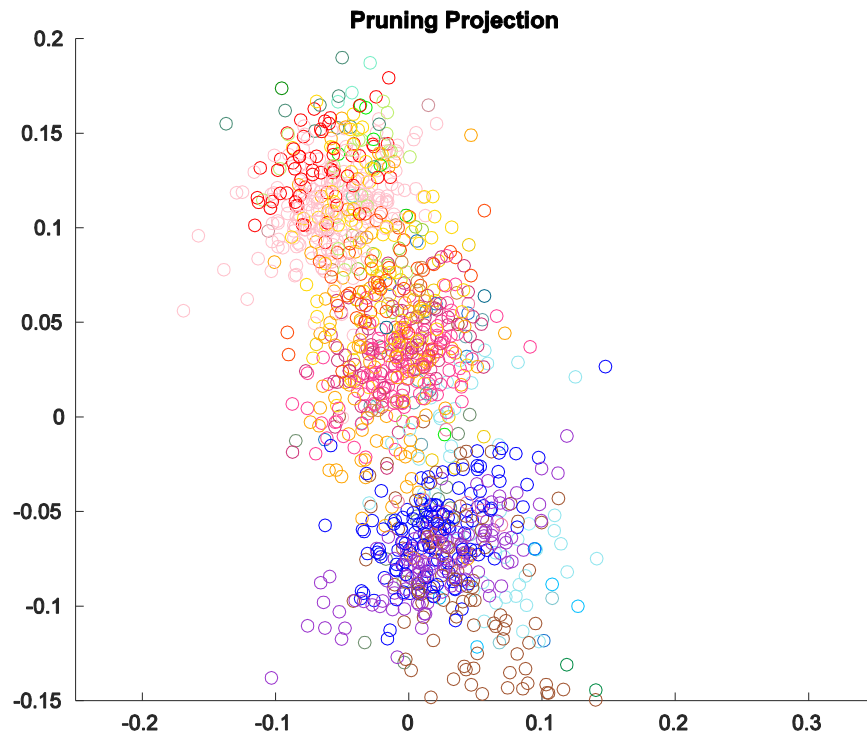
[Novembre et al. '08]: Take top two singular vectors of people x SNP matrix (POPRES)



“Genes Mirror Geography in Europe”

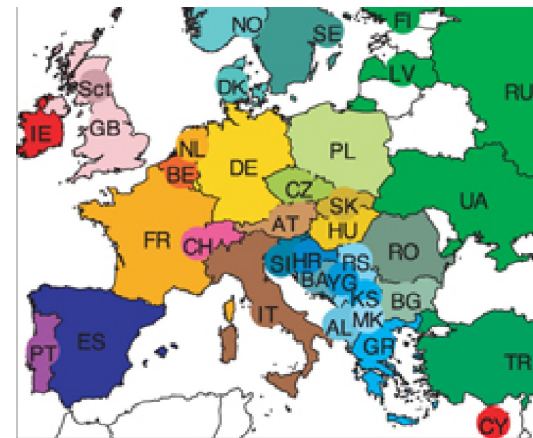
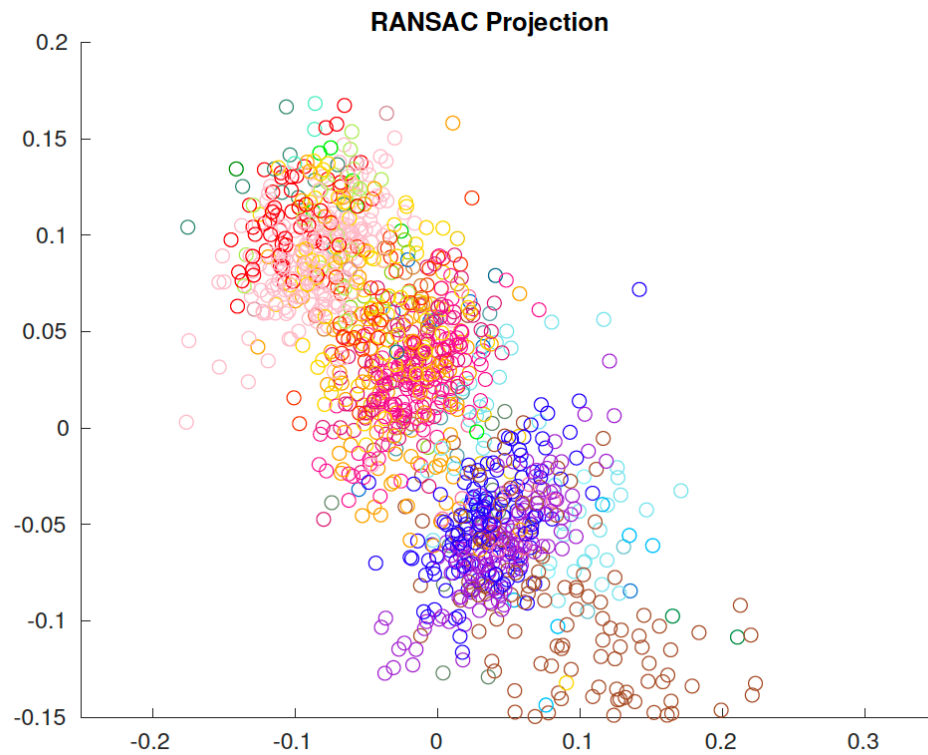
EXPERIMENTS: PRUNING PROJECTION

A comparison of error rate on semi-synthetic data:



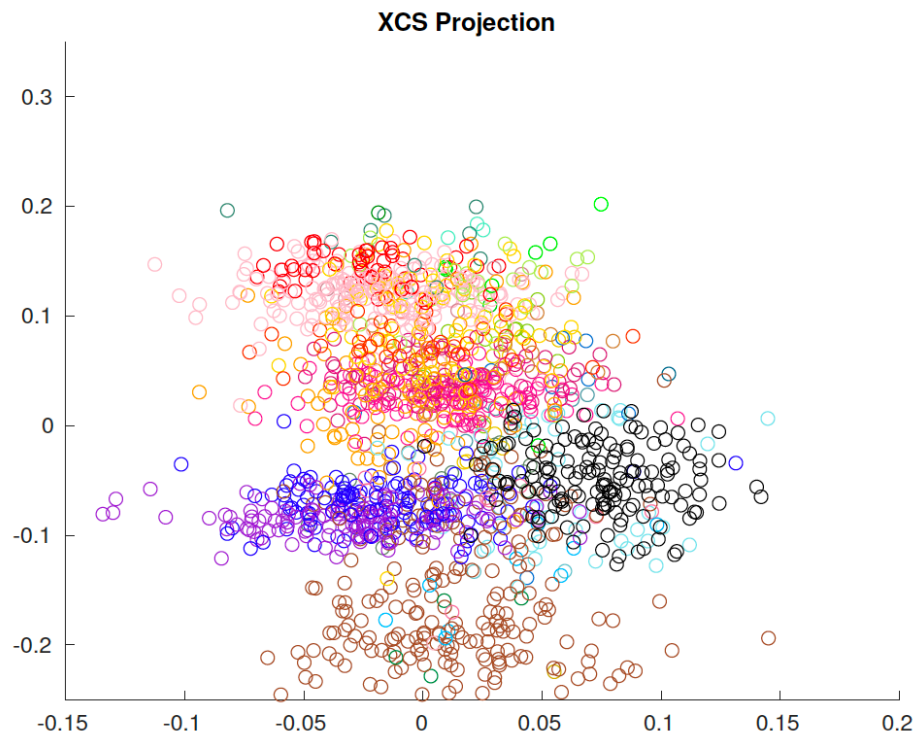
EXPERIMENTS: RANSAC PROJECTION

A comparison of error rate on semi-synthetic data:

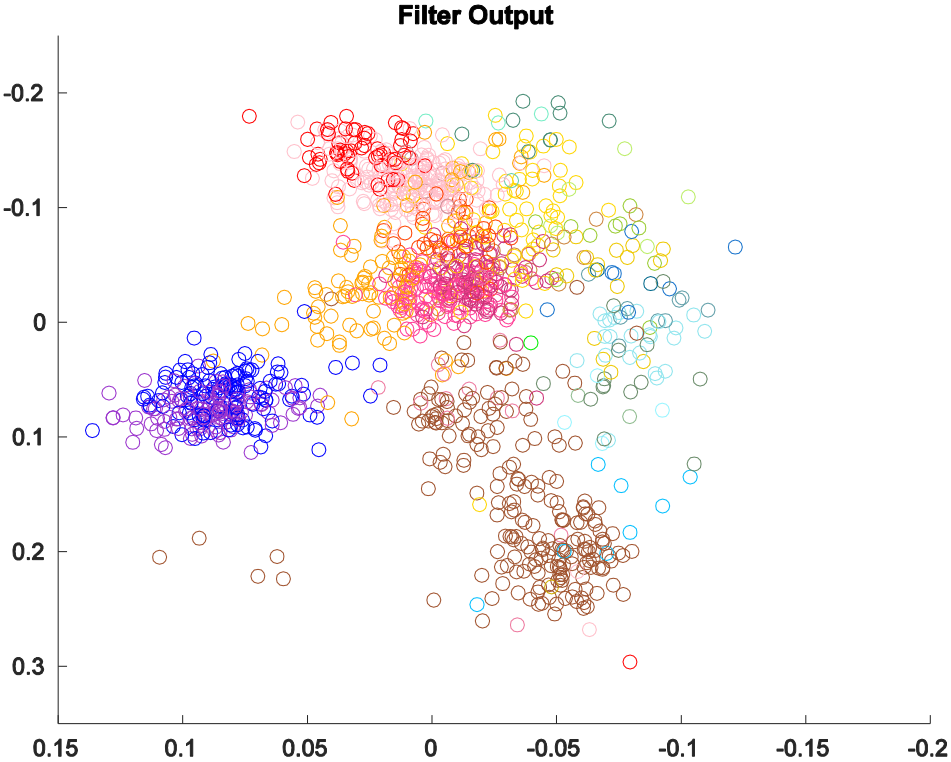


EXPERIMENTS: ROBUST PCA [XCS]

A comparison of error rate on semi-synthetic data:



EXPERIMENTS: FILTER PROJECTION



OUTLINE

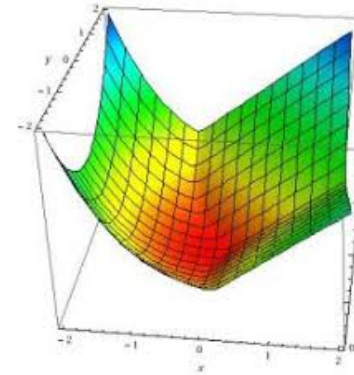
Part III

- General Framework for Robust Mean Estimation
- Experiments
- **Robust Stochastic Optimization**
- Learning with Majority of Outliers

ROBUST STOCHASTIC OPTIMIZATION

Sever: A Robust Meta-Algorithm for Stochastic Optimization.

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



ROBUST STOCHASTIC CONVEX OPTIMIZATION

Problem: Given loss function $\ell(X, w)$ and ϵ -corrupted samples from a distribution \mathcal{D} over X , minimize $f(w) = \mathbb{E}_{X \sim \mathcal{D}}[\ell(X, w)]$

Difficulty: Corrupted data can move the gradients.

Theorem: Suppose ℓ is convex and $\text{Cov}_{X \sim \mathcal{D}}[\nabla \ell(X, w)] \preceq \sigma^2 \cdot I$. Under mild assumptions on \mathcal{D} , can recover a point such that

$$f(\hat{w}) - \min_w f(w) \leq O(\sigma \sqrt{\epsilon}) .$$

Main Idea: Filter at minimizer of empirical risk.

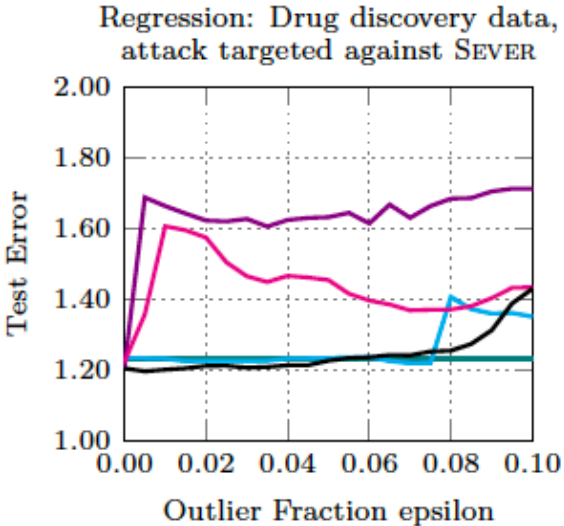
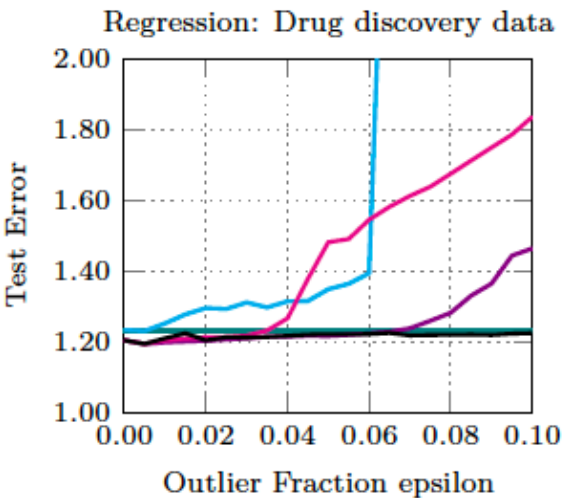
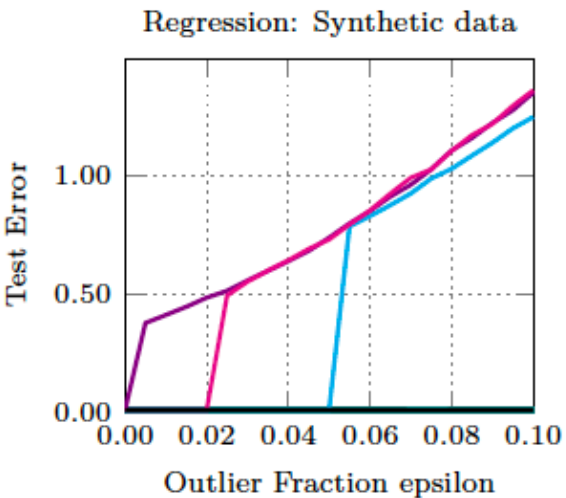
SPECIFIC APPLICATIONS

Corollary: Outlier-robust learning algorithms with dimension-independent error guarantees for:

- SVMs
- Linear Regression
- Logistic Regression
- GLMs
- **Experimental Performance Against Data Poisoning Attacks.**

Concurrent works obtained tighter guarantees in terms of either sample complexity or error, by focusing on specific tasks and distributional assumptions [Klivans-Kothari-Meka'18, Diakonikolas-Kong-Stewart'18, ...].

EXPERIMENTS: RIDGE REGRESSION



— uncorrupted — l2 — loss — gradientCentered — SEVER

OUTLINE

Part III

- General Framework for Robust Mean Estimation
- Experiments
- Robust Stochastic Optimization
- **Learning with Majority of Outliers**

LEARNING WITH A MAJORITY OF OUTLIERS

- So far focused on setting where $\epsilon < 1/2$.
- What can we learn from a dataset in which the **majority** of points are corrupted?

Problem: Given a set of points $x_1, \dots, x_N \in \mathbb{R}^d$ and $0 < \alpha \leq 1/2$ such that:

- An unknown subset of αN points are drawn from an unknown $D \in \mathcal{F}$, and
- The remaining $(1 - \alpha)N$ points are arbitrary,
approximate the mean μ of D .



Which is the “real” D ?

LIST-DECODABLE LEARNING

- Return *several hypotheses* with the guarantee that at least one is close.

List-Decodable Mean Estimation:

Given a set of points $x_1, \dots, x_N \in \mathbb{R}^d$ and $0 < \alpha \leq 1/2$ such that:

- An unknown subset of αN points are drawn from an unknown $D \in \mathcal{F}$, and
- The remaining $(1 - \alpha)N$ points are arbitrary,

output a small list of s hypotheses vectors such that one is close to the mean μ of D .

- Model defined in [Balcan-Blum-Vempala'08]
- First studied for mean estimation [Charikar-Steinhardt-Valiant'17]
- Application: Learning Mixture Models

LIST-DECODABLE MEAN ESTIMATION

Theorem [Charikar-Steinhardt-Valiant'17]: Let $0 < \alpha \leq 1/2$. If D has covariance $\Sigma \preceq I$ there is an efficient algorithm that uses $N \geq d/\alpha$ corrupted points, and outputs a list of $s = O(1/\alpha)$ vectors $\hat{\mu}_1, \dots, \hat{\mu}_s$ such that with high probability

$$\min_i \|\hat{\mu}_i - \mu\|_2 = \tilde{O}(1/\sqrt{\alpha}).$$

Theorem [Diakonikolas-Kane-Stewart'18] Any list-decodable mean estimator for bounded covariance distributions must have error $\Omega(1/\sqrt{\alpha})$ as long as the list size is any function of α .

- Initial algorithm [CSV'17] based on ellipsoid method.
- Generalization of filtering (“multi-filtering”) works for list-decodable setting [DKS'18].