High-Dimensional Robust Mean Estimation via Gradient Descent









Yu Cheng UIC

Ilias Diakonikolas UW Madison

Rong Ge Mahdi Soltanolkotabi Duke USC

Mean Estimation

- Input: \mathbb{N} samples $\{X_1, \dots, X_N\}$ drawn from $\mathcal{N}(\mu^*, I)$ on \mathbb{R}^d .
- Goal: Learn μ^* .



Mean Estimation

- Input: \mathbb{N} samples $\{X_1, \dots, X_N\}$ drawn from $\mathcal{N}(\mu^*, I)$ on \mathbb{R}^d .
- Goal: Learn μ^* .

Empirical mean
$$\widehat{\mu} = \frac{1}{N} \sum_{i=1}^{N} X_i$$
 works:
 $\|\widehat{\mu} - \mu^*\|_2 \le \epsilon$ when $N = \Omega(d/\epsilon^2)$.

Robust Mean Estimation



Robust Mean Estimation



replaces ϵN samples with arbitrary points.

Goal: Learn μ^* given an ϵ -corrupted set of N samples.

Previous Work

Algorithm	Error Guarantee	Poly-Time?
Coordinate-wise Median	$O(\epsilon\sqrt{d})$	Yes
Geometric Median	$O(\epsilon\sqrt{d})$	Yes
Tukey Median	$O(\epsilon)$	No
Tournament	$O(\epsilon)$	No
Pruning	$O(\epsilon\sqrt{d})$	Yes

Previous Work

Algorithm	Error Guarantee	Runtime
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	Polynomial
[Diakonikolas+ '16]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	

Previous Work

Algorithm	Error Guarantee	Runtime
[Lai+ '16]	$O(\epsilon \sqrt{\log d})$	Dolumomial
[Diakonikolas+ '16]		Polynoimai
[C Diakonikolas Ge '19]	$O(\epsilon \sqrt{\log(1/\epsilon)})$	$\tilde{O}(Nd)/\text{poly}(\epsilon)$
[Dong Hopkins Li '19]		$\tilde{O}(Nd)$

These algorithms have near-optimal sample complexity.

Motivation

Existing algorithms are fairly sophisticated (e.g., ellipsoid method, iterative spectral methods, matrix multiplicative weight update) and they are not parameter free.

Is it possible to solve robust estimation tasks by standard first-order methods?

Our Results

- A natural non-convex formulation of robust mean estimation.
- Any approximate stationary point of this non-convex objective gives a good solution for mean estimation.
- Gradient descent converges to an approximate stationary point in a polynomial number of iterations.

Non-Convex Formulation

$$\mu_w = \sum_i w_i X_i$$
 and $\Sigma_w = \sum_i w_i (X_i - \mu_w) (X_i - \mu_w)^\top$

[Diakonikolas+'16]:

If Σ_w has small spectral norm, then μ_w is close to the true mean.

minimize $\|\Sigma_w\|_2$ subject to $w \in \Delta_{N,2\epsilon}$

$$\Delta_{N,2\epsilon} = \{ w \in \mathbb{R}^N : ||w||_1 = 1 \text{ and } 0 \le w_i \le \frac{1}{(1-2\epsilon)N} \}$$

Yu Cheng (UIC)

Our Results

minimize $\|\Sigma_w\|_2$ subject to $w \in \Delta_{N,2\epsilon}$

Despite its non-convexity, we can show that any (approximate) stationary point *w* defines a μ_w that is $O(\epsilon \sqrt{\log(1/\epsilon)})$ -close to μ^* .

Our Results

 $\min \|\Sigma_w\|_2 \quad \text{s.t. } w \in \Delta_{N,2\epsilon}$

$\|\Sigma_w\|_2$ may not be differentiable w.r.t. W.

- Sub-gradient: use $\frac{\partial (v^{\mathsf{T}} \Sigma_w v)}{\partial w}$ where v is any top eigenvector of Σ_w .
- Softmax: minimize $\frac{1}{\rho} \operatorname{tr} \exp(\rho \Sigma_w)$, which is differentiable.

We prove structural and algorithmic results for both approaches.

Our Algorithms
Sub-gradientSoftmaxStart with any
$$w_0 \in \mathcal{K} = \Delta_{N,2\epsilon}$$
.
For $t = 0 \dots T - 1$
Let $v \in \operatorname{argmax}_{\|v\|_2=1} v^\top \Sigma_w v$.
 $w_{t+1} \leftarrow \mathcal{P}_{\mathcal{K}} \left(w_t - \eta \frac{\partial (v^\top \Sigma_w v)}{\partial w} \right)$.
end for
of iterations: $T = \tilde{O}(N^2 d^4)$...
For ...
For ...

 $-\eta \frac{\partial smax(\Sigma_w)}{\partial w} \Big).$

MATLAB Implementation

Projected Sub-gradient Descent

Our Results

- A natural non-convex formulation of robust mean estimation.
- Any approximate stationary point of our non-convex objective gives a near-optimal solution for mean estimation.
- Gradient descent converges to an approximate stationary point in a polynomial number of iterations.

W is a bad solution.

 $\Rightarrow v^{\mathsf{T}} \Sigma_w v$ is much larger than it should be.

- \Rightarrow We can find *i* and *j* such that
 - it is feasible to increase w_i and decrease w_j .
 - $v^{\mathsf{T}}\Sigma_w v$ becomes smaller after the change.

 \Rightarrow *w* is not a first-order stationary point.

 $v^{\mathsf{T}}\Sigma_w v$ = variance in the direction v.



Simple case: $\mu_w = 0$ and Σ_w has a unique largest eigenvector v.

We have
$$\Sigma_w = \sum_i w_i X_i X_i^{\mathsf{T}}$$
 and $v^{\mathsf{T}} \Sigma_w v = \sum_i w_i y_i^2$ where $y_i = X_i^{\mathsf{T}} v$.
$$\frac{\partial (v^{\mathsf{T}} \Sigma_w v)}{\partial w_i} = y_i^2$$

Simple case: $\mu_w = 0$ and Σ_w has a unique largest eigenvector v.

We have
$$\Sigma_w = \sum_i w_i X_i X_i^{\mathsf{T}}$$
 and $v^{\mathsf{T}} \Sigma_w v = \sum_i w_i y_i^2$ where $y_i = X_i^{\mathsf{T}} v$.
$$\frac{\partial (v^{\mathsf{T}} \Sigma_w v)}{\partial w_i} = y_i^2$$

 $\sum_{i \in \text{bad}} w_i y_i^2 \text{ is very large } \Rightarrow \exists i \text{ s.t. } w_i > 0 \text{ and } y_i^2 \text{ is large.}$

Challenges in bounding # of iterations:

•
$$w_{t+1} \leftarrow \mathcal{P}_{\mathcal{K}}\left(w_t - \eta \frac{\partial (v^{\mathsf{T}} \Sigma_w v)}{\partial w}\right)$$
 can make $\|\Sigma_w\|_2$ larger.

• Non-convex + Non-smooth + Constraints

Our Contributions

High-Dimensional Robust Statistics







Open Problems

- Faster convergence rate.
- More general robust estimation tasks:
 - Covariance estimation.
 - Sparse PCA.
 - Robust regression.