

Being Robust (in High Dimensions) Can Be Practical

Ilias Diakonikolas¹ Gautam Kamath² Daniel M. Kane³ Jerry Li²
Ankur Moitra⁴ Alistair Stewart¹

¹CS, USC

²EECS & CSAIL, MIT

³CSE & Math, UCSD

⁴Math & CSAIL, MIT

ICML, 2017

Presenter: Bargav Jayaraman

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

Introduction to Robust Estimation

- Given that samples come from a nice distribution, but adversary has corrupted a constant fraction of samples, goal is to robustly estimate the statistics - mean and covariance
- In one-dimension, robust alternatives to mean and covariance exist - *median* and *interquantile range*
- In high dimensions, there is a trade-off between robustness and computational efficiency
 - **Tukey median** - hard to compute; heuristics based computation does not scale with dimensions
 - **Minimum volume ellipsoid** - hard to compute; heuristics based computation scale poorly with dimensions

Outline

1 Introduction

- Introduction to Robust Estimation
- **Related Work**
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

- **Robust mean and covariance estimation** - [DKK⁺16] gave an algorithm for agnostically learning the parameters of a Gaussian $\mathcal{N}(\mu, \Sigma)$ that satisfy $d_{TV}(\mathcal{N}, \mathcal{N}') \leq O(\epsilon)$ where ϵ samples are corrupted, and the computational complexity of the algorithm is polynomial in dimensionality d and $1/\epsilon$. [LRV16] proposed unknown mean estimation where $d_{TV}(\mathcal{N}, \mathcal{N}') \leq O(\epsilon\sqrt{\log d})$.
- **Robust PCA** - [CLMW11] proposed robust PCA with semidefinite programming which can tolerate a constant fraction of corruptions. [XCS10] used semidefinite programming for robust PCA with outliers.

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- **Contributions**

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

Contributions

- 1 Modification to the algorithm of [DKK⁺16] with the definition of *good sets* to estimate the mean with $O(d/\epsilon^2)$ samples and covariance with $O(d^2/\epsilon^2)$ samples.
- 2 Improvement to the number of corruptions that can be tolerated by empirically tuning the threshold for filtering of corrupt points.
- 3 Same bounds are shown to be valid even for weaker distributional assumptions of the underlying data.
- 4 Comparison of models via visual representation of genetic data that encodes the map of Europe.

Outline

- 1 Introduction
 - Introduction to Robust Estimation
 - Related Work
 - Contributions
- 2 Proposed Method
 - **Formal Framework**
 - Nearly Sample-Optimal Efficient Robust Learning
 - Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation
- 3 Experiments
 - Experimental Settings
 - Synthetic Data
 - Semi-Synthetic Data
- 4 Conclusion
- 5 Extra Slide

Formal Framework

Notation: For a vector v , $\|v\|_2$ is the Euclidean norm, and for a matrix M , $\|M\|_2$ is the spectral norm and $\|M\|_F$ is the Frobenius norm. $X \in_u S$ means sample X is drawn from the empirical distribution defined by S .

Definition (ϵ -corruption)

Given $\epsilon > 0$ and a distribution family \mathcal{D} , the algorithm specifies the number of samples m and the adversary generates m samples X_1, X_2, \dots, X_m from some $D \in \mathcal{D}$. It then draws $m' \sim \text{Bin}(\epsilon, m)$ from an appropriate distribution and replaces m' of the input samples with arbitrary points. The altered samples are given to the algorithm.

Goal of the algorithm is to return the parameters of \hat{D} that are close to true distribution D . For mean, Euclidean distance is used and for covariance, Mahalanobis distance is used $\|\Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2} - I\|_F$.

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- **Nearly Sample-Optimal Efficient Robust Learning**
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

Nearly Sample-Optimal Efficient Robust Learning

The overall filtering procedure for robust estimation is the following iterative procedure:

- 1 find some univariate test (via spectral methods) that is violated by the corrupted points
- 2 find some concrete tail bound violated by the corrupted set of points
- 3 throw away all the points which violate this tail bound

Theorem (1)

Let G be a sub-gaussian distribution on \mathbb{R}^d with parameter $\nu = \Theta(1)$, mean μ^G , covariance matrix I , and $\epsilon > 0$. Let S be an ϵ -corrupted set of samples from G of size $\Omega((d/\epsilon^2)\text{polylog}(d/\epsilon))$. There exists an efficient algorithm that, on input S and $\epsilon > 0$, returns a mean vector $\hat{\mu}$ so that with probability at least $9/10$ we have $\|\hat{\mu} - \mu^G\|_2 = O(\epsilon\sqrt{\log 1/\epsilon})$.

Theorem (2)

Let P be a distribution on \mathbb{R}^d with unknown mean vector μ^P and unknown covariance matrix $\Sigma_P \leq \sigma^2 I$. Let S be an ϵ -corrupted set of samples from P of size $\Theta((d/\epsilon) \log d)$. There exists an efficient algorithm that, on input S and $\epsilon > 0$, with probability $9/10$ outputs $\|\hat{\mu} - \mu^P\|_2 = O(\sqrt{\epsilon}\sigma)$.

The main difference between this theorem and the previous one is the choice of filtering threshold. Instead of looking for a violation of a concentration inequality, here threshold is chosen at random.

Caution: This method may throw away even some uncorrupted points but it only rejects $O(\epsilon)$ samples with high probability.

Theorem (3)

Let $G \sim \mathcal{N}(0, \Sigma)$ be a Gaussian in d dimensions, and let $\epsilon > 0$. Let S be an ϵ -corrupted set of samples from G of size $\Omega((d^2/\epsilon^2)\text{polylog}(d/\epsilon))$. There exists an efficient algorithm that, given S and ϵ , returns the parameters of a Gaussian distribution $G' \sim \mathcal{N}(0, \Sigma)$ so that with probability at least $9/10$, it holds $\|I - \Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2}\|_F = O(\epsilon \log(1/\epsilon))$.

Outline

- 1 Introduction
 - Introduction to Robust Estimation
 - Related Work
 - Contributions
- 2 Proposed Method
 - Formal Framework
 - Nearly Sample-Optimal Efficient Robust Learning
 - **Filtering**
 - Robust Mean Estimation
 - Robust Covariance Estimation
- 3 Experiments
 - Experimental Settings
 - Synthetic Data
 - Semi-Synthetic Data
- 4 Conclusion
- 5 Extra Slide

Robust Mean Estimation

Algorithms which achieve Theorems 1 and 2 have common template, with three parameters:

- $\text{Thresh}(\epsilon)$ - threshold function to terminate if covariance has spectral norm bounded by $\text{Thresh}(\epsilon)$.
- $\text{Tail}(T, d, \epsilon, \delta, \tau)$ - univariate tail bound violated by only τ fraction of uncorrupted points, but more of corrupted points.
- $\delta(\epsilon, s)$ - slack function (required for technical reasons).

Robust Mean Estimation Algorithm Template

Algorithm 1 Filter-based algorithm template for robust mean estimation

- 1: **Input:** An ε -corrupted set of samples S , $\text{Thres}(\varepsilon)$, $\text{Tail}(T, d, \varepsilon, \delta, \tau)$, $\delta(\varepsilon, s)$
- 2: Compute the sample mean $\mu^{S'} = \mathbb{E}_{X \in_u S'}[X]$
- 3: Compute the sample covariance matrix Σ
- 4: Compute approximations for the largest absolute eigenvalue of Σ , $\lambda^* := \|\Sigma\|_2$, and the associated unit eigenvector v^* .
- 5: **if** $\|\Sigma\|_2 \leq \text{Thres}(\varepsilon)$ **then**
- 6: **return** $\mu^{S'}$.
- 7: Let $\delta = \delta(\varepsilon, \|\Sigma\|_2)$.
- 8: Find $T > 0$ such that

$$\Pr_{X \in_u S'} \left[|v^* \cdot (X - \mu^{S'})| > T + \delta \right] > \text{Tail}(T, d, \varepsilon, \delta, \tau).$$

- 9: **return** $\{x \in S' : |v^* \cdot (x - \mu^{S'})| \leq T + \delta\}$.
-

Algorithm for Sub-Gaussian (Theorem 1)

Algorithm 2 Filter algorithm for a sub-gaussian with unknown mean and identity covariance

1: **procedure** FILTER-SUB-GAUSSIAN-UNKNOWN-MEAN(S', ε, τ)

input: A multiset S' such that there exists an (ε, τ) -good S with $\Delta(S, S') \leq 2\varepsilon$

output: Multiset S'' or mean vector $\hat{\mu}$ satisfying Proposition A.7

- 2: Compute the sample mean $\mu^{S'} = \mathbb{E}_{X \in_u S'}[X]$ and the sample covariance matrix Σ , i.e., $\Sigma = (\Sigma_{i,j})_{1 \leq i,j \leq d}$ with $\Sigma_{i,j} = \mathbb{E}_{X \in_u S'}[(X_i - \mu_i^{S'})(X_j - \mu_j^{S'})]$.
- 3: Compute approximations for the largest absolute eigenvalue of $\Sigma - I$, $\lambda^* := \|\Sigma - I\|_2$, and the associated unit eigenvector v^* .
- 4: **if** $\|\Sigma - I\|_2 \leq O(\varepsilon \log(1/\varepsilon))$, **then return** $\mu^{S'}$.
- 5: Let $\delta := 3\sqrt{\varepsilon\|\Sigma - I\|_2}$. Find $T > 0$ such that

$$\Pr_{X \in_u S'} \left[|v^* \cdot (X - \mu^{S'})| > T + \delta \right] > 8 \exp(-T^2/2\nu) + 8 \frac{\varepsilon}{T^2 \log(d \log(\frac{d}{\varepsilon\tau}))}.$$

- 6: **return** the multiset $S'' = \{x \in S' : |v^* \cdot (x - \mu^{S'})| \leq T + \delta\}$.
-

Algorithm with Bounded Second Moment (Theorem 2)

Algorithm 3 Filter under second moment assumptions

- 1: **function** FILTERUNDER2NDMOMENT(S)
 - 2: Compute μ^S , Σ^S , the mean and covariance matrix of S .
 - 3: Find the eigenvector v^* with highest eigenvalue λ^* of Σ^S .
 - 4: **if** $\lambda^* \leq 9$ **then**
 - 5: **return** μ^S
 - 6: **else**
 - 7: Draw Z from the distribution on $[0, 1]$ with probability density function $2x$.
 - 8: Let $T = Z \max\{|v^* \cdot x - \mu^S| : x \in S\}$.
 - 9: Return the set $S' = \{x \in S : |v^* \cdot (X - \mu^S)| < T\}$.
-

Robust Covariance Estimation (Theorem 3)

Algorithm 4 Filter algorithm for a Gaussian with unknown covariance matrix.

- 1: **procedure** FILTER-GAUSSIAN-UNKNOWN-COVARIANCE(S', ε, τ)
- input:** A multiset S' such that there exists an (ε, τ) -good set S with $\Delta(S, S') \leq 2\varepsilon$
- output:** Either a set S'' with $\Delta(S, S'') < \Delta(S, S')$ or the parameters of a Gaussian G' with $d_{TV}(G, G') = O(\varepsilon \log(1/\varepsilon))$.
- Let $C > 0$ be a sufficiently large universal constant.
- 2: Let Σ' be the matrix $\mathbb{E}_{X \in S'}[XX^T]$ and let G' be the mean 0 Gaussian with covariance matrix Σ' .
- 3: **if** there is any $x \in S'$ so that $x^T(\Sigma')^{-1}x \geq Cd \log(|S'|/\tau)$ **then**
- 4: **return** $S'' = S' - \{x : x^T(\Sigma')^{-1}x \geq Cd \log(|S'|/\tau)\}$.
- 5: Compute an approximate eigendecomposition of Σ' and use it to compute $\Sigma'^{-1/2}$
- 6: Let $x_{(1)}, \dots, x_{(|S'|)}$ be the elements of S' .
- 7: For $i = 1, \dots, |S'|$, let $y_{(i)} = \Sigma'^{-1/2}x_{(i)}$ and $z_{(i)} = y_{(i)}^{\otimes 2}$.
- 8: Let $T_{S'} = -I^b I^{bT} + (1/|S'|) \sum_{i=1}^{|S'|} z_{(i)} z_{(i)}^T$.
- 9: Approximate the top eigenvalue λ^* and corresponding unit eigenvector v^* of $T_{S'}$.
- 10: Let $p^*(x) = \frac{1}{\sqrt{2}}((\Sigma'^{-1/2}x)^T v^* (\Sigma'^{-1/2}x) - \text{tr}(v^* v^{*T}))$
- 11: **if** $\lambda^* \leq (1 + C\varepsilon \log^2(1/\varepsilon))Q_{G'}(p^*)$ **then**
- 12: **return** G'
- 13: Let μ be the median value of $p^*(X)$ over $X \in S'$.
- 14: Find a $T \geq C'$ so that

$$\Pr_{X \in S'}(|p^*(X) - \mu| \geq T + 4/3) \geq \text{Tail}(T, d, \varepsilon, \tau)$$

- 15: **return** $S'' = \{X \in S' : |p^*(X) - \mu| < T\}$.

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- **Experimental Settings**
- Synthetic Data
- Semi-Synthetic Data

4 Conclusion

5 Extra Slide

Experimental Settings

Experiments were performed over:

- **Synthetic dataset unknown mean** - $\epsilon = 0.1$,
 $d = [100, 150, \dots, 400]$, $n = 10d/\epsilon^2$ samples are generated where
($1 - \epsilon$) fraction come from $\mathcal{N}(\mu, I)$
- **Synthetic dataset unknown covariance** - $\epsilon = 0.1$,
 $d = [10, 20, \dots, 100]$, $n = 0.5d/\epsilon^2$ samples are generated where
($1 - \epsilon$) fraction come from $\mathcal{N}(0, \Sigma)$
- **Semi-synthetic dataset** - Genotype of thousands of individuals.
PCA is used to project into two dimensions, which have a striking
resemblance to the map of Europe.

Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

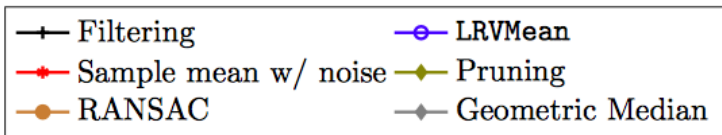
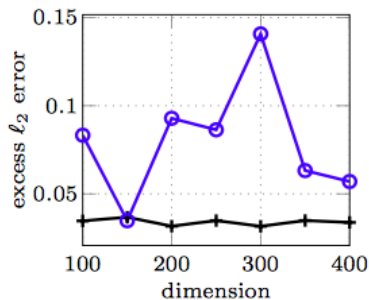
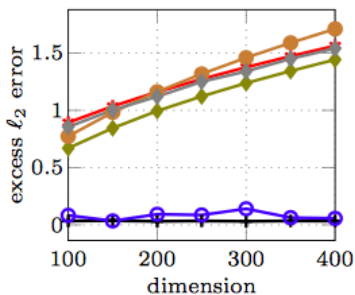
3 Experiments

- Experimental Settings
- **Synthetic Data**
- Semi-Synthetic Data

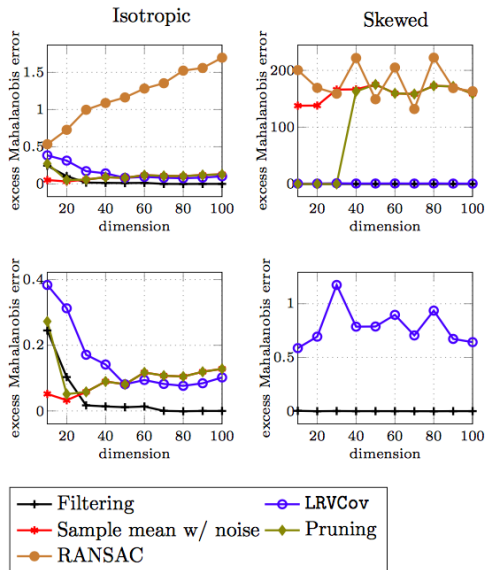
4 Conclusion

5 Extra Slide

Synthetic Data - Mean Estimation



Synthetic Data - Covariance Estimation



Outline

1 Introduction

- Introduction to Robust Estimation
- Related Work
- Contributions

2 Proposed Method

- Formal Framework
- Nearly Sample-Optimal Efficient Robust Learning
- Filtering
 - Robust Mean Estimation
 - Robust Covariance Estimation

3 Experiments

- Experimental Settings
- Synthetic Data
- **Semi-Synthetic Data**

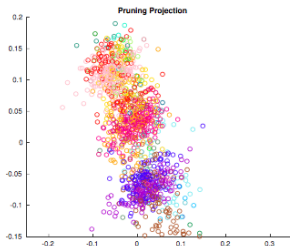
4 Conclusion

5 Extra Slide

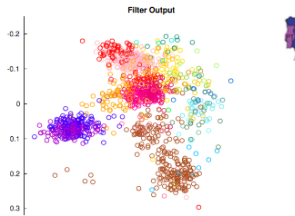
Semi-Synthetic Data



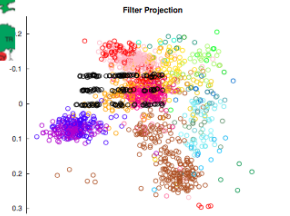
The data projected onto the top two directions of the original data set without noise



The data projected onto the top two directions of the noisy data set after pruning

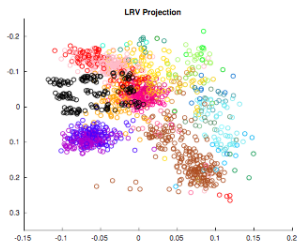
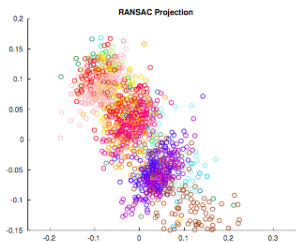
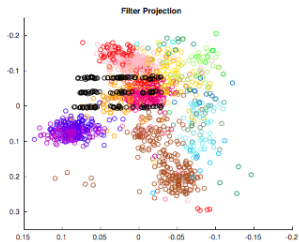


The filtered set of points projected onto the top two directions returned by the filter

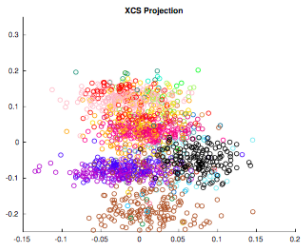
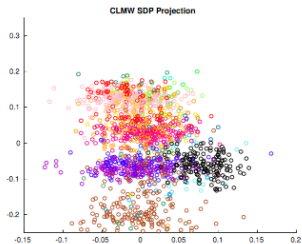
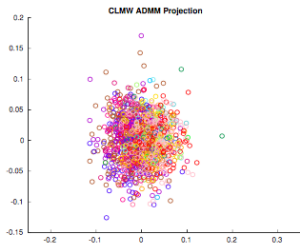


The data projected onto the top two directions returned by the filter

Semi-Synthetic Data



Semi-Synthetic Data



Conclusion

- Robust estimation of mean and covariance in high-dimensional data
- Experiments over synthetic data comparing the statistical accuracy of proposed method over existing approaches
- Experiments over Genotype dataset to visually represent the performance of various approaches in recovering the map of Europe
- Theoretical proof of correctness of the proposed approach

Roadmap of Proof of Theorem 1

Theorem A.3 $|S'| = \Omega(d/\epsilon^2 \text{poly}(\log(d/\epsilon)))$

$\exists \text{ Algo}(S', \epsilon) \rightarrow \hat{\mu}$ s.t. with prob $1-\tau$

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

← proved by

Fact A.2
Tail Bound
of
Sub-Gaussian

$$\Pr_{x \sim G} [|\langle x, \mu \rangle| > T] \leq \exp(-T^2/2\nu)$$

using

Def A.4
Goodness of
Set of Samples

Given G sub-gaussian with μ^G & $\Sigma = I$

Set S is (ϵ, τ) -good w.r.t. G if:

1. $\forall x \in S, \|x - \mu^G\|_2 \leq O(\sqrt{d \log(|S|/\tau)})$

2. $L(x) = \nu \cdot (x - \mu^G)^T, \nu \nu^T = I$
 $|\Pr_{x \in S} [L(x) > \sigma] - \Pr_{x \sim G} [L(x) > \sigma]| \leq \frac{\epsilon}{T \log(d \log(d/\epsilon))}$

3. $\|M_S - \mu^G\|_2 \leq \epsilon$

4. $\|M_S - I\|_2 \leq \epsilon$

Fact A.8

$$\Pr_{x \sim G} [|\langle x, \mu^G \rangle| > T] \leq 2 \exp(-T^2/2\nu)$$

$$\Pr_{x \in S} [|\langle x, \mu^G \rangle| > T] \leq \frac{2 \exp(-T^2/2\nu)}{1 + \frac{\epsilon}{T^2 \log(d \log(d/\epsilon))}}$$

Claim A.9

$$\Pr_{x \sim G} [|\langle x, \mu^S \rangle| > T + \|M_S - \mu^G\|_2] \leq 2 \exp(-T^2/2\nu)$$

$$\Pr_{x \in S} [|\langle x, \mu^S \rangle| > T + \|M_S - \mu^G\|_2] \leq \frac{2 \exp(-T^2/2\nu)}{1 + \frac{\epsilon}{T^2 \log(d \log(d/\epsilon))}}$$

Prop. A.7 Let G be sub-gaussian with μ^G & I

Let S be (ϵ, τ) -good w.r.t. G

Let S' be any multiset, $\Delta(S, S') \leq 2\epsilon$ & for any $x, y \in S'$

$$\|x - y\|_2 \leq O(\sqrt{d \log(d/\epsilon)})$$

$\exists \text{ Algo}(S', \epsilon)$

$$\hat{\mu} \text{ s.t. } \|\hat{\mu} - \mu^G\|_2 = O(\epsilon \sqrt{\log(d/\epsilon)})$$

$S'' \subseteq S'$ s.t.

$$\Delta(S, S'') \leq \Delta(S, S') - \epsilon/d$$

$$\alpha = d \log(d/\epsilon) \log(d \log(d/\epsilon))$$

← proved by

Corollary A.13

$$\Sigma - I = \frac{\epsilon I}{|S'|} M_\epsilon + O(\epsilon \log(d/\epsilon)) + O(\epsilon |E|/|S'|)^2 \|M_\epsilon\|_2$$

Lemma A.12

$$\mu^S - \mu^G = \frac{\epsilon I}{|S'|} (M_\epsilon - \mu^G) + O(\epsilon \sqrt{\log(d/\epsilon)})$$

Corollary A.11

$$M_S - I = \frac{\epsilon I}{|S'|} M_\epsilon + O(\epsilon \sqrt{\log(d/\epsilon)})$$

Lemma A.10

$$\|M_S\|_2 = O(\log(|S'|/|\mathcal{I}|) + \frac{\epsilon |S'|}{|\mathcal{I}|})$$