Learning Gaussian Covariance Robustly

Daniel M. Kane

Departments of CS/Math University of California, San Diego dakane@ucsd.edu

June 23rd, 2019

D, Kane (UCSD)

June, 2019 1 / 16

• • = • • = •

Outline

- Problem Setup
- Rough Estimates
- Refined Estimates
- Unknown Mean

< 3 >

< 行

Basic Problem

- Consider $G = N(0, \Sigma) \subset \mathbb{R}^n$.
- Given N samples, ϵ -fraction adversarially corrupted.
- Learn approximation to Σ .

★ Ξ →

How closely can we expect to learn Σ ?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

How closely can we expect to learn Σ ? Can't learn *G* to better than ϵ total variation.

< ∃ ►

How closely can we expect to learn Σ ? Can't learn *G* to better than ϵ total variation.

$$d_{TV}(N(0,\Sigma),N(0,\Sigma')) = \Theta(\min(1,\|\Sigma^{-1/2}\Sigma'\Sigma^{-1/2}-I\|_F)),$$

where

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2.$$

< ∃ >

How closely can we expect to learn Σ ? Can't learn *G* to better than ϵ total variation.

$$d_{\mathcal{T}\mathcal{V}}(\mathcal{N}(0,\Sigma),\mathcal{N}(0,\Sigma')) = \Theta(\min(1,\|\Sigma^{-1/2}\Sigma'\Sigma^{-1/2}-I\|_{\mathcal{F}})),$$

where

$$\|A\|_F^2 = \sum_{i,j} A_{i,j}^2.$$

Hope get estimate $\hat{\Sigma}$ so that:

$$\|\Sigma^{-1/2}\hat{\Sigma}\Sigma^{-1/2}-I\|_F=\tilde{O}(\epsilon).$$

∃ >

Basic Technique

Learning the mean of a Gaussian is equivalent to

- Learning $\mathbb{E}[L(G)]$ for degree-1 polynomials *L*.
- Learning the first moments of G.

< ∃ ►

Basic Technique

Learning the mean of a Gaussian is equivalent to

- Learning $\mathbb{E}[L(G)]$ for degree-1 polynomials *L*.
- Learning the first moments of G.

Learning the covariance of a mean 0 Gaussian is equivalent to:

- Learning $\mathbb{E}[p(G)]$ for even, degree-2 polynomials p.
- Learning the second moments of G.
- Learning $\mathbb{E}[GG^T]$.

★ ∃ ► < ∃ ►</p>

Basic Technique

Learning the mean of a Gaussian is equivalent to

- Learning $\mathbb{E}[L(G)]$ for degree-1 polynomials *L*.
- Learning the first moments of G.

Learning the covariance of a mean 0 Gaussian is equivalent to:

- Learning $\mathbb{E}[p(G)]$ for even, degree-2 polynomials p.
- Learning the second moments of G.
- Learning $\mathbb{E}[GG^T]$.

We will use the last of these formulations.

We have reduced the problem to robustly estimating the mean of the n^2 -dimensional random variable $X = GG^T$. Since $Cov(G) = \Sigma = \mathbb{E}[X]$.

We have reduced the problem to robustly estimating the mean of the n^2 -dimensional random variable $X = GG^T$. Since $Cov(G) = \Sigma = \mathbb{E}[X]$.

Let $\Sigma = \text{Cov}(X)$. If $\Sigma \ll I_{n^2}$, can learn $\mathbb{E}[X]$ to L^2 error (and thus, Σ to Frobenius error) $O(\sqrt{\epsilon})$.

We have reduced the problem to robustly estimating the mean of the n^2 -dimensional random variable $X = GG^T$. Since $Cov(G) = \Sigma = \mathbb{E}[X]$.

Let $\Sigma = \text{Cov}(X)$. If $\Sigma \ll I_{n^2}$, can learn $\mathbb{E}[X]$ to L^2 error (and thus, Σ to Frobenius error) $O(\sqrt{\epsilon})$.

So, what is Σ ?

• Suppose that y_i are an orthonormal basis of linear functions of G.

•
$$\operatorname{Cov}(y_i, y_j) = \delta_{i,j}$$

A D N A B N A B N A B N

- Suppose that y_i are an orthonormal basis of linear functions of G.
 Cov(y_i, y_j) = δ_{i,j}
- y_iy_j(i ≠ j) and (y_i² − 1)/√2 form an orthonormal basis for even degree-2 polynomials of G.

・ 何 ト ・ ヨ ト ・ ヨ ト

- Suppose that y_i are an orthonormal basis of linear functions of G.
 Cov(y_i, y_j) = δ_{i,j}
- y_iy_j(i ≠ j) and (y_i² − 1)/√2 form an orthonormal basis for even degree-2 polynomials of G.
- For matrix A,

$$\begin{aligned} A^{\textit{flat}} \mathbf{\Sigma} A^{\textit{flat}} &= \operatorname{Var}(A \cdot X) = \operatorname{Var}(G^{\mathsf{T}} A G) \\ &= 2 \left\| \Sigma^{1/2} \left(\frac{A + A^{\mathsf{T}}}{2} \right) \Sigma^{1/2} \right\|_{\mathsf{F}}^2 \end{aligned}$$

.

・ 何 ト ・ ヨ ト ・ ヨ ト

- Suppose that y_i are an orthonormal basis of linear functions of G.
 Cov(y_i, y_j) = δ_{i,j}
- y_iy_j(i ≠ j) and (y_i² − 1)/√2 form an orthonormal basis for even degree-2 polynomials of G.
- For matrix A,

$$\begin{aligned} A^{\textit{flat}} \mathbf{\Sigma} A^{\textit{flat}} &= \operatorname{Var}(A \cdot X) = \operatorname{Var}(G^{\mathsf{T}} A G) \\ &= 2 \left\| \Sigma^{1/2} \left(\frac{A + A^{\mathsf{T}}}{2} \right) \Sigma^{1/2} \right\|_{\mathsf{F}}^2 \end{aligned}$$

So, for example, if $\Sigma \leq I$, $\Sigma \ll I$.

.

Bootstrapping

- To learn Σ , need to learn $\mathbb{E}[X]$ robustly.
- Can learn $\mathbb{E}[X]$ robustly, if we have an upper bound on Σ .
- Can find Σ if we know Σ .

< (日) × (日) × (4)

Bootstrapping

- To learn Σ , need to learn $\mathbb{E}[X]$ robustly.
- Can learn $\mathbb{E}[X]$ robustly, if we have an upper bound on Σ .
- Can find Σ if we know Σ .

Bootstrap better and better approximations to Σ !

★ Ξ >

Upper Bounds

Critical Point: If $\Sigma \leq \Sigma_0$, then $\Sigma \leq \Sigma_0$, i.e.

$$2\left\|\boldsymbol{\Sigma}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}^{1/2}\right\|_{F}^{2} \leq 2\left\|\boldsymbol{\Sigma}_{0}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}_{0}^{1/2}\right\|_{F}^{2}$$

for all A.

イロト イヨト イヨト イヨト

Upper Bounds

Critical Point: If $\Sigma \leq \Sigma_0$, then $\Sigma \leq \Sigma_0$, i.e.

$$2\left\|\boldsymbol{\Sigma}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}^{1/2}\right\|_{F}^{2} \leq 2\left\|\boldsymbol{\Sigma}_{0}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}_{0}^{1/2}\right\|_{F}^{2}$$

for all A.

So if $\Sigma \leq \Sigma_0$, then $\operatorname{Cov}(\Sigma_0^{-1/2} X \Sigma_0^{-1/2}) \ll I_{n^2}$. Can get estimate $\hat{\Sigma}$ with $\left\| \Sigma_0^{-1/2} \left(\hat{\Sigma} - \Sigma \right) \Sigma_0^{-1/2} \right\|_F = O(\sqrt{\epsilon}).$

イロト イポト イヨト イヨト

Upper Bounds

Critical Point: If $\Sigma \leq \Sigma_0$, then $\Sigma \leq \Sigma_0$, i.e.

$$2\left\|\boldsymbol{\Sigma}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}^{1/2}\right\|_{F}^{2} \leq 2\left\|\boldsymbol{\Sigma}_{0}^{1/2}\left(\frac{\boldsymbol{A}+\boldsymbol{A}^{T}}{2}\right)\boldsymbol{\Sigma}_{0}^{1/2}\right\|_{F}^{2}$$

for all A.

So if
$$\Sigma \leq \Sigma_0$$
, then $\operatorname{Cov}(\Sigma_0^{-1/2} X \Sigma_0^{-1/2}) \ll I_{n^2}$. Can get estimate $\hat{\Sigma}$ with
 $\left\| \Sigma_0^{-1/2} \left(\hat{\Sigma} - \Sigma \right) \Sigma_0^{-1/2} \right\|_F = O(\sqrt{\epsilon}).$

So $\hat{\Sigma} = \Sigma + O(\sqrt{\epsilon})\Sigma_0$.

< □ > < 同 > < 回 > < 回 > < 回 >

- Start with some upper bound $\Sigma_0 \ge \Sigma$ (twice the sample covariance works with high probability).
- Get approximation $\hat{\Sigma}_0$.

→ Ξ →

- Start with some upper bound $\Sigma_0 \ge \Sigma$ (twice the sample covariance works with high probability).
- Get approximation $\hat{\Sigma}_0$.
- Use $\Sigma_1 = \hat{\Sigma}_0 + C \sqrt{\epsilon} \Sigma_0$ as new upper bound.
- Get approximation $\hat{\Sigma}_1$.

• . . .

- Start with some upper bound $\Sigma_0 \ge \Sigma$ (twice the sample covariance works with high probability).
- Get approximation $\hat{\Sigma}_0$.
- Use $\Sigma_1 = \hat{\Sigma}_0 + C \sqrt{\epsilon} \Sigma_0$ as new upper bound.
- Get approximation $\hat{\Sigma}_1$.
- Use $\Sigma_2 = \hat{\Sigma}_1 + C \sqrt{\epsilon} \Sigma_1$ as new upper bound.

★ ∃ ▶

- Start with some upper bound $\Sigma_0 \ge \Sigma$ (twice the sample covariance works with high probability).
- Get approximation $\hat{\Sigma}_0$.
- Use $\Sigma_1 = \hat{\Sigma}_0 + C \sqrt{\epsilon} \Sigma_0$ as new upper bound.
- Get approximation $\hat{\Sigma}_1$.
- Use $\Sigma_2 = \hat{\Sigma}_1 + C \sqrt{\epsilon} \Sigma_1$ as new upper bound.

• . . .

Have $\Sigma_{i+1} \leq \Sigma + O(\sqrt{\epsilon})\Sigma_i$. Eventually get $\Sigma_{\infty} \leq \Sigma(1 + O(\sqrt{\epsilon}))$, and $\hat{\Sigma}$ with

$$\|\Sigma^{-1/2}\hat{\Sigma}\Sigma^{-1/2}-I\|_{F}=O(\sqrt{\epsilon}).$$

• • = • • = •

Error Idea

- Have error $O(\sqrt{\epsilon})$.
 - Best possible using only bounds on Cov(X).

< □ > < 同 > < 回 > < 回 > < 回 >

Error Idea

- Have error $O(\sqrt{\epsilon})$.
 - Best possible using only bounds on Cov(X).
- Hope to do better given:
 - Accurate approximation to Cov(X).
 - Tail bounds for X.

• • = • • =

Error Idea

- Have error $O(\sqrt{\epsilon})$.
 - Best possible using only bounds on Cov(X).
- Hope to do better given:
 - Accurate approximation to Cov(X).
 - Tail bounds for X.

Simplifying Assumption: $\Sigma \approx I$.

→ 3 → 4 3

Concentration

Standard Result: If p is a degree-2 polynomial with Var(p(G)) = O(1), then

 $\Pr(|p(G) - \mathbb{E}[p(G)]| > T) = O(\exp(-\Omega(T))).$

イロト 不得下 イヨト イヨト 二日

Concentration

Standard Result: If *p* is a degree-2 polynomial with Var(p(G)) = O(1), then

$$\Pr(|p(G) - \mathbb{E}[p(G)]| > T) = O(\exp(-\Omega(T))).$$

Therefore, X has exponential concentration about its mean in any direction.

• Know Σ to error $O(\delta)$ in Frobenius norm

-

・ロト ・ 日 ト ・ 目 ト ・

- Know Σ to error $O(\delta)$ in Frobenius norm
- Compute $\boldsymbol{\Sigma}$ to error $O(\delta)$ in operator norm

→ Ξ →

- Know Σ to error $O(\delta)$ in Frobenius norm
- Compute $\boldsymbol{\Sigma}$ to error $O(\delta)$ in operator norm
- Compute Σ to error $O(\sqrt{\epsilon\delta} + \epsilon \log(1/\epsilon))$ in Frobenius norm

★ ∃ ▶

Setup

- Know Σ to error $O(\delta)$ in Frobenius norm
- Compute ${f \Sigma}$ to error ${\cal O}(\delta)$ in operator norm
- Compute Σ to error $O(\sqrt{\epsilon\delta} + \epsilon \log(1/\epsilon))$ in Frobenius norm
- Iterate to get error $O(\epsilon \log(1/\epsilon))$

Combine

- Last Talk: If $\Sigma = I$, robustly learn μ .
- This Talk: If $\mu = 0$, robustly learn Σ .

• • • • • • • • • • • •

Combine

- Last Talk: If $\Sigma = I$, robustly learn μ .
- This Talk: If $\mu = 0$, robustly learn Σ .
- Question: What if neither Σ nor μ is known?

→ Ξ →



• Consider differences of pairs of samples $G_{2i} - G_{2i+1}$.

< □ > < □ > < □ > < □ > < □ >



- Consider differences of pairs of samples $G_{2i} G_{2i+1}$.
- Distributed as $N(0, 2\Sigma)$, with 2ϵ error.

< □ > < 同 > < 回 > < 回 > < 回 >

Trick

- Consider differences of pairs of samples $G_{2i} G_{2i+1}$.
- Distributed as $N(0, 2\Sigma)$, with 2ϵ error.
- Use to learn $\hat{\Sigma}$, an approximation to Σ with $O(\epsilon \log(1/\epsilon))$ error.

▲ □ ▶ ▲ □ ▶ ▲ □

Trick

- Consider differences of pairs of samples $G_{2i} G_{2i+1}$.
- Distributed as $N(0, 2\Sigma)$, with 2ϵ error.
- Use to learn $\hat{\Sigma}$, an approximation to Σ with $O(\epsilon \log(1/\epsilon))$ error.
- $\hat{\Sigma}^{-1/2}G \approx N(\hat{\Sigma}^{-1/2}\mu, I)$
 - Treat difference as $O(\epsilon \log(1/\epsilon))$ adversarial error.
 - Use to learn approximation to μ .

• • = • • = •

Trick

- Consider differences of pairs of samples $G_{2i} G_{2i+1}$.
- Distributed as $N(0, 2\Sigma)$, with 2ϵ error.
- Use to learn $\hat{\Sigma}$, an approximation to Σ with $O(\epsilon \log(1/\epsilon))$ error.

•
$$\hat{\Sigma}^{-1/2} G pprox \textit{N}(\hat{\Sigma}^{-1/2} \mu, \textit{I})$$

- Treat difference as $O(\epsilon \log(1/\epsilon))$ adversarial error.
- Use to learn approximation to μ .

Final result: Learn distribution for G to $\tilde{O}(\epsilon)$ error in total variational distance.

Conclusion

We can learn the mean and covariance of an unknown Gaussian robustly. In order to do so, we need to consider the 2nd and 4th moments of the distribution in question. Later we will look into cases where even higher moments are useful.