Iterative Methods for Efficient Statistics in High Dimensions

Manolis Zampetakis MIT

Constantinos DaskalakisThemis GouleakisVasilis KontonisChristos TzamosMITUSCUW MadisonUW Madison

NTUA 2019

Given *n* samples from a Gaussian distribution in *d* dimensions estimate:

- the **mean** and
- the covariance matrix.

Given *n* samples from a Gaussian distribution in *d* dimensions estimate:

- the **mean** and
- the covariance matrix.

Answer:

$$\hat{\boldsymbol{\mu}} = rac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\boldsymbol{\Sigma}} = rac{1}{n} \sum_{i=1}^{n} x_i x_i^T$$

Given *n* samples from a Gaussian distribution in *d* dimensions estimate:

- the **mean** and
- the covariance matrix.

Answer:

$$\hat{\boldsymbol{\mu}} = rac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\boldsymbol{\Sigma}} = rac{1}{n} \sum_{i=1}^{n} x_i x_i^T$$

How large n should be?

Theorem ([Folklore]). Let $x_1, x_2, ..., x_n$ i.i.d. samples from $\mathcal{N}(\mu^*, \Sigma^*)$, then the estimates

$$\hat{\mu} = rac{1}{n}\sum_{i=1}^n x_i$$
 and $\hat{\Sigma} = rac{1}{n}\sum_{i=1}^n x_i x_i^T$

satisfy with probability at least 99%:

 $d_{TV}(\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}), \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)) \leq \varepsilon$

for $n = \Theta(d^2/\varepsilon^2)$.

Theorem ([Folklore]). Let $x_1, x_2, ..., x_n$ i.i.d. samples from $\mathcal{N}(\mu^*, \Sigma^*)$, then the estimates

$$\hat{\boldsymbol{\mu}} = rac{1}{n}\sum_{i=1}^n \boldsymbol{x}_i$$
 and $\hat{\boldsymbol{\Sigma}} = rac{1}{n}\sum_{i=1}^n \boldsymbol{x}_i \boldsymbol{x}_i^T$

satisfy with probability at least 99%:

$$d_{TV}(\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}), \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)) \leq \varepsilon$$

for $n = \Theta(d^2/\varepsilon^2)$.

Time and Sample Efficient in High Dimensions.

Estimation Obstacles

Truncated Samples Inhomogeneous Population

Global Convergence of Iterative Methods

Estimation Obstacles

o Truncated Sampleso Inhomogeneous Population

Optimization Problems

Global Convergence of Iterative Methods

Estimation Obstacles

o Truncated Sampleso Inhomogeneous Population

Global Convergence of Iterative Methods

Optimization Problems

Estimation Obstacles

o Truncated Sampleso Inhomogeneous Population

Optimization Problems

Convex Landscape

Non-Convex Landscape

Iterative Methods

Global Convergence of Iterative Methods

Estimation Obstacles

• Truncated Samples
• Inhomogeneous Population

Optimization Problems

Convex Landscape

Non-Convex Landscape

Iterative Methods

Global Convergence of Iterative Methods

Anecdotal Story [Gil Kalai's Blog Post]



"My friend the baker," said Poincaré, "I weighed every loaf of bread that I bought from you in the last year and the distribution is Gaussian with mean 950 grams. How can you claim that your average loaf is 1 kilogram?"



Anecdotal Story (continued) [Gil Kalai's Blog Post]



A year later the two pals meet again.

"How are you doing dear Henri" asked the baker "are my bread loaves heavy enough for you?"

"Yes, for me they are," answered Poincaré "but when I weighed all the loaves last year I discovered that your mean value is still 950 grams."

"*How is this possible?*" asked the baker.

A year later th

"How are you a heavy enough f

"Yes, for me the last year I disco

"How is this po



Example [Hausman, Wise (Econometrica) 1976]

Data Set:

(wage rate, income, education level – I.Q.)

Study: intelligence vs income for low wage rate workers

Example [Hausman, Wise (Econometrica) 1976]

Data Set:

(wage rate, income, education level – I.Q.)

Study: intelligence vs income for low wage rate workers

Truncation!: samples collected only if income was less than 1.5 times the poverty level.

Example [Hausman, Wise (Econometrica) 1976]

Data Set:

(wage rate, income, education level – I.Q.)

Findings (e.g. [Hause 1971]):

Intelligence and education level have mostly **no effect** on the income and sometimes even **negative effect**! (?)

Example [Hausman, Wise (Econometrica) 1976]

Data Set:

(wage rate, income, education level – I.Q.)

Findings (e.g. [Hause 1971]):

Intelligence and education level have mostly **no effect** on the income and sometimes even **negative effect**! (?)

Explanation: Because of the **truncation bias**. Workers of the same rate excluded if they work more hours and have higher income.

Example [Hausman, Wise (Econometrica) 1976]

Explanation: Because of the **truncation bias**. Workers of the same rate excluded if they work more hours and have higher income.



Example [Hausman, Wise (Econometrica) 1976]

Explanation: Because of the **truncation bias**. Workers of the same rate excluded if they work more hours and have higher income.



Multidimensional Example









Since Galton there has been an extensive literature on estimation from truncated samples, motivated by

measurement limitations,



Since Galton there has been an extensive literature on estimation from truncated samples, motivated by

measurement limitations,

ethical considerations,



Since Galton there has been an extensive literature on estimation from truncated samples, motivated by

measurement limitations,

- ethical considerations,
- privacy considerations,



Since Galton there has been an extensive literature on estimation from truncated samples, motivated by

measurement limitations,

ethical considerations,

privacy considerations,

For this reason truncated samples appear in:
physics,
economics,
social sciences,
psychological studies.

□ Known Survival Set

Gaussian Estimation [Daskalakis, Gouleakis, Tzamos, Z. '18]
 Linear Regression [Daskalakis, Gouleakis, Tzamos, Z. '19]

Unknown Survival Set [Kontonis, Tzamos, Z. '19]

Generation Known Survival Set

Gaussian Estimation [Daskalakis, Gouleakis, Tzamos, Z. '18]
 Linear Regression [Daskalakis, Gouleakis, Tzamos, Z. '19]

Unknown Survival Set [Kontonis, Tzamos, Z. '19]

Generation Known Survival Set

Gaussian Estimation [Daskalakis, Gouleakis, Tzamos, Z. '18]
 Linear Regression [Daskalakis, Gouleakis, Tzamos, Z. '19]

Unknown Survival Set [Kontonis, Tzamos, Z. '19]

Truncated Statistics



TRUNCATED AND CENSORED SAMPLES

STATISTICS: textbooks and monograph

Theory and Applications

A. CLIFFORD COHEN

Colume 119

REGRESSION MODELS Censored, Sample Selected, or Truncated Data *RICHARD BREEN* Series: Quantitative Applications in the Social Sciences in the Social Sciences

Truncated Statistics





STATISTICS, textbacks and monoon

Theory and Applications

A. CLIFFORD COHEN

REGRESSION Censored, Sample Selected, or Truncated Data *RICHARD BREEN* Series: Quantitative Applications in the Social Sciences 111 (*) a SAGE UNIVERSITY PAPER

1. Only for known axes aligned box!

Truncated Statistics





STATISTICS: textbooks and mono

Theory and Applications

A. CLIFFORD COHEN

REGRESSION MODELS Censored, Sample Selected, or Truncated Data *RICHARD BREEN* Series: Quantitative Applications in the Social Sciences Only for known axes aligned box!
 No sample complexity analysis.
Truncated Statistics





STATISTICS: textbooks and mo

Theory and Applications

A. CLIFFORD COHEN

REGRESSION MODELS Censored, sample Selected, or Truncated Data *RICHARD BREEN* Series: Quantitative Applications in the Social Sciences Only for known axes aligned box!
No sample complexity analysis.
No efficient algorithm.

Truncated Statistics



TRUNCATED AND CENSORED SAMPLES

Theory and Applications

A. CLIFFORD COHEN



- 1. Only for known axes aligned box!
- 2. No sample complexity analysis.
- 3. No efficient algorithm.

Even for the case of estimating a multi-normal distribution.

We provide efficient estimation for **truncated multi-normal distribution**.

- 1. Only for known axes aligned box!
- 2. No sample complexity analysis.
- 3. No efficient algorithm.

We provide efficient estimation for **truncated multi-normal distribution**.

1. Only for known axes aligned box!

- 2. No sample complexity analysis.
- 3. No efficient algorithm.

The survival set should only have non-trivial mass!

We provide efficient estimation for **truncated multi-normal distribution**.

1. Only for known axes aligned box!

- 2. No sample complexity analysis.
- 3. No efficient algorithm.

The survival set should only have non-trivial mass!

Can even be **unknown** under

stronger assumptions.

We provide efficient estimation for truncated multi-normal distribution.

- 1. Only for known axes aligned box!
- 2. No sample complexity analysis.
- 3. No efficient algorithm.

We get almost optimal # of samples.

We provide efficient estimation for **truncated multi-normal distribution**.

- 1. Only for known axes aligned box!
- 2. No sample complexity analysis.

3. No efficient algorithm.

We provide a **convex** programming formulation which can be solved efficiently.

Multivariate Gaussian Distribution $\mathcal{N}(\mu, \Sigma)$

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right).$$

Multivariate Gaussian Distribution $\mathcal{N}(\mu, \Sigma)$

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right).$$

Let $S \subseteq \mathbb{R}^d$, we define the mass of *S* with respect to measure $\mathcal{N}(\mu, \Sigma)$

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};S) = \int_{S} \mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};\boldsymbol{y}) d\boldsymbol{y}.$$

Let $S \subseteq \mathbb{R}^d$, we define the *truncated normal distribution* $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, S)$ as

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma},S;\boldsymbol{x}) = \begin{cases} \frac{1}{\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};S)} \cdot \mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};\boldsymbol{x}) & \boldsymbol{x} \in S \\ & 0 & \boldsymbol{x} \notin S \end{cases}.$$

Let $S \subseteq \mathbb{R}^d$, we define the *truncated normal distribution* $\mathcal{N}(\mu, \Sigma, S)$ as

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma},S;\boldsymbol{x}) = \begin{cases} \frac{1}{\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};S)} \cdot \mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma};\boldsymbol{x}) & \boldsymbol{x} \in S \\ & 0 & \boldsymbol{x} \notin S \end{cases}.$$

Given samples from $\mathcal{N}(\mu^*, \Sigma^*, S)$ compute estimates $\hat{\mu}$, $\hat{\Sigma}$ such that $d_{TV}(\mathcal{N}(\hat{\mu}, \hat{\Sigma}), \mathcal{N}(\mu^*, \Sigma^*)) \leq \varepsilon.$

Theorem ([Daskalakis Gouleakis Tzamos Z. '18]). Let $S \subseteq \mathbb{R}^d$ that satisfies Assumptions 1 and 2 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, \Sigma^*, S)$, then we can efficiently compute estimates $\hat{\mu}$ and $\hat{\Sigma}$ that satisfy with probability at least 99%:

$$d_{TV}(\mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}), \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)) \leq \varepsilon$$

for $n = \tilde{\Theta}(d^2 / \varepsilon^2)$.

Assumption 1 (CONSTANT MASS)

$$\mathcal{N}(\boldsymbol{\mu^*}, \boldsymbol{\Sigma^*}; S) \geq \alpha$$

Assumption 1 (CONSTANT MASS)

 $\mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*; S) \ge \alpha$ If $\alpha \to 0$ then # of sample $\to \infty$.

Assumption 1 (CONSTANT MASS)

 $\mathcal{N}(\boldsymbol{\mu^*}, \boldsymbol{\Sigma^*}; S) \geq \alpha$

Assumption 2 (ORACLE ACCESS)

Given $x \in \mathbb{R}^d$ we can answer if $x \in S$ or not.

Assumption 1 (CONSTANT MASS)

 $\mathcal{N}(\boldsymbol{\mu^*}, \boldsymbol{\Sigma^*}; S) \geq \alpha$

Assumption 2 (ORACLE ACCESS)

Given $x \in \mathbb{R}^d$ we can answer if $x \in S$ or not.

Impossible for unknown and **arbitrary** S!

In this talk

Theorem ([Daskalakis Gouleakis Tzamos Z. '18]). Let $S \subseteq \mathbb{R}^d$ that satisfies Assumptions 1 and 2 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, I, S)$, then we can efficiently compute an estimate $\hat{\mu}$ that satisfies with probability at least 99%:

 $d_{TV}(\mathcal{N}(\hat{\boldsymbol{\mu}}, \boldsymbol{I}), \mathcal{N}(\boldsymbol{\mu^*}, \boldsymbol{I})) \leq \varepsilon$

for $n = \tilde{\Theta}(d/\varepsilon^2)$.

In this talk

Theorem ([Daskalakis Gouleakis Tzamos Z. '18]). Let $S \subseteq \mathbb{R}^d$ that satisfies Assumptions 1 and 2 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, I, S)$, then we can efficiently compute an estimate $\hat{\mu}$ that satisfies with probability at least 99%:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2 \le \varepsilon$$

for $n = \tilde{\Theta}(d/\varepsilon^2)$.

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$$

4.
$$\hat{r} \leftarrow \text{Sample } \mathcal{N}(\mu^{(t-1)}, I, S)$$

5.
$$\boldsymbol{v}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do
- 3. $r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$
- 4. $\hat{r} \leftarrow \text{Sample } \mathcal{N}(\mu^{(t-1)}, I, S)$

5.
$$\boldsymbol{v}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6. $\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do
- 3. $r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$
- 4. $\hat{\mathbf{r}} \leftarrow \text{Sample } \mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \mathbf{I}, S)$

5.
$$\boldsymbol{v}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6. $\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$$

4.
$$\hat{r} \leftarrow \text{Sample } \mathcal{N}(\mu^{(t-1)}, I, S)$$

5.
$$\boldsymbol{\nu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$$

4.
$$\hat{\mathbf{r}} \leftarrow \text{Sample } \mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \mathbf{I}, S)$$

5.
$$\boldsymbol{v}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$$

4.
$$\hat{\mathbf{r}} \leftarrow \text{Sample } \mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \mathbf{I}, S)$$

5.
$$\boldsymbol{\nu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, S),$$

4.
$$\hat{\mathbf{r}} \leftarrow \text{Sample } \mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \mathbf{I}, S)$$

5.
$$\boldsymbol{\nu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \boldsymbol{\hat{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{ project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

7. output $\frac{1}{T+1} \sum_{t=0}^{T} \mu^{(t)}$

Algorithm: Stochastic Gradient Descent in the population negative log-likelihood function.

Algorithm: Stochastic Gradient Descent in the population negative log-likelihood function.

Conditions for Fast Convergence of SGD

1.
$$\mathbb{E}[\mathbf{r}-\hat{\mathbf{r}}] = \nabla \bar{\ell}(\boldsymbol{\mu}),$$

2. $\bar{\ell}(\mu)$ is *strongly* convex,

3.
$$\mathbb{E}\left[\|\boldsymbol{r}-\hat{\boldsymbol{r}}\|_{2}^{2}\right]$$
 is bounded.

Algorithm: Stochastic Gradient Descent in the population negative log-likelihood function.

Conditions for Fast Convergence of SGD

1.
$$\mathbb{E}[\boldsymbol{r}-\hat{\boldsymbol{r}}] = \nabla \bar{\ell}(\boldsymbol{\mu}),$$

2. $\bar{\ell}(\mu)$ is *strongly* convex,

3.
$$\mathbb{E}\left[\|\boldsymbol{r}-\hat{\boldsymbol{r}}\|_{2}^{2}\right]$$
 is bounded.

 $\overline{\ell}(\mu)$ is strongly convex.

Does not hold for any $\mu \in \mathbb{R}^d$.





- 1. estimate $\mu^{(0)} = \frac{1}{n} \sum_{i=1}^{n} x_i$,
- 2. for t = 1 ... T do

3.
$$r \leftarrow \text{Sample } \mathcal{N}(\mu^*, I, C),$$

4.
$$\hat{\mathbf{r}} \leftarrow \text{Sample } \mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \mathbf{I}, C)$$

5.
$$\boldsymbol{\nu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t-1)} - \eta_t (\boldsymbol{r} - \hat{\boldsymbol{r}})$$

6.
$$\mu^{(t)} \leftarrow \text{project } \nu^{(t)} \text{ to the ball } B = \left\{ x \mid \left\| x - \mu^{(0)} \right\| \le R \right\}$$

7. output $\frac{1}{T+1} \sum_{t=0}^{T} \mu^{(t)}$

Algorithm: Stochastic Gradient Descent in the population negative log-likelihood function.

Conditions for Fast Convergence of SGD

1.
$$\mathbb{E}[\mathbf{r}-\hat{\mathbf{r}}] = \nabla \bar{\ell}(\boldsymbol{\mu}),$$

2. $\bar{\ell}(\mu)$ is *strongly* convex,

3.
$$\mathbb{E}\left[\|\boldsymbol{r}-\hat{\boldsymbol{r}}\|_{2}^{2}\right]$$
 is bounded.



□ Known Survival Set

Gaussian Estimation [Daskalakis, Gouleakis, Tzamos, Z. '18]
Linear Regression [Daskalakis, Gouleakis, Tzamos, Z. '19]

Unknown Survival Set [Kontonis, Tzamos, Z. '19]



□ Known Survival Set

Gaussian Estimation [Daskalakis, Gouleakis, Tzamos, Z. '18]
Linear Regression [Daskalakis, Gouleakis, Tzamos, Z. '19]

Unknown Survival Set [Kontonis, Tzamos, Z. '19]

Goal: Learn the parameters of the Gaussian distribution and the set!

What learning the set means?


What learning the set means?



What learning the set means?



Gaussian Surface Area

$\mathcal{N}_0 = \mathcal{N}(0, I)$

Gaussian Surface Area

Gaussian Surface Area: $\Gamma(S)$ "Surface of *S* with respect to the Gaussian measure."

$$\Gamma(S) = \lim_{\delta \to 0} \frac{\mathcal{N}_0(S + B(\mathbf{0}, \delta)) - \mathcal{N}_0(S)}{\delta}$$

$\mathcal{N}_0 = \mathcal{N}(0, I)$

Gaussian Surface Area

Gaussian Surface Area: $\Gamma(S)$ "Surface of *S* with respect to the Gaussian measure."

$$\Gamma(S) = \lim_{\delta \to 0} \frac{\mathcal{N}_0(S + B(\mathbf{0}, \delta)) - \mathcal{N}_0(S)}{\delta}$$

$$\Gamma(\mathcal{S}) = \max_{S \in \mathcal{S}} \Gamma(S)$$

$\mathcal{N}_0 = \mathcal{N}(0, I)$

Gaussian Surface Area

Gaussian Surface Area: $\Gamma(S)$ "Surface of *S* with respect to the Gaussian measure."

$$\Gamma(S) = \lim_{\delta \to 0} \frac{\mathcal{N}_0(S + B(\mathbf{0}, \delta)) - \mathcal{N}_0(S)}{\delta}$$

$$\Gamma(\mathcal{S}) = \max_{S \in \mathcal{S}} \Gamma(S)$$





In this talk

Theorem ([Kontonis, Tzamos, Z. '19]). Let $S \subseteq \mathbb{R}^d$ unknown and $S \in S$ with Gaussian Surface Area $\Gamma(S)$ that satisfies Assumptions 1 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, I, S)$, then we can efficiently compute an estimate $\hat{\mu}$ that satisfies with probability at least 99%:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2 \leq \varepsilon$$

for $n = d^{\tilde{\Theta}(\Gamma^2(S)/\varepsilon^8)}$.

In this talk

Theorem ([Kontonis, Tzamos, Z. '19]). Let $S \subseteq \mathbb{R}^d$ unknown and $S \in S$ with Gaussian Surface Area $\Gamma(S)$ that satisfies Assumptions 1 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, I, S)$, then we can efficiently compute an estimate $\hat{\mu}$ that satisfies with probability at least 99%:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2 \leq \varepsilon$$

for $n = d^{\tilde{\Theta}(\Gamma^2(S)/\varepsilon^8)}$.

Concept Class	Surface Area	Sample Complexity
PTFs of degree <i>k</i>	<i>O</i> (<i>k</i>) [Kane '11]	$d^{O(k^2)}$
Intersections of <i>k</i> halfspaces	$O(\sqrt{\log k})$ [KOS '08]	$d^{O(\log k)}$
General convex sets	$O(d^{1/4})$ [Ball 1993]	$d^{O(\sqrt{d})}$



- ▷ **Univariate** Hermite Polynomials $h_0(x) = 1$, $h_1(x) = x$, $h_2(x) = \frac{x^2-1}{\sqrt{2}}$, ...,
- ▷ Multivariate: For $V = (v_1, ..., v_d) H_V(x) = \prod_{i=1}^d h_{v_i}(x_i)$



- ▷ **Univariate** Hermite Polynomials $h_0(x) = 1$, $h_1(x) = x$, $h_2(x) = \frac{x^2-1}{\sqrt{2}}$, ...,
- ▷ **Multivariate**: For $V = (v_1, ..., v_d)$ $H_V(x) = \prod_{i=1}^d h_{v_i}(x_i)$

Orthonormal basis under the Gaussian measure!



- ▷ Univariate Hermite Polynomials $h_0(x) = 1$, $h_1(x) = x$, $h_2(x) = \frac{x^2-1}{\sqrt{2}}$, ...,
- ▷ **Multivariate**: For $V = (v_1, \ldots, v_d)$ $H_V(x) = \prod_{i=1}^d h_{v_i}(x_i)$
- ▷ Approximate a function *f* using Hermite polynomials

 $f_m(x) = \sum_{V:|V| \le m} \hat{f}(V) \cdot H_V(x) \qquad \hat{f}(V) = \mathbb{E}_{x \sim \mathcal{N}_0}[f(x)H_V(x)]$



- ▷ Univariate Hermite Polynomials $h_0(x) = 1$, $h_1(x) = x$, $h_2(x) = \frac{x^2 1}{\sqrt{2}}$,...
- Multivariate
 Approximate



Unknown Set – Unknown Gaussian

 $\psi(x) = 1_S(x)w(x)$ weighted characteristic function of *S* $\Gamma(S)$ = Gaussian surface area of *S* α = Gaussian mass of *S*

Unknown Set – Unknown Gaussian

 $\psi(x) = 1_S(x)w(x)$ weighted characteristic function of *S* $\Gamma(S)$ = Gaussian surface area of *S* α = Gaussian mass of *S*

Theorem ([Kontonis, Tzamos, Z. '19]). If $m \ge O\left(\frac{\Gamma(S)^2}{\varepsilon^4}\right)$ then $\mathbb{E}_{x \sim \mathcal{N}_0}\left[(\psi(x) - \psi_m(x))^2\right] \le \varepsilon.$

Unknown Set – Unknown Gaussian

 $\psi(x) = 1_S(x)w(x)$ weighted characteristic function of *S* $\Gamma(S)$ = Gaussian surface area of *S* α = Gaussian mass of *S*

Theorem ([Kontonis, Tzamos, Z. '19]). If $m \ge O\left(\frac{\Gamma(S)^2}{\varepsilon^4}\right)$ then $\mathbb{E}_{x \sim \mathcal{N}_0}\left[(\psi(x) - \psi_m(x))^2\right] \le \varepsilon.$

We need #samples $n = d^{O(m)}$.

Optimization Problem

$$M_{\psi,h}(\mu) = \mathbb{E}_{x \sim \mathcal{N}_S^*} \left[e^{h(\mu;x)} \psi(x) \mathcal{N}_0(x) \right]$$

Optimization Problem

$$M_{\psi,h}(\mu) = \mathbb{E}_{x \sim \mathcal{N}_S^*} \left[e^{h(\mu;x)} \psi(x) \mathcal{N}_0(x) \right]$$

Theorem ([Kontonis, Tzamos, Z. '19]).

 $\triangleright M_{\psi,h}(\mu)$ is a convex function of μ .

$$\triangleright \ \nabla_{\mu} M_{\psi,h}(\mu) = 0 \Leftrightarrow \mu = \mu^*.$$

 \triangleright We have sample access to $\nabla_{\mu} M_{\psi,h}(\mu).(*)$

\Rightarrow We can run SGD!

Main Result

Theorem ([Kontonis, Tzamos, Z. '19]). Let $S \subseteq \mathbb{R}^d$ unknown and $S \in S$ with Gaussian Surface Area $\Gamma(S)$ that satisfies Assumptions 1 and x_1, x_2, \ldots, x_n i.i.d. samples from $\mathcal{N}(\mu^*, I, S)$, then we can efficiently compute an estimate $\hat{\mu}$ that satisfies with probability at least 99%:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}^*\|_2 \leq \varepsilon$$

for $n = d^{\tilde{\Theta}(\Gamma^2(S)/\varepsilon^8)}$.

Contents

Estimation Obstacles

Truncated Samples
 Inhomogeneous Population

Optimization Problems

Convex Landscape

Non-Convex Landscape

Global Convergence of Iterative Methods



In 1894 Karl Pearson analyzed a data set of the lengths of 1000 crabs.













How we can explain this asymmetry?





Inhomogeneous population!





Pearson [Pearson 1894] used the theoretically rigorous **method of moments** in order to estimate the means and variances of the mixture of normal distributions but:



Pearson [Pearson 1894] used the theoretically rigorous **method of moments** in order to estimate the means and variances of the mixture of normal distributions but:

- 1. only asymptotic results,
- 2. only single dimensional.

Expectation – Maximization Algorithm

EM algorithm was proposed by [Dempster Lair Rubin 1977] as a general technique to solve estimation problems when likelihood is **non-convex**.

Expectation – Maximization Algorithm

EM algorithm was proposed by [Dempster Lair Rubin 1977] as a general technique to solve estimation problems when likelihood is **non-convex**.



...with more than 50,000 citations according to Google Scholar!

We assume that we have access to samples from the following mixture of two multidimensional normal distributions

$$p_{\mu_1,\mu_2}(x) = \frac{1}{2}\mathcal{N}(\mu_1, I; x) + \frac{1}{2}\mathcal{N}(\mu_2, I; x).$$

Goal: Estimate the means μ_1, μ_2 .

Numerical Example We are given the following samples				Crab Species		Leng	th	
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

Nume We ar	erical Exar e given the	n <u>ple</u> e following	g samples		Crab Species		Leng	th
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

We assume that length follows a Gaussian distribution with unit variance.

Nume We ar	erical Exar e given the	n <u>ple</u> e following	g samples		Crab Species		Leng	th
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

and we are looking for an estimation of the mean length of species 1, μ_1 and the mean length of species 2 μ_2 . Initially we set $\mu_1 = 1.6$, $\mu_2 = 1.7$.

Nume We ar	erical Exar e given the	n <u>ple</u> e following	g samples		Crab Species		Leng	;th
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

and we are looking for an estimation of the mean length of species 1, μ_1 and the mean length of species 2 μ_2 . Initially we set $\mu_1 = 1.6$, $\mu_2 = 1.7$.

(E-step) Species 1 Posterior

0.505	0.504	0.502	0.496	0.492	0.499	0.495	
pecies 2 Posterior							

0.495	0.496	0.498	0.504	0.508	0.501	0.505
-------	-------	-------	-------	-------	-------	-------

Numerical Example We are given the following samples					Crab Species		Leng	th
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

(M-step) Species 1 Mean Length

0.505.1.5+0.504.1.52+0.502.1.57+0.496.1.72+0.492.1.83+0.499.1.64+0.495 1.75

0.505+0.504+0.502+0.496+0.492+0.499+0.495

Numerical Example We are given the following samples					Crab Species		Leng	th
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

(M-step) Species 2 Mean Length

0.495·1.5+0.496·1.52+0.498·1.57+0.504·1.72+0.508·1.83+0.501·1.64+0.505·1.75 0.495+0.496+0.498+0.504+0.508+0.501+0.505

Numerical Example We are given the following samples				Crab Species		Leng	th	
	1.5	1.52	1.57	1.72	1.83	1.64	1.75	

(M-step) Species 2 Mean Length

0.495·1.5+0.496·1.52+0.498·1.57+0.504·1.72+0.508·1.83+0.501·1.64+0.505·1.75 0.495+0.496+0.498+0.504+0.508+0.501+0.505




The Estimation Algorithm

1. pick $\mu^{(0)}$ at random,

2. for t = 1 ... T do

3.
$$p_i \leftarrow \text{evaluate } \frac{\mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{I}; \boldsymbol{x}_i)}{\mathcal{N}(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{I}; \boldsymbol{x}_i) + \mathcal{N}(-\boldsymbol{\mu}^{(t-1)}, \boldsymbol{I}; \boldsymbol{x}_i)}$$

4.
$$\boldsymbol{\mu}^{(t)} \leftarrow \frac{\sum_i p_i \boldsymbol{x}_i}{\sum_i p_i}$$

5. output $\mu^{(T)}$

Our Result

Theorem ([Daskalakis Tzamos Z. '17]). Given access to $n = \tilde{\Theta}(d/\epsilon^2)$ samples from p_{μ} the EM algorithm gives an estimation $\hat{\mu}$ s.t. with probability 99 %

 $d_{TV}(p_{\mu},p_{\hat{\mu}})\leq \varepsilon.$

Our Result

Theorem ([Daskalakis Tzamos Z. '17]). Given access to $n = \tilde{\Theta}(d/\epsilon^2)$ samples from p_{μ} the EM algorithm gives an estimation $\hat{\mu}$ s.t. with probability 99 %

 $d_{TV}(p_{\mu},p_{\hat{\mu}})\leq \varepsilon.$

First global convergence guarantees for a non-trivial instance of EM algorithm since its definition in 1977!

Concurrent and independent work by [Xu Hsu Maleki '16].











































Iterative Methods



Iterative Methods



Iterative Methods



Solution Concept

\varepsilon-Approximate Fixed Point Let d^* be a distance metric,

 $d^*(\hat{x}, f(\hat{x})) \leq \varepsilon.$

Solution Concept

ε -Approximate Fixed Point Let d^* be a distance metric,

 $d^*(\hat{x}, f(\hat{x})) \leq \varepsilon.$

How to prove **global convergence** to approximate fixed points?

A function $f : \mathcal{D} \to \mathcal{D}$ is a **contraction map** with respect to *d* if $d(f(x), f(y)) \le c \cdot d(x, y) \quad \forall x, y \in \mathcal{D}.$

and *c* < 1.

A function $f : \mathcal{D} \to \mathcal{D}$ is a **contraction map** with respect to *d* if $d(f(x), f(y)) \leq c \cdot d(x, y) \quad \forall x, y \in \mathcal{D}.$

and c < 1.



Banach's Fixed Point Theorem

Banach's Fixed Point Theorem



Banach's Fixed Point Theorem



Banach's Fixed Point Theorem



Banach's Fixed Point Theorem



Banach's Fixed Point Theorem

Let $f : \mathcal{D} \to \mathcal{D}$ be a contraction map, then f has a unique fixed point x^* . Moreover f convergences globally and geometrically to x^* .

How general is Banach's Theorem?







A function $f : \mathcal{D} \to \mathcal{D}$ is a **contraction map** with respect to *d* if

$$d(f(x), f(y)) \leq c \cdot d(x, y) \quad \forall x, y \in \mathcal{D}.$$

and *c* < 1.

We are free to choose any distance metric d!





Converse Fixed Point Theorems!

Bessaga's and Meyers's Converse Fixed Point Theorems Let $f : \mathcal{D} \to \mathcal{D}$ with a unique fixed point x^* , then there exists a distance metric d' such that f is a contraction map with respect to d'.

Our Converse Fixed Point Theorems

Theorem ([Daskalakis Tzamos Z. '18]). Let $f : \mathcal{D} \to \mathcal{D}$ with a unique fixed point x^* . If f converges globally and geometrically to x^* then this can be shown through a contraction map argument.

Our Converse Fixed Point Theorems

Theorem ([Daskalakis Tzamos Z. '18]). Let $f : \mathcal{D} \to \mathcal{D}$ with a unique fixed point x^* . If f converges globally and geometrically to x^* then this can be shown through a contraction map argument.

Fast Convergence
Our Converse Fixed Point Theorems

Theorem ([Daskalakis Tzamos Z. '18]). Let $f : \mathcal{D} \to \mathcal{D}$ with a unique fixed point x^* . If f converges globally and geometrically to x^* then this can be shown through a contraction map argument.

Fast Convergence Hardness

Our Converse Fixed Point Theorems

Theorem ([Daskalakis Tzamos Z. '18]). Let $f : \mathcal{D} \to \mathcal{D}$ with a unique fixed point x^* . If f converges globally and geometrically to x^* then this can be shown through a contraction map argument.



Complexity of Total Search Problems

FNP: class of search problems whose decision version is in NP.

TFNP: class of total search problems of FNP, i.e. a solution always exists [Megiddo Papadimitriou '91]

Subclasses of TFNP introduced by [Johnson Papadimitriou Yannakakis '88], [**Papadimitriou '94]**, [Daskalakis Papadimitriou '11], [Jerabek '16]



Complexity of Total Search Problems

FNP: class of search problems whose decision version is in NP.

TFNP: class of total search problems of FNP, i.e. a solution always exists [Megiddo Papadimitriou '91]

Subclasses of TFNP introduced by [Johnson Papadimitriou Yannakakis '88], [**Papadimitriou '94]**, [Daskalakis Papadimitriou '11], [Jerabek '16]





Estimation Obstacles

O Truncated SamplesO Inhomogeneous Population

Optimization Problems

Convex Landscape

Non-Convex Landscape

Iterative Methods

Global Convergence of Iterative Methods

• Truncated estimation in other models? [Daskalakis Ilyas Z. '20]

- Truncated estimation in other models? [Daskalakis Ilyas Z. '20]
- Estimation in truncated + Inhomogeneous Population?

- Truncated estimation in other models? [Daskalakis Ilyas Z. '20]
- Estimation in truncated + Inhomogeneous Population?
- (known estimation problem) + (censoring)

- Truncated estimation in other models? [Daskalakis Ilyas Z. '20]
- Estimation in truncated + Inhomogeneous Population?
- (known estimation problem) + (censoring)
- Better understanding of complexity of Fixed Points, e.g. relations to cryptography? [Bitanski Paneth Rosen '15], [Hubacek Yogev '17] [Sotiraki Zirdelis Z. '18]

- Truncated estimation in other models? [Daskalakis Ilyas Z. '20]
- Estimation in truncated + Inhomogeneous Population?
- (known estimation problem) + (censoring)
- Better understanding of complexity of Fixed Points, e.g. relations to cryptography? [Bitanski Paneth Rosen '15], [Hubacek Yogev '17] [Sotiraki Zirdelis Z. '18]

Thank you! 😊