

Computational Barriers in Statistical Estimation

Andrea Montanari

Stanford University

July 9, 2017

Statistical estimation/Statistical learning

Class of models ($\Theta \subseteq \mathbb{R}^d$)

$$\mathcal{C}_\Theta \equiv \{ \mathbb{P}_\theta : \theta \in \Theta \}$$

Data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim_{iid} \mathbb{P}_{\theta_0}(\cdot)$$

Estimate θ_0 from data $\mathbf{x}_1^n = (\mathbf{x}_1, \dots, \mathbf{x}_n)$

Statistical estimation/Statistical learning

Class of models ($\Theta \subseteq \mathbb{R}^d$)

$$\mathcal{C}_\Theta \equiv \{ \mathbb{P}_\theta : \theta \in \Theta \}$$

Data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim_{iid} \mathbb{P}_{\theta_0}(\cdot)$$

Estimate θ_0 from data $\mathbf{x}_1^n = (\mathbf{x}_1, \dots, \mathbf{x}_n)$

Minimax theory

Loss function

$$\begin{aligned} L : \mathbb{R}^p \times \mathbb{R}^p &\rightarrow \mathbb{R} \\ (\hat{\theta}, \theta_0) &\mapsto L(\hat{\theta}, \theta_0) \end{aligned}$$

Minimax risk

$$R_n^*(\Theta) = \inf_{\hat{\theta}(\cdot)} \sup_{\theta_0 \in \Theta} \mathbb{E}_{\theta_0} L(\hat{\theta}(x_1^n), \theta_0)$$

[Wald, 1950]

Kindergarten example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \mathbb{R}^d.$$

Theorem

$$R_n^*(\Theta) = \frac{d}{n}.$$

- Foundation of the least squares, maximum likelihood, ...

Kindergarten example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \mathbb{R}^d.$$

Theorem

$$R_n^*(\Theta) = \frac{d}{n}.$$

- ▶ Foundation of the least squares, maximum likelihood, ...

Kindergarten example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \mathbb{R}^d.$$

Theorem

$$R_n^*(\Theta) = \frac{d}{n}.$$

- Foundation of the least squares, maximum likelihood, ...

A more sophisticated example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_0 \leq s_0\}.$$

Theorem (Donoho, Johnstone 1990s)

If $s_0/d \rightarrow 0$, then

$$R_n^*(\Theta) = \frac{2s_0}{n} \log(d/s_0) \cdot (1 + o(1)).$$

- ▶ Key role in compressed sensing, sparse learning, ...

A more sophisticated example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_0 \leq s_0\}.$$

Theorem (Donoho, Johnstone 1990s)

If $s_0/d \rightarrow 0$, then

$$R_n^*(\Theta) = \frac{2s_0}{n} \log(d/s_0) \cdot (1 + o(1)).$$

► Key role in compressed sensing, sparse learning, ...

A more sophisticated example

$$\mathbb{P}_{\boldsymbol{\theta}}(\cdot) = \mathbf{N}(\boldsymbol{\theta}, \mathbf{I}_d)$$

$$L(\boldsymbol{\theta}_0, \hat{\boldsymbol{\theta}}) = \|\boldsymbol{\theta}_0 - \hat{\boldsymbol{\theta}}\|_2^2$$

$$\Theta = \{\boldsymbol{\theta} \in \mathbb{R}^d : \|\boldsymbol{\theta}\|_0 \leq s_0\}.$$

Theorem (Donoho, Johnstone 1990s)

If $s_0/d \rightarrow 0$, then

$$R_n^*(\Theta) = \frac{2s_0}{n} \log(d/s_0) \cdot (1 + o(1)).$$

- ▶ Key role in compressed sensing, sparse learning, ...

What is this talk about?

$$R_n^{\text{Poly}}(\Theta) = \inf_{\hat{\theta}(\cdot) \in \text{Poly}} \sup_{\theta_0 \in \Theta} \mathbb{E}_{\theta_0} L(\hat{\theta}(x_1^n), \theta_0)$$

Developments

- ▶ Often we expect $R_n^{\text{Poly}}(\Theta) \gtrsim R_n(\Theta)$
- ▶ Accurate predictions
- ▶ Convergence of fields (Statistics, CS Theory, Physics)
- ▶ New algorithms?

What is this talk about?

$$R_n^{\text{Poly}}(\Theta) = \inf_{\hat{\theta}(\cdot) \in \text{Poly}} \sup_{\theta_0 \in \Theta} \mathbb{E}_{\theta_0} L(\hat{\theta}(x_1^n), \theta_0)$$

Developments

- ▶ Often we expect $R_n^{\text{Poly}}(\Theta) \gtrsim R_n(\Theta)$
- ▶ Accurate predictions
- ▶ Convergence of fields (Statistics, CS Theory, Physics)
- ▶ New algorithms?

Outline

- 1 Local algorithms
- 2 Landscapes
- 3 SDP relaxations
- 4 Conclusion

Local and message passing algorithms

Example #1: Sparse low-rank matrix

Unknowns:

$$\Theta(s_0, d) = \left\{ \boldsymbol{\theta} \in \{0, 1\}^d : \|\boldsymbol{\theta}\|_0 = s_0 \right\}.$$

Loss:

$$L(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}) = \frac{1}{d} \text{Hamming}(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta})$$

Example #1: Sparse low-rank matrix

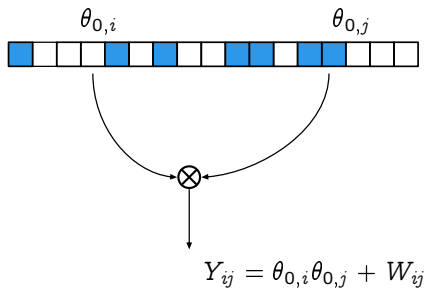
Data: $(\mathbf{x}_1, \dots, \mathbf{x}_n)$

$$\mathbf{x}_\ell = (i_\ell, j_\ell, Y_{i_\ell, j_\ell}) \in [d] \times [d] \times \mathbb{R},$$

$$i_\ell, j_\ell \sim_{iid} \text{Unif}([d]),$$

$$Y_{i_\ell, j_\ell} |_{i_\ell, j_\ell} \sim N(\theta_{0, i_\ell} \theta_{0, j_\ell}, \sigma^2)$$

Pictorially



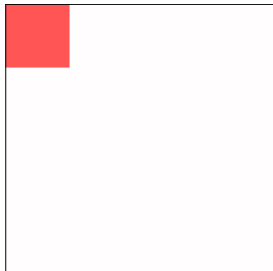
Example #1: A different formulation

Data $Y \in \mathbb{R}^{n \times n}$

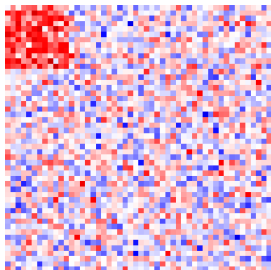
$$Y = \mathcal{P}_E(\theta_0 \theta_0^\top + W)$$

- ▶ $E \subseteq \binom{[d]}{2}$ uniformly random s.t. $|E| = n$
- ▶ $\mathcal{P}_E : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$: projector that zeroes entries not in E
- ▶ $(W_{ij})_{i < j} \sim_{iid} N(0, \sigma^2)$, $W = W^\top$

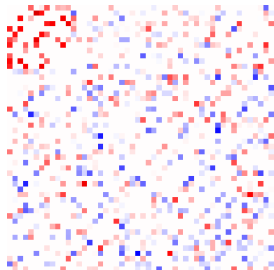
$$\theta_0 \theta_0^\top$$



$$\theta_0 \theta_0^\top + W$$



$$Y = \mathcal{P}_E(\theta_0 \theta_0^\top + W)$$



Example #1: Yet another formulation

- ▶ $G = (V, E) \sim \mathcal{G}(d, n)$
(uniform random with d vertices, n edges)

- ▶ For each $(i, j) \in E$

$$Y_{ij} = \theta_i \theta_j + W_{ij}, \quad W_{ij} \sim \mathcal{N}(0, \sigma^2)$$

Problem parameters

$$\delta = \frac{n}{d} \quad \text{(half) average graph degree}$$

$$\varepsilon = \frac{s_0}{d} \quad \text{sparsity}$$

Asymptotics

- ▶ Dense graph: $n \asymp d^2$, $s_0 \asymp \sqrt{d}$ (\sim Planted clique problem)
- ▶ Sparse graph: $n \asymp d$, $s_0 \asymp d$ (δ, ε fixed)

Problem parameters

$$\delta = \frac{n}{d} \quad \text{(half) average graph degree}$$

$$\varepsilon = \frac{s_0}{d} \quad \text{sparsity}$$

Asymptotics

- ▶ Dense graph: $n \asymp d^2$, $s_0 \asymp \sqrt{d}$ (\sim Planted clique problem)
- ▶ Sparse graph: $n \asymp d$, $s_0 \asymp d$ (δ, ε fixed)

$$R^{\text{Poly}}(\delta, \varepsilon; d) \equiv R_{n=d\delta}^{\text{Poly}}(\Theta(s_0 = d\varepsilon, d))$$

$$\lim_{d \rightarrow \infty} R^{\text{Poly}}(\delta, \varepsilon; d) = ?$$

We do not know, but ...

$$R^{\text{Poly}}(\delta, \varepsilon; d) \equiv R_{n=d\delta}^{\text{Poly}}(\Theta(s_0 = d\varepsilon, d))$$

$$\lim_{d \rightarrow \infty} R^{\text{Poly}}(\delta, \varepsilon; d) = ?$$

We do not know, but ...

First simplifications

- ▶ Worst case prior

$$\theta \sim \text{Unif}(\Theta(s_0, d))$$

- ▶ Roughly ($\varepsilon = s_0/d$)

$$(\theta_i)_{i \leq d} \sim_{iid} \text{Bern}(\varepsilon)$$

First simplifications

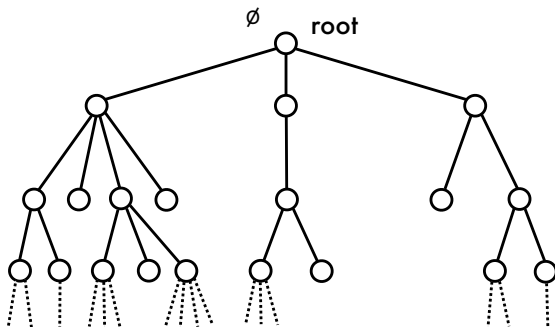
- ▶ Worst case prior

$$\theta \sim \text{Unif}(\Theta(s_0, d))$$

- ▶ Roughly ($\varepsilon = s_0/d$)

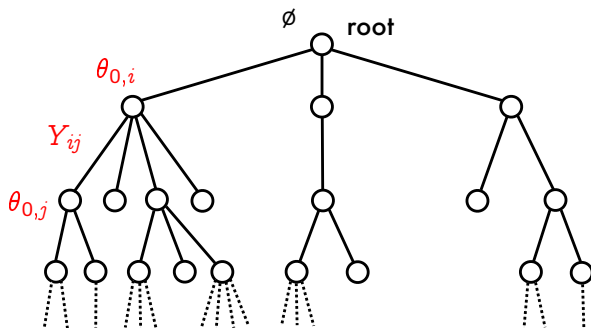
$$(\theta_i)_{i \leq d} \sim_{iid} \text{Bern}(\varepsilon)$$

Local weak limit: $d \rightarrow \infty$



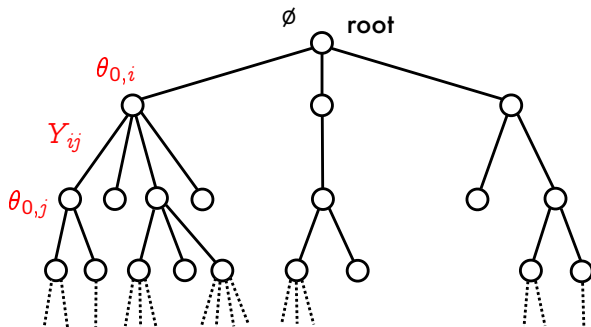
$$G \xrightarrow{lwc} \text{GW}(\text{Pois}(2\delta))$$

Local weak limit: $d \rightarrow \infty$



$$(G, \theta, \mathbf{Y}) \xrightarrow{lwc} (\text{GW}(\text{Pois}(2\delta)), \theta, \mathbf{Y})$$

First guess



$$\lim_{d \rightarrow \infty} R^{\text{Poly}}(\delta, \varepsilon, d) \stackrel{?}{=} \inf_{\hat{\theta}} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(\mathbf{Y}) \neq \theta_{0,\emptyset})$$

It gets interesting

$$\lim_{d \rightarrow \infty} R^{\text{Poly}}(\delta, \varepsilon, d) \stackrel{?}{=} \inf_{\hat{\theta}} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(\mathbf{Y}) \neq \theta_{0, \emptyset})$$

How do you define the r.h.s.?

A natural idea:

$$Y_{\ell}^0 = (Y_{ij} : d(\emptyset, i) \leq \ell),$$
$$\hat{\theta}_{\emptyset}(Y_{\ell}^0) = \arg \max_{\sigma \in \{0,1\}} \mathbb{P}(\theta_{\emptyset} = \sigma | Y_{\ell}^0).$$

It gets interesting

$$\lim_{d \rightarrow \infty} R^{\text{Poly}}(\delta, \varepsilon, d) \stackrel{?}{=} \inf_{\hat{\theta}} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(\mathbf{Y}) \neq \theta_{0, \emptyset})$$

How do you define the r.h.s.?

A natural idea:

$$\begin{aligned} \mathbf{Y}_{\ell}^0 &= (Y_{ij} : d(\emptyset, i) \leq \ell), \\ \hat{\theta}_{\emptyset}(\mathbf{Y}_{\ell}^0) &= \arg \max_{\sigma \in \{0,1\}} \mathbb{P}(\theta_{\emptyset} = \sigma \mid \mathbf{Y}_{\ell}^0). \end{aligned}$$

Risk of local algorithms

$$R^{\text{loc}}(\delta, \varepsilon) = \lim_{\ell \rightarrow \infty} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(\mathbf{Y}_{\ell}^0) \neq \theta_{0,\emptyset})$$

Remark

Belief propagation achieves $R^{\text{loc}}(\delta, \varepsilon)$

Risk of local algorithms

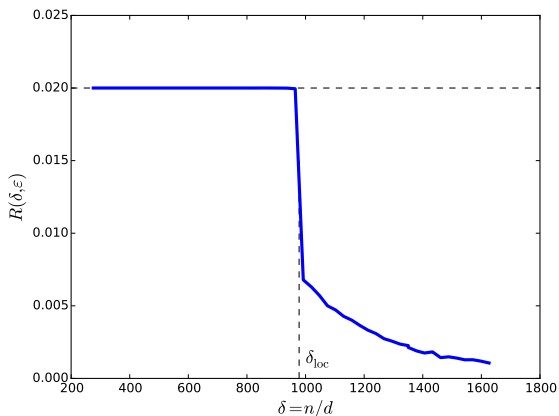
$$R^{\text{loc}}(\delta, \varepsilon) = \lim_{\ell \rightarrow \infty} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(\mathbf{Y}_{\ell}^0) \neq \theta_{0, \emptyset})$$

Remark

Belief propagation achieves $R^{\text{loc}}(\delta, \varepsilon)$

Can be computed

$(\varepsilon = 0.02, \sigma = 1.5)$



$$R^{loc}(\delta, \varepsilon) = \lim_{\ell \rightarrow \infty} \mathbb{P}_{\text{Tree}}(\hat{\theta}_{\emptyset}(Y_{\ell}^0) \neq \theta_{0, \emptyset})$$

Phase transition for local algorithms

Theorem (\sim Deshpande, Montanari, 2015)

As $\varepsilon \rightarrow 0$ (with $\delta \rightarrow \infty$, σ fixed)

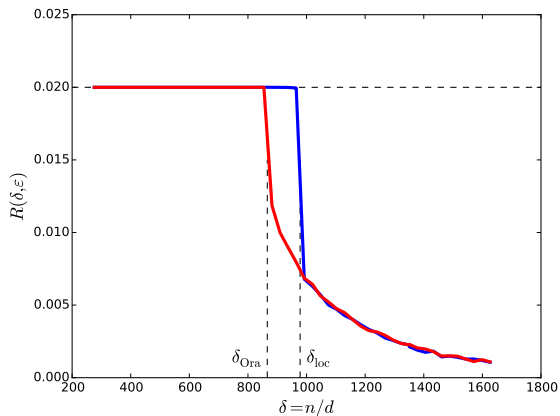
$$\frac{1}{\varepsilon} R^{\text{loc}}(\delta, \varepsilon) \rightarrow \begin{cases} 1 & \text{if } \delta \leq (1 - o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon\varepsilon^2}, \\ 0 & \text{if } \delta \geq (1 + o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon\varepsilon^2}, \end{cases}$$

A different definition of the limit

$$\begin{aligned} \mathbf{Y}_\ell^+ &= \{(Y_{ij} : d(\emptyset, i) \leq \ell); (\theta_i : d(\emptyset, i) > \ell)\} \\ \hat{\theta}_\emptyset(\mathbf{Y}_\ell^+) &= \arg \max_{\sigma \in \{0,1\}} \mathbb{P}(\theta_\emptyset = \sigma \mid \mathbf{Y}_\ell^+) \end{aligned}$$

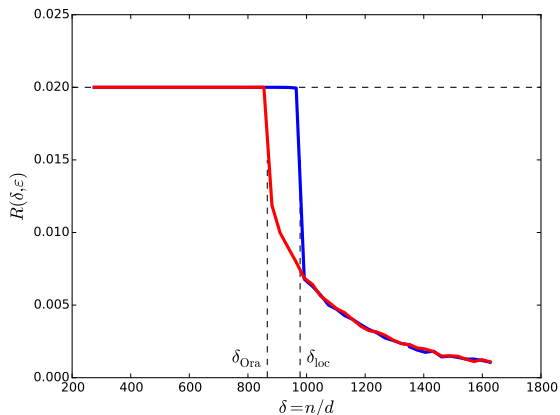
$$R^{\text{Ora}}(\delta, \varepsilon) = \lim_{\ell \rightarrow \infty} \mathbb{P}_{\text{Tree}}(\hat{\theta}_\emptyset(\mathbf{Y}_\ell^0) \neq \theta_{0,\emptyset})$$

It also can be computed



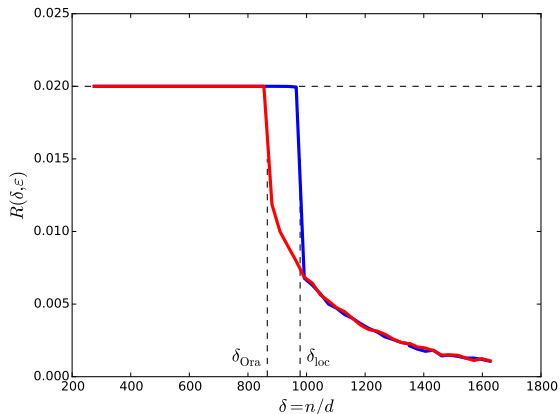
$$R^{\text{Ora}}(\delta, \epsilon) \leq R^*(\delta, \epsilon) \leq R^{\text{Poly}}(\delta, \epsilon) \leq R^{\text{loc}}(\delta, \epsilon)$$

It also can be computed



$$R^{\text{Ora}}(\delta, \epsilon) \leq R^*(\delta, \epsilon) \leq R^{\text{Poly}}(\delta, \epsilon) \leq R^{\text{loc}}(\delta, \epsilon)$$

It also can be computed



$$R^{\text{Ora}}(\delta, \epsilon) \leq R^*(\delta, \epsilon) \leq R^{\text{Poly}}(\delta, \epsilon) \leq R^{\text{loc}}(\delta, \epsilon)$$

Minimax risk

Theorem (\sim Montanari, 2015)

As $\varepsilon \rightarrow 0$ (with $\delta \rightarrow \infty$, σ fixed)

$$\frac{1}{\varepsilon} R^*(\delta, \varepsilon) \rightarrow \begin{cases} 1 & \text{if } \delta \leq (1 - o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \\ 0 & \text{if } \delta \geq (1 + o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \end{cases}$$

- ▶ Sharp threshold (dense): Lelarge, Miolane 2017; Barbier et al. 2017
- ▶ Do not know of any polytime algorithm working for

$$1.01 \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon) < \delta < 0.99 \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon).$$

Minimax risk

Theorem (\sim Montanari, 2015)

As $\varepsilon \rightarrow 0$ (with $\delta \rightarrow \infty$, σ fixed)

$$\frac{1}{\varepsilon} R^*(\delta, \varepsilon) \rightarrow \begin{cases} 1 & \text{if } \delta \leq (1 - o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \\ 0 & \text{if } \delta \geq (1 + o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \end{cases}$$

- ▶ Sharp threshold (dense): Lelarge, Miolane 2017; Barbier et al. 2017
- ▶ Do not know of any polytime algorithm working for

$$1.01 \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon) < \delta < 0.99 \cdot \frac{\sigma^2}{2\varepsilon \varepsilon^2}.$$

Minimax risk

Theorem (\sim Montanari, 2015)

As $\varepsilon \rightarrow 0$ (with $\delta \rightarrow \infty$, σ fixed)

$$\frac{1}{\varepsilon} R^*(\delta, \varepsilon) \rightarrow \begin{cases} 1 & \text{if } \delta \leq (1 - o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \\ 0 & \text{if } \delta \geq (1 + o_\varepsilon(1)) \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon), \end{cases}$$

- ▶ Sharp threshold (dense): Lelarge, Miolane 2017; Barbier et al. 2017
- ▶ Do not know of any polytime algorithm working for

$$1.01 \cdot \frac{\sigma^2}{2\varepsilon} \log(1/\varepsilon) < \delta < 0.99 \cdot \frac{\sigma^2}{2\varepsilon \varepsilon^2}.$$

Open problems

- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by Gibbs sampling?
- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by convex optimization?
- ▶ Is $R^{\text{Poly}}(\delta, \varepsilon) = R^{\text{loc}}(\delta, \varepsilon)$?

Open problems

- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by Gibbs sampling?
- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by convex optimization?
- ▶ Is $R^{\text{Poly}}(\delta, \varepsilon) = R^{\text{loc}}(\delta, \varepsilon)$?

Open problems

- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by Gibbs sampling?
- ▶ Can we beat $R^{\text{loc}}(\delta, \varepsilon)$ by convex optimization?
- ▶ Is $R^{\text{Poly}}(\delta, \varepsilon) = R^{\text{loc}}(\delta, \varepsilon)$?

Ubiquitous

Sparse principal component analysis

- ▶ Berthet, Rigollet, 2013
- ▶ Deshpande, Montanari, 2014
- ▶ Barbier et al 2016; Miolane 2017

Hidden clique problem

(2 asymmetric communities)

- ▶ Jerrum, 1992
- ▶ Feige, Krauthgamer, 200
- ▶ Deshpande, Montanari, 2015; Montanari 2015

Community detection ($k \geq 5$ symmetric communities)

- ▶ Decelle, Krzakala, Moore, Zdeborova 2011
- ▶ Bordenave, Lelarge, Massoulié 2015
- ▶ Abbe, Sandon 2015

Tensor PCA

- ▶ See below

Landscapes

Empirical risk minimization/M-estimation

Class of models ($\Theta \subseteq \mathbb{R}^d$)

$$\mathcal{C}_\Theta \equiv \{ \mathbb{P}_\theta : \theta \in \Theta \}$$

Data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim_{iid} \mathbb{P}_{\theta_0}(\cdot)$$

$$\text{minimize } \hat{\mathcal{L}}_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n \ell(\theta; \mathbf{x}_i)$$

Empirical risk minimization/M-estimation

Class of models ($\Theta \subseteq \mathbb{R}^d$)

$$\mathcal{C}_\Theta \equiv \{ \mathbb{P}_\theta : \theta \in \Theta \}$$

Data

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \sim_{iid} \mathbb{P}_{\theta_0}(\cdot)$$

$$\text{minimize } \hat{\mathcal{L}}_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n \ell(\theta; \mathbf{x}_i)$$

Empirical risk minimization/M-estimation

Rationale

$$\theta_0 = \arg \min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta), \quad \mathcal{L}(\theta) = \mathbb{E} \widehat{\mathcal{L}}_n(\theta)$$

- ▶ What can we say generically?
- ▶ How does complexity show up?

- ▶ What can we say generically?
- ▶ How does complexity show up?

Uniform convergence

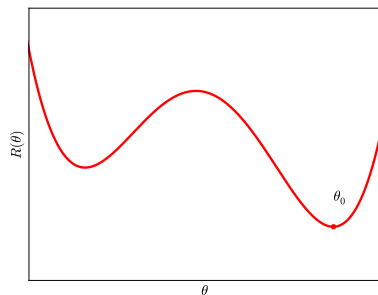
Theorem (Vapnik, Chervonenkis, 1968; ...)

Under conditions [omitted], with high probability

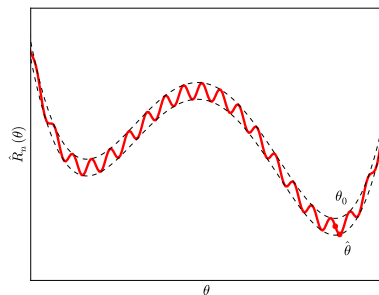
$$\sup_{\theta \in \Theta} |\hat{\mathcal{L}}_n(\theta) - \mathcal{L}(\theta)| \leq C \sqrt{\frac{d_*}{n}}.$$

(d_ = VC dimension; ...).*

Uniform convergence



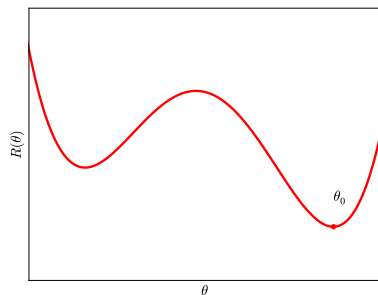
Population risk



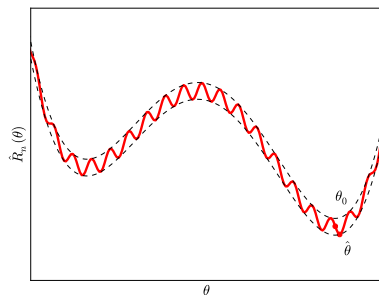
Empirical risk

Will optimization algorithms get stuck in local minima?
Landscape analysis

Uniform convergence



Population risk

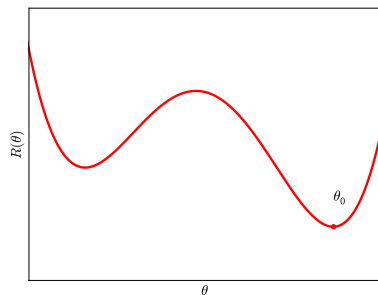


Empirical risk

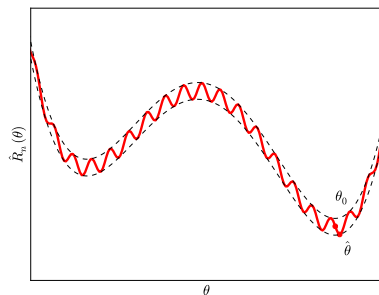
Will optimization algorithms get stuck in local minima?

Landscape analysis

Uniform convergence



Population risk



Empirical risk

Will optimization algorithms get stuck in local minima?
Landscape analysis

Assumptions

$\theta \in \mathbb{B}^p(r)$ = Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

'Under mild assumptions'

Assumptions

$\theta \in \mathbb{B}^p(r)$ = Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

'Under mild assumptions'

Assumptions

$\theta \in \mathbb{B}^p(r) =$ Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

'Under mild assumptions'

Assumptions

$\theta \in \mathbb{B}^p(r)$ = Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

'Under mild assumptions'

Assumptions

$\theta \in \mathbb{B}^p(r)$ = Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

'Under mild assumptions'

Assumptions

$\theta \in \mathbb{B}^p(r)$ = Ball of radius r in \mathbb{R}^p

Data: $\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_n$ iid

A1 $\nabla_{\theta} \ell(\theta; \mathbf{Z})$ is τ^2 -sub-Gaussian

A2 For any $\lambda \in \mathbb{B}^p(1)$, $\mathcal{Z}_{\lambda} \equiv \langle \lambda, \nabla^2 \ell(\theta; \mathbf{Z}) \lambda \rangle$ is τ^2 -sub-Exponential.

A3 The Hessian of the population risk at $\mathbf{0}$ is bounded by a polynomial

$$\|\nabla^2 \mathcal{L}(\mathbf{0})\|_{\text{op}} \leq \tau^2 d^C.$$

A4 The Hessian of the loss is Lipschitz continuous with integrable constant

$$\mathbb{E}\{\|\nabla^2 \ell(\cdot; \mathbf{Z})\|_{\text{Lip}}\} \leq \tau^3 d^C.$$

‘Under mild assumptions’

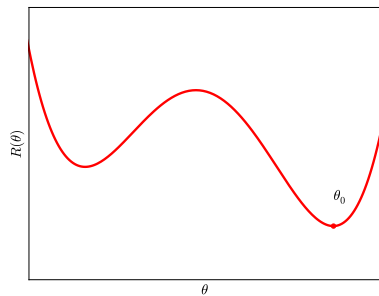
Lemma

Under assumptions A1, A2, A3, A4, if $n \geq Cp \log p$, then with probability at least $1 - \delta$, the following hold:

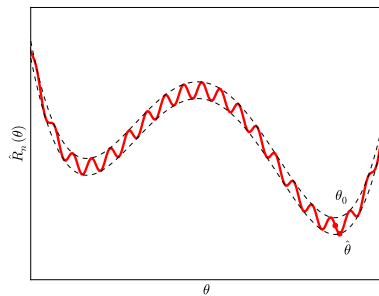
$$\sup_{\boldsymbol{\theta} \in \mathbb{B}^p(\tau)} \left\| \nabla \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) - \nabla \mathcal{L}(\boldsymbol{\theta}) \right\|_2 \leq \tau \sqrt{\frac{Cd \log n}{n}},$$

$$\sup_{\boldsymbol{\theta} \in \mathbb{B}^p(\tau)} \left\| \nabla^2 \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) - \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \right\|_{\text{op}} \leq \tau^2 \sqrt{\frac{Cd \log n}{n}}.$$

This cannot happen!

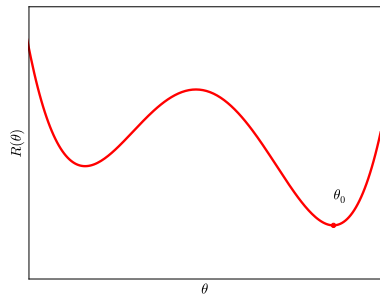


Population risk

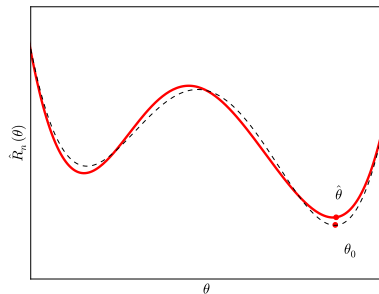


Empirical risk

This can happen!



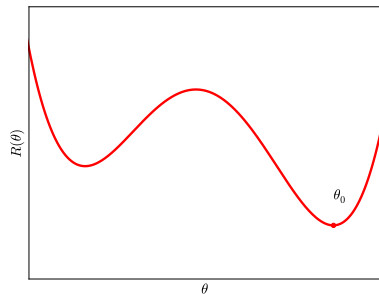
Population risk



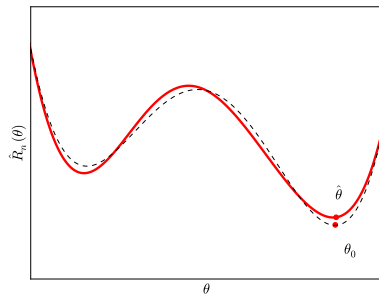
Empirical risk

Nice population risk \Rightarrow Nice empirical risk

This can happen!



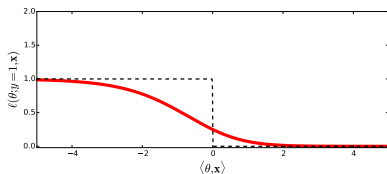
Population risk



Empirical risk

Nice population risk \Rightarrow Nice empirical risk

Example: Binary classification



$$\mathbf{z}_i = (y_i, \mathbf{x}_i), \quad y_i \in \{0, 1\}, \quad \mathbf{x}_i \in \mathbb{R}^d$$

$$\mathbb{P}(y_i = 1 | \mathbf{x}_i) = \sigma(\langle \boldsymbol{\theta}_0, \mathbf{x}_i \rangle)$$

$$\hat{\mathcal{L}}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle))^2.$$

- ▶ Rosenblatt 1958 (Perceptron); ... many extensions
- ▶ More robust than logistic regression

Sample application: Binary classification

Theorem (Mei, Bai, Montanari 2017)

Assume \mathbf{X}_i to be centered, sub-Gaussian, with $\mathbb{E}\{\mathbf{X}\mathbf{X}^\top\} \succeq \delta\mathbf{I}_d$.
For nice^a functions σ , whp:

1. The population risk has a unique critical point $\hat{\theta}_n$.
2. Gradient descent converges exponentially fast to $\hat{\theta}_n$.
3. The estimation error is $\|\hat{\theta}_n - \theta_0\|_2 \leq C\sqrt{(d \log n)/n}$.

$$^a \sigma'(x) > 0, \|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty \leq C.$$

Similar results for robust regression, one-bit compressed sensing, ...
[see paper]

Sample application: Binary classification

Theorem (Mei, Bai, Montanari 2017)

Assume X_i to be centered, sub-Gaussian, with $\mathbb{E}\{X X^\top\} \succeq \delta I_d$.
For nice^a functions σ , whp:

1. The population risk has a unique critical point $\hat{\theta}_n$.
2. Gradient descent converges exponentially fast to $\hat{\theta}_n$.
3. The estimation error is $\|\hat{\theta}_n - \theta_0\|_2 \leq C \sqrt{(d \log n)/n}$.

$${}^a \sigma'(x) > 0, \|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty \leq C.$$

Similar results for robust regression, one-bit compressed sensing, ...
[see paper]

Sample application: Binary classification

Theorem (Mei, Bai, Montanari 2017)

Assume \mathbf{X}_i to be centered, sub-Gaussian, with $\mathbb{E}\{\mathbf{X}\mathbf{X}^\top\} \succeq \delta\mathbf{I}_d$.
For nice^a functions σ , whp:

1. The population risk has a unique critical point $\hat{\boldsymbol{\theta}}_n$.
2. Gradient descent converges exponentially fast to $\hat{\boldsymbol{\theta}}_n$.
3. The estimation error is $\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|_2 \leq C\sqrt{(d \log n)/n}$.

$$^a \sigma'(x) > 0, \|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty \leq C.$$

Similar results for robust regression, one-bit compressed sensing, ...
[see paper]

Sample application: Binary classification

Theorem (Mei, Bai, Montanari 2017)

Assume \mathbf{X}_i to be centered, sub-Gaussian, with $\mathbb{E}\{\mathbf{X}\mathbf{X}^\top\} \succeq \delta\mathbf{I}_d$.
For nice^a functions σ , whp:

1. The population risk has a unique critical point $\hat{\boldsymbol{\theta}}_n$.
2. Gradient descent converges exponentially fast to $\hat{\boldsymbol{\theta}}_n$.
3. The estimation error is $\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|_2 \leq C\sqrt{(d \log n)/n}$.

$${}^a \sigma'(x) > 0, \|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty \leq C.$$

Similar results for robust regression, one-bit compressed sensing, ...
[see paper]

Sample application: Binary classification

Theorem (Mei, Bai, Montanari 2017)

Assume \mathbf{X}_i to be centered, sub-Gaussian, with $\mathbb{E}\{\mathbf{X}\mathbf{X}^\top\} \succeq \delta\mathbf{I}_d$.
For nice^a functions σ , whp:

1. The population risk has a unique critical point $\hat{\boldsymbol{\theta}}_n$.
2. Gradient descent converges exponentially fast to $\hat{\boldsymbol{\theta}}_n$.
3. The estimation error is $\|\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0\|_2 \leq C\sqrt{(d \log n)/n}$.

$$^a \sigma'(x) > 0, \|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty \leq C.$$

Similar results for robust regression, one-bit compressed sensing, ...
[see paper]

Non-convex literature

Convergence to ‘statistical neighborhood’

- ▶ Loh, Wainwright, 2012
- ▶ Loh Wainwright, 2013
- ▶ Yang, Wang, Liu, Eldar, Zhang, 2015
- ▶ ...

Smart initialization

- ▶ Keshavan, Montanari, Oh, 2009
- ▶ Chen, Candés, 2015
- ▶ Anandkumar, Ge, Jenzamin, 2015
- ▶ ...

Unique local minimum

- ▶ Loh, 2015 [High-dim. regression, $n \gtrsim s_0^2$]
- ▶ Sun, Qu, Wright, 2016 [Phase retrieval]
- ▶ Ge, Lee, Ma, 2016 [Matrix completion]

Non-convex literature

Convergence to ‘statistical neighborhood’

- ▶ Loh, Wainwright, 2012
- ▶ Loh Wainwright, 2013
- ▶ Yang, Wang, Liu, Eldar, Zhang, 2015
- ▶ ...

Smart initialization

- ▶ Keshavan, Montanari, Oh, 2009
- ▶ Chen, Candés, 2015
- ▶ Anandkumar, Ge, Jenzamin, 2015
- ▶ ...

Unique local minimum

- ▶ Loh, 2015 [High-dim. regression, $n \gtrsim s_0^2$]
- ▶ Sun, Qu, Wright, 2016 [Phase retrieval]
- ▶ Ge, Lee, Ma, 2016 [Matrix completion]

- ▶ What can we say generically?
- ▶ How does complexity show up?

Intuition

Population risk very flat \Leftrightarrow Many local minima

[Close to where local algorithms fail?]

Simplest example: Spiked tensor model

- ▶ Unknown parameter $\theta_0 \in \mathbb{R}^n$, $\|\theta_0\|_2 = 1$

- ▶ Data¹

$$Y = \lambda \theta_0^{\otimes k} + W$$

- ▶ $W =$ symmetric Gaussian noise tensor.
 $(W_{i_1, \dots, i_k})_{i_1 < \dots < i_k} \sim_{iid} N(0, 1/n)$

¹Equivalently, $Y_{i_1, \dots, i_k} = \lambda \theta_{0, i_1} \cdots \theta_{0, i_k} + W_{i_1, \dots, i_k}$.

Spiked tensor model: What do we know?

$$Y = \lambda \theta_0^{\otimes k} + W$$

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

For any $\varepsilon > 0$, there exist constants λ_{IT} , $\lambda_{ML}(\varepsilon)$, $C(\varepsilon)$, such that:

- ▶ If $\lambda > \lambda_{ML}(\varepsilon)$, then $\mathbb{E}\{|\langle \hat{\theta}^{ML}, \theta_0 \rangle|\} \geq 1 - \varepsilon$.
- ▶ No estimator can achieve $\mathbb{E}\{|\langle \hat{\theta}, \theta_0 \rangle|\} \geq \varepsilon$ unless $\lambda > \lambda_{IT}$.
- ▶ There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{Poly}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided $\lambda \geq C(\varepsilon)n^{(k-2)/4}$.

No efficient estimator is known for $1 \ll \lambda \ll n^{(k-2)/4}$!

Spiked tensor model: What do we know?

$$Y = \lambda \theta_0^{\otimes k} + W$$

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

For any $\varepsilon > 0$, there exist constants λ_{IT} , $\lambda_{\text{ML}}(\varepsilon)$, $C(\varepsilon)$, such that:

- ▶ *If $\lambda > \lambda_{\text{ML}}(\varepsilon)$, then $\mathbb{E}\{|\langle \hat{\theta}^{\text{ML}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$.*
- ▶ *No estimator can achieve $\mathbb{E}\{|\langle \hat{\theta}, \theta_0 \rangle|\} \geq \varepsilon$ unless $\lambda > \lambda_{\text{IT}}$.*
- ▶ *There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided $\lambda \geq C(\varepsilon)n^{(k-2)/4}$.*

No efficient estimator is known for $1 \ll \lambda \ll n^{(k-2)/4}$!

Spiked tensor model: What do we know?

$$Y = \lambda \theta_0^{\otimes k} + W$$

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

For any $\varepsilon > 0$, there exist constants λ_{IT} , $\lambda_{\text{ML}}(\varepsilon)$, $C(\varepsilon)$, such that:

- ▶ *If $\lambda > \lambda_{\text{ML}}(\varepsilon)$, then $\mathbb{E}\{|\langle \hat{\theta}^{\text{ML}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$.*
- ▶ *No estimator can achieve $\mathbb{E}\{|\langle \hat{\theta}, \theta_0 \rangle|\} \geq \varepsilon$ unless $\lambda > \lambda_{\text{IT}}$.*
- ▶ *There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided $\lambda \geq C(\varepsilon)n^{(k-2)/4}$.*

No efficient estimator is known for $1 \ll \lambda \ll n^{(k-2)/4}$!

Spiked tensor model: What do we know?

$$Y = \lambda \theta_0^{\otimes k} + W$$

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

For any $\varepsilon > 0$, there exist constants λ_{IT} , $\lambda_{\text{ML}}(\varepsilon)$, $C(\varepsilon)$, such that:

- ▶ *If $\lambda > \lambda_{\text{ML}}(\varepsilon)$, then $\mathbb{E}\{|\langle \hat{\theta}^{\text{ML}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$.*
- ▶ *No estimator can achieve $\mathbb{E}\{|\langle \hat{\theta}, \theta_0 \rangle|\} \geq \varepsilon$ unless $\lambda > \lambda_{\text{IT}}$.*
- ▶ *There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided $\lambda \geq C(\varepsilon)n^{(k-2)/4}$.*

No efficient estimator is known for $1 \ll \lambda \ll n^{(k-2)/4}$!

Spiked tensor model: What do we know?

$$Y = \lambda \theta_0^{\otimes k} + W$$

Theorem (Montanari, Richard, 2014; Hopkins, Shi, Steurer, 2015)

For any $\varepsilon > 0$, there exist constants λ_{IT} , $\lambda_{\text{ML}}(\varepsilon)$, $C(\varepsilon)$, such that:

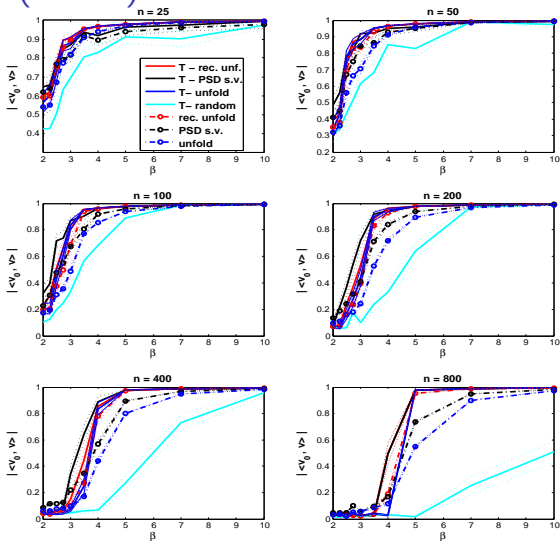
- ▶ *If $\lambda > \lambda_{\text{ML}}(\varepsilon)$, then $\mathbb{E}\{|\langle \hat{\theta}^{\text{ML}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$.*
- ▶ *No estimator can achieve $\mathbb{E}\{|\langle \hat{\theta}, \theta_0 \rangle|\} \geq \varepsilon$ unless $\lambda > \lambda_{\text{IT}}$.*
- ▶ *There exists poly-time estimator achieving $\mathbb{E}\{|\langle \hat{\theta}^{\text{Poly}}, \theta_0 \rangle|\} \geq 1 - \varepsilon$, provided $\lambda \geq C(\varepsilon)n^{(k-2)/4}$.*

No efficient estimator is known for $1 \ll \lambda \ll n^{(k-2)/4}$!

More precise results

- ▶ Montanari, Reichman, Zeitouni, 2015
- ▶ Bandeira, Perry, Wein, 2017
- ▶ Krzakala, Lelarge, Miolane, Zdeborova, 2017
- ▶ ...

In practice ($k = 3$)



What does the landscape look like?

Maximum likelihood

$$\begin{aligned} & \text{minimize} && \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) = \| \mathbf{Y} - \lambda \boldsymbol{\theta}^{\otimes k} \|_F \\ & \text{subject to} && \| \boldsymbol{\theta} \|_2 = 1 \end{aligned}$$

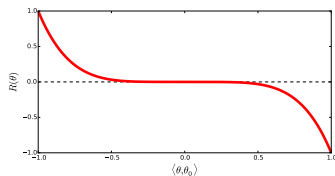
$$\widehat{\mathcal{L}}_n(\boldsymbol{\theta}) = \text{const.} - 2\lambda \langle \mathbf{Y}, \boldsymbol{\theta}^{\otimes k} \rangle$$

Maximum likelihood

$$\begin{aligned} & \text{minimize} && \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) = \| \mathbf{Y} - \lambda \boldsymbol{\theta}^{\otimes k} \|_F \\ & \text{subject to} && \| \boldsymbol{\theta} \|_2 = 1 \end{aligned}$$

$$\widehat{\mathcal{L}}_n(\boldsymbol{\theta}) = \text{const.} - 2\lambda \langle \mathbf{Y}, \boldsymbol{\theta}^{\otimes k} \rangle$$

Risk



Maximum likelihood

$$\begin{aligned} &\text{minimize} && \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) = -\langle \mathbf{Y}, \boldsymbol{\theta}^{\otimes k} \rangle \\ &\text{subject to} && \|\boldsymbol{\theta}\|_2 = 1 \end{aligned}$$

‘Population’ risk

$$\mathcal{L}(\boldsymbol{\theta}) = -\lambda \langle \boldsymbol{\theta}_0, \boldsymbol{\theta} \rangle^k$$

Back-of-the-envelope

Expected gradient

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} \boldsymbol{\theta}_0$$

Random initialization $\langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle = \Theta(n^{-1/2})$:

$$\begin{aligned} \langle \boldsymbol{\theta}_0, \nabla \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) \rangle &= -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} - k \langle \mathbf{W}, \boldsymbol{\theta}_0 \otimes \boldsymbol{\theta}^{\otimes(k-1)} \rangle \\ &= -\lambda \Theta(n^{-(k-1)/2}) + \Theta(n^{-1/2}) \end{aligned}$$

► Convergence: $\lambda \gg n^{(k-2)/2}$

Back-of-the-envelope

Expected gradient

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} \boldsymbol{\theta}_0$$

Random initialization $\langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle = \Theta(n^{-1/2})$:

$$\begin{aligned} \langle \boldsymbol{\theta}_0, \nabla \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) \rangle &= -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} - k \langle \mathbf{W}, \boldsymbol{\theta}_0 \otimes \boldsymbol{\theta}^{\otimes(k-1)} \rangle \\ &= -\lambda \Theta(n^{-(k-1)/2}) + \Theta(n^{-1/2}) \end{aligned}$$

► Convergence: $\lambda \gg n^{(k-2)/2}$

Back-of-the-envelope

Expected gradient

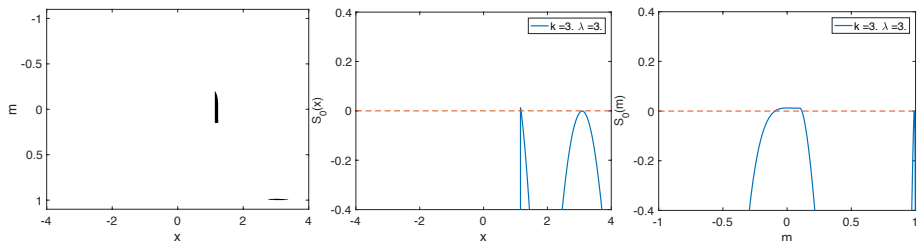
$$\nabla \mathcal{L}(\boldsymbol{\theta}) = -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} \boldsymbol{\theta}_0$$

Random initialization $\langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle = \Theta(n^{-1/2})$:

$$\begin{aligned} \langle \boldsymbol{\theta}_0, \nabla \widehat{\mathcal{L}}_n(\boldsymbol{\theta}) \rangle &= -k\lambda \langle \boldsymbol{\theta}, \boldsymbol{\theta}_0 \rangle^{k-1} - k \langle \mathbf{W}, \boldsymbol{\theta}_0 \otimes \boldsymbol{\theta}^{\otimes(k-1)} \rangle \\ &= -\lambda \Theta(n^{-(k-1)/2}) + \Theta(n^{-1/2}) \end{aligned}$$

► Convergence: $\lambda \gg n^{(k-2)/2}$

Expected number of local minima: $k = 3, \lambda = 3$



- ▶ Exponential in black region ($m = \langle \theta, \theta_0 \rangle, x = \langle Y, \theta^{\otimes k} \rangle$)
- ▶ N = number of local minima

$$\mathbb{E}N(m, x) = e^{nS_0(m, x) + o(n)}$$

[Ben Arous, Mei, Montanari, Nica, 2017]

Complexity of landscape $\overset{?}{\leftrightarrow}$ Complexity for local algorithms

SDP relaxations

An emerging dichotomy

In several statistical estimation problems

- ▶ Either local (or message-passing) algorithms work...
- ▶ ...or SDP hierarchies do not work

Why?

An emerging dichotomy

In several statistical estimation problems

- ▶ Either local (or message-passing) algorithms work...
- ▶ ...or SDP hierarchies do not work

Why?

An emerging dichotomy

In several statistical estimation problems

- ▶ Either local (or message-passing) algorithms work...
- ▶ ...or SDP hierarchies do not work

Why?

A possible explanation

Perhaps SDPs on random instances can be solved by local algorithms...

Simplest example

Centered adjacency matrix of $G = (V, E)$ ($d =$ average degree)

$$A_{ij}^{\text{cen}} = \begin{cases} 1 - \frac{d}{n} & \text{if } (i, j) \in E, \\ -\frac{d}{n} & \text{otherwise.} \end{cases}$$

SDP(A^{cen}):

$$\begin{aligned} & \text{maximize} && \langle A^{\text{cen}}, \mathbf{X} \rangle, \\ & \text{subject to} && \mathbf{X} \in \mathbb{R}^{n \times n}, \mathbf{X} \succeq 0, \\ & && X_{ii} = 1. \end{aligned}$$

- ▶ Graph clustering, embedding, testing latent structure,...

What does it mean?

Input: Graph $G_n = (V_n, E_n)$

1. Generate $z = (z(i))_{i \in V} \sim_{iid} \mathcal{N}(0, 1)$
2. Compute, for each $v \in V_n$, $\xi_v = F(\mathcal{B}_\ell(v; G_n), z|_{\mathcal{B}_\ell(v; G_n)})$
3. Output $\mathbf{X} = \mathbb{E}_z\{\xi\xi^T\} + \dots \in \mathbb{R}^{n \times n}$

Can this achieve $\langle \mathbf{A}^{\text{cen}}, \mathbf{X} \rangle \geq (1 - o_n(1))\text{SDP}(\mathbf{A}^{\text{cen}})$?

Erdős-Renyi random graphs

Theorem (Fan, Montanari, 2016)

Let $G \sim \mathcal{G}(n, d/n)$ and $A^{\text{cen}} = A_G^{\text{cen}}$. Then, a.s.,

$$\begin{aligned} 2\sqrt{d} \left(1 - \frac{1}{d+1}\right) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq 2\sqrt{d} \left(1 - \frac{1}{2d}\right). \end{aligned}$$

Further, the lower bound is achieved by local algorithms.

- ▶ A local algorithm achieves $8/9$ of $\text{SDP}(A^{\text{cen}})$.

Erdős-Renyi random graphs

Theorem (Fan, Montanari, 2016)

Let $G \sim \mathcal{G}(n, d/n)$ and $A^{\text{cen}} = A_G^{\text{cen}}$. Then, a.s.,

$$\begin{aligned} 2\sqrt{d} \left(1 - \frac{1}{d+1}\right) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq 2\sqrt{d} \left(1 - \frac{1}{2d}\right). \end{aligned}$$

Further, the lower bound is achieved by local algorithms.

- ▶ A local algorithm achieves $8/9$ of $\text{SDP}(A^{\text{cen}})$.

Erdős-Renyi random graphs

Theorem (Fan, Montanari, 2016)

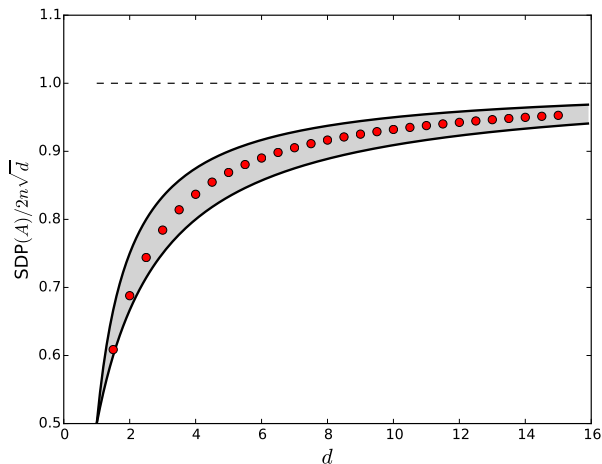
Let $G \sim \mathcal{G}(n, d/n)$ and $A^{\text{cen}} = A_G^{\text{cen}}$. Then, a.s.,

$$\begin{aligned} 2\sqrt{d} \left(1 - \frac{1}{d+1}\right) &\leq \liminf_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq \\ &\leq \limsup_{n \rightarrow \infty} \frac{1}{n} \text{SDP}(A^{\text{cen}}) \leq 2\sqrt{d} \left(1 - \frac{1}{2d}\right). \end{aligned}$$

Further, the lower bound is achieved by local algorithms.

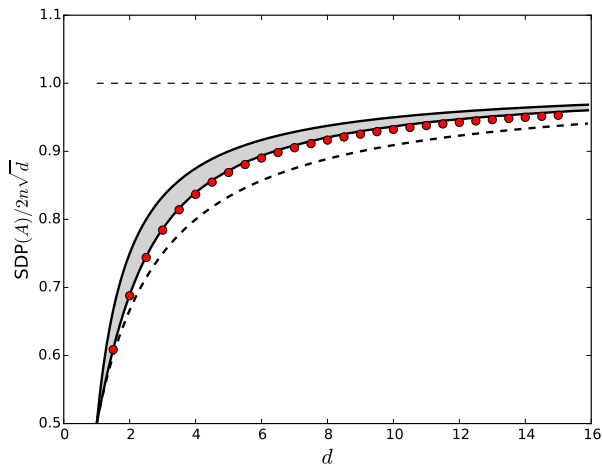
- ▶ A local algorithm achieves $8/9$ of $\text{SDP}(A^{\text{cen}})$.

Bounds vs numerical simulations



[Related results for planted partition model, regular graphs, ...]

A better local algorithm (see paper)



[Related results for planted partition model, regular graphs, ...]

Conclusion

Conclusion

- ▶ Which statistical problems are tractable?
- ▶ Multiple points of view:
 - ▶ Local-message passing algorithms
 - ▶ Landscape analysis
 - ▶ SDP hierarchies

Thanks!

Conclusion

- ▶ Which statistical problems are tractable?
- ▶ Multiple points of view:
 - ▶ Local-message passing algorithms
 - ▶ Landscape analysis
 - ▶ SDP hierarchies

Thanks!