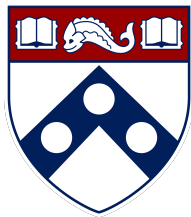


## **Nonconvex Optimization for High-Dimensional Estimation (Part 3)**



Yuxin Chen

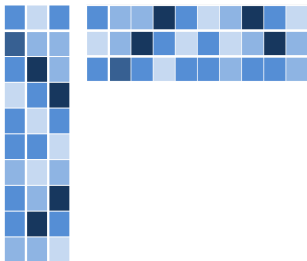
Wharton Statistics & Data Science, Spring 2022

*Bridging convex and nonconvex optimization in  
estimation and inference*

# Noisy low-rank matrix completion

observations:  $M_{i,j} = M_{i,j}^* + \text{noise}, \quad (i, j) \in \Omega$

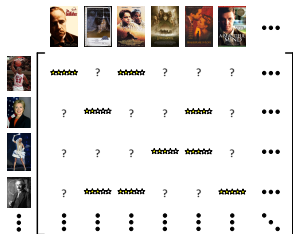
goal: estimate  $M^*$



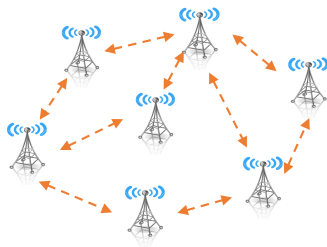
unknown rank- $r$  matrix  $M^* \in \mathbb{R}^{n \times n}$



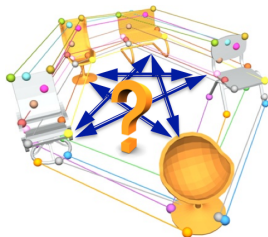
sampling set  $\Omega$



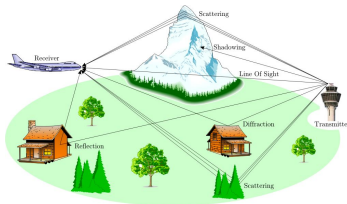
## recommendation systems



localization



shape matching



channel estimation



# Noisy low-rank matrix completion

---

observations:  $M_{i,j} = M_{i,j}^* + \text{noise}, \quad (i,j) \in \Omega$

goal: estimate  $M^*$

**convex relaxation:**

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \underbrace{\sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2}_{\text{squared loss}} + \lambda \|\mathbf{Z}\|_*$$

$$- \quad \|\mathbf{Z}\|_* = \sum_{i=1}^n \sigma_i(\mathbf{Z})$$

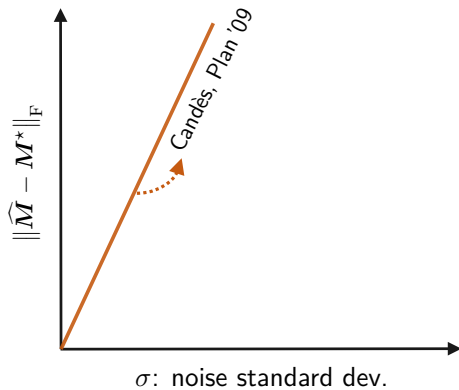
# Prior statistical guarantees for convex relaxation

---

- **random sampling:** each  $(i, j) \in \Omega$  indep. with prob.  $p$
- **random noise:** i.i.d. sub-Gaussian noise with variance  $\sigma^2$
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ : rank  $r = O(1)$ , well-conditioned,...

Candès, Plan '09

$\sigma n^{1.5}$

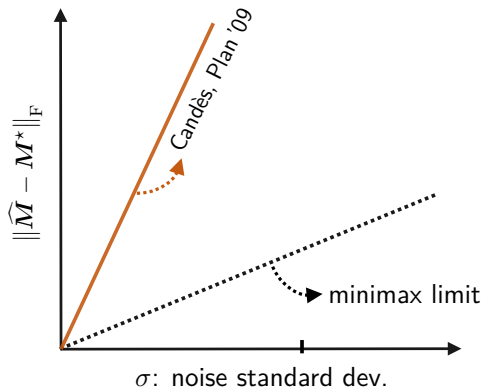


minimax limit

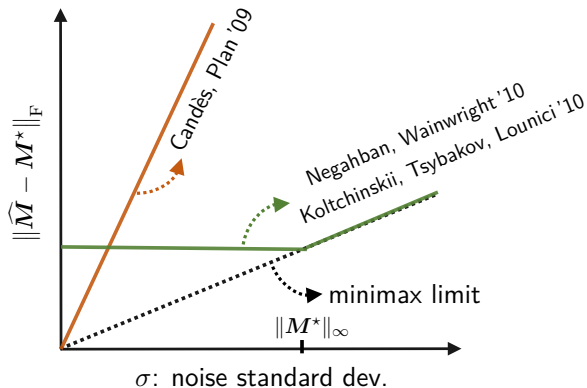
Candès, Plan '09

$$\sigma \sqrt{n/p}$$

$$\sigma n^{1.5}$$

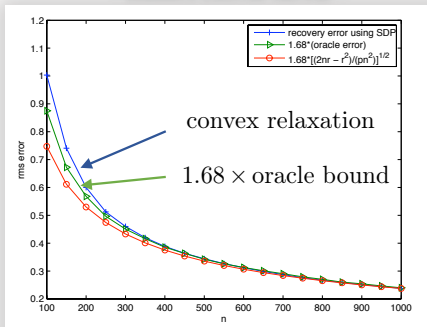


minimax limit	$\sigma \sqrt{n/p}$
Candès, Plan '09	$\sigma n^{1.5}$
Negahban, Wainwright '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$
Koltchinskii, Tsybakov, Lounici '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$



# Matrix Completion with Noise

Emmanuel J. Candès and Yaniv Plan



*Existing theory for convex relaxation does not match practice . . .*

# Matrix Completion with Noise

Emmanuel J. Candès and Yaniv Plan

with adversarial noise. Consequently, our analysis loses a  $\sqrt{n}$  factor vis a vis an optimal bound that is achievable via the help of an oracle.

*Existing theory for convex relaxation does not match practice . . .*

# What are the roadblocks?

---

Strategy:  $M^{\text{cvx}}$  is optimizer if  $\underbrace{\text{there exists } W}_{\text{dual certificate}}$  s.t.

$(M^{\text{cvx}}, W)$  obeys KKT optimality condition



# What are the roadblocks?

---

Strategy:  $M^{\text{cvx}}$  is optimizer if  $\underbrace{\text{there exists } W}_{\text{dual certificate}}$  s.t.

$(M^{\text{cvx}}, W)$  obeys KKT optimality condition



David Gross

- **noiseless case:**  $\underbrace{M^{\text{cvx}} \leftarrow M^*}_{\text{exact recovery}}; W \leftarrow \text{golfing scheme}$

# What are the roadblocks?

---

Strategy:  $M^{\text{cvx}}$  is optimizer if  $\underbrace{\text{there exists } W}_{\text{dual certificate}}$  s.t.

$(M^{\text{cvx}}, W)$  obeys KKT optimality condition



David Gross

- **noiseless case:**  $\underbrace{M^{\text{cvx}} \leftarrow M^*}_{\text{exact recovery}}; W \leftarrow \text{golfing scheme}$
- **noisy case:**  $M^{\text{cvx}}$  is very complicated, hard to construct  $W \dots$

dual certification (golfing scheme)



dual certification (golfing scheme)

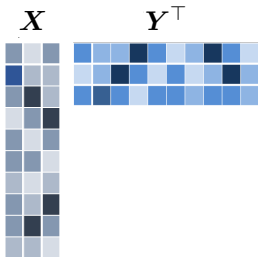


nonconvex optimization

## A detour: nonconvex optimization

---

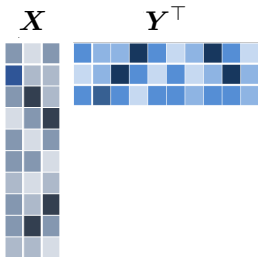
**Burer–Monteiro:** represent  $Z$  by  $XY^\top$  with  $\underbrace{X, Y \in \mathbb{R}^{n \times r}}_{\text{low-rank factors}}$



# A detour: nonconvex optimization

---

**Burer–Monteiro:** represent  $Z$  by  $XY^\top$  with  $\underbrace{X, Y \in \mathbb{R}^{n \times r}}_{\text{low-rank factors}}$



$$\underset{X, Y \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(X, Y) = \underbrace{\sum_{(i,j) \in \Omega} \left[ (XY^\top)_{i,j} - M_{i,j} \right]^2}_{\text{squared loss}} + \text{reg}(X, Y)$$

# A detour: nonconvex optimization

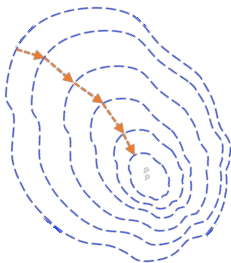
---

- Burer, Monteiro '03
- Rennie, Srebro '05
- Keshavan, Montanari, Oh '09 '10
- Jain, Netrapalli, Sanghavi '12
- Hardt '13
- Sun, Luo '14
- Chen, Wainwright '15
- Tu, Boczar, Simchowitz, Soltanolkotabi, Recht '15
- Zhao, Wang, Liu '15
- Zheng, Lafferty '16
- Yi, Park, Chen, Caramanis '16
- Ge, Lee, Ma '16
- Ge, Jin, Zheng '17
- Ma, Wang, Chi, Chen '17
- Chen, Li '18
- Chen, Liu, Li '19
- ...

# A detour: nonconvex optimization

---

$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}, \mathbf{Y}) = \sum_{(i,j) \in \Omega} \left[ (\mathbf{X} \mathbf{Y}^\top)_{i,j} - M_{i,j} \right]^2 + \text{reg}(\mathbf{X}, \mathbf{Y})$$

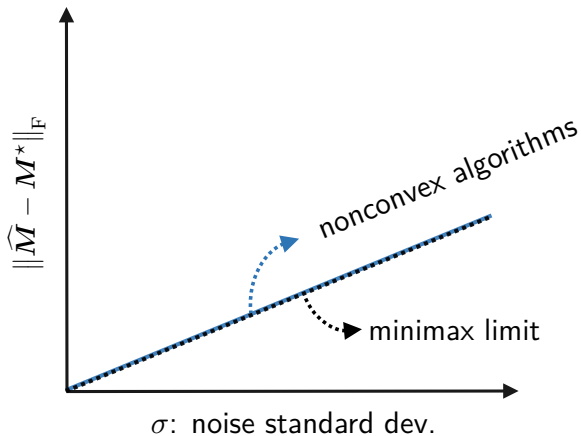


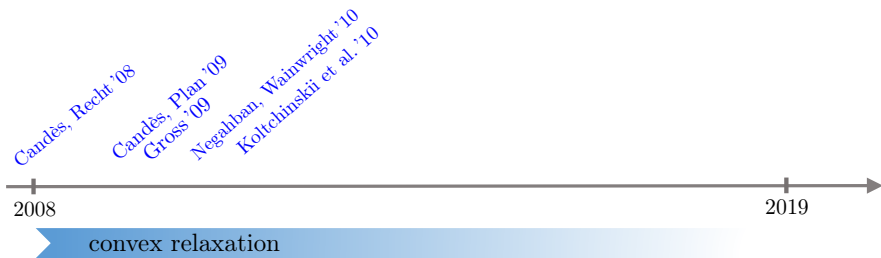
- **suitable initialization:**  $(\mathbf{X}^0, \mathbf{Y}^0)$
- **gradient descent:** for  $t = 0, 1, \dots$

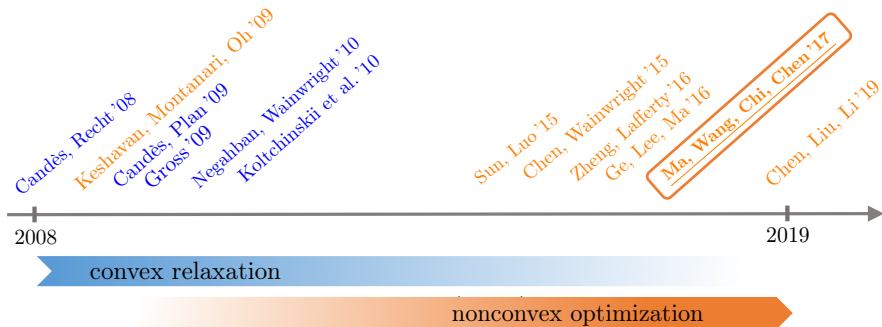
$$\begin{aligned} \mathbf{X}^{t+1} &= \mathbf{X}^t - \eta_t \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta_t \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t) \end{aligned}$$

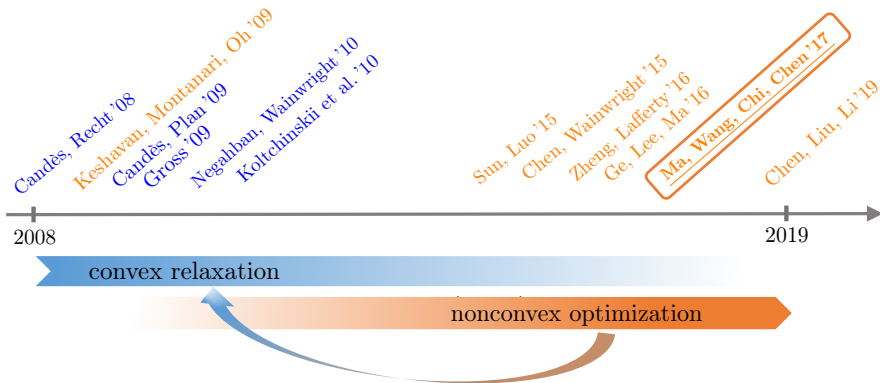


minimax limit	$\sigma \sqrt{n/p}$
nonconvex algorithms	$\sigma \sqrt{n/p}$ (optimal!)









# An interesting experiment

---

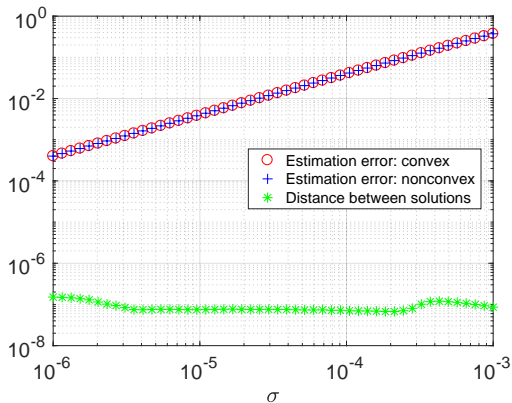
**convex:** 
$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

**nonconvex:** 
$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \sum_{(i,j) \in \Omega} \left[ (\mathbf{X}\mathbf{Y}^\top)_{i,j} - M_{i,j} \right]^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{X}\|_F^2 + \frac{\lambda}{2} \|\mathbf{Y}\|_F^2}_{\text{reg}(\mathbf{X}, \mathbf{Y})}$$

$$- \quad \|\mathbf{Z}\|_* = \min_{\mathbf{Z} = \mathbf{X}\mathbf{Y}^\top} \frac{1}{2} \|\mathbf{X}\|_F^2 + \frac{1}{2} \|\mathbf{Y}\|_F^2$$

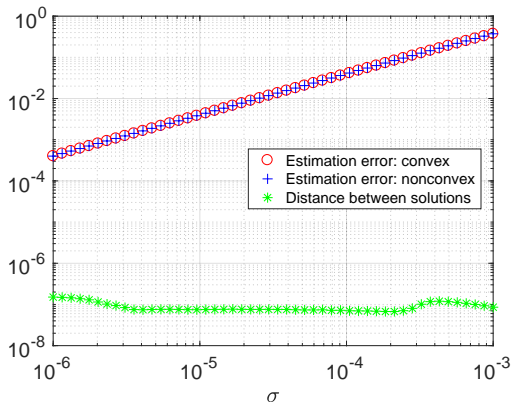
# An interesting experiment

$$n = 1000, r = 5, p = 0.2, \lambda = 5\sigma\sqrt{np}$$



# An interesting experiment

$$n = 1000, r = 5, p = 0.2, \lambda = 5\sigma\sqrt{np}$$



Convex and nonconvex solutions are exceedingly close!

convex



nonconvex



$$\text{stability}\left(\begin{array}{|c|} \hline \text{convex} \\ \hline \end{array}\right) \approx \text{stability}\left(\begin{array}{|c|} \hline \text{nonconvex} \\ \hline \end{array}\right)$$



## Main results: $r = O(1)$

---

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent

## Main results: $r = O(1)$

---

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (\lambda \asymp \sigma \sqrt{np})$$

## Main results: $r = O(1)$

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (\lambda \asymp \sigma \sqrt{np})$$

### Theorem 1 (Chen, Chi, Fan, Ma, Yan '19)

With high prob., any minimizer  $M^{\text{cvx}}$  of convex program obeys

1.  $M^{\text{cvx}}$  is nearly rank- $r$

$$\|M^{\text{cvx}} - \text{proj}_{\text{rank-}r}(M^{\text{cvx}})\|_{\text{F}} \ll \frac{1}{n^5} \cdot \sigma \sqrt{\frac{n}{p}}$$

## Main results: $r = O(1)$

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (\lambda \asymp \sigma \sqrt{np})$$

### Theorem 1 (Chen, Chi, Fan, Ma, Yan '19)

With high prob., any minimizer  $\mathbf{M}^{\text{cvx}}$  of convex program obeys

1.  $\mathbf{M}^{\text{cvx}}$  is nearly rank- $r$

2. 
$$\|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \sigma \sqrt{\frac{n}{p}}$$

## Main results: $r = O(1)$

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent

$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (\lambda \asymp \sigma \sqrt{np})$$

### Theorem 1 (Chen, Chi, Fan, Ma, Yan '19)

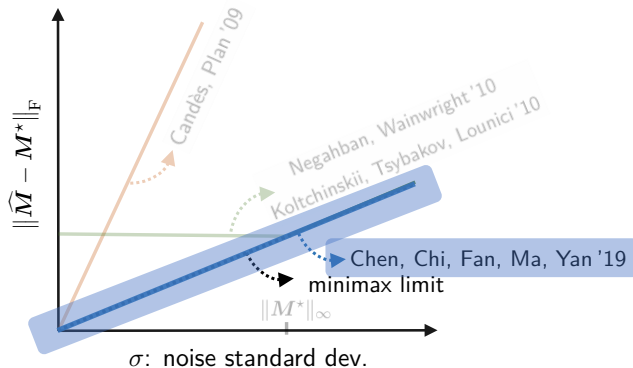
With high prob., any minimizer  $M^{\text{cvx}}$  of convex program obeys

1.  $M^{\text{cvx}}$  is nearly rank- $r$

2.

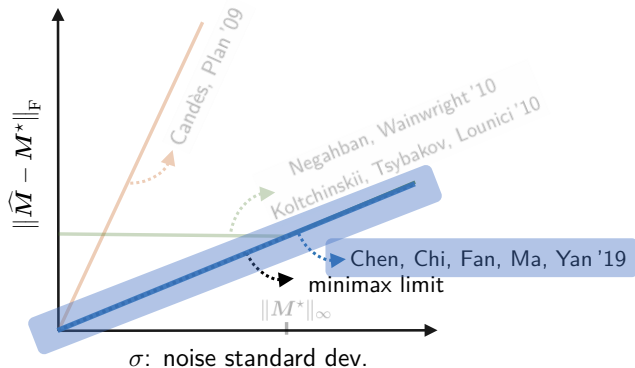
$$\|M^{\text{cvx}} - M^*\|_{\text{F}} \lesssim \sigma \sqrt{\frac{n}{p}}$$
$$\|M^{\text{cvx}} - M^*\|_{\infty} \lesssim \sigma \sqrt{\frac{n \log n}{p}} \cdot \frac{1}{n}$$

$$\|M^{\text{cvx}} - M^*\|_F \lesssim \sigma \sqrt{\frac{n}{p}}$$



- minimax optimal when  $r = O(1)$

$$\|M^{\text{cvx}} - M^*\|_F \lesssim \sigma \sqrt{\frac{n}{p}} \qquad \|M^{\text{cvx}} - M^*\|_\infty \lesssim \sigma \sqrt{\frac{n \log n}{p}} \cdot \frac{1}{n}$$



- minimax optimal when  $r = O(1)$
- estimation errors are spread out across all entries

# Implicit regularization

---

No need to enforce spikiness constraint as in Negahban & Wainwright

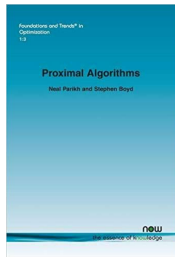
$$\underset{\|\mathbf{Z}\|_\infty \leq \alpha}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

- convex relaxation automatically controls spikiness of solutions



# Statistical guarantees for iterative algorithms

---

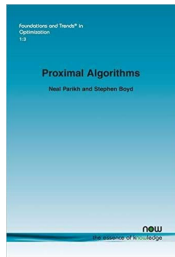


$$\underset{\mathbf{Z}}{\text{minimize}} \quad g(\mathbf{Z}) := \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (1)$$

Many algorithms (e.g. SVT, SOFT-IMPUTE, FPC, FISTA) have been proposed to solve (1), typically without statistical guarantees

# Statistical guarantees for iterative algorithms

---



$$\underset{\mathbf{Z}}{\text{minimize}} \quad g(\mathbf{Z}) := \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_* \quad (1)$$

Many algorithms (e.g. SVT, SOFT-IMPUTE, FPC, FISTA) have been proposed to solve (1), typically without statistical guarantees

We provide statistical guarantees for any  $\mathbf{Z}$  with  $g(\mathbf{Z}) \leq g(\mathbf{Z}_{\text{opt}}) + \varepsilon$  for some sufficiently small  $\varepsilon > 0$

## Main results: general case

---

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{r^2 \log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ : well-conditioned, incoherent

## Main results: general case

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{r^2 \log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ : well-conditioned, incoherent

### Theorem 2 (Chen, Chi, Fan, Ma, Yan '19)

With high prob., any minimizer  $\mathbf{M}^{\text{cvx}}$  of convex program obeys

1.  $\mathbf{M}^{\text{cvx}}$  is nearly rank- $r$

$$2. \quad \|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \frac{\sigma}{\sigma_{\min}(\mathbf{M}^*)} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|_{\text{F}}$$

$$\|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\|_{\infty} \lesssim \sqrt{r} \frac{\sigma}{\sigma_{\min}(\mathbf{M}^*)} \sqrt{\frac{n \log n}{p}} \|\mathbf{M}^*\|_{\infty}$$

$$\|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\| \lesssim \frac{\sigma}{\sigma_{\min}(\mathbf{M}^*)} \sqrt{\frac{n}{p}} \|\mathbf{M}^*\|$$

## Main results: general case

---

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{r^2 \log^3 n}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $\mathbf{M}^* \in \mathbb{R}^{n \times n}$ : well-conditioned, incoherent

sample complexity bound  $O(nr^2 \log^3 n)$  is suboptimal in  $r$

*A little analysis:  
connection between convex and nonconvex solutions*

# Link between convex and nonconvex optimizers

---

$(X, Y)$  is nonconvex optimizer

# Link between convex and nonconvex optimizers

---

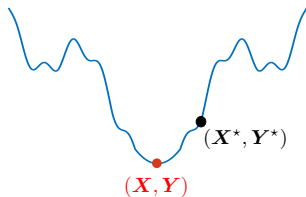
$(\mathbf{X}, \mathbf{Y})$  is nonconvex optimizer  $\overset{?}{\implies} \mathbf{X}\mathbf{Y}^\top$  is convex solution



# Link between convex and nonconvex optimizers

---

- $\lambda$  is properly chosen
- $(\mathbf{X}, \mathbf{Y})$  is close to truth (in  $\ell_{2,\infty}$  sense)

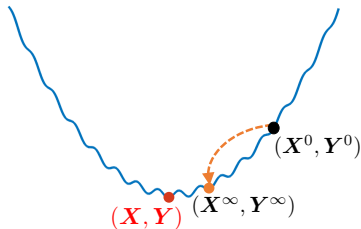


$(\mathbf{X}, \mathbf{Y})$  is nonconvex optimizer  $\overset{\checkmark}{\implies} \mathbf{X}\mathbf{Y}^\top$  is convex solution

i.e.  $\text{dist}(\text{convex solution}, \text{nonconvex solution}) = 0$

# Approximate nonconvex optimizers

---



**Issue:** we do NOT know statistical properties of nonconvex optimizers

- It is unclear whether nonconvex algorithms converge to optimizers

# Approximate nonconvex optimizers

---

**Strategy:** resort to “approximate stationary points” instead  
 $\nabla f(\mathbf{X}, \mathbf{Y}) \approx \mathbf{0}$

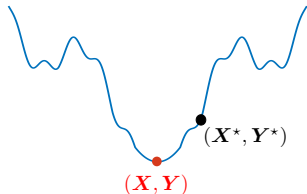
# Approximate nonconvex optimizers

---

**Strategy:** resort to “approximate stationary points” instead

$$\nabla f(\mathbf{X}, \mathbf{Y}) \approx \mathbf{0}$$

- $\lambda$  is properly chosen
- $(\mathbf{X}, \mathbf{Y})$  is close to truth (in  $\ell_{2,\infty}$  sense)



$$\nabla f(\mathbf{X}, \mathbf{Y}) \approx \mathbf{0} \quad \overset{\checkmark}{\implies} \quad \text{dist}(\mathbf{X}\mathbf{Y}^\top, \text{convex solutions}) \approx 0$$

# Construct approximate nonconvex optimizers via GD

---

starting from  $(\mathbf{X}^0, \mathbf{Y}^0) = \text{truth}$  or spectral initialization:

$$\begin{aligned}\mathbf{X}^{t+1} &= \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)\end{aligned}\quad t = 0, 1, \dots, T$$

# Construct approximate nonconvex optimizers via GD

---

starting from  $(\mathbf{X}^0, \mathbf{Y}^0) = \text{truth}$  or spectral initialization:

$$\begin{aligned}\mathbf{X}^{t+1} &= \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)\end{aligned}\quad t = 0, 1, \dots, T$$

- when  $T$  is large: there exists point with very small gradient

$$\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}} \lesssim \frac{1}{\sqrt{\eta T}}$$

# Construct approximate nonconvex optimizers via GD

---

starting from  $(\mathbf{X}^0, \mathbf{Y}^0) = \text{truth}$  or spectral initialization:

$$\begin{aligned}\mathbf{X}^{t+1} &= \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)\end{aligned}\quad t = 0, 1, \dots, T$$

- when  $T$  is large: there exists point with very small gradient

$$\|\nabla f(\mathbf{X}, \mathbf{Y})\|_{\text{F}} \lesssim \frac{1}{\sqrt{\eta T}}$$

- hopefully not far from  $(\mathbf{X}^*, \mathbf{Y}^*)$  (in  $\ell_{2,\infty}$  sense in particular)

# Analyzing nonconvex GD: leave-one-out analysis

---

Leave out a small amount of information from data and run GD



# Analyzing nonconvex GD: leave-one-out analysis

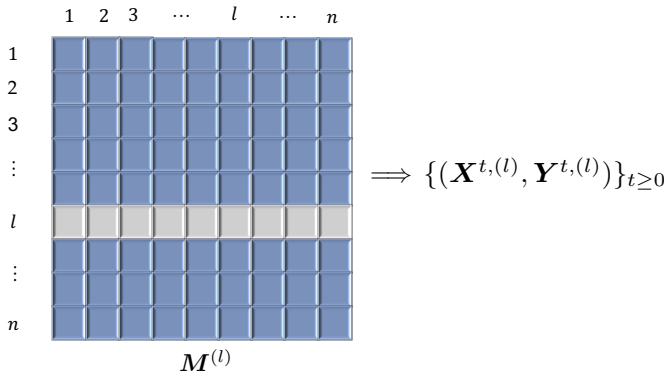
---

Leave out a small amount of information from data and run GD

- Stein '72
- El Karoui, Bean, Bickel, Lim, Yu '13
- El Karoui '15
- Javanmard, Montanari '15
- Zhong, Boumal '17
- Lei, Bickel, El Karoui '17
- Sur, Chen, Candès '17
- Abbe, Fan, Wang, Zhong '17
- Chen, Fan, Ma, Wang '17
- Ma, Wang, Chi, Chen '17
- Chen, Chi, Fan, Ma '18
- Ding, Chen '18
- Dong, Shi '18
- Chen, Liu, Li '19

# Analyzing nonconvex GD: leave-one-out analysis

For each  $1 \leq l \leq n$ , introduce leave-one-out iterates  $\{(\mathbf{X}^{t,(l)}, \mathbf{Y}^{t,(l)})\}$  by replacing  $l^{\text{th}}$  row (or column) with true values



- exploit partial statistical independence
- exploit leave-one-out stability

*Inference and uncertainty quantification*

# Reasoning about uncertainty

---

	2		2	
		6		
3	1		4	
		4		1
	0			

# Reasoning about uncertainty

---

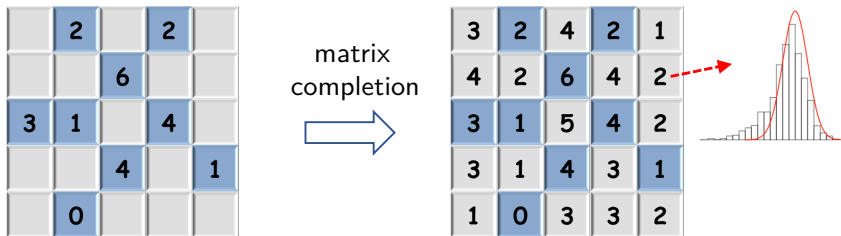
	2		2	
		6		
3	1		4	
		4		1
	0			

matrix  
completion



3	2	4	2	1
4	2	6	4	2
3	1	5	4	2
3	1	4	3	1
1	0	3	3	2

# Reasoning about uncertainty

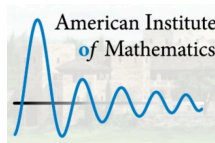


How to assess uncertainty, or “confidence”, of obtained estimates?

## INFERENCE IN HIGH DIMENSIONAL REGRESSION

organized by

Peter Buehlmann, Andrea Montanari, and Jonathan Taylor



- (3) Confidence intervals for matrix completion. In matrix completion, the data analyst is given a large data matrix with a number of missing entries. In many interesting applications (e.g. to collaborative filtering) it is indeed the case that the vast majority of entries is missing. In order to fill the missing entries, the assumption is made that the underlying –unknown– matrix has a low-rank structure.

Substantial work has been devoted to methods for computing point estimates of the missing entries. In applications, it would be very interesting to compute confidence intervals as well. This requires developing distributional characterizations of standard matrix completion methods.

# Challenges

---

$$\mathbf{M}^{\text{cvx}} \triangleq \arg \min_{\mathbf{Z} \in \mathbb{R}^{n \times n}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

- convex estimate  $\mathbf{M}^{\text{cvx}}$  is biased towards small norm



# Challenges

---

$$\mathbf{M}^{\text{cvx}} \triangleq \arg \min_{\mathbf{Z} \in \mathbb{R}^{n \times n}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

- convex estimate  $\mathbf{M}^{\text{cvx}}$  is biased towards small norm
- very challenging to pin down distributions of obtained estimates

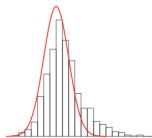
# Challenges

---

$$\mathbf{M}^{\text{cvx}} \triangleq \arg \min_{\mathbf{Z} \in \mathbb{R}^{n \times n}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

- convex estimate  $\mathbf{M}^{\text{cvx}}$  is biased towards small norm
- very challenging to pin down distributions of obtained estimates
- existing orderwise bounds come with unspecified (but huge) pre-constants
  - overly wide confidence intervals

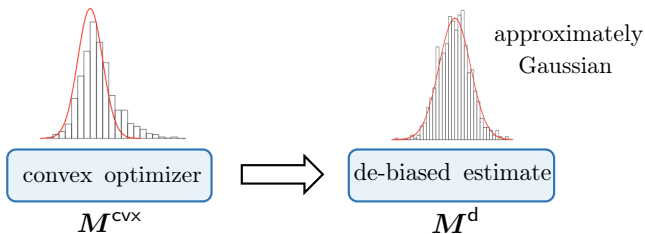
— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



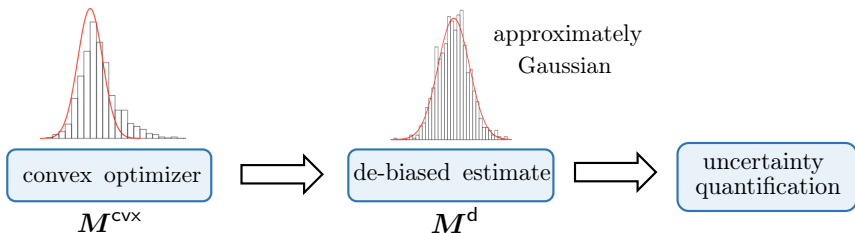
convex optimizer

$M^{\text{cvx}}$

— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



# De-biasing convex estimate

---

$$M^{\text{cvx}} \xrightarrow{\text{de-biasing}} \underbrace{M^{\text{cvx}} + \frac{1}{p} \mathcal{P}_{\Omega}(M^{\star} + E - M^{\text{cvx}})}_{\text{(nearly) unbiased estimate of } M^{\star}}$$

# De-biasing convex estimate

---

$$M^{\text{cvx}} \xrightarrow{\text{de-biasing}} \underbrace{M^{\text{cvx}} + \frac{1}{p} \mathcal{P}_{\Omega}(M^{\star} + E - M^{\text{cvx}})}_{\text{(nearly) unbiased estimate of } M^{\star}}$$

- **issue:** high-rank after de-biasing; statistical accuracy suffers

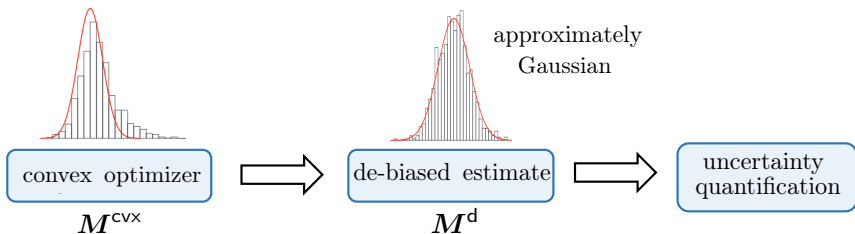
# De-biasing convex estimate

---

$$M^{\text{cvx}} \xrightarrow{\text{de-biasing}} \underbrace{\text{proj}_{\text{rank-}r} \left( M^{\text{cvx}} + \frac{1}{p} \mathcal{P}_{\Omega}(M^{\star} + E - M^{\text{cvx}}) \right)}_{\text{1 iteration of singular value projection (Jain, Meka, Dhillon '10)}} =: M^{\text{d}}$$

- **issue:** high-rank after de-biasing; statistical accuracy suffers
- **solution:** low-rank projection





# Distributional guarantees for low-rank factors

---

$$\begin{aligned} \mathbf{X}^d \mathbf{Y}^{d\top} &\leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{d\top} \mathbf{X}^d = \mathbf{Y}^{d\top} \mathbf{Y}^d} \text{rank-}r \text{ decomp. of } M^d \\ \mathbf{X}^* \mathbf{Y}^{*\top} &\leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{*\top} \mathbf{X}^* = \mathbf{Y}^{*\top} \mathbf{Y}^*} \text{rank-}r \text{ decomp. of } M^* \end{aligned}$$

# Distributional guarantees for low-rank factors

---

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{\log^3 n}{n}$
- **random noise:** i.i.d.  $\mathcal{N}(0, \sigma^2)$  (not too large)
- true matrix  $M^\star \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , well-conditioned, incoherent
- regularization parameter:  $\lambda \asymp \sigma \sqrt{np}$

$$\begin{aligned} \mathbf{X}^d \mathbf{Y}^{d\top} &\leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{d\top} \mathbf{X}^d = \mathbf{Y}^{d\top} \mathbf{Y}^d} \text{rank-}r \text{ decomp. of } M^d \\ \mathbf{X}^\star \mathbf{Y}^{\star\top} &\leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{\star\top} \mathbf{X}^\star = \mathbf{Y}^{\star\top} \mathbf{Y}^\star} \text{rank-}r \text{ decomp. of } M^\star \end{aligned}$$

# Distributional guarantees for low-rank factors

$$\mathbf{X}^d \mathbf{Y}^{d\top} \leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{d\top} \mathbf{X}^d = \mathbf{Y}^{d\top} \mathbf{Y}^d} \text{ rank-}r \text{ approx. of } \mathbf{M}^d$$

$$\mathbf{X}^* \mathbf{Y}^{*\top} \leftarrow \underbrace{\text{balanced}}_{\mathbf{X}^{*\top} \mathbf{X}^* = \mathbf{Y}^{*\top} \mathbf{Y}^*} \text{ rank-}r \text{ decomp. of } \mathbf{M}^*$$

## Theorem 3 (Chen, Fan, Ma, Yan '19)

With high prob., there exists global rotation matrix  $\mathbf{H} \in \mathbb{R}^{r \times r}$  s.t.

$$\mathbf{X}^d \mathbf{H} - \mathbf{X}^* \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,\cdot}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{*\top} \mathbf{Y}^*)^{-1})$$

$$\mathbf{Y}^d \mathbf{H} - \mathbf{Y}^* \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,\cdot}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{*\top} \mathbf{X}^*)^{-1})$$

# Implications

---

$$\begin{aligned} \mathbf{X}^{\text{d}} \mathbf{H} - \mathbf{X}^{\star} &\approx \mathbf{Z}^X, & \mathbf{Z}_{i,\cdot}^X &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{\star\top} \mathbf{Y}^{\star})^{-1}) \\ \mathbf{Y}^{\text{d}} \mathbf{H} - \mathbf{Y}^{\star} &\approx \mathbf{Z}^Y, & \mathbf{Z}_{i,\cdot}^Y &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{\star\top} \mathbf{X}^{\star})^{-1}) \end{aligned}$$

- estimation errors for different rows of  $\mathbf{X}^{\star}$  are nearly independent

$$\mathbf{X}_{i,\cdot}^{\text{d}} \mathbf{H} - \mathbf{X}_{i,\cdot}^{\star} \quad \text{nearly ind. of} \quad \mathbf{X}_{j,\cdot}^{\text{d}} \mathbf{H} - \mathbf{X}_{j,\cdot}^{\star}.$$

# Implications

---

$$\mathbf{X}^{\text{d}}\mathbf{H} - \mathbf{X}^{\star} \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,\cdot}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{Y}^{\star\top}\mathbf{Y}^{\star})^{-1})$$

$$\mathbf{Y}^{\text{d}}\mathbf{H} - \mathbf{Y}^{\star} \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,\cdot}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{X}^{\star\top}\mathbf{X}^{\star})^{-1})$$

- accurate uncertainty quantification for low-rank factors, e.g.

$$\mathbf{X}_{i,\cdot}^{\text{d}}\mathbf{H} - \mathbf{X}_{i,\cdot}^{\star} \sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{Y}^{\star\top}\mathbf{Y}^{\star})^{-1}) + \text{negligible term}$$

$$\mathbf{Y}_{i,\cdot}^{\text{d}}\mathbf{H} - \mathbf{Y}_{i,\cdot}^{\star} \sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{X}^{\star\top}\mathbf{X}^{\star})^{-1}) + \text{negligible term}$$

# Implications

---

$$\mathbf{X}^d \mathbf{H} - \mathbf{X}^\star \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,\cdot}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{\star\top} \mathbf{Y}^\star)^{-1})$$

$$\mathbf{Y}^d \mathbf{H} - \mathbf{Y}^\star \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,\cdot}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{\star\top} \mathbf{X}^\star)^{-1})$$

- accurate uncertainty quantification for low-rank factors, e.g.

$$\mathbf{X}_{i,\cdot}^d - \mathbf{X}_{i,\cdot}^\star \mathbf{H}^\top \sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{d\top} \mathbf{Y}^d)^{-1}) + \text{negligible term}$$

$$\mathbf{Y}_{i,\cdot}^d - \mathbf{Y}_{i,\cdot}^\star \mathbf{H}^\top \sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{d\top} \mathbf{X}^d)^{-1}) + \text{negligible term}$$

# Implications

---

$$\begin{aligned} \mathbf{X}^d \mathbf{H} - \mathbf{X}^\star &\approx \mathbf{Z}^X, & \mathbf{Z}_{i,\cdot}^X &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{\star\top} \mathbf{Y}^\star)^{-1}) \\ \mathbf{Y}^d \mathbf{H} - \mathbf{Y}^\star &\approx \mathbf{Z}^Y, & \mathbf{Z}_{i,\cdot}^Y &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{\star\top} \mathbf{X}^\star)^{-1}) \end{aligned}$$

- accurate uncertainty quantification for low-rank factors, e.g.

$$\begin{aligned} \mathbf{X}_{i,\cdot}^d - \mathbf{X}_{i,\cdot}^\star \mathbf{H}^\top &\sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{d\top} \mathbf{Y}^d)^{-1}) + \text{negligible term} \\ \mathbf{Y}_{i,\cdot}^d - \mathbf{Y}_{i,\cdot}^\star \mathbf{H}^\top &\sim \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{d\top} \mathbf{X}^d)^{-1}) + \text{negligible term} \end{aligned}$$

— *asymptotically optimal*



# Implications

---

$$\begin{aligned}\mathbf{X}^d \mathbf{H} - \mathbf{X}^\star &\approx \mathbf{Z}^X, & \mathbf{Z}_{i,\cdot}^X &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{\star\top} \mathbf{Y}^\star)^{-1}) \\ \mathbf{Y}^d \mathbf{H} - \mathbf{Y}^\star &\approx \mathbf{Z}^Y, & \mathbf{Z}_{i,\cdot}^Y &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{\star\top} \mathbf{X}^\star)^{-1})\end{aligned}$$

- accurate uncertainty quantification for matrix entries: if  $\|\mathbf{X}_{i,\cdot}^\star\|_2 + \|\mathbf{Y}_{j,\cdot}^\star\|_2$  is not too small, then

$$M_{i,j}^d - M_{i,j}^\star \sim \mathcal{N}(0, v_{i,j}^\star) + \text{negligible term}$$

$$\text{where } v_{i,j}^\star \triangleq \frac{\sigma^2}{p} \left\{ \mathbf{X}_{i,\cdot}^\star (\mathbf{X}^{\star\top} \mathbf{X}^\star)^{-1} \mathbf{X}_{i,\cdot}^{\star\top} + \mathbf{Y}_{j,\cdot}^\star (\mathbf{Y}^{\star\top} \mathbf{Y}^\star)^{-1} \mathbf{Y}_{j,\cdot}^{\star\top} \right\}$$

# Implications

---

$$\begin{aligned}\mathbf{X}^d \mathbf{H} - \mathbf{X}^\star &\approx \mathbf{Z}^X, & \mathbf{Z}_{i,\cdot}^X &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{Y}^{\star\top} \mathbf{Y}^\star)^{-1}) \\ \mathbf{Y}^d \mathbf{H} - \mathbf{Y}^\star &\approx \mathbf{Z}^Y, & \mathbf{Z}_{i,\cdot}^Y &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{\star\top} \mathbf{X}^\star)^{-1})\end{aligned}$$

- accurate uncertainty quantification for matrix entries: if  $\|\mathbf{X}_{i,\cdot}^\star\|_2 + \|\mathbf{Y}_{j,\cdot}^\star\|_2$  is not too small, then

$$M_{i,j}^d - M_{i,j}^\star \sim \mathcal{N}(0, \hat{v}_{i,j}) + \text{negligible term}$$

$$\text{where } \hat{v}_{i,j} \triangleq \frac{\sigma^2}{p} \left\{ \mathbf{X}_{i,\cdot}^d (\mathbf{X}^{d\top} \mathbf{X}^d)^{-1} \mathbf{X}_{i,\cdot}^{d\top} + \mathbf{Y}_{j,\cdot}^d (\mathbf{Y}^{d\top} \mathbf{Y}^d)^{-1} \mathbf{Y}_{j,\cdot}^{d\top} \right\}$$

# Implications

---

$$\begin{aligned}\mathbf{X}^{\text{d}}\mathbf{H} - \mathbf{X}^{\star} &\approx \mathbf{Z}^X, & \mathbf{Z}_{i,\cdot}^X &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{Y}^{\star\top}\mathbf{Y}^{\star})^{-1}) \\ \mathbf{Y}^{\text{d}}\mathbf{H} - \mathbf{Y}^{\star} &\approx \mathbf{Z}^Y, & \mathbf{Z}_{i,\cdot}^Y &\stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p}(\mathbf{X}^{\star\top}\mathbf{X}^{\star})^{-1})\end{aligned}$$

- accurate uncertainty quantification for matrix entries: if  $\|\mathbf{X}_{i,\cdot}^{\star}\|_2 + \|\mathbf{Y}_{j,\cdot}^{\star}\|_2$  is not too small, then

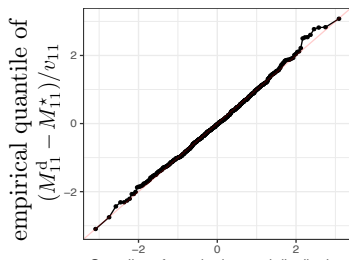
$$M_{i,j}^{\text{d}} - M_{i,j}^{\star} \sim \mathcal{N}(0, \hat{v}_{i,j}) + \text{negligible term}$$

$$\text{where } \hat{v}_{i,j} \triangleq \frac{\sigma^2}{p} \left\{ \mathbf{X}_{i,\cdot}^{\text{d}} (\mathbf{X}^{\text{d}\top} \mathbf{X}^{\text{d}})^{-1} \mathbf{X}_{i,\cdot}^{\text{d}\top} + \mathbf{Y}_{j,\cdot}^{\text{d}} (\mathbf{Y}^{\text{d}\top} \mathbf{Y}^{\text{d}})^{-1} \mathbf{Y}_{j,\cdot}^{\text{d}\top} \right\}$$

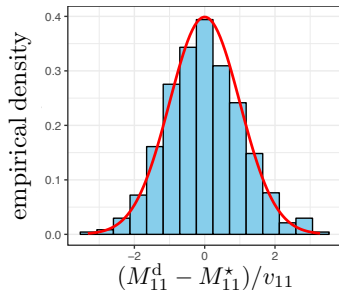
*— asymptotically optimal*

# Numerical experiments

---



quantile of a standard Gaussian



$$n = 1000, p = 0.2, r = 5, \|\mathbf{M}^*\| = 1, \kappa = 1, \sigma = 10^{-3}$$

convex



nonconvex

convex



nonconvex



$$\text{inference} \left( \boxed{\text{convex}} \right) \approx \text{inference} \left( \boxed{\text{nonconvex}} \right)$$

Same inference procedures work for both cvx & noncvx estimates!

## A bit of intuition

---

Consider rank-1 PSD case  $\mathbf{M}^\star = \mathbf{x}^\star \mathbf{x}^{\star\top}$ ,  $p = 1$  (no missing data)

$$\text{minimize}_{\mathbf{x}} \quad \frac{1}{2} \|\mathbf{x} \mathbf{x}^\top - \mathbf{x}^\star \mathbf{x}^{\star\top} - \mathbf{E}\|_{\text{F}}^2 + \lambda \|\mathbf{x}\|_2^2$$

## A bit of intuition

---

Consider rank-1 PSD case  $M^* = x^* x^{*\top}$ ,  $p = 1$  (no missing data)

$$\text{minimize}_x \quad \frac{1}{2} \|xx^\top - x^* x^{*\top} - E\|_F^2 + \lambda \|x\|_2^2$$

- first-order optimality condition

$$(xx^\top - x^* x^{*\top} - E)x + \lambda x = 0$$



# A bit of intuition

---

$$(xx^\top - x^*x^{*\top} - E)x \underbrace{+\lambda x}_{\text{causes bias}} = 0$$

## A bit of intuition

---

$$(xx^\top - x^*x^{*\top} - E)x \underbrace{+ \lambda x}_{\text{causes bias}} = 0$$



$$(x^d x^{d\top} - x^*x^{*\top} - E)x^d = 0, \quad x^d = \sqrt{\frac{\lambda + \|x\|_2^2}{\|x\|_2^2}} x$$

# A bit of intuition

---

$$(xx^\top - x^*x^{*\top} - E)x \underbrace{+ \lambda x}_{\text{causes bias}} = 0$$



$$(x^d x^{d\top} - x^* x^{*\top} - E)x^d = 0, \quad x^d = \sqrt{\frac{\lambda + \|x\|_2^2}{\|x\|_2^2}} x$$



$$x^d - x^* = \underbrace{\frac{1}{\|x^d\|_2^2} E x^d}_{\text{nearly Gaussian}} + \underbrace{\frac{(x^* - x^d)^\top x^d}{\|x^d\|_2^2} x^*}_{\text{hopefully small}}$$

## **Back to estimation: de-biased estimator is optimal**

---

Distributional theory in turn allows us to track estimation accuracy

# Back to estimation: de-biased estimator is optimal

Distributional theory in turn allows us to track estimation accuracy

## Theorem 4 (Chen, Fan, Ma, Yan '19)

$$\frac{\|M^d - M^*\|_F^2}{n^2} = \underbrace{\frac{(2 + o(1)) \textcolor{red}{nr}\sigma^2}{\textcolor{red}{n^2p}}}_{\text{Oracle lower bound}} \quad \text{with high prob.}$$

# Back to estimation: de-biased estimator is optimal

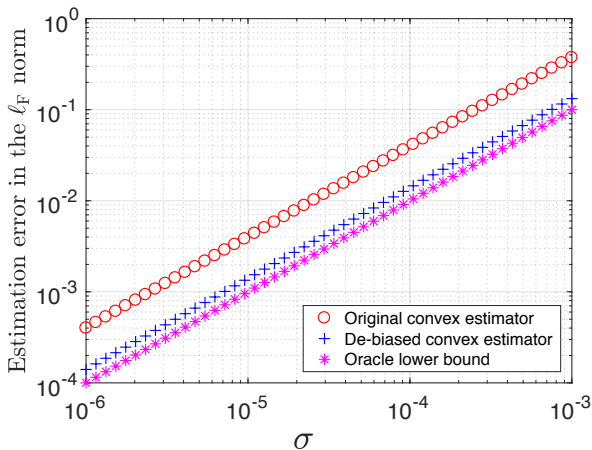
Distributional theory in turn allows us to track estimation accuracy

## Theorem 4 (Chen, Fan, Ma, Yan '19)

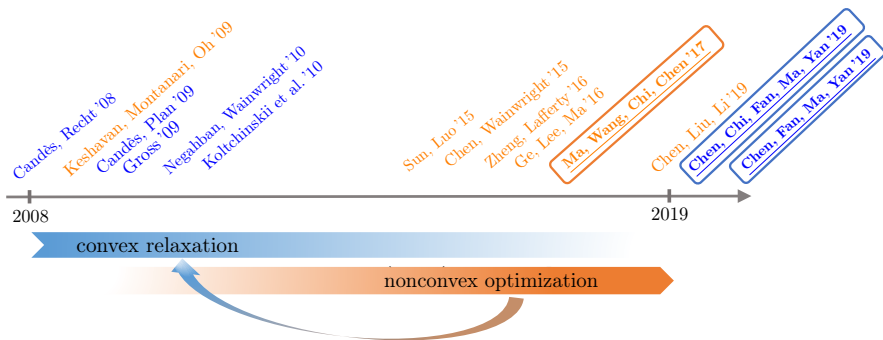
$$\frac{\|M^d - M^*\|_F^2}{n^2} = \underbrace{\frac{(2 + o(1)) \textcolor{red}{n} \textcolor{red}{r} \sigma^2}{\textcolor{red}{n}^2 \textcolor{red}{p}}}_{\text{Oracle lower bound}} \quad \text{with high prob.}$$

- precise characterization of estimation accuracy
- achieves full statistical efficiency (including **pre-constant**)

## Numerical evidence ( $r = 5$ , $p = 0.2$ , $n = 1000$ )



Euclidean estimation error vs. noise standard deviation  $\sigma$





# Reference

---

- “*Matrix completion with noise*,” E. J. Candès, Y. Plan, *Proceedings of the IEEE*, vol. 98, no. 6, 2010.
- “*Restricted strong convexity and weighted matrix completion: Optimal bounds with noise*,” S. Negahban, M. Wainwright, *Journal of Machine Learning Research*, vol. 13, no. 1, 2012.
- “*Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion*,” V. Koltchinskii, K. Lounici, A. B. Tsybakov, *Annals of Statistics*, vol. 39, no. 5, 2011.
- “*Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization*,” Y. Chen, Y. Chi, J. Fan, C. Ma, Y. Yan, *SIAM Journal on Optimization*, vol. 30, no. 4, 2020.

# Reference

---

- “*Inference and uncertainty quantification for noisy matrix completion*,” Y. Chen, J. Fan, C. Ma, Y. Yan, *Proceedings of the National Academy of Sciences*, vol. 116, no. 46, 2019.
- “*Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution*,” C. Ma, K. Wang, Y. Chi, Y. Chen, *Foundations of Computational Mathematics*, vol. 20, no. 3, 2020.
- “*Bridging convex and nonconvex optimization in robust PCA: Noise, outliers and missing data*,” Y. Chen, J. Fan, C. Ma, Y. Yan, *Annals of Statistics*, vol. 49, no. 5, 2021.