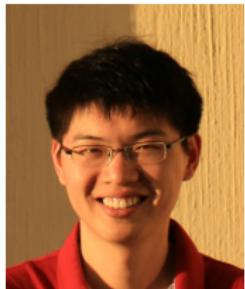


# Inference and Uncertainty Quantification for Noisy Matrix Completion



Yuxin Chen

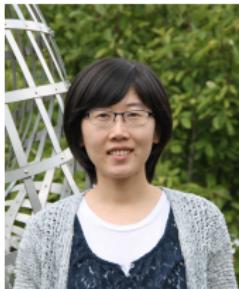
EE, Princeton University



Cong Ma  
Princeton ORFE



Yuling Yan  
Princeton ORFE



Yuejie Chi  
CMU ECE



Jianqing Fan  
Princeton ORFE

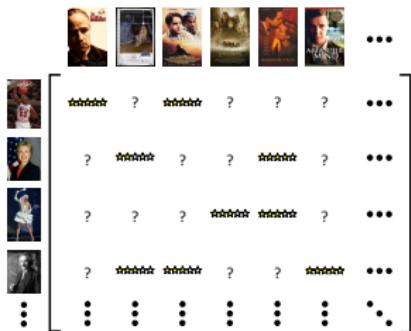
# Low-rank matrix completion

$$\begin{bmatrix} \checkmark & ? & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \\ \checkmark & ? & ? & \checkmark & ? & ? \\ ? & ? & \checkmark & ? & ? & \checkmark \\ \checkmark & ? & ? & ? & ? & ? \\ ? & \checkmark & ? & ? & \checkmark & ? \\ ? & ? & \checkmark & \checkmark & ? & ? \end{bmatrix}$$

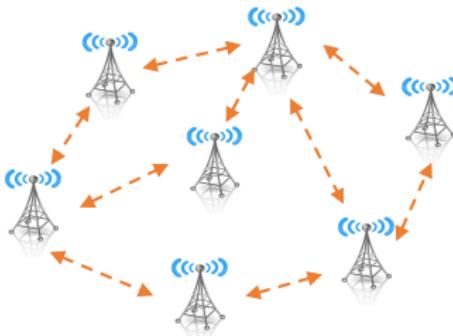


figure credit: E. J. Candès

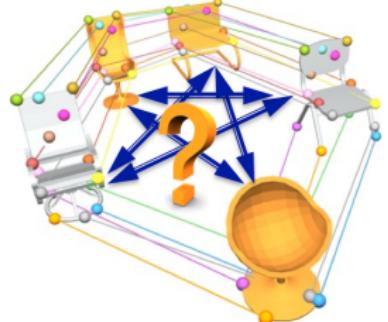
Given partial samples of a low-rank matrix  $M^*$ , fill in missing entries



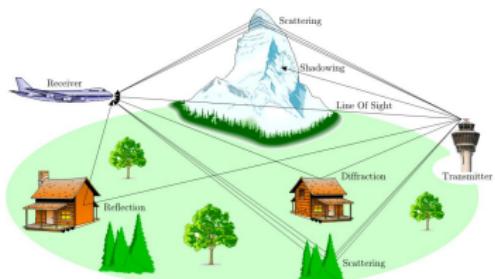
recommendation systems



localization



shape matching



channel estimation

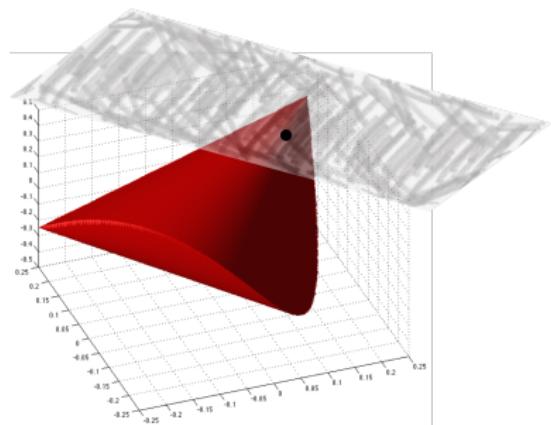
# Estimation via convex relaxation

$$\underset{Z}{\text{minimize}} \quad \underbrace{f(Z; \text{data}) + \lambda \|Z\|_*}_{\text{empirical loss}}$$



low-rank matrix

*Composition C* by Piet Mondrian



semidefinite relaxation

# One step further: reasoning about uncertainty?

---

	2		2	
		6		
3	1		4	
	4		4	1
	0			

# One step further: reasoning about uncertainty?

---

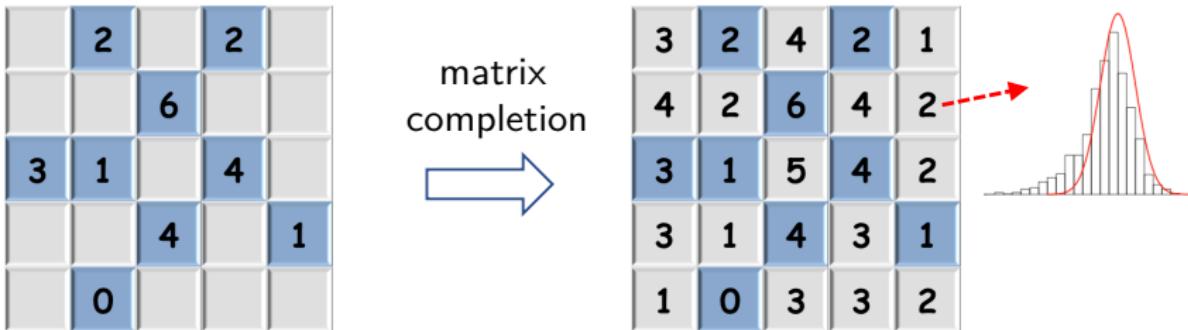
	2	2		
	6			
3	1	4		
	4		1	
	0			

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

# One step further: reasoning about uncertainty?



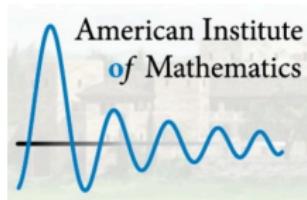
How to assess uncertainty, or “confidence”, of obtained estimates due to imperfect data acquisition?

- noise
- incomplete measurements
- ...

## 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

---

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{f(\boldsymbol{Z}; \text{data})}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

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

# Challenges

---

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{f(\boldsymbol{Z}; \text{data})}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

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

# Challenges

---

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{f(\boldsymbol{Z}; \text{data})}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

- convex estimate  $\boldsymbol{M}^{\text{cvx}}$  is biased towards small norm
- very challenging to pin down distributions of obtained estimates
- existing estimation error bounds are highly sub-optimal  
→ overly wide confidence intervals

# This talk

---

1. **Stability:** order-wise optimal  $\ell_2$  estimation error bounds

# This talk

---

1. **Stability:** order-wise optimal  $\ell_2$  estimation error bounds
2. **Uncertainty quantification:** optimal entrywise confidence intervals (including both rates & pre-constants)

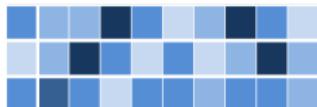
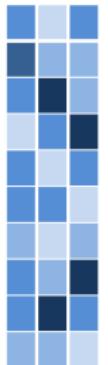
*Stability analysis*

# 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  $\Omega$

# Noisy low-rank matrix completion

---

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

goal: estimate  $M^*$

**convex relaxation:**

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

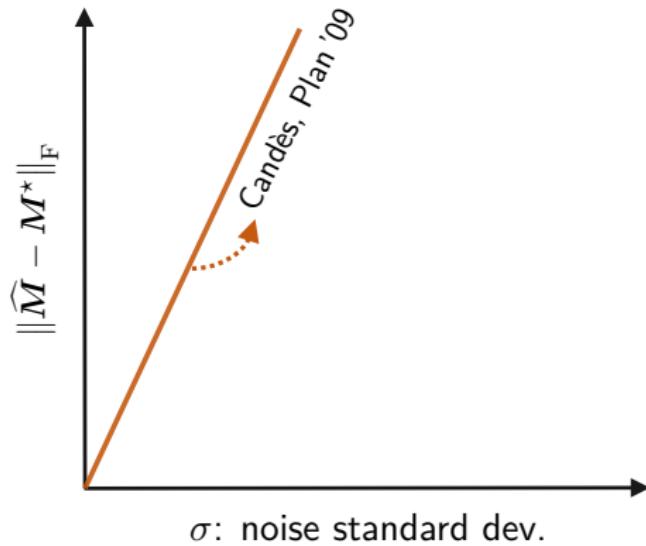
# Prior statistical guarantees for convex relaxation

---

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

Candès, Plan '09

$\sigma n^{1.5}$

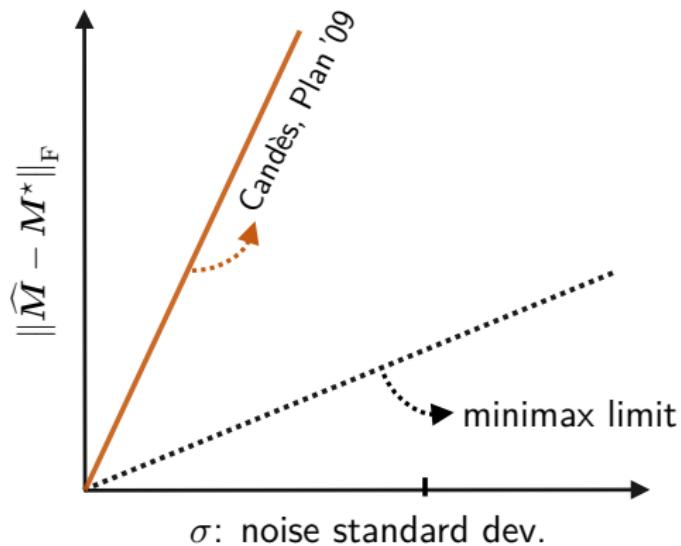


minimax limit

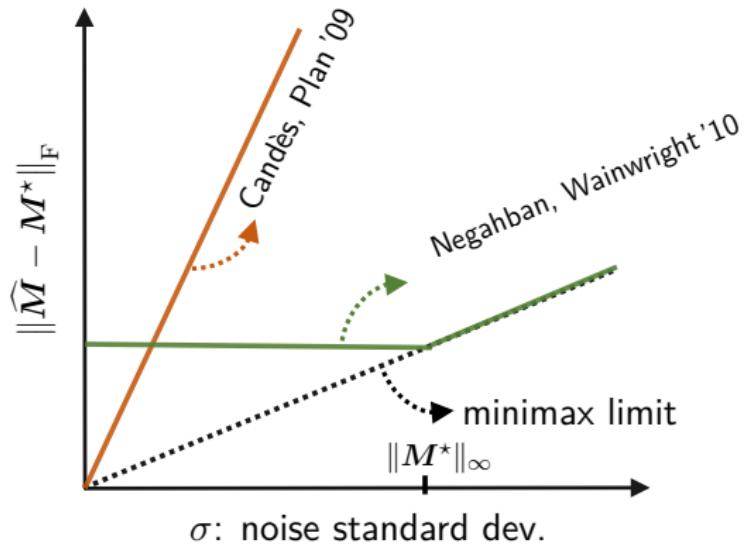
$$\sigma\sqrt{n/p}$$

Candès, Plan '09

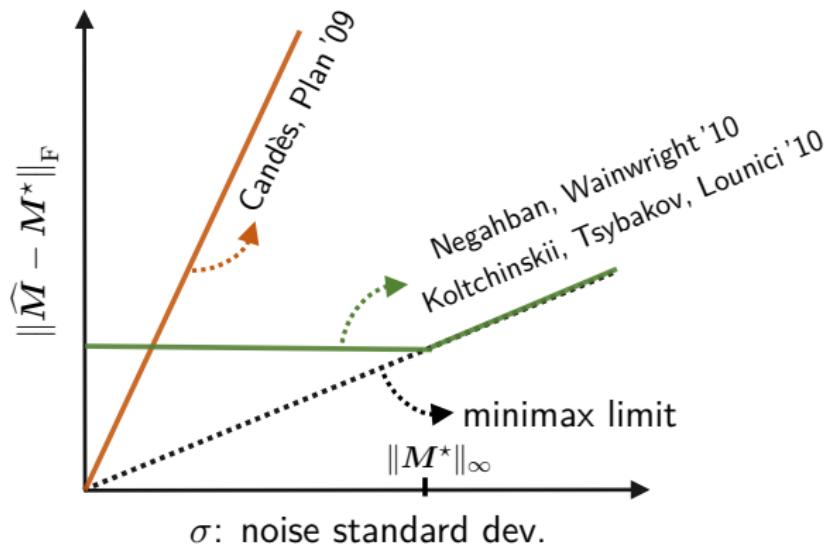
$$\sigma n^{1.5}$$



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

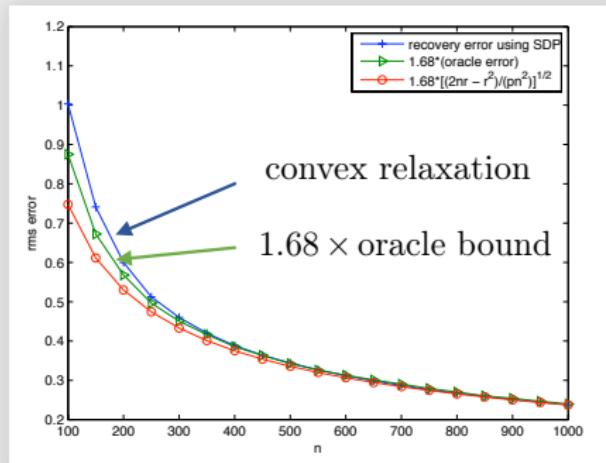


minimax limit	$\sigma\sqrt{n/p}$
Candès, Plan '09	$\sigma n^{1.5}$
Negahban, Wainwright '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$
Koltchinskii, Tsybakov, Lounici '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$



# Matrix Completion with Noise

Emmanuel J. Candès and Yannick Plan



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

# Matrix Completion with Noise

Emmanuel J. Candès and Yannick 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)

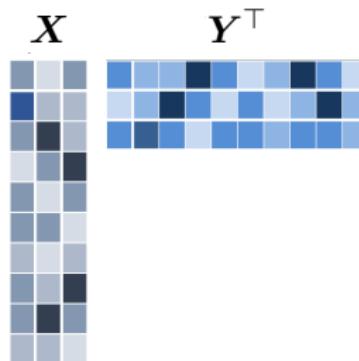


nonconvex optimization

# A detour: nonconvex optimization

---

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

$$\mathbf{X} \quad \mathbf{Y}^\top$$


# A detour: nonconvex optimization

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

$$\begin{array}{c} \mathbf{X} \qquad \mathbf{Y}^\top \\ \begin{matrix} \text{---} & \text{---} \\ \text{---} & \text{---} \end{matrix} \end{array}$$

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

# A detour: nonconvex optimization

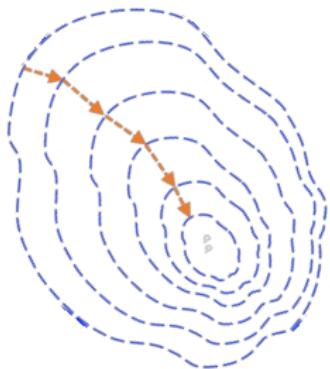
---

- 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  $t = 0, 1, \dots$

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

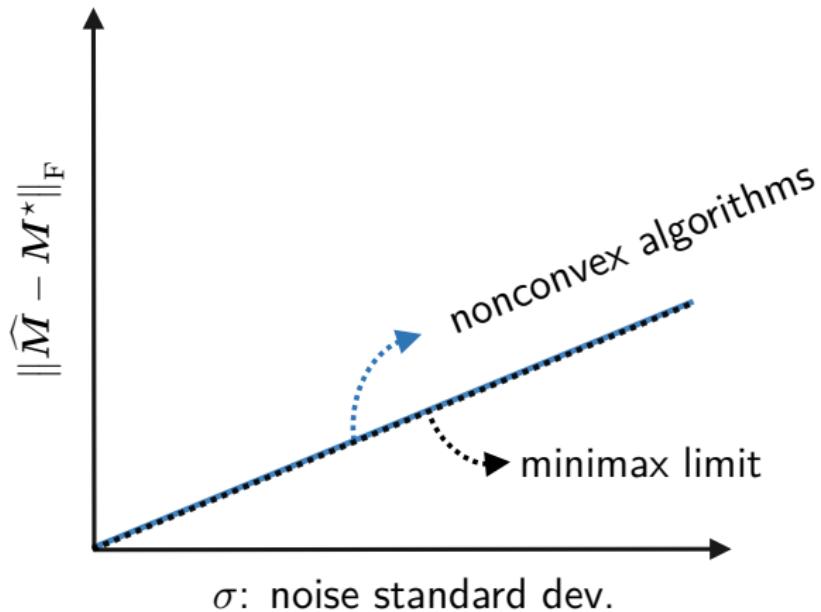
$$\mathbf{Y}^{t+1} = \mathbf{Y}^t - \eta_t \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)$$

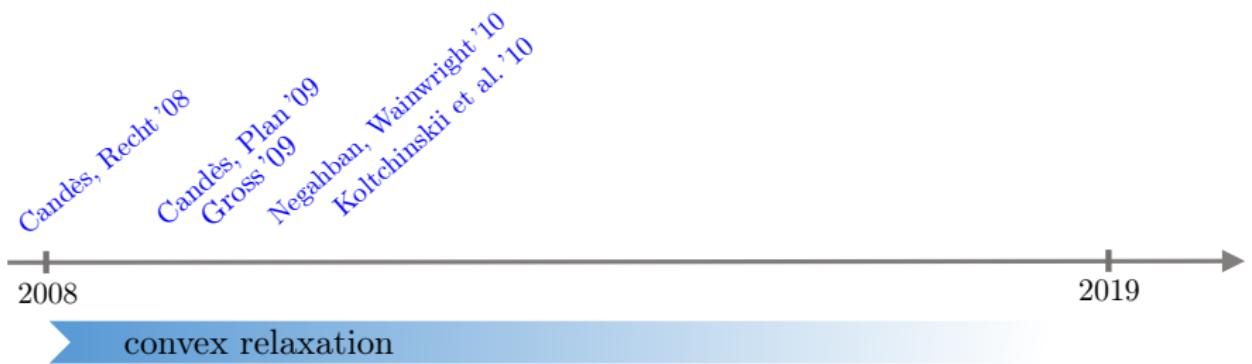
## A detour: nonconvex optimization

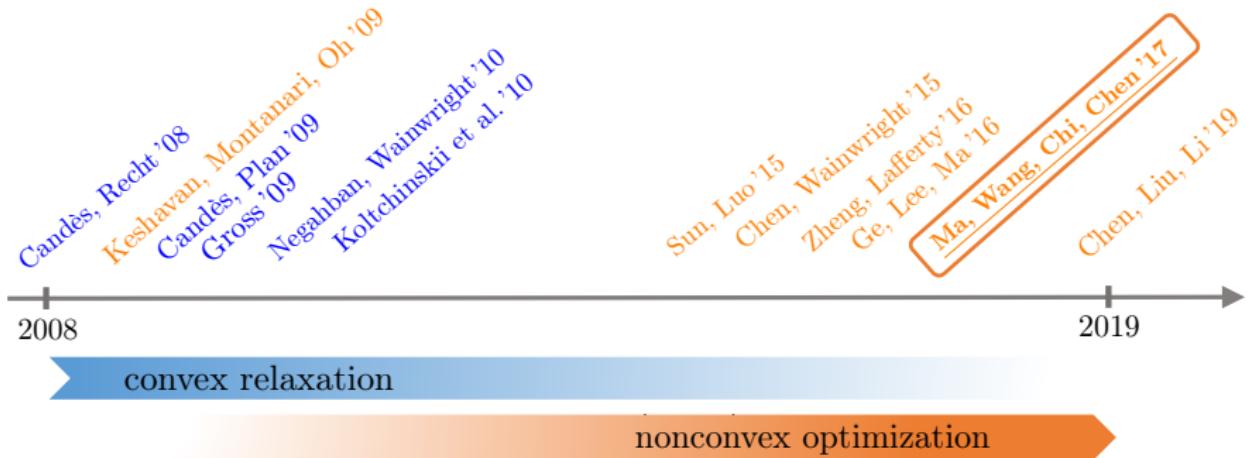
---

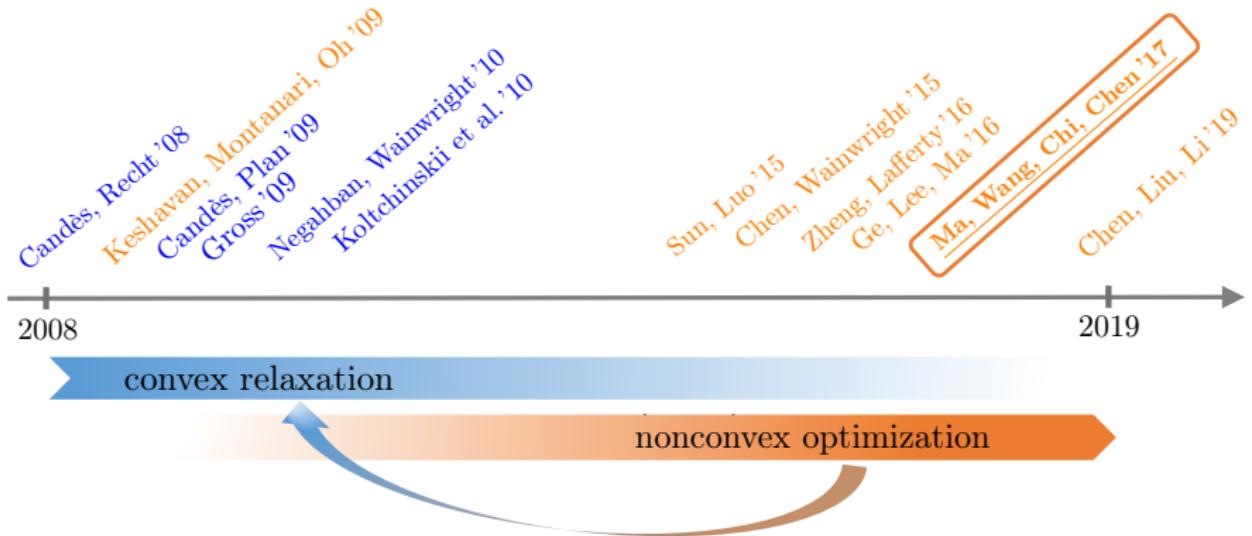
- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  (not too large)
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, ...

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









# A motivating experiment

---

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

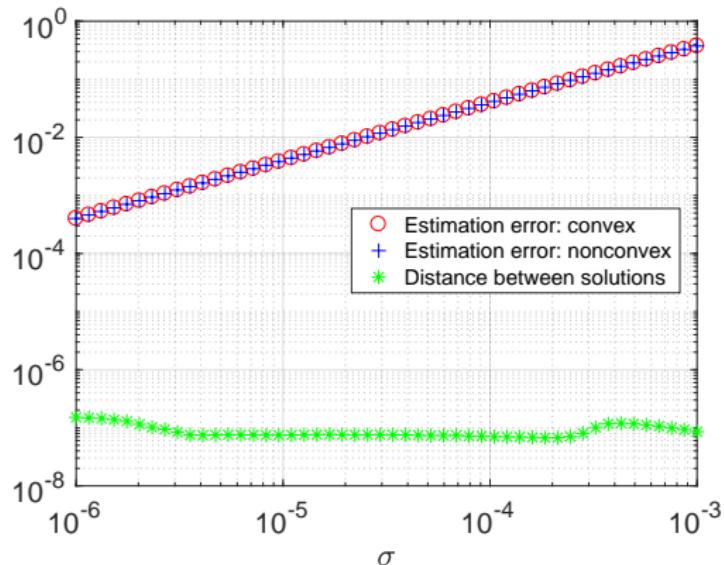
**nonconvex:**  $\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \sum_{(i,j) \in \Omega} \left[ (\mathbf{XY}^\top)_{i,j} - M_{i,j} \right]^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{X}\|_\text{F}^2 + \frac{\lambda}{2} \|\mathbf{Y}\|_\text{F}^2}_{\text{reg}(\mathbf{X},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

—  $\|\mathbf{Z}\|_* = \min_{\mathbf{Z} = \mathbf{XY}^\top} \frac{1}{2} \|\mathbf{X}\|_\text{F}^2 + \frac{1}{2} \|\mathbf{Y}\|_\text{F}^2$

# A motivating experiment

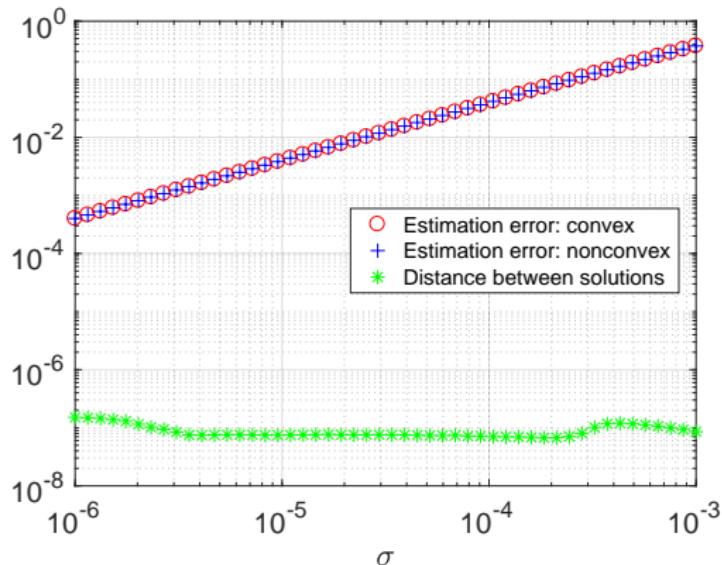
---

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



# A motivating experiment

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



Convex and nonconvex solutions are exceedingly close!

convex



nonconvex

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

## 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  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

## 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  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

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

## 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  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

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

### 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  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

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

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

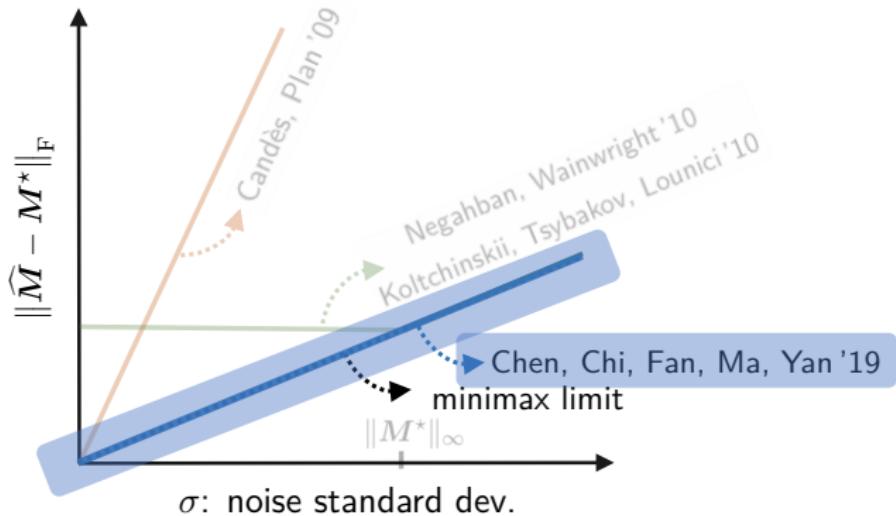
With high prob., any minimizer  $M^{\text{cvx}}$  of convex program obeys

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

2.

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

$$\|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \sigma \sqrt{\frac{n}{p}} : \quad \text{minimax optimal when } r = O(1)$$



# Implicit regularization

---

No need to enforce spikiness constraint as in Negahban & Wainwright

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

- convex programming automatically controls spikiness of solutions

## 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  $M^* \in \mathbb{R}^{n \times n}$ : incoherent, well-conditioned

## 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  $M^* \in \mathbb{R}^{n \times n}$ : incoherent, well-conditioned

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

## Main results: general case

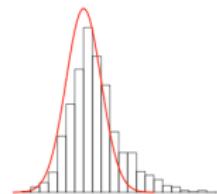
---

- **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  $M^* \in \mathbb{R}^{n \times n}$ : incoherent, well-conditioned

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

*Inference and uncertainty quantification*

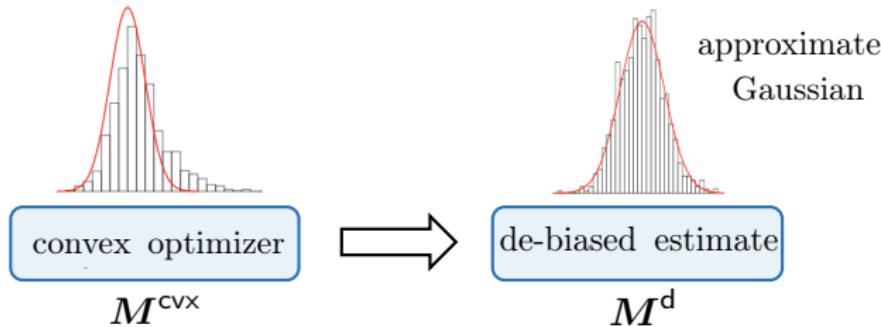
— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



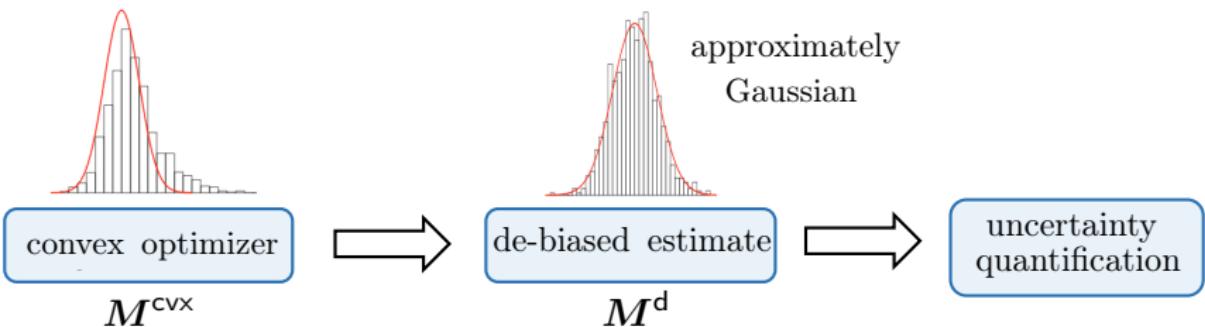
convex optimizer

$M^{\text{CVX}}$

— 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

---

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

## De-biasing convex estimate

---

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

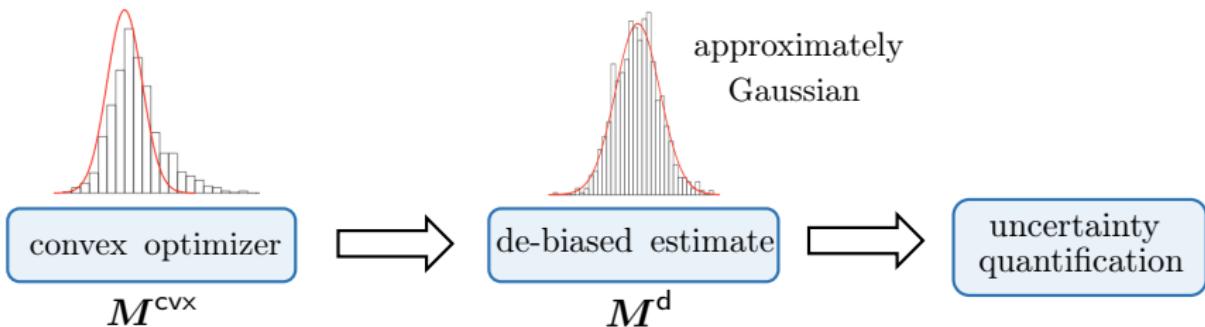
- **issue:** high-rank after de-biasing; statistical accuracy suffers

## De-biasing convex estimate

---

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

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



# 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  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned
- regularization parameter:  $\lambda \asymp \sigma \sqrt{np}$

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

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

# Distributional guarantees for low-rank factors

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

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

## 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

---

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

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

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

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

# Implications

---

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

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

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

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

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

# Implications

---

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

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

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

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

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

# Implications

---

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

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

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

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

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

— *asymptotically optimal*

# Implications

---

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

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

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

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

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

# Implications

---

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

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

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

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

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

---

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

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

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

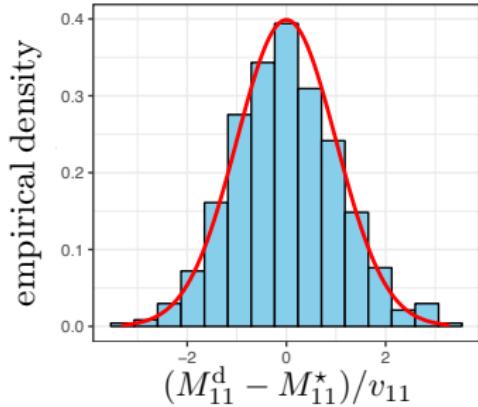
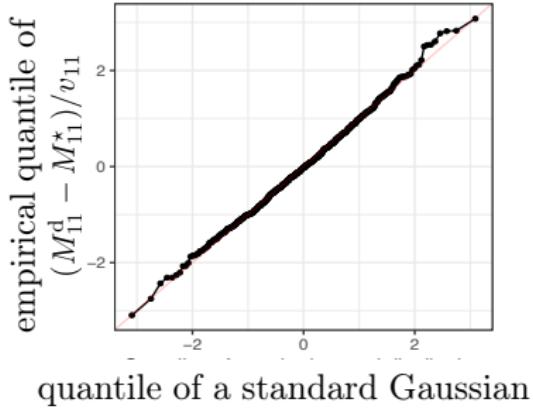
$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}(0, \hat{v}_{i,j}) + \text{negligible term}$$

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

— *asymptotically optimal*

# Numerical experiments

---



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

convex



nonconvex

convex



nonconvex



inference  $(\text{convex})$



inference  $(\text{nonconvex})$

Same inference procedures work for both cvx & noncvx estimates!

## A bit of intuition

---

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

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

## A bit of intuition

---

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

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

- first-order optimality condition

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

# A bit of intuition

---

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

# A bit of intuition

---

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

$\Updownarrow$

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

# A bit of intuition

---

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

$\Updownarrow$

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

$\Updownarrow$

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

# An alternative de-biased estimator

---

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

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z} \in \mathbb{R}^{n \times n}} \|\mathcal{P}_{\Omega}(\boldsymbol{Z} - \boldsymbol{M})\|_{\text{F}}^2 + \lambda \|\boldsymbol{Z}\|_{*}$$

# An alternative de-biased estimator

---

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

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z} \in \mathbb{R}^{n \times n}} \|\mathcal{P}_{\Omega}(\boldsymbol{Z} - \boldsymbol{M})\|_{\text{F}}^2 + \lambda \|\boldsymbol{Z}\|_{*}$$



$$\widehat{\boldsymbol{M}}^{\text{d}} = \boldsymbol{M}^{\text{cvx}} + \text{linear-map} \left( \mathcal{P}_{\Omega}(\boldsymbol{M} - \boldsymbol{M}^{\text{cvx}}) \right)$$

# An alternative de-biased estimator

---

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

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z} \in \mathbb{R}^{n \times n}} \|\mathcal{P}_{\Omega}(\boldsymbol{Z} - \boldsymbol{M})\|_{\text{F}}^2 + \lambda \|\boldsymbol{Z}\|_{*}$$



$$\widehat{\boldsymbol{M}}^{\text{d}} = \boldsymbol{M}^{\text{cvx}} + \frac{1}{p} \mathcal{P}_T \left( \mathcal{P}_{\Omega} (\boldsymbol{M} - \boldsymbol{M}^{\text{cvx}}) \right)$$

—  $\mathcal{P}_T$ : projection onto  $T$  (tangent space at  $\text{proj}_{\text{rank-}r}(\boldsymbol{M}^{\text{cvx}})$ )

$\widehat{\boldsymbol{M}}^{\text{d}} \approx \boldsymbol{M}^{\text{d}}$  !

## 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)

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

# Back to estimation: de-biased estimator is optimal

Distributional theory in turn allows us to track estimation accuracy

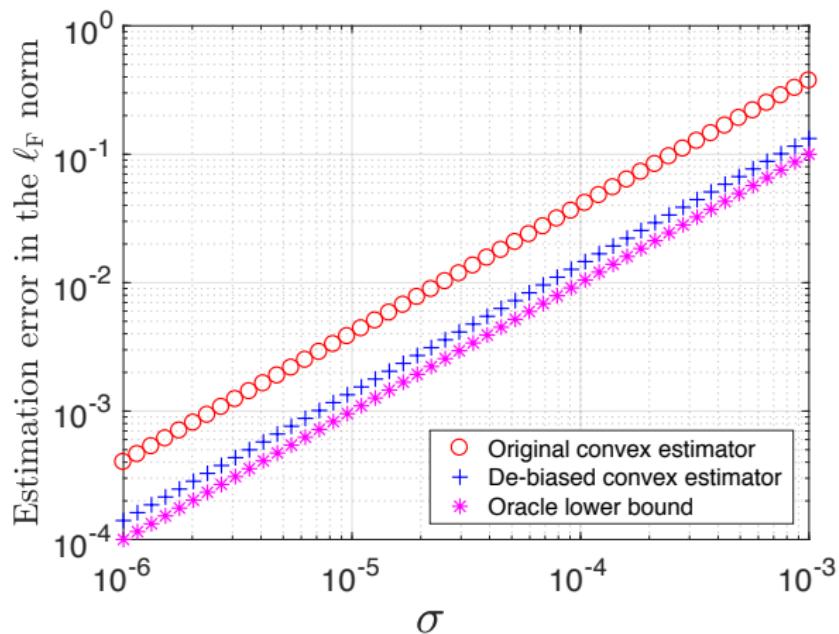
## Theorem 4 (Chen, Fan, Ma, Yan '19)

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

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

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

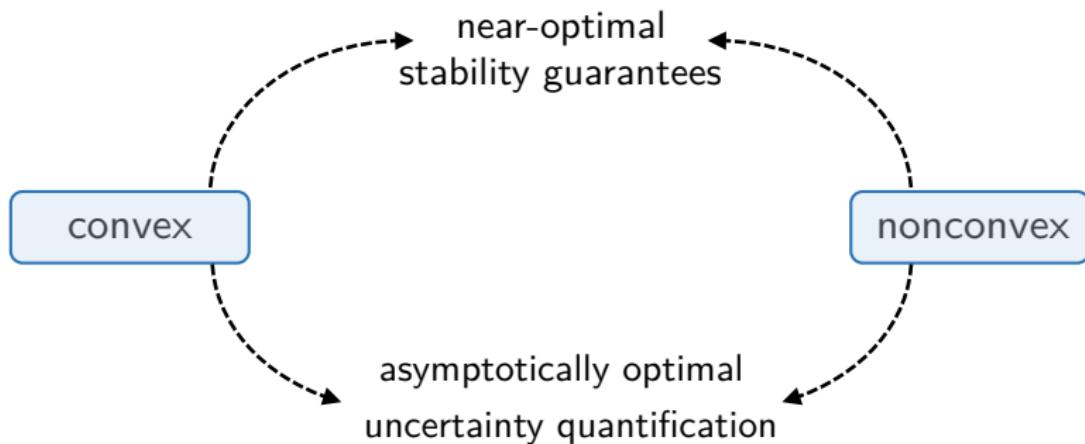
---



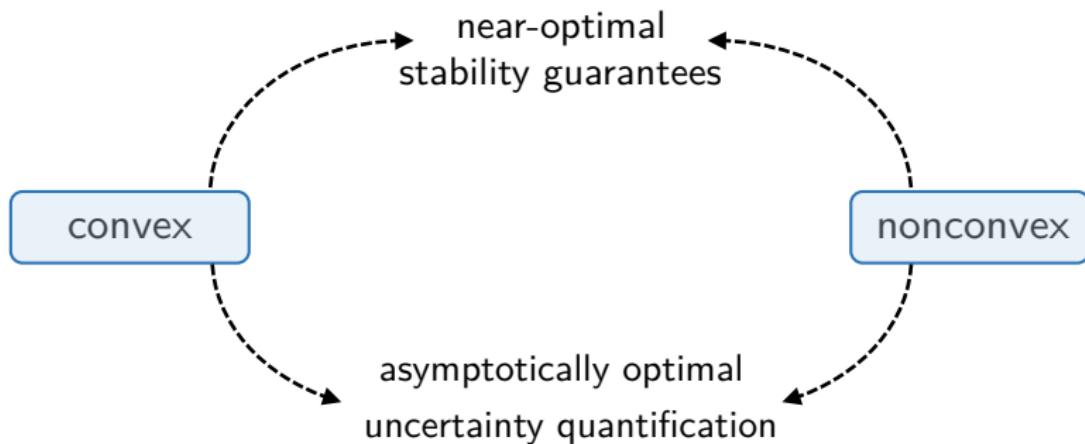
Euclidean estimation error vs. noise standard deviation  $\sigma$

# Concluding remarks

---

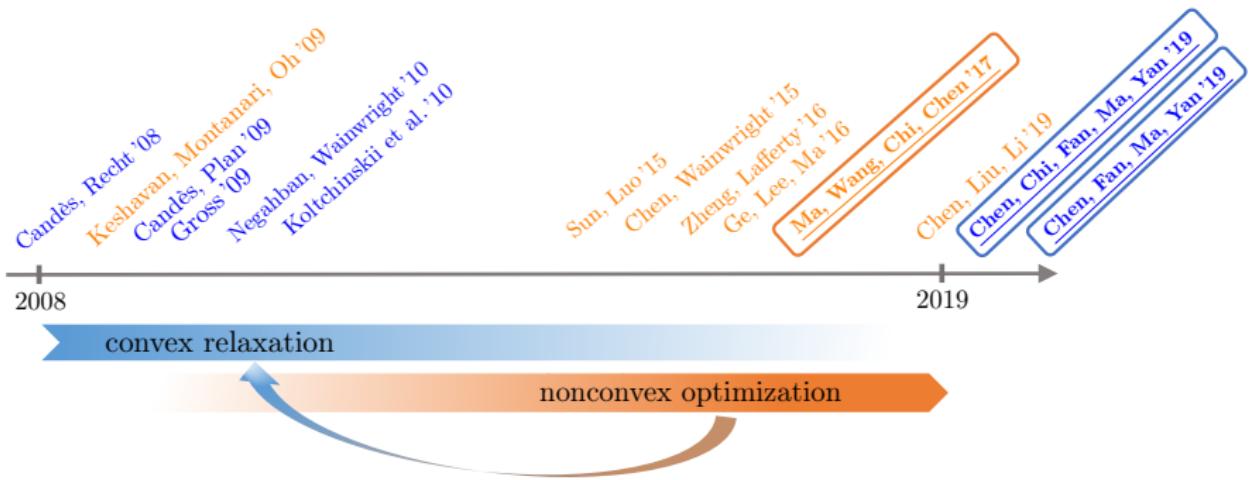


# Concluding remarks



- improve dependency on rank & cond. number
- what if  $M^*$  is only approximately low-rank?
- more general sampling patterns / graphs





"Noisy matrix completion: understanding statistical guarantees for convex relaxation via nonconvex optimization," Y. Chen, Y. Chi, J. Fan, C. Ma, Y. Yan, 2019

"Inference and uncertainty quantification for noisy matrix completion," Y. Chen, J. Fan, C. Ma, Y. Yan, 2019

*Backup slides — 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

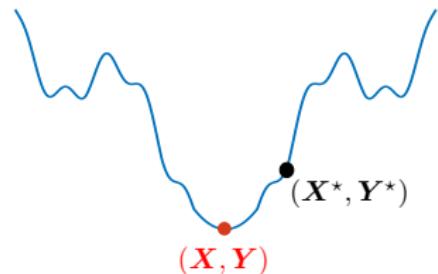
---

$(X, Y)$  is nonconvex optimizer  $\xrightarrow{?}$   $XY^\top$  is convex solution

# Link between convex and nonconvex optimizers

---

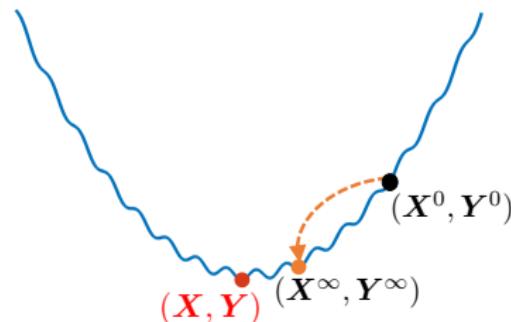
- $(\mathbf{X}, \mathbf{Y})$  is close to truth (in  $\ell_{2,\infty}$  sense)
- a little condition on noise size



$(\mathbf{X}, \mathbf{Y})$  is nonconvex optimizer  $\xrightarrow{\text{✓}}$   $\mathbf{XY}^\top$  is convex solution

# Approximate nonconvex optimizers

---



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

- It is unclear whether nonconvex algorithms converge to optimizers (due to lack of strong convexity)

# Approximate nonconvex optimizers

---

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

# Approximate nonconvex optimizers

---

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

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

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

# Approximate nonconvex optimizers

---

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

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  
 $\|\nabla f(\mathbf{X}, \mathbf{Y})\|_F \lesssim \frac{1}{\sqrt{\eta T}}$
- hopefully not far from  $(\mathbf{X}^*, \mathbf{Y}^*)$

## 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

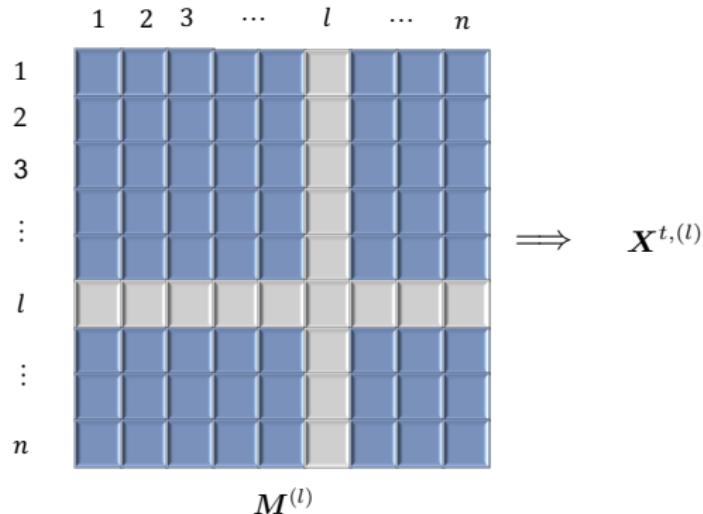
---

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  $\mathbf{X}^{t,(l)}$  by replacing  $l^{\text{th}}$  row and column with true values

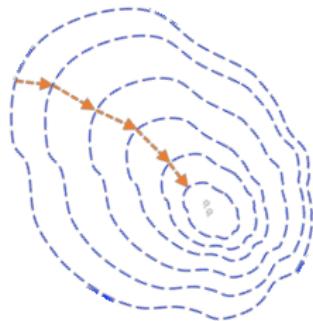


- exploit partial statistical independence
- exploit leave-one-out stability

*Backup slides: gradient descent for nonconvex matrix completion*

# Gradient descent for nonconvex matrix completion

---



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

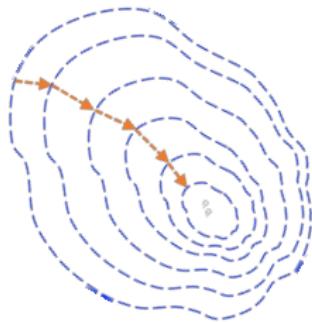
Prior works analyze regularized GD

- not guaranteed to return small-gradient solutions
- no  $\ell_{2,\infty}$  error control

— Keshavan et al. '09, Sun, Luo '15, Chen, Wainwright '15, Zheng, Lafferty '16

# Gradient descent for nonconvex matrix completion

---



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

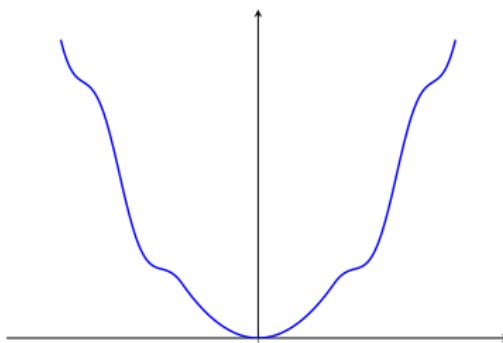
Our work and Chen et al. analyze **vanilla** GD

- regularization-free
- optimal  $\ell_{2,\infty}$  error control

— Ma, Wang, Chi, Chen '17, Chen, Liu, Li '19

# Gradient descent theory revisited

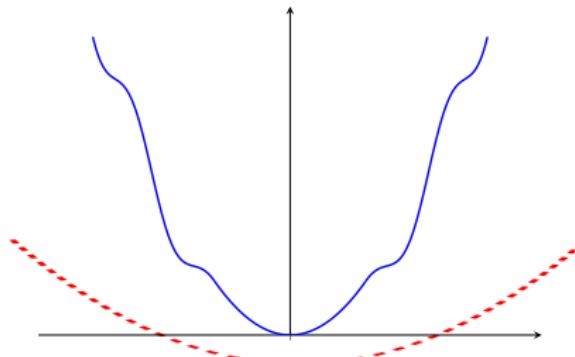
---



Two standard conditions that enable geometric convergence of GD

# Gradient descent theory revisited

---

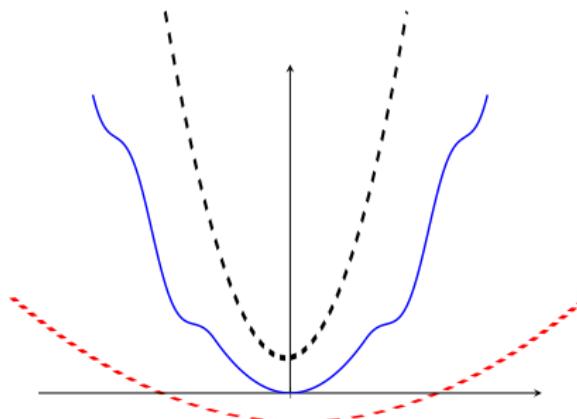


Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity

# Gradient descent theory revisited

---



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity
- (local) smoothness

# Gradient descent theory revisited

---

$f$  is said to be  $\alpha$ -strongly convex and  $\beta$ -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{X}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{X}$$

**$\ell_2$  error contraction:** GD with  $\eta = 1/\beta$  obeys

$$\|\mathbf{X}^{t+1} - \mathbf{X}^\star\|_{\text{F}} \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{X}^t - \mathbf{X}^\star\|_{\text{F}}$$

## Incoherence region

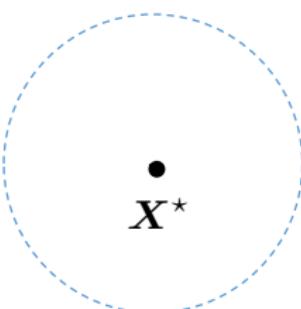
---

Which region enjoys both restricted strong convexity and smoothness?

## Incoherence region

---

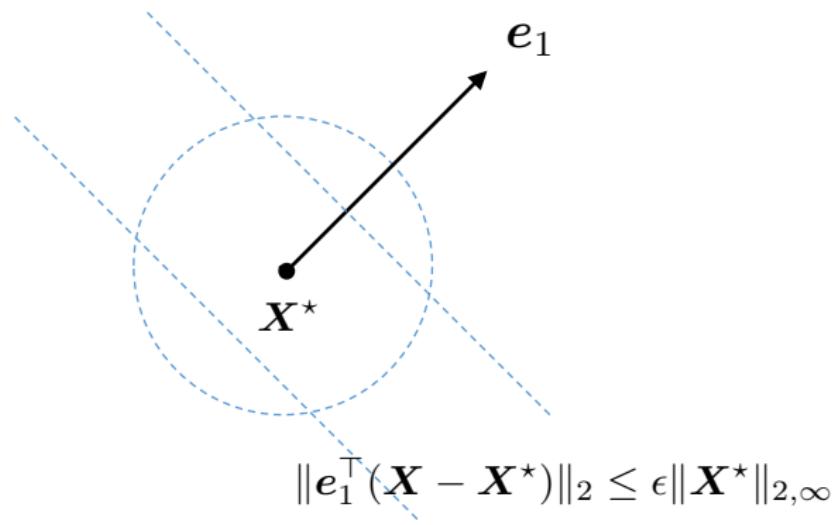
Which region enjoys both restricted strong convexity and smoothness?



- $X$  is not far away from  $X^*$

# Incoherence region

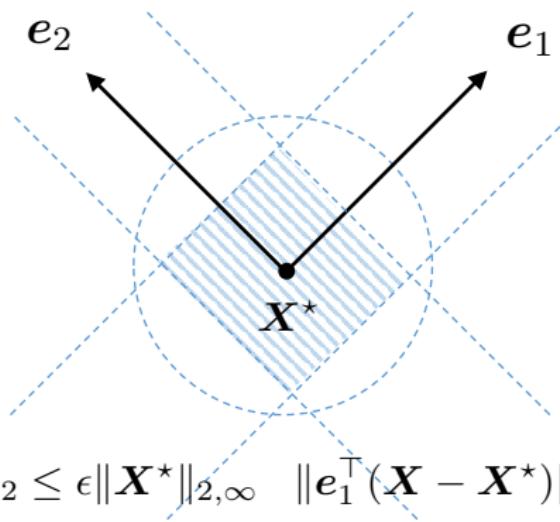
Which region enjoys both restricted strong convexity and smoothness?



- $X$  is not far away from  $X^*$
- $X$  is incoherent w.r.t. standard basis vectors (**incoherence region**)

# Incoherence region

Which region enjoys both restricted strong convexity and smoothness?



$$\|e_2^\top (X - X^*)\|_2 \leq \epsilon \|X^*\|_{2,\infty} \quad \|e_1^\top (X - X^*)\|_2 \leq \epsilon \|X^*\|_{2,\infty}$$

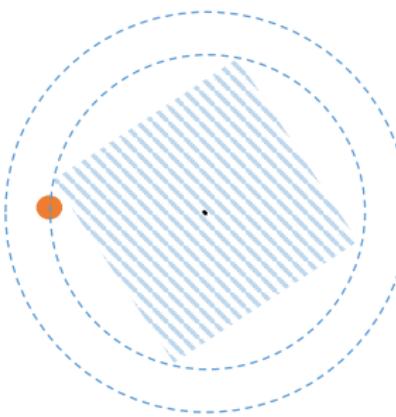
- $X$  is not far away from  $X^*$
- $X$  is incoherent w.r.t. standard basis vectors (**incoherence region**)

# Inadequacy of generic gradient descent theory

---



region of local strong convexity + smoothness



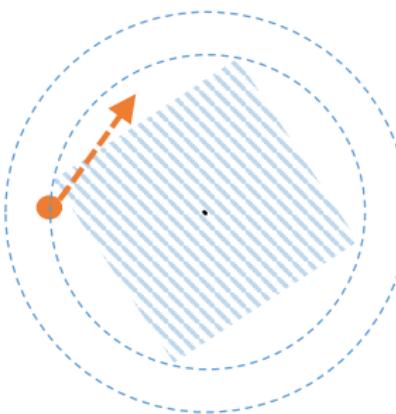
- Generic optimization theory does NOT ensure GD stays in incoherence region

# Inadequacy of generic gradient descent theory

---



region of local strong convexity + smoothness



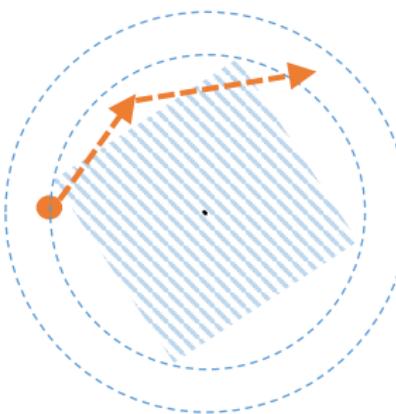
- Generic optimization theory does NOT ensure GD stays in incoherence region

# Inadequacy of generic gradient descent theory

---



region of local strong convexity + smoothness



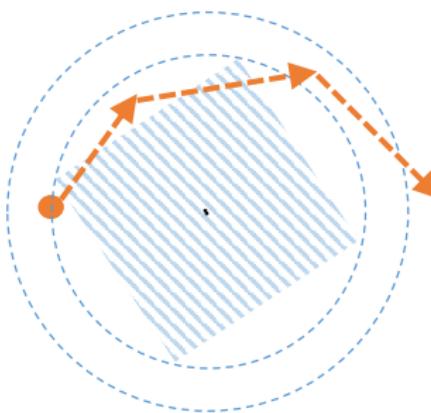
- Generic optimization theory does NOT ensure GD stays in incoherence region

# Inadequacy of generic gradient descent theory

---



region of local strong convexity + smoothness



- Generic optimization theory does NOT ensure GD stays in incoherence region
- Calls for new analysis tools

## **Key proof idea: leave-one-out analysis**

---

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

## Key proof idea: 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

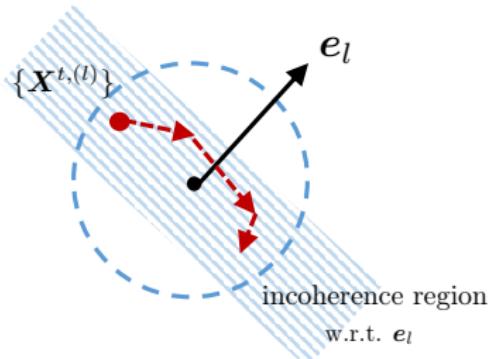
## Key proof idea: leave-one-out analysis

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

$$\begin{array}{ccccccc} & 1 & 2 & 3 & \cdots & l & \cdots & n \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ l \\ \vdots \\ n \end{matrix} & \begin{array}{|c|c|c|c|c|c|c|c|} \hline & \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{grey} & \text{blue} & \text{blue} \\ \hline \text{blue} & & & & & & & \\ \text{blue} & & & & & & & \\ \text{blue} & & & & & & & \\ \vdots & & & & & & & \\ \text{grey} & & \text{grey} & \text{grey} & \text{grey} & \text{grey} & \text{grey} & \text{grey} \\ \vdots & & & & & & & \\ \text{blue} & & \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} & \implies & \mathbf{X}^{t,(l)} \\ & \mathbf{M}^{(l)} & & & & & \end{array}$$

## Key proof idea: leave-one-out analysis

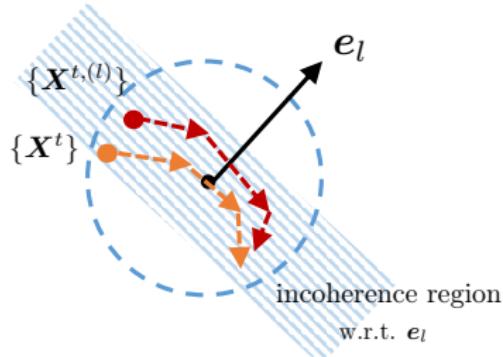
---



- Leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  contains more information of  $l^{\text{th}}$  row of truth; indep. of randomness in  $l^{\text{th}}$  row

## Key proof idea: leave-one-out analysis

---



- Leave-one-out iterates  $\{X^{t,(l)}\}$  contains more information of  $l^{\text{th}}$  row of truth; indep. of randomness in  $l^{\text{th}}$  row
- Leave-one-out iterates  $\{X^{t,(l)}\} \approx$  true iterates  $\{X^t\}$