

Stability, nonconvex optimization & asymmetry in low-rank matrix estimation



Yuxin Chen

EE, Princeton University

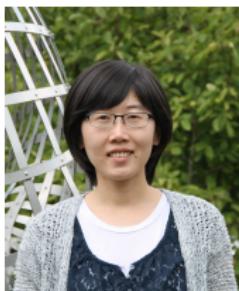
Noisy matrix completion: understanding stability of convex relaxation via nonconvex optimization



Cong Ma
Princeton ORFE



Yuling Yan
Princeton ORFE



Yuejie Chi
CMU ECE

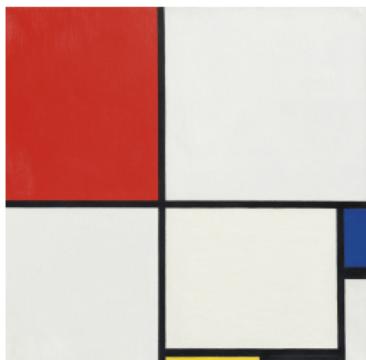


Jianqing Fan
Princeton ORFE

Convex relaxation for low-rank structure

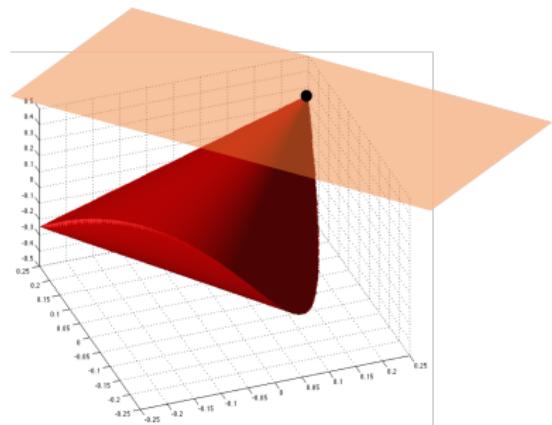
$$\underset{Z}{\text{minimize}} \quad \|Z\|_*$$

subj. to noiseless data constraints



low-rank matrix

figure credit: Piet Mondrian



semidefinite relaxation

Convex relaxation for low-rank structure

$$\underset{Z}{\text{minimize}} \quad \|Z\|_*$$

subj. to **noiseless** data constraints

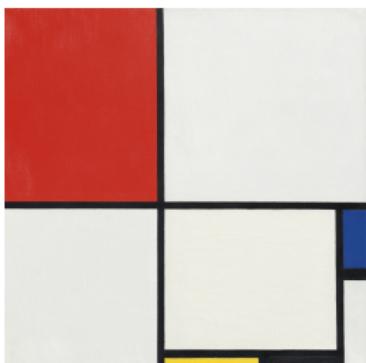
- ✓ matrix sensing (Recht, Fazel, Parrilo '07)
- ✓ phase retrieval (Candès, Strohmer, Voroninski '11, Candès, Li '12)
- ✓ matrix completion (Candès, Recht '08, Candès, Tao '08, Gross '09)
- ✓ robust PCA (Chandrasekaran et al. '09, Candès et al. '09)
- ✓ Hankel matrix completion (Fazel et al. '13, Chen, Chi '13, Cai et al. '15)
- ✓ blind deconvolution (Ahmed, Recht, Romberg '12, Ling, Strohmer '15)
- ✓ joint alignment / matching (Chen, Huang, Guibas '14)

...

Stability of convex relaxation against noise

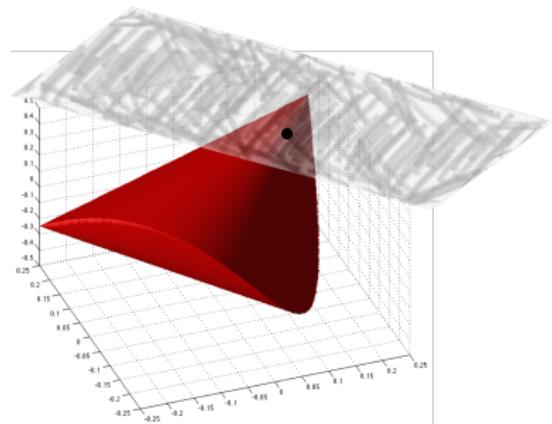
minimize _{Z} $\|Z\|_*$

subj. to **noisy** data constraints



low-rank matrix

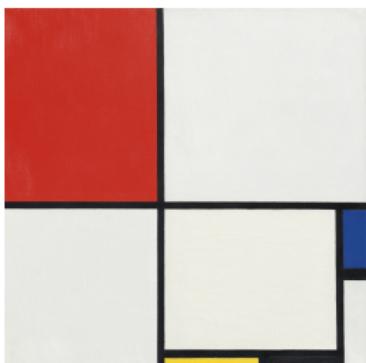
figure credit: Piet Mondrian



semidefinite relaxation

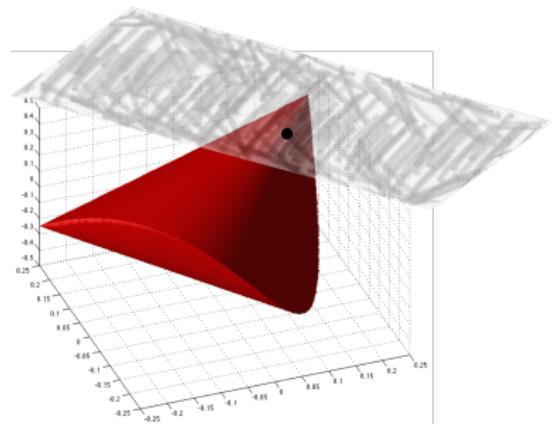
Stability of convex relaxation against noise

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



low-rank matrix

figure credit: Piet Mondrian



semidefinite relaxation

Stability of convex relaxation against noise

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

- ✓ matrix sensing (RIP measurements) (Candès, Plan '10)
- ✓ phase retrieval (Gaussian measurements) (Candès et al. '11)
- ? matrix completion
(Candès, Plan '09, Negahban, Wainwright '10, Koltchinskii et al. '10)
- ? robust PCA (Zhou, Li, Wright, Candès, Ma '10)
- ? Hankel matrix completion (Chen, Chi '13)
- ? blind deconvolution (Ahmed, Recht, Romberg '12, Ling, Strohmer '15)
- ? joint alignment / matching
- ...

Stability of convex relaxation against noise

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

- ✓ matrix sensing (RIP measurements) (Candès, Plan '10)
- ✓ phase retrieval (Gaussian measurements) (Candès et al. '11)
- ? **this talk: matrix completion** (Candès, Plan '09, Negahban, Wainwright '10, Koltchinskii et al. '10)
- ? robust PCA (Zhou, Li, Wright, Candès, Ma '10)
- ? Hankel matrix completion (Chen, Chi '13)
- ? blind deconvolution (Ahmed, Recht, Romberg '12, Ling, Strohmer '15)
- ? joint alignment / matching

...

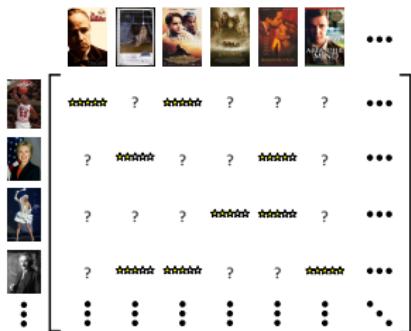
Low-rank matrix completion

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

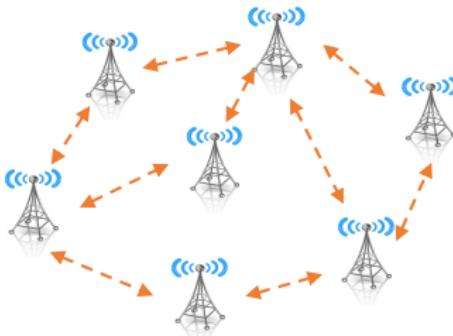


figure credit: E. J. Candès

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



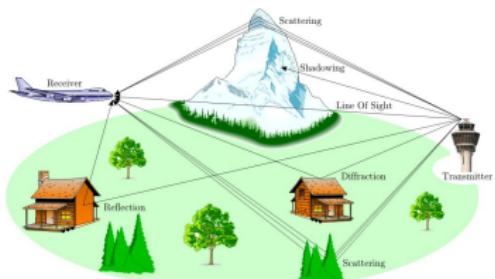
recommendation systems



localization



shape matching

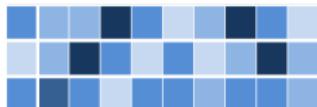
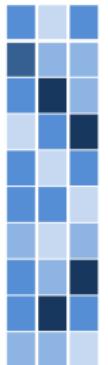


channel estimation

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 Ω

Noisy low-rank matrix completion

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

goal: estimate M^*

convex relaxation:

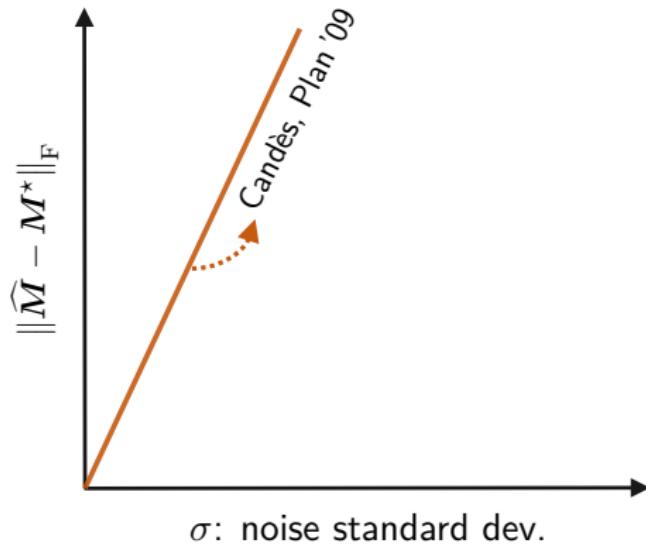
$$\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad \underbrace{\sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2}_{\text{squared loss}} + \lambda \|\mathbf{Z}\|_*$$

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 σ^2
- true matrix $M^* \in \mathbb{R}^{n \times n}$: rank $r = O(1)$, incoherent, ...

Candès, Plan '09

$\sigma n^{1.5}$

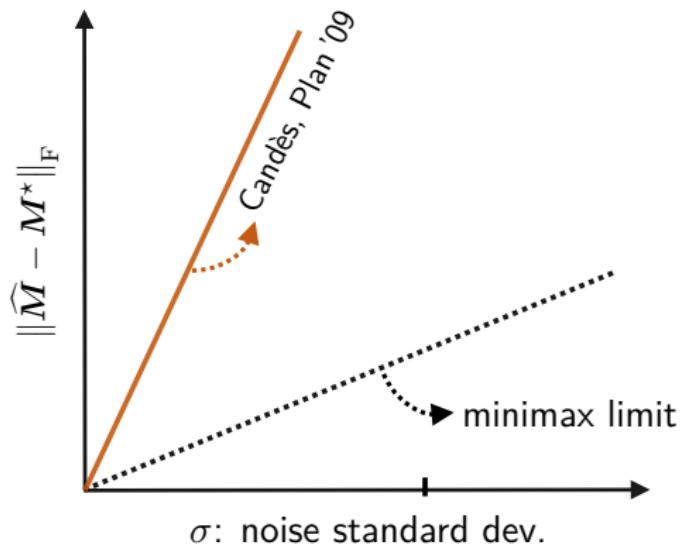


minimax limit

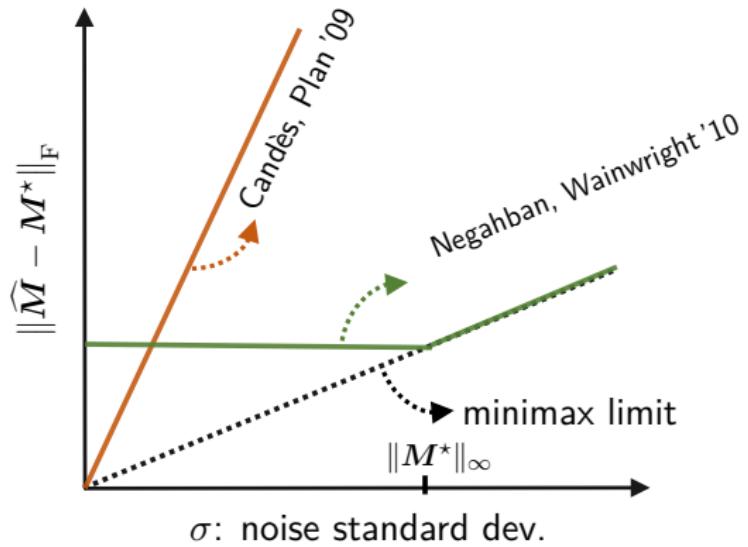
$$\sigma\sqrt{n/p}$$

Candès, Plan '09

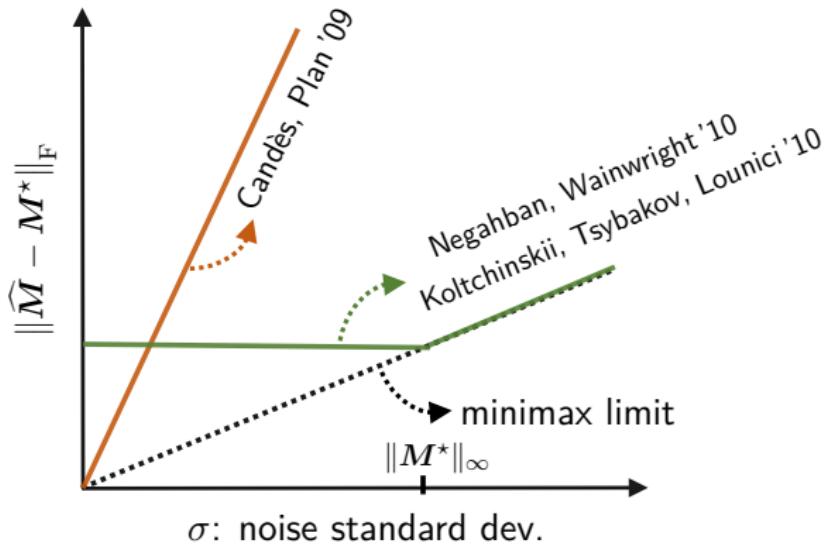
$$\sigma n^{1.5}$$



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

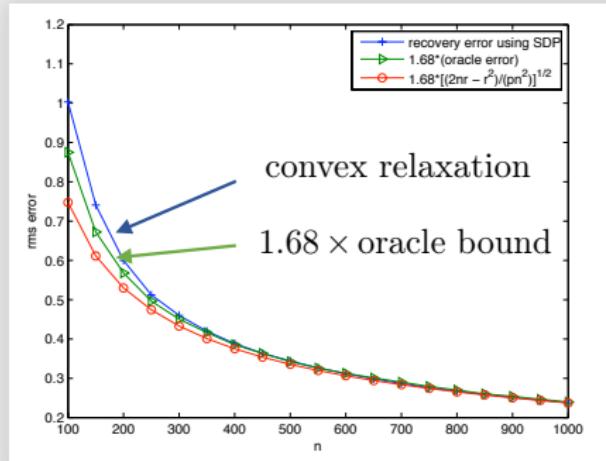


minimax limit	$\sigma\sqrt{n/p}$
Candès, Plan '09	$\sigma n^{1.5}$
Negahban, Wainwright '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$
Koltchinskii, Tsybakov, Lounici '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$



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

$(\widehat{M}_{\text{cvx}}, \mathbf{W})$ obeys KKT optimality condition

What are the roadblocks?

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

$(\widehat{\mathbf{M}}_{\text{cvx}}, \mathbf{W})$ obeys KKT optimality condition



David Gross

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

What are the roadblocks?

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

$(\widehat{\mathbf{M}}_{\text{cvx}}, \mathbf{W})$ obeys KKT optimality condition



David Gross

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

dual certification (golfing scheme)



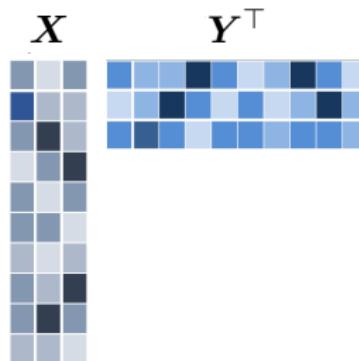
dual certification (golfing scheme)



nonconvex optimization

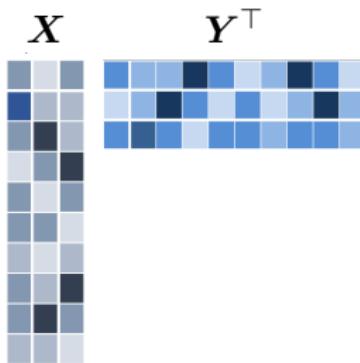
A detour: nonconvex optimization

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 $\underbrace{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}_{\text{low-rank factors}}$



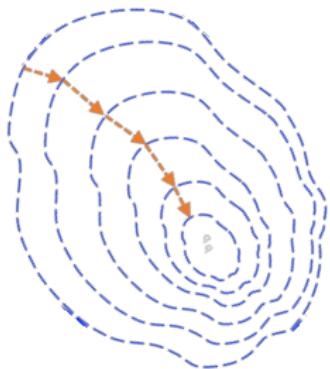
$$\underset{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}, \mathbf{Y}) = \underbrace{\sum_{(i,j) \in \Omega} \left[(\mathbf{X}\mathbf{Y}^\top)_{i,j} - M_{i,j} \right]^2}_{\text{squared loss}} + \text{reg}(\mathbf{X}, \mathbf{Y})$$

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 $t = 0, 1, \dots$

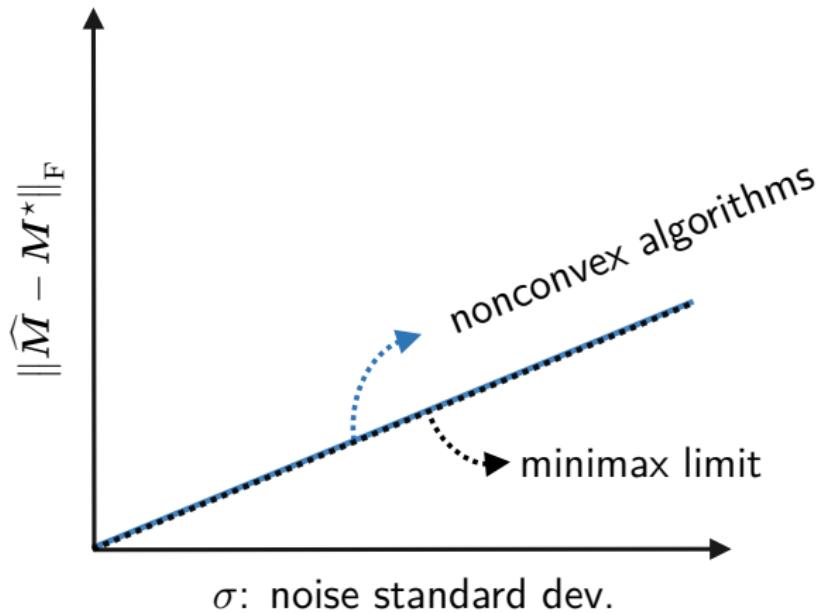
$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta_t \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t)$$

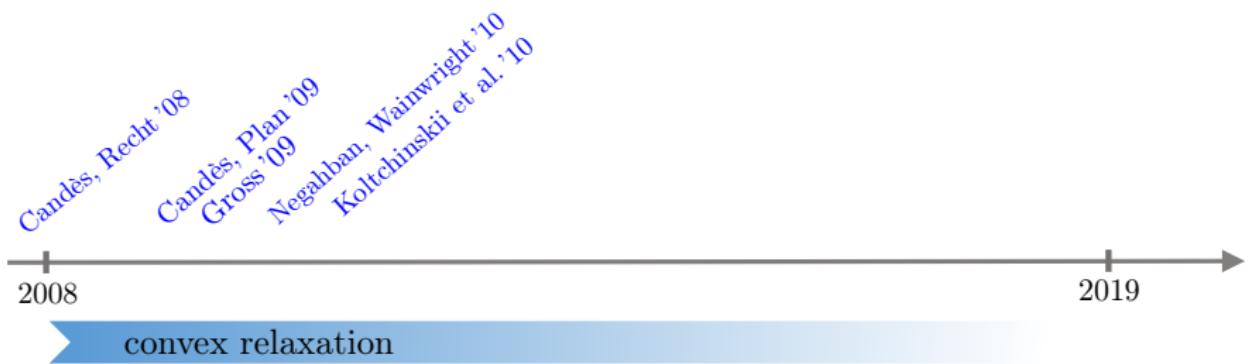
$$\mathbf{Y}^{t+1} = \mathbf{Y}^t - \eta_t \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)$$

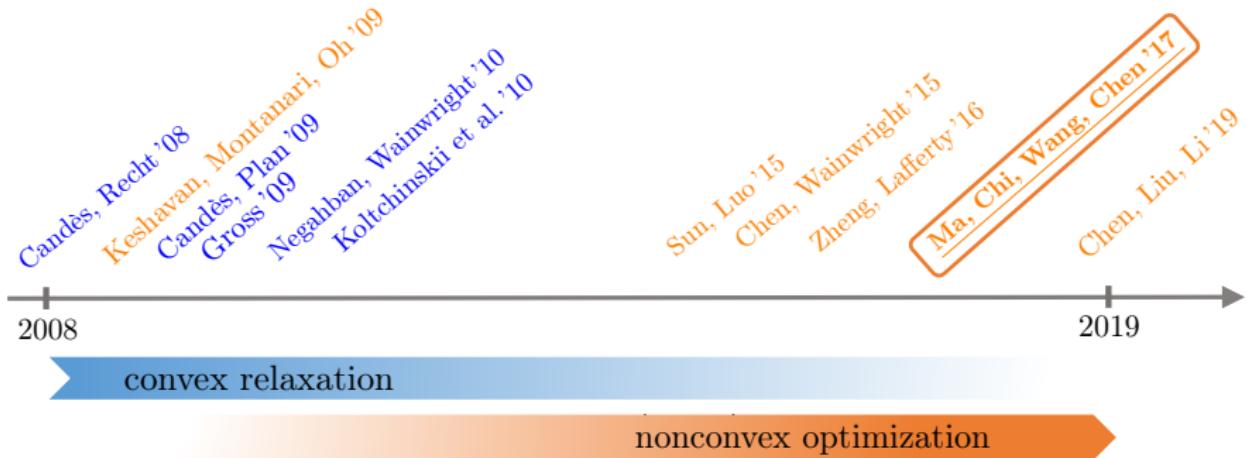
A detour: nonconvex optimization

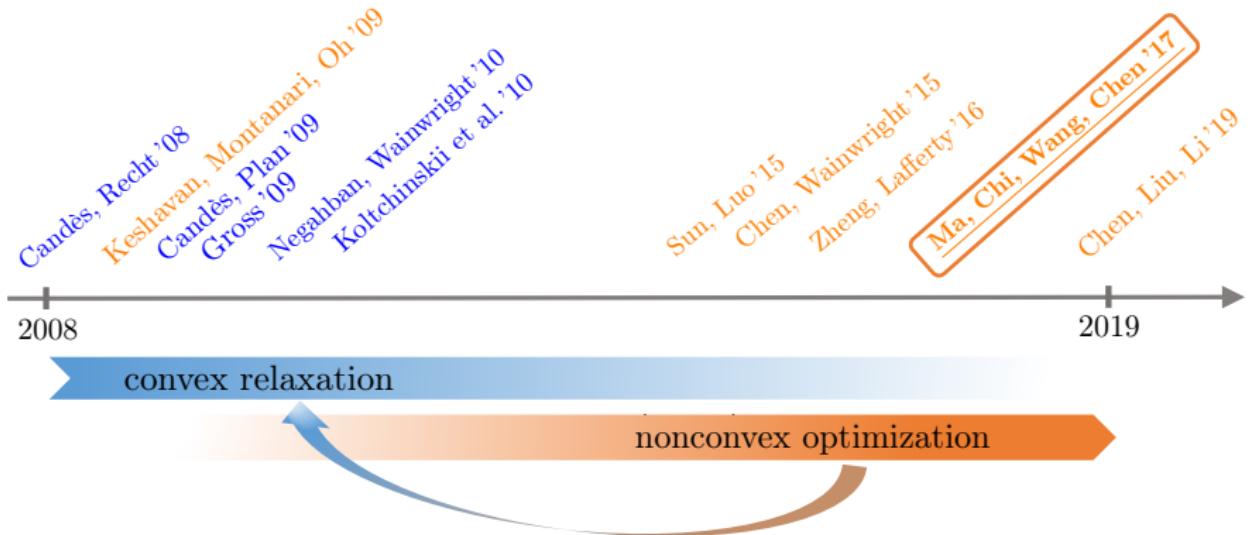
- **random sampling:** each $(i, j) \in \Omega$ with prob. p
- **random noise:** i.i.d. sub-Gaussian noise with variance σ^2
- true matrix $M^* \in \mathbb{R}^{n \times n}$: $r = O(1)$, incoherent, ...

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









A motivating experiment

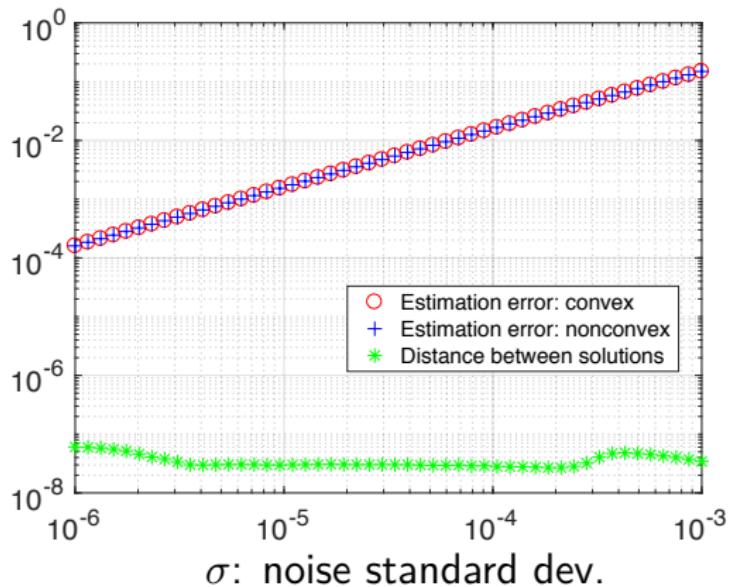
convex: $\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$

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

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

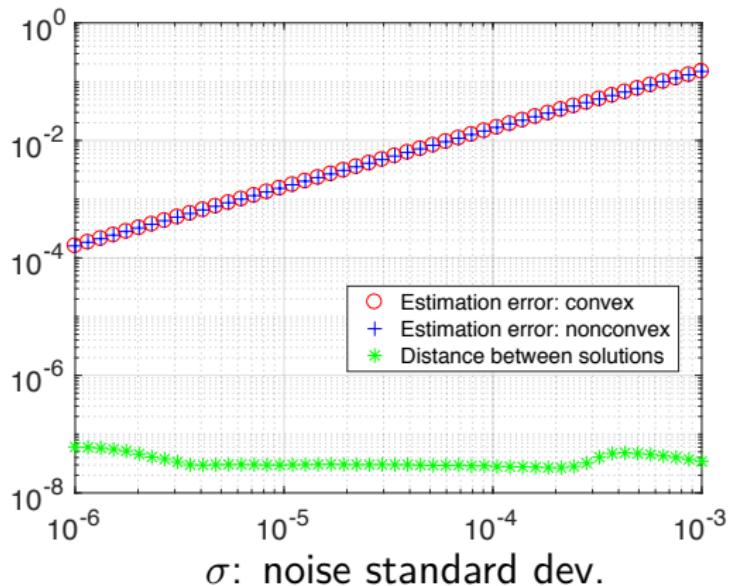
A motivating experiment

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



A motivating experiment

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



Convex and nonconvex solutions are exceedingly close!

convex



nonconvex



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

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 noise with variance σ^2
- 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 noise with variance σ^2
- 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 noise with variance σ^2
- 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 $\widehat{\mathbf{M}}_{\text{cvx}}$ of convex program obeys

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

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

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 noise with variance σ^2
- 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 \widehat{M}_{cvx} of convex program obeys

1. \widehat{M}_{cvx} is nearly rank- r

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

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 noise with variance σ^2
- true matrix $M^* \in \mathbb{R}^{n \times n}$: $r = O(1)$, incoherent, well-conditioned

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

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

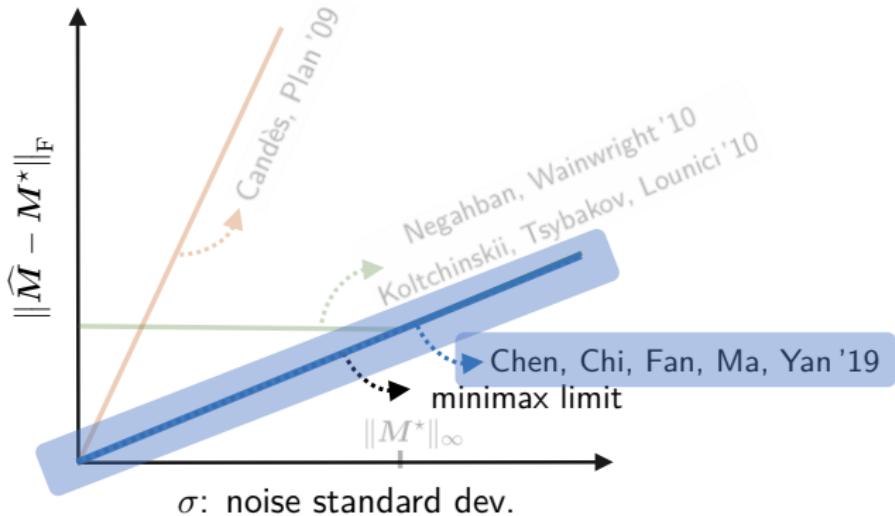
With high prob., any minimizer \widehat{M}_{cvx} of convex program obeys

1. \widehat{M}_{cvx} is nearly rank- r

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

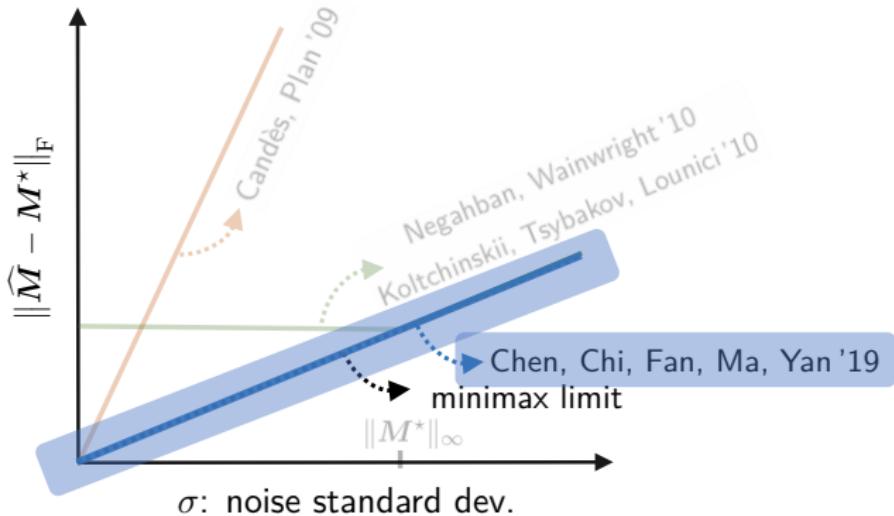
$$\|\widehat{M}_{\text{cvx}} - M^*\|_{\infty} \lesssim \sigma \sqrt{\frac{n \log n}{p}} \cdot \frac{1}{n}$$

$$\|\widehat{\boldsymbol{M}}_{\text{cvx}} - \boldsymbol{M}^*\|_{\text{F}} \lesssim \sigma \sqrt{\frac{n}{p}}$$



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

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



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

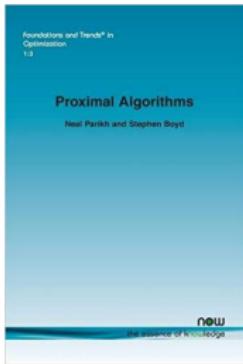
Implicit regularization

No need to enforce spikiness constraint as in Negahban & Wainwright

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

- convex programming automatically controls spikiness of solutions

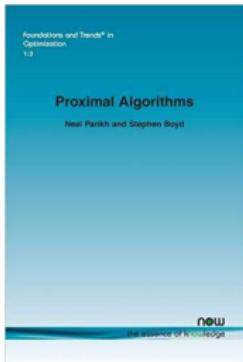
Statistical guarantees for iterative algorithms



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

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 noise with variance σ^2
- 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 noise with variance σ^2
- 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 \widehat{M}_{cvx} of convex program obeys

1. \widehat{M}_{cvx} is nearly rank- r

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

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

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

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 noise with variance σ^2
- 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 !

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

Link between convex and nonconvex optimizers

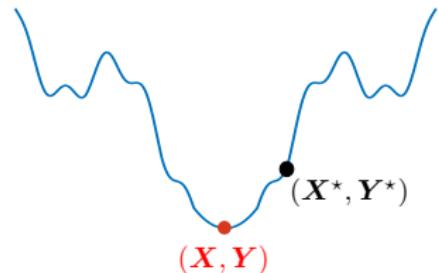
(X, Y) is nonconvex optimizer

Link between convex and nonconvex optimizers

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

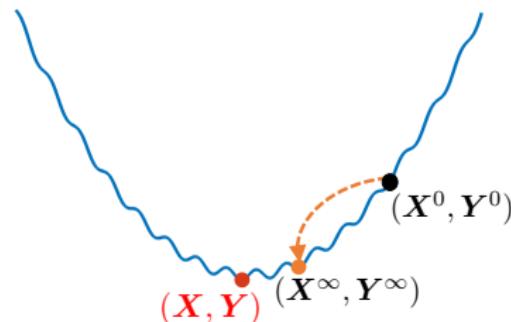
Link between convex and nonconvex optimizers

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



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

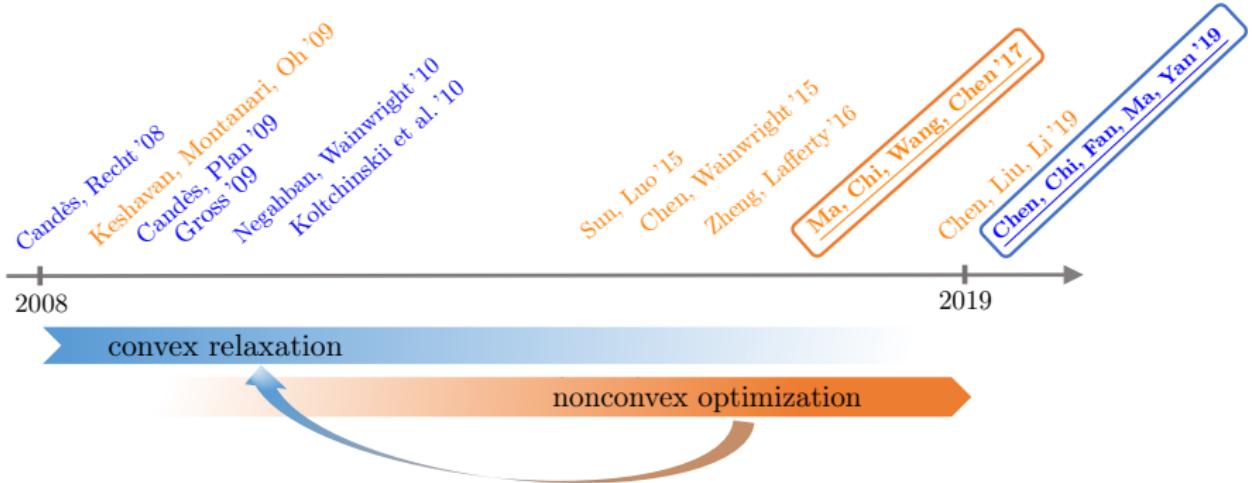
Leave out a small amount of information from data and run GD

- 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^{th} row and column with true values

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



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

Asymmetry helps: eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices

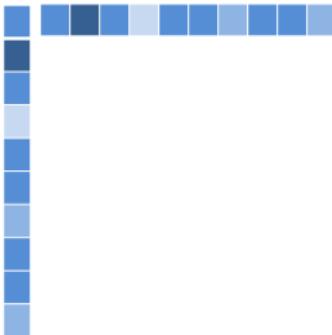


Chen Cheng
PKU Math & Stanford Stats



Jianqing Fan
Princeton ORFE

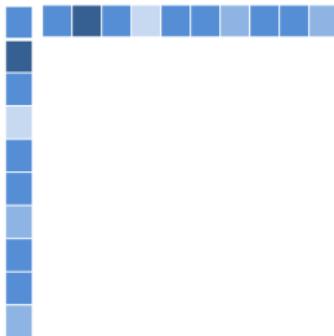
Eigenvalue / eigenvector estimation



M^* : truth

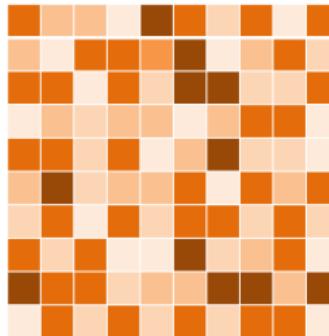
- A rank-1 matrix: $M^* = \lambda^* u^* u^{*\top} \in \mathbb{R}^{n \times n}$

Eigenvalue / eigenvector estimation



M^* : truth

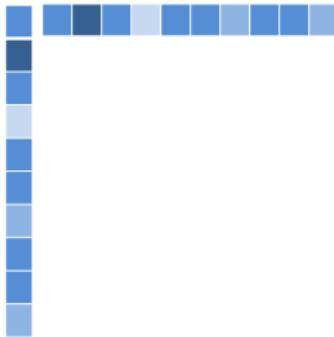
+



H : noise

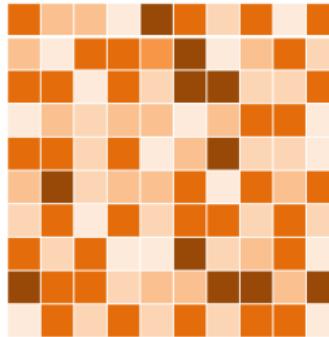
- A rank-1 matrix: $M^* = \lambda^* u^* u^{*\top} \in \mathbb{R}^{n \times n}$
- Observed noisy data: $M = M^* + H$

Eigenvalue / eigenvector estimation



M^* : truth

+



H : noise

- A rank-1 matrix: $M^* = \lambda^* u^* u^{*\top} \in \mathbb{R}^{n \times n}$
- Observed noisy data: $M = M^* + H$
- **Goal:** estimate eigenvalue λ^* and eigenvector u^*

Non-symmetric noise matrix

$$M = \begin{array}{c} \text{A vertical column of blue squares followed by a horizontal row of blue squares.} \\ + \end{array} \boxed{\begin{array}{c} \text{A 10x10 grid of orange and brown squares. The pattern is asymmetric, with many zero entries and varying intensities of orange and brown.} \\ H: \text{asymmetric matrix} \end{array}}$$
$$M^* = \lambda^* u^* u^{*\top}$$

This may arise when, e.g., we have 2 samples for each entry of M^* and arrange them in an asymmetric manner

A natural estimation strategy: SVD

$$M = \begin{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \end{matrix} + \boxed{\begin{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \end{matrix}}$$

$M^* = \lambda^* u^* u^{*\top}$

H : *asymmetric* matrix

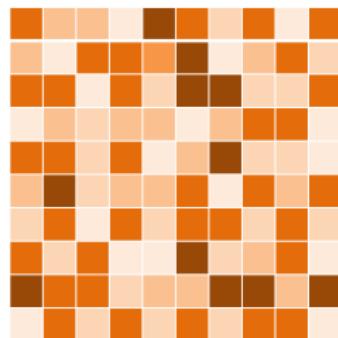
- Use leading singular value λ^{svd} of M to estimate λ^*
- Use leading left singular vector of M to estimate u^*

A less popular strategy: eigen-decomposition

$$M = \begin{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ \begin{matrix} M^* \\ + \end{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \end{matrix}$$

$$M^* = \lambda^* u^* u^{*\top}$$

+



H : *asymmetric* matrix

- Use leading singular value λ^{svd} eigenvalue λ^{eigs} of M to estimate λ^*
- Use leading singular vector eigenvector of M to estimate u^*

SVD vs. eigen-decomposition

For *asymmetric* matrices:

- Numerical stability

SVD > eigen-decomposition

SVD vs. eigen-decomposition

For *asymmetric* matrices:

- Numerical stability

$$\text{SVD} \quad > \quad \text{eigen-decomposition}$$

- **(Folklore?)** Statistical accuracy

$$\text{SVD} \quad \asymp \quad \text{eigen-decomposition}$$

SVD vs. eigen-decomposition

For *asymmetric* matrices:

- Numerical stability

$$\text{SVD} \quad > \quad \text{eigen-decomposition}$$

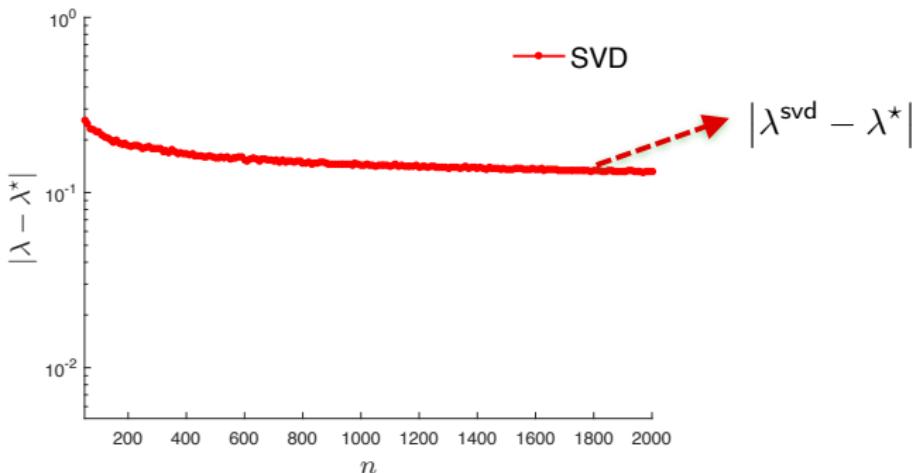
- **(Folklore?)** Statistical accuracy

$$\text{SVD} \quad \asymp \quad \text{eigen-decomposition}$$

Shall we always prefer SVD over eigen-decomposition?

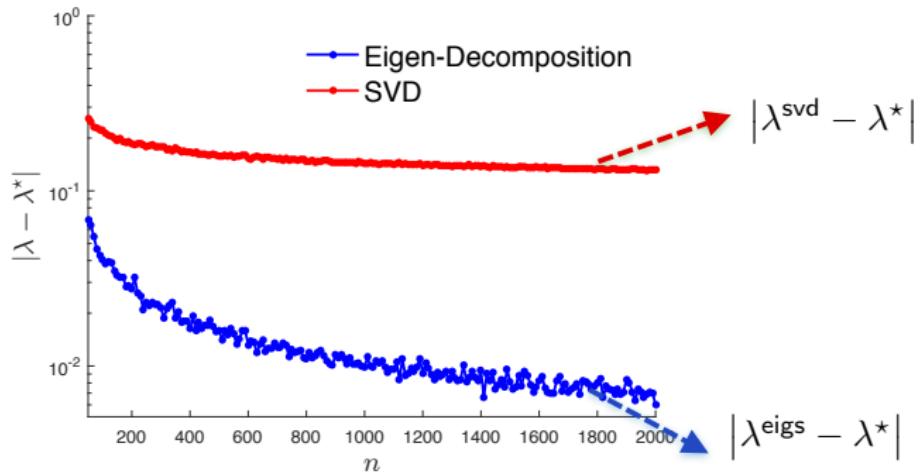
A curious numerical experiment: Gaussian noise

$$M = \underbrace{u^* u^{*\top}}_{M^*} + H; \quad \{H_{i,j}\} : \text{i.i.d. } \mathcal{N}(0, \sigma^2), \sigma = \frac{1}{\sqrt{n \log n}}$$



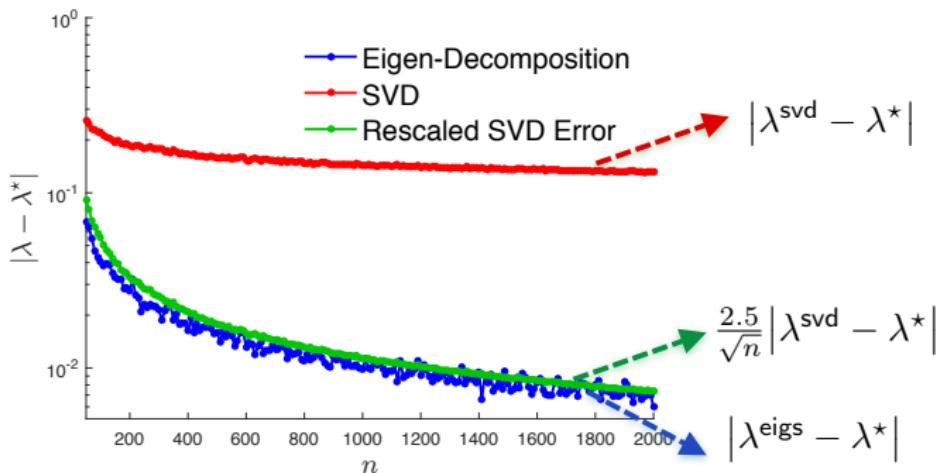
A curious numerical experiment: Gaussian noise

$$M = \underbrace{u^* u^{*\top}}_{M^*} + H; \quad \{H_{i,j}\} : \text{i.i.d. } \mathcal{N}(0, \sigma^2), \sigma = \frac{1}{\sqrt{n \log n}}$$



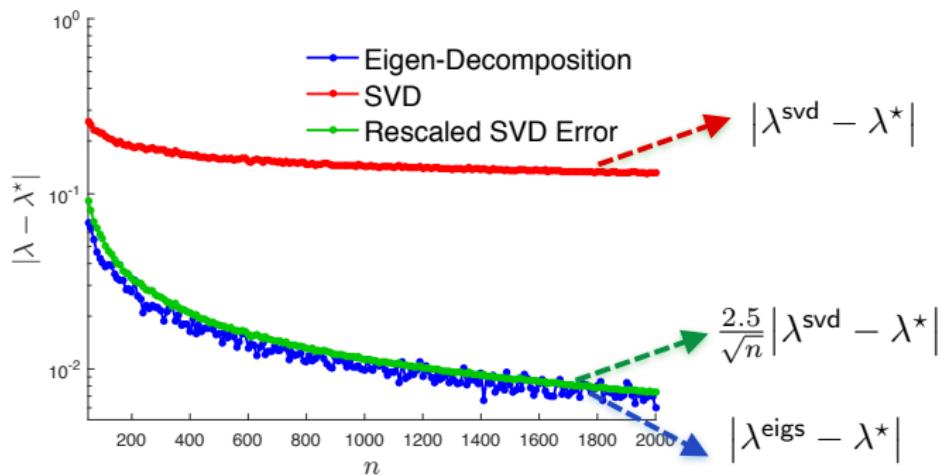
A curious numerical experiment: Gaussian noise

$$M = \underbrace{u^* u^{*\top}}_{M^*} + H; \quad \{H_{i,j}\} : \text{i.i.d. } \mathcal{N}(0, \sigma^2), \sigma = \frac{1}{\sqrt{n \log n}}$$



A curious numerical experiment: Gaussian noise

$$M = \underbrace{u^* u^{*\top}}_{M^*} + H; \quad \{H_{i,j}\} : \text{i.i.d. } \mathcal{N}(0, \sigma^2), \sigma = \frac{1}{\sqrt{n \log n}}$$



empirically, $|\lambda^{\text{eigs}} - \lambda^*| \approx \frac{2.5}{\sqrt{n}} |\lambda^{\text{svd}} - \lambda^*|$

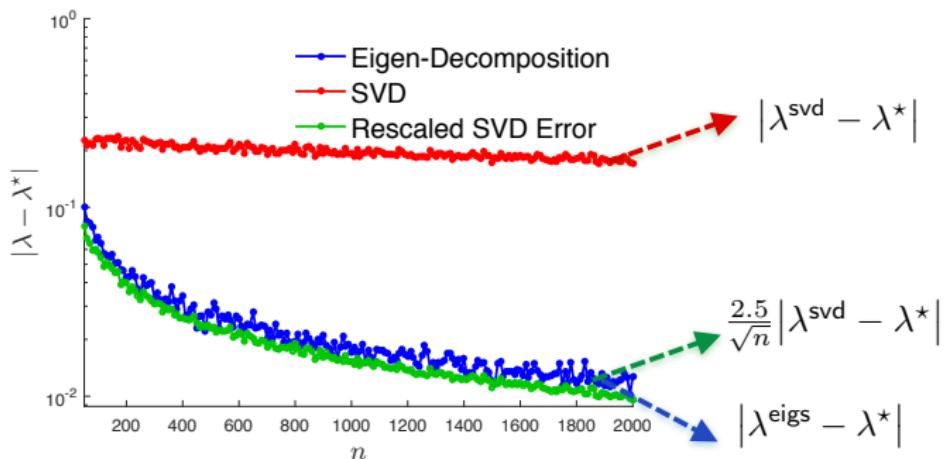
Another numerical experiment: matrix completion

$$M^* = \mathbf{u}^* \mathbf{u}^{*\top}; \quad M_{i,j} = \begin{cases} \frac{1}{p} M_{i,j}^* & \text{with prob. } p, \\ 0, & \text{else,} \end{cases} \quad p = \frac{3 \log n}{n}$$

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

Another numerical experiment: matrix completion

$$M^* = \mathbf{u}^* \mathbf{u}^{*\top}; \quad M_{i,j} = \begin{cases} \frac{1}{p} M_{i,j}^* & \text{with prob. } p, \\ 0, & \text{else,} \end{cases} \quad p = \frac{3 \log n}{n}$$



empirically, $|\lambda^{\text{eigs}} - \lambda^*| \approx \frac{2.5}{\sqrt{n}} |\lambda^{\text{svd}} - \lambda^*|$

Why does eigen-decomposition work so much better than SVD?

Problem setup

$$\mathbf{M} = \underbrace{\mathbf{u}^* \mathbf{u}^{*\top}}_{\mathbf{M}^*} + \mathbf{H} \in \mathbb{R}^{n \times n}$$

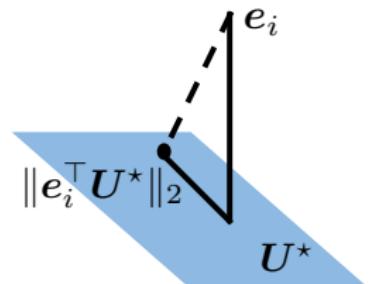
- \mathbf{H} : noise matrix
 - **independent entries:** $\{H_{i,j}\}$ are independent
 - **zero mean:** $\mathbb{E}[H_{i,j}] = 0$
 - **variance:** $\text{Var}(H_{i,j}) \leq \sigma^2$
 - **magnitudes:** $\mathbb{P}\{|H_{i,j}| \geq B\} \lesssim n^{-12}$

Problem setup

$$\mathbf{M} = \underbrace{\mathbf{u}^* \mathbf{u}^{*\top}}_{\mathbf{M}^*} + \mathbf{H} \in \mathbb{R}^{n \times n}$$

- \mathbf{H} : noise matrix
 - **independent entries**: $\{H_{i,j}\}$ are independent
 - **zero mean**: $\mathbb{E}[H_{i,j}] = 0$
 - **variance**: $\text{Var}(H_{i,j}) \leq \sigma^2$
 - **magnitudes**: $\mathbb{P}\{|H_{i,j}| \geq B\} \lesssim n^{-12}$
- \mathbf{M}^* obeys incoherence condition

$$\max_{1 \leq i \leq n} |\mathbf{e}_i^\top \mathbf{u}^*| \leq \sqrt{\frac{\mu}{n}}$$



Classical linear algebra results

$$|\lambda^{\text{svd}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Weyl})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Bauer-Fike})$$

Classical linear algebra results

$$|\lambda^{\text{svd}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Weyl})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Bauer-Fike})$$

\Downarrow matrix Bernstein inequality

$$|\lambda^{\text{svd}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n$$

$$|\lambda^{\text{eigs}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n$$

Classical linear algebra results

$$|\lambda^{\text{svd}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Weyl})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Bauer-Fike})$$

\Downarrow matrix Bernstein inequality

$$|\lambda^{\text{svd}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n \quad (\text{reasonably tight if } \|\mathbf{H}\| \text{ is large})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n$$

Classical linear algebra results

$$|\lambda^{\text{svd}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Weyl})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \leq \|\mathbf{H}\| \quad (\text{Bauer-Fike})$$

\Downarrow matrix Bernstein inequality

$$|\lambda^{\text{svd}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n \quad (\text{reasonably tight if } \|\mathbf{H}\| \text{ is large})$$

$$|\lambda^{\text{eigs}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n \quad (\text{can be significantly improved})$$

Main results: eigenvalue perturbation

Theorem 3 (Chen, Cheng, Fan '18)

With high prob., leading eigenvalue λ^{eigs} of M obeys

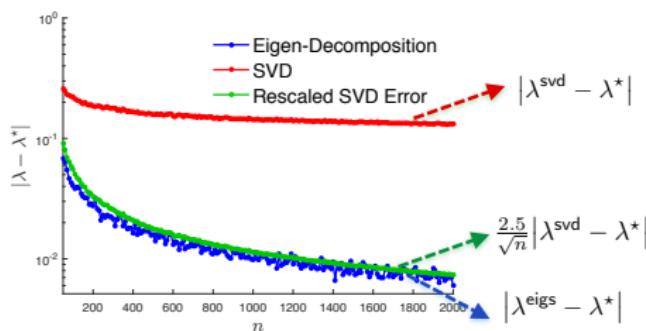
$$|\lambda^{\text{eigs}} - \lambda^*| \lesssim \sqrt{\frac{\mu}{n} (\sigma \sqrt{n \log n} + B \log n)}$$

Main results: eigenvalue perturbation

Theorem 3 (Chen, Cheng, Fan '18)

With high prob., leading eigenvalue λ^{eigs} of M obeys

$$|\lambda^{\text{eigs}} - \lambda^*| \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$



- Eigen-decomposition is $\sqrt{\frac{n}{\mu}}$ times better than SVD!

— recall $|\lambda^{\text{svd}} - \lambda^*| \lesssim \sigma \sqrt{n \log n} + B \log n$

Main results: entrywise eigenvector perturbation

Theorem 4 (Chen, Cheng, Fan '18)

With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

Main results: entrywise eigenvector perturbation

Theorem 4 (Chen, Cheng, Fan '18)

With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

- if $\|\mathbf{H}\| \ll |\lambda^*|$, then

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_2 \ll \|\mathbf{u}^*\|_2 \quad \text{(classical bound)}$$

Main results: entrywise eigenvector perturbation

Theorem 4 (Chen, Cheng, Fan '18)

With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

- if $\|\mathbf{H}\| \ll |\lambda^*|$, then

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_2 \ll \|\mathbf{u}^*\|_2 \quad (\text{classical bound})$$

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \ll \|\mathbf{u}^*\|_\infty \quad (\text{our bound})$$

Main results: entrywise eigenvector perturbation

Theorem 4 (Chen, Cheng, Fan '18)

With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

- if $\|\mathbf{H}\| \ll |\lambda^*|$, then

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_2 \ll \|\mathbf{u}^*\|_2 \quad (\text{classical bound})$$

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \ll \|\mathbf{u}^*\|_\infty \quad (\text{our bound})$$

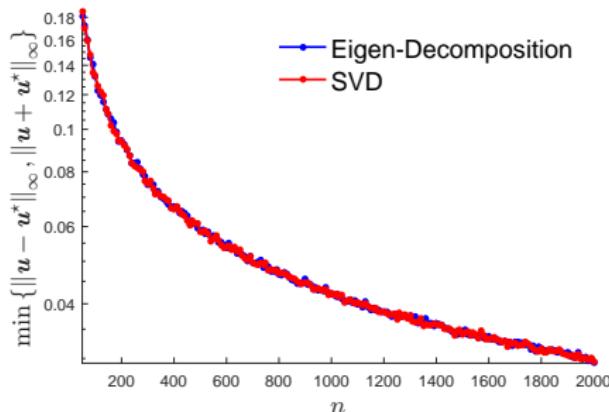
- entrywise eigenvector perturbation is well-controlled

Main results: entrywise eigenvector perturbation

Theorem 4 (Chen, Cheng, Fan '18)

With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \|\mathbf{u} \pm \mathbf{u}^*\|_\infty \lesssim \sqrt{\frac{\mu}{n}} (\sigma \sqrt{n \log n} + B \log n)$$



$$\{H_{i,j}\} : \text{i.i.d. } \mathcal{N}(0, \sigma^2); \sigma^2 = \frac{1}{n \log n}$$

Main results: perturbation of linear forms of eigenvectors

Theorem 5 (Chen, Cheng, Fan '18)

Fix any unit vector \mathbf{a} . With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \{|\mathbf{a}^\top (\mathbf{u} \pm \mathbf{u}^*)|\} \lesssim \max \left\{ |\mathbf{a}^\top \mathbf{u}^*|, \sqrt{\frac{\mu}{n}} \right\} (\sigma \sqrt{n \log n} + B \log n)$$

Main results: perturbation of linear forms of eigenvectors

Theorem 5 (Chen, Cheng, Fan '18)

Fix any unit vector \mathbf{a} . With high prob., leading eigenvector \mathbf{u} of M obeys

$$\min \{|\mathbf{a}^\top (\mathbf{u} \pm \mathbf{u}^*)|\} \lesssim \max \left\{ |\mathbf{a}^\top \mathbf{u}^*|, \sqrt{\frac{\mu}{n}} \right\} (\sigma \sqrt{n \log n} + B \log n)$$

- if $\|\mathbf{H}\| \ll |\lambda^*|$, then

$$\min \{|\mathbf{a}^\top (\mathbf{u} \pm \mathbf{u}^*)|\} \ll \max \left\{ |\mathbf{a}^\top \mathbf{u}^*|, \|\mathbf{u}^*\|_\infty \right\}$$

- perturbation of an *arbitrary* linear form of leading eigenvector is well-controlled

Intuition: asymmetry reduces bias

From Neumann series one can verify

some sort of Taylor expansion

$$|\lambda - \lambda^*| \asymp \left| \frac{u^{*\top} H u^*}{\lambda} + \frac{u^{*\top} H^2 u^*}{\lambda^2} + \frac{u^{*\top} H^3 u^*}{\lambda^3} + \dots \right|$$

Intuition: asymmetry reduces bias

From Neumann series one can verify
some sort of Taylor expansion

$$|\lambda - \lambda^*| \asymp \left| \frac{u^{*\top} H u^*}{\lambda} + \frac{u^{*\top} H^2 u^*}{\lambda^2} + \frac{u^{*\top} H^3 u^*}{\lambda^3} + \dots \right|$$

To develop some intuition, let's look at 2nd order term

Intuition: asymmetry reduces bias

From Neumann series one can verify
some sort of Taylor expansion

$$|\lambda - \lambda^*| \asymp \left| \frac{\mathbf{u}^{*\top} \mathbf{H} \mathbf{u}^*}{\lambda} + \boxed{\frac{\mathbf{u}^{*\top} \mathbf{H}^2 \mathbf{u}^*}{\lambda^2}} + \frac{\mathbf{u}^{*\top} \mathbf{H}^3 \mathbf{u}^*}{\lambda^3} + \dots \right|$$

To develop some intuition, let's look at 2nd order term

- if \mathbf{H} is symmetric,

$$\mathbb{E}[\mathbf{u}^{*\top} \mathbf{H}^2 \mathbf{u}^*] = \mathbb{E}[\|\mathbf{H} \mathbf{u}^*\|_2^2] = n\sigma^2$$

Intuition: asymmetry reduces bias

From Neumann series one can verify
some sort of Taylor expansion

$$|\lambda - \lambda^*| \asymp \left| \frac{\mathbf{u}^{*\top} \mathbf{H} \mathbf{u}^*}{\lambda} + \boxed{\frac{\mathbf{u}^{*\top} \mathbf{H}^2 \mathbf{u}^*}{\lambda^2}} + \frac{\mathbf{u}^{*\top} \mathbf{H}^3 \mathbf{u}^*}{\lambda^3} + \dots \right|$$

To develop some intuition, let's look at 2nd order term

- if \mathbf{H} is symmetric,

$$\mathbb{E}[\mathbf{u}^{*\top} \mathbf{H}^2 \mathbf{u}^*] = \mathbb{E}[\|\mathbf{H} \mathbf{u}^*\|_2^2] = n\sigma^2$$

- if \mathbf{H} is asymmetric,

$$\underbrace{\mathbb{E}[\mathbf{u}^{*\top} \mathbf{H}^2 \mathbf{u}^*] = \mathbb{E}[\langle \mathbf{H}^\top \mathbf{u}^*, \mathbf{H} \mathbf{u}^* \rangle]}_{\text{much smaller than symmetric case}} = \sigma^2$$

What happens if M^* is also not symmetric?

- A rank-1 matrix: $M^* = \lambda^* u^* v^{*\top} \in \mathbb{R}^{n_1 \times n_2}$
- Suppose we observe 2 independent noisy copies

$$M_1 = M^* + H_1, \quad M_2 = M^* + H_2$$

- **Goal:** estimate λ^* , u^* and v^*

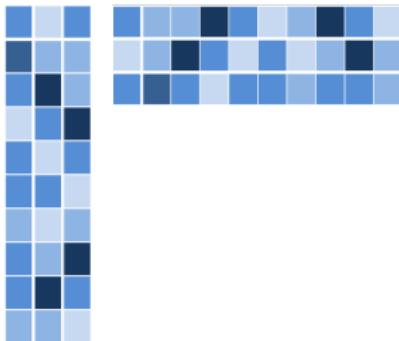
Asymmetrization + dilation

Compute leading eigenvalue / eigenvector of

$$\begin{bmatrix} \mathbf{0} & M_1 \\ M_2^\top & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & M^* + \mathbf{H}_1 \\ M^{*\top} + \mathbf{H}_2^\top & \mathbf{0} \end{bmatrix}$$

- Our findings (eigenvalue / eigenvector perturbation) continue to hold for this case!

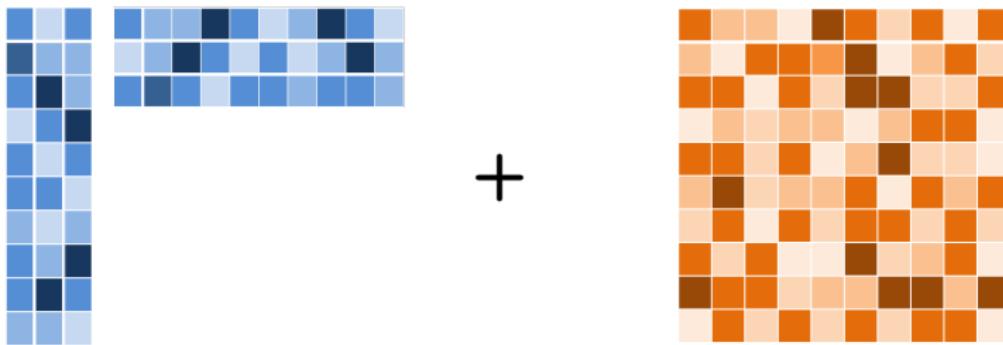
Rank- r case



M^* : truth

- A rank- r and well-conditioned matrix: $M^* = \sum_{i=1}^r \lambda_i^* u_i^* u_i^{*\top}$
- Observed noisy data: $M = M^* + H$, where $\{H_{i,j}\}$ are independent
- **Goal:** estimate λ^*

Rank- r case



M^* : truth

H : noise

- A rank- r and well-conditioned matrix: $M^* = \sum_{i=1}^r \lambda_i^* u_i^* u_i^{*\top}$
- Observed noisy data: $M = M^* + H$, where $\{H_{i,j}\}$ are independent
- **Goal:** estimate λ^*

Eigenvalue perturbation: rank- r case

Theorem 6 (Chen, Cheng, Fan '18)

With high prob., i th largest eigenvalue λ_i ($1 \leq i \leq r$) of M obeys

$$|\lambda_i - \lambda_j^*| \lesssim \sqrt{\frac{\mu r^2}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

for some $1 \leq j \leq r$

Eigenvalue perturbation: rank- r case

Theorem 6 (Chen, Cheng, Fan '18)

With high prob., i th largest eigenvalue λ_i ($1 \leq i \leq r$) of M obeys

$$|\lambda_i - \lambda_j^*| \lesssim \sqrt{\frac{\mu r^2}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

for some $1 \leq j \leq r$

- Eigen-decomposition is $\sqrt{\frac{n}{\mu r^2}}$ times better than SVD!

Eigenvalue perturbation: rank- r case

Theorem 6 (Chen, Cheng, Fan '18)

With high prob., i th largest eigenvalue λ_i ($1 \leq i \leq r$) of M obeys

$$|\lambda_i - \lambda_j^*| \lesssim \sqrt{\frac{\mu r^2}{n}} (\sigma \sqrt{n \log n} + B \log n)$$

for some $1 \leq j \leq r$

- Eigen-decomposition is $\sqrt{\frac{n}{\mu r^2}}$ times better than SVD!
- Might be improvable to $\sqrt{\frac{\mu r}{n}} (\sigma \sqrt{n \log n} + B \log n)$?

Concluding remarks

Eigen-decomposition could be much more powerful than SVD
when dealing with non-symmetric data matrices

Concluding remarks

Eigen-decomposition could be much more powerful than SVD
when dealing with non-symmetric data matrices

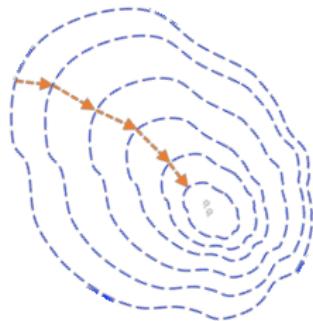
Future directions:

- Eigenvector perturbation for rank- r case
- Beyond i.i.d. noise

Y. Chen, C. Cheng, J. Fan, "Asymmetry helps: Eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices", [arXiv:1811.12804](https://arxiv.org/abs/1811.12804), 2018

Backup slides: gradient descent for nonconvex matrix completion

Gradient descent for nonconvex matrix completion



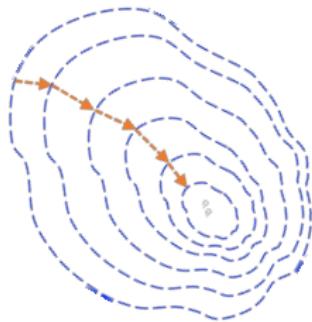
$$\begin{aligned}\mathbf{X}^{t+1} &= \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)\end{aligned}$$

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



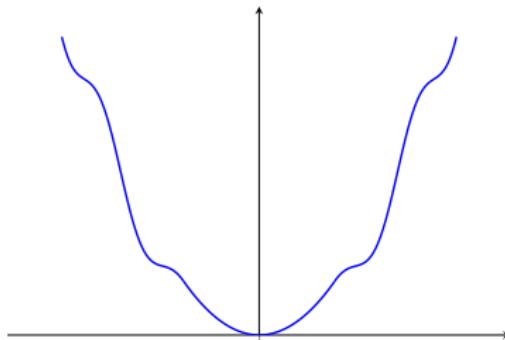
$$\begin{aligned}\mathbf{X}^{t+1} &= \mathbf{X}^t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^t, \mathbf{Y}^t) \\ \mathbf{Y}^{t+1} &= \mathbf{Y}^t - \eta \nabla_{\mathbf{Y}} f(\mathbf{X}^t, \mathbf{Y}^t)\end{aligned}$$

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

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

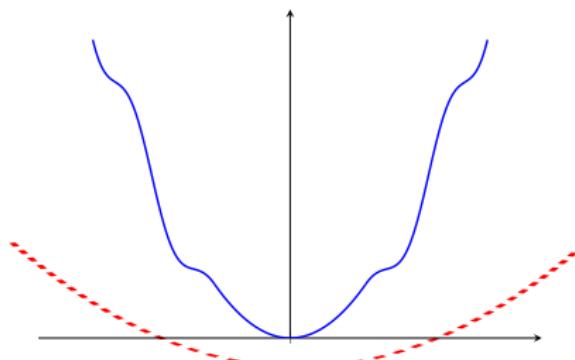
— Ma, Wang, Chi, Chen '17, Chen, Liu, Li '19

Gradient descent theory revisited



Two standard conditions that enable geometric convergence of GD

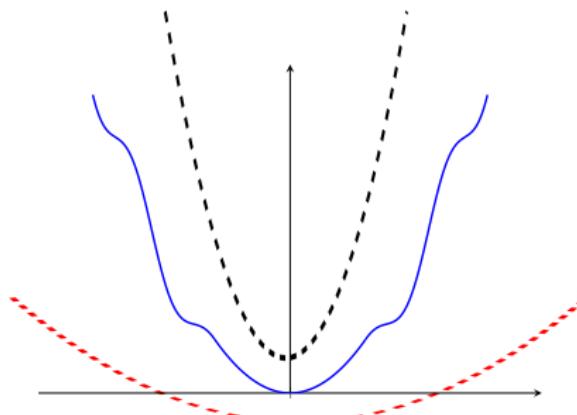
Gradient descent theory revisited



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 α -strongly convex and β -smooth if

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

ℓ_2 error contraction: GD with $\eta = 1/\beta$ obeys

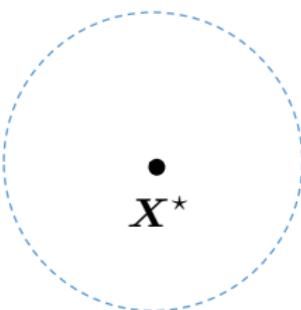
$$\|\mathbf{X}^{t+1} - \mathbf{X}^\star\|_{\text{F}} \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{X}^t - \mathbf{X}^\star\|_{\text{F}}$$

Incoherence region

Which region enjoys both restricted strong convexity and smoothness?

Incoherence region

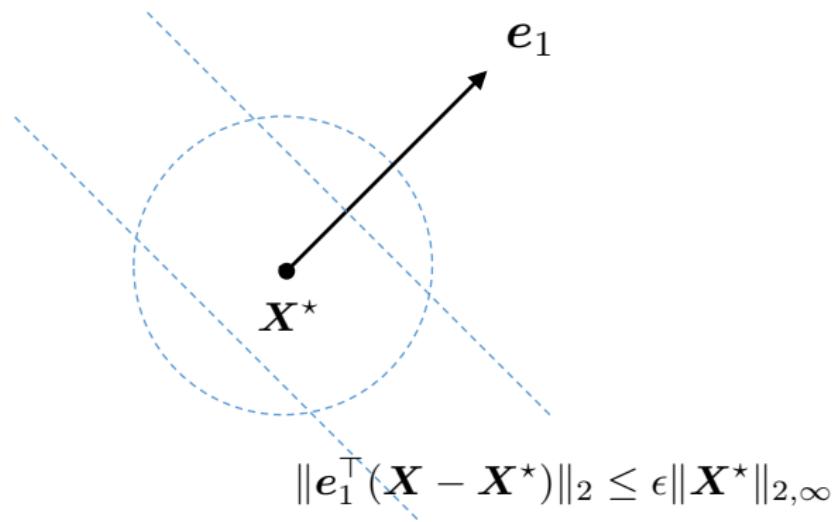
Which region enjoys both restricted strong convexity and smoothness?



- X is not far away from X^*

Incoherence region

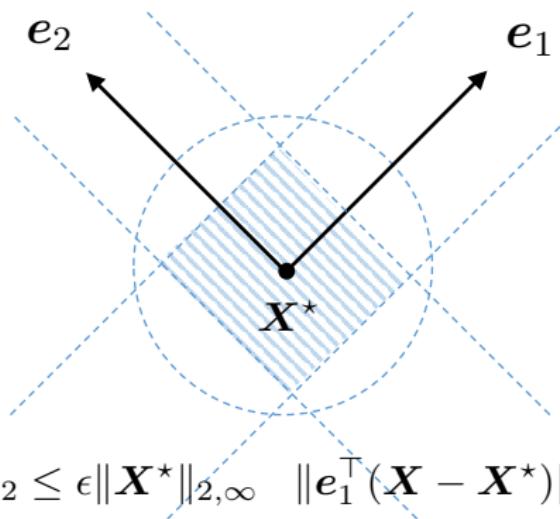
Which region enjoys both restricted strong convexity and smoothness?



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

Incoherence region

Which region enjoys both restricted strong convexity and smoothness?



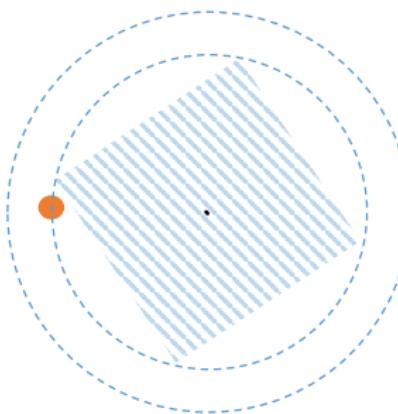
$$\|e_2^\top (X - X^*)\|_2 \leq \epsilon \|X^*\|_{2,\infty} \quad \|e_1^\top (X - X^*)\|_2 \leq \epsilon \|X^*\|_{2,\infty}$$

- 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

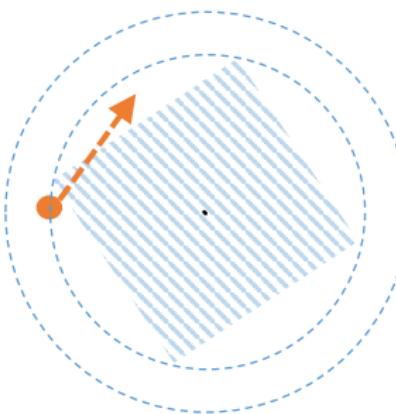


- Generic optimization theory does NOT ensure GD stays in incoherence region

Inadequacy of generic gradient descent theory



region of local strong convexity + smoothness

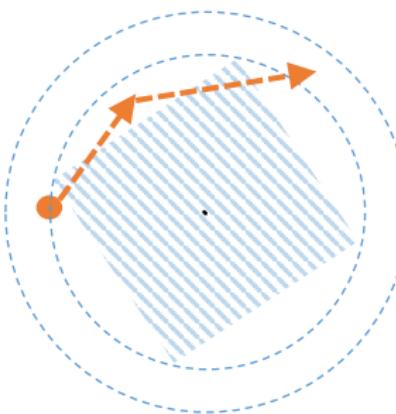


- Generic optimization theory does NOT ensure GD stays in incoherence region

Inadequacy of generic gradient descent theory



region of local strong convexity + smoothness

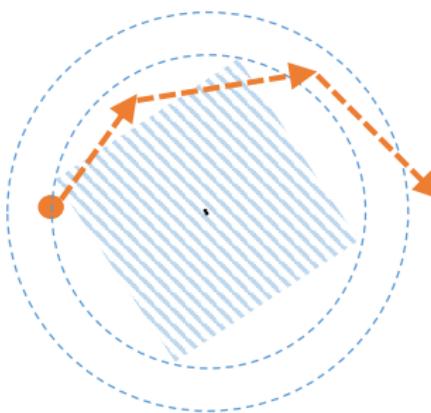


- Generic optimization theory does NOT ensure GD stays in incoherence region

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

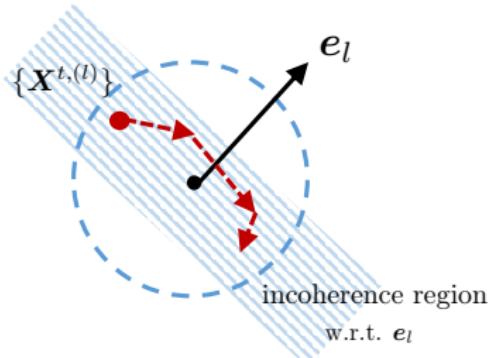
- Stein '72
- El Karoui, Bean, Bickel, Lim, Yu '13
- El Karoui '15
- Javanmard, Montanari '15
- Zhong, Boumal '17
- Lei, Bickel, El Karoui '17
- Sur, Chen, Candès '17
- Abbe, Fan, Wang, Zhong '17
- Chen, Fan, Ma, Wang '17
- Ma, Wang, Chi, Chen '17
- Chen, Chi, Fan, Ma '18
- Ding, Chen '18
- Dong, Shi '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^{th} row and column with true values

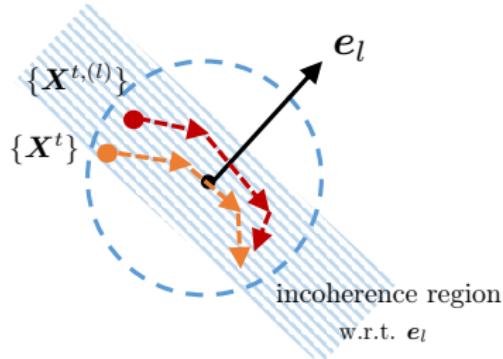
$$\begin{array}{ccccccc} & 1 & 2 & 3 & \cdots & l & \cdots & n \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \\ l \\ \vdots \\ n \end{matrix} & \begin{array}{|c|c|c|c|c|c|c|c|} \hline & \text{blue} & \text{blue} & \text{blue} & \text{blue} & \text{grey} & \text{blue} & \text{blue} \\ \hline \text{blue} & &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

Key proof idea: leave-one-out analysis



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

Key proof idea: leave-one-out analysis



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