

# **Nonconvex Optimization Meets Statistics: A Few Recent Stories**



Yuxin Chen

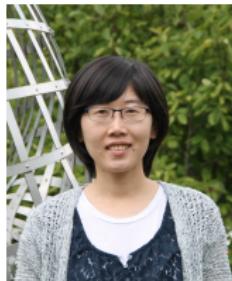
Electrical Engineering, Princeton University



Cong Ma  
Princeton ORFE



Yuling Yan  
Princeton ORFE



Yuejie Chi  
CMU ECE



Jianqing Fan  
Princeton ORFE

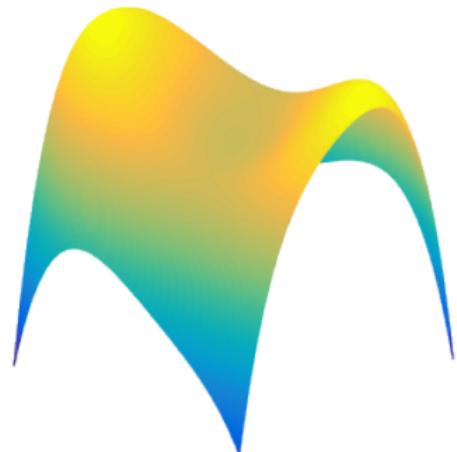
# Nonconvex problems are everywhere

---

Empirical risk minimization is usually nonconvex

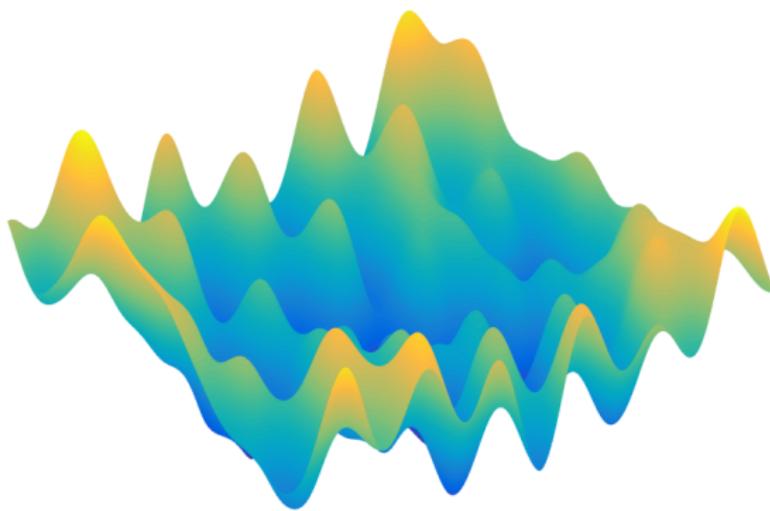
$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}; \text{data})$$

- low-rank matrix completion
- blind deconvolution
- dictionary learning
- mixture models
- deep neural nets
- ...



# Nonconvex optimization may be super scary

---

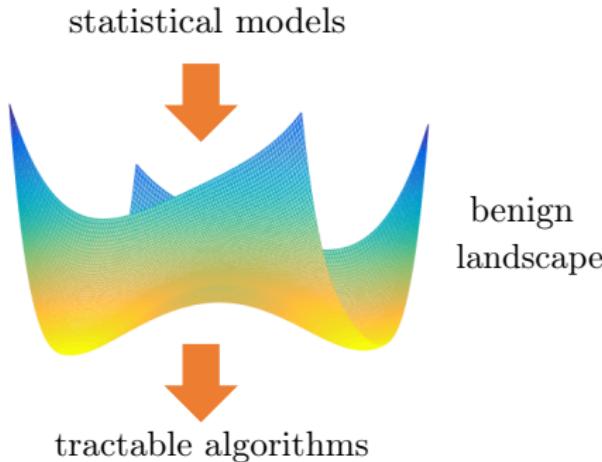


There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

# Statistical models come to rescue

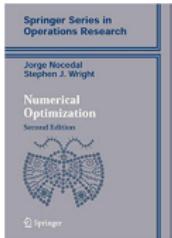
---



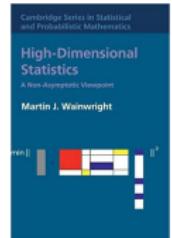
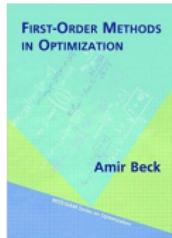
When data are generated by certain statistical models, problems are often much nicer than worst-case instances

— *Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview*

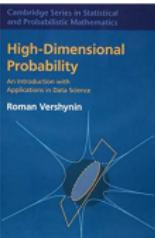
Chi, Lu, Chen '18



nonconvex optimization



(high-dimensional) statistics



1. Random initialization when solving random quadratic systems  
— *optimal computational efficiency*
  
2. Inference and uncertainty quantification for matrix completion  
— *a distributional theory*
  
3. Bridging convex & nonconvex optimization in matrix completion  
— *an implicit gift*

*Random initialization when  
solving random quadratic systems of equations*

# Solving quadratic systems of equations

$$\begin{array}{c} A \\ \left. \right\} m \\ \left. \right\} n \\ x \\ = \\ Ax \\ \longrightarrow \\ y = |Ax|^2 \end{array}$$

The diagram illustrates the computation of quadratic measurements. On the left, a matrix  $A$  of size  $m \times n$  is shown as a grid of colored squares. To its right is a vector  $x$  represented by vertical bars of varying shades of blue. An equals sign follows. To the right of that is the product  $Ax$ , shown as a vertical column of 10 blue numbers. A large arrow points from  $Ax$  to the final result  $y = |Ax|^2$ , which is also a vertical column of 10 blue numbers.

$Ax$	$y =  Ax ^2$
1	1
-3	9
2	4
-1	1
4	16
2	4
-2	4
-1	1
3	9
4	16

Estimate  $\boldsymbol{x}^* \in \mathbb{R}^n$  from  $m$  random quadratic measurements

$$y_k = (\boldsymbol{a}_k^\top \boldsymbol{x}^*)^2, \quad k = 1, \dots, m$$

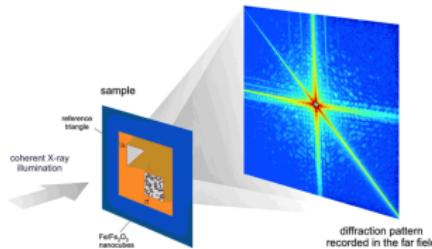
assume w.l.o.g.  $\|\boldsymbol{x}^*\|_2 = 1$

# Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

- electric field  $x(t_1, t_2) \longrightarrow$  Fourier transform  $\hat{x}(f_1, f_2)$

*Fig credit: Stanford SLAC*



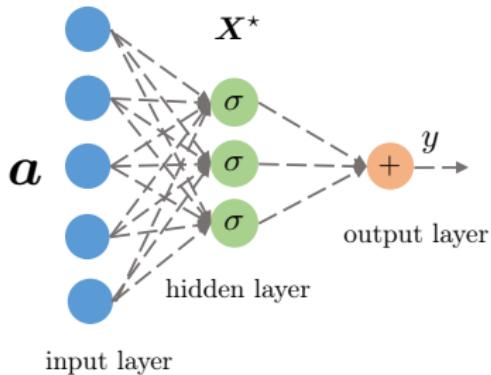
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

**Phase retrieval:** recover signal  $x(t_1, t_2)$  from intensity  $|\hat{x}(f_1, f_2)|^2$

# Motivation: learning neural nets with quadratic activation

---

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17



input features:  $a$ ; weights:  $\mathbf{X}^* = [x_1^*, \dots, x_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(a^\top x_i^*) \stackrel{\sigma(z)=z^2}{=} \sum_{i=1}^r (a^\top x_i^*)^2$$

## Wirtinger flow (Candès, Li, Soltanolkotabi '14)

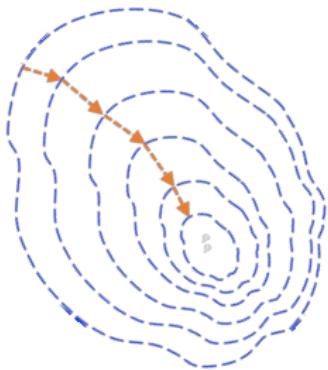
---

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[ (\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$

# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

---

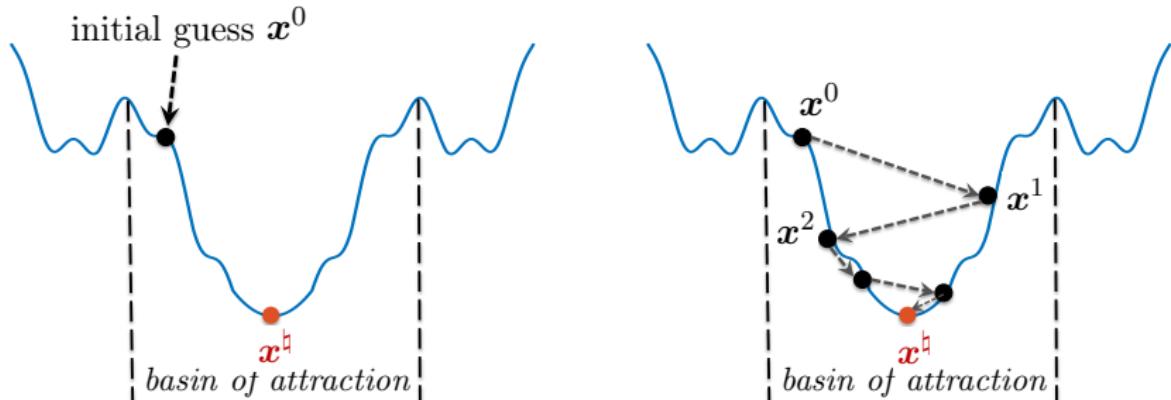
$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[ (\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$



- **spectral initialization:**  $\boldsymbol{x}^0 \leftarrow$  leading eigenvector of certain data matrix
- **gradient descent:**

$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \eta_t \nabla f(\boldsymbol{x}^t), \quad t = 0, 1, \dots$$

# Rationale of two-stage approach



1. initialize within  $\underbrace{\text{local basin sufficiently close to } x^*}_{\text{(restricted) strongly convex; no saddles / spurious local mins}}$
2. iterative refinement

# A highly incomplete list of two-stage methods

---

## phase retrieval:

- Netrapalli, Jain, Sanghavi '13
- Candès, Li, Soltanolkotabi '14
- Chen, Candès '15
- Cai, Li, Ma '15
- Wang, Giannakis, Eldar '16
- Zhang, Zhou, Liang, Chi '16
- Kolte, Ozgur '16
- Zhang, Chi, Liang '16
- Soltanolkotabi '17
- Vaswani, Nayer, Eldar '16
- Chi, Lu '16
- Wang, Zhang, Giannakis, Akcakaya, Chen '16
- Tan, Vershynin '17
- Ma, Wang, Chi, Chen '17
- Duchi, Ruan '17
- Jeong, Gunturk '17
- Yang, Yang, Fang, Zhao, Wang, Neykov '17
- Qu, Zhang, Wright '17
- Goldstein, Studer '16
- Bahmani, Romberg '16
- Hand, Voroninski '16
- Wang, Giannakis, Saad, Chen '17
- Barmherzig, Sun '17
- ...

## other problems:

- Keshavan, Montanari, Oh '09
- Sun, Luo '14
- Chen, Wainwright '15
- Tu, Boczar, Simchowitz, Soltanolkotabi, Recht '15
- Zheng, Lafferty '15
- Balakrishnan, Wainwright, Yu '14
- Chen, Suh '15
- Chen, Candès '16
- Li, Ling, Strohmer, Wei '16
- Yi, Park, Chen, Caramanis '16
- Jin, Kakade, Netrapalli '16
- Huang, Kakade, Kong, Valiant '16
- Ling, Strohmer '17
- Aghasi, Ahmed, Hand '17
- Lee, Tian, Romberg '17
- Li, Chi, Zhang, Liang '17
- Cai, Wang, Wei '17
- Abbe, Bandeira, Hall '14
- Chen, Kamath, Suh, Tse '16
- Zhang, Zhou '17
- Boumal '16
- Zhong, Boumal '17
- Li, Ma, Chen, Chi '18
- Chen, Liu, Li '19
- Charisopoulos, Davis, Diaz, Drusvyatskiy '19
- Charisopoulos, Chen, Davis, Diaz, Ding, Drusvyatskiy '19
- ...

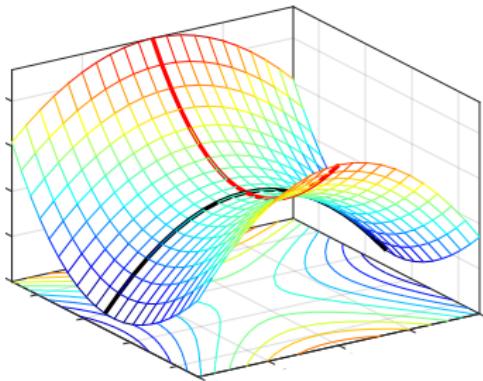
*Is carefully-designed initialization necessary  
for fast convergence?*

*Is carefully-designed initialization necessary  
for fast convergence?*

*Can we initialize GD randomly, which is **simpler** and  
**model-agnostic**?*

# What does prior theory say?

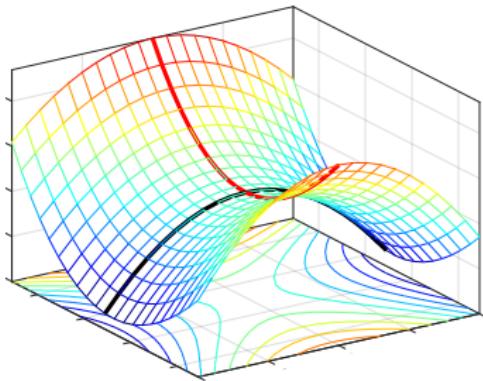
---



- **landscape:** no spurious local mins (Sun, Qu, Wright '16)

# What does prior theory say?

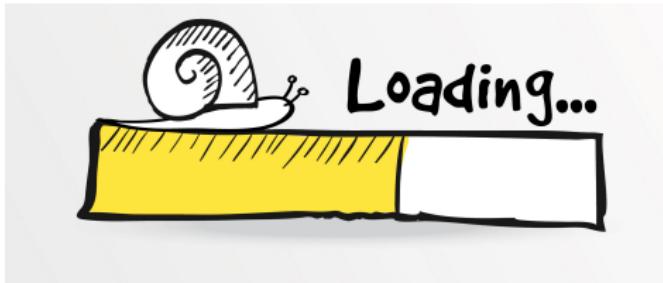
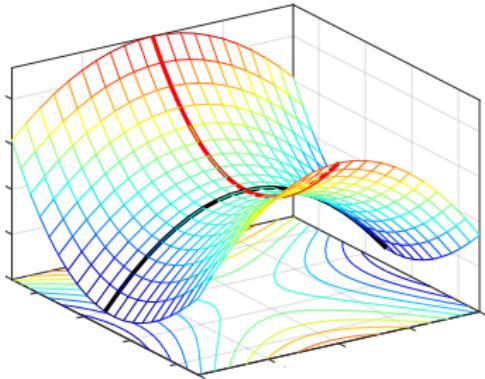
---



- **landscape:** no spurious local mins (Sun, Qu, Wright '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

# What does prior theory say?

---

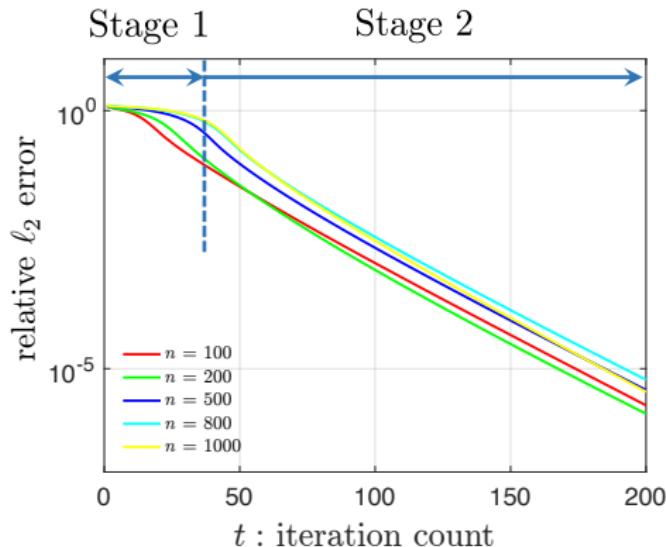


- **landscape:** no spurious local mins (Sun, Qu, Wright '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

“almost surely” might mean “take forever”

# Numerical efficiency of randomly initialized GD

$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **tens of iterations**

## Our theory

---

These numerical findings can be formalized when  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ :

# Our theory

---

These numerical findings can be formalized when  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ :

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

## Theorem 1 (Chen, Chi, Fan, Ma '18)

Under i.i.d. Gaussian design, GD with  $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$  achieves

# Our theory

---

These numerical findings can be formalized when  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ :

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min\{\|\mathbf{x}^t \pm \mathbf{x}^*\|_2\}$$

## Theorem 1 (Chen, Chi, Fan, Ma '18)

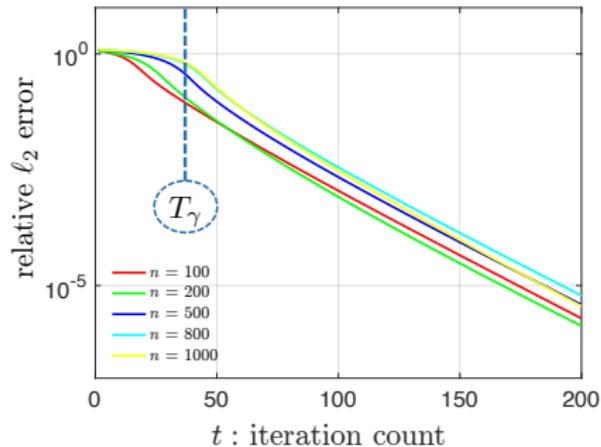
Under i.i.d. Gaussian design, GD with  $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$  achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma$$

with high prob. for  $T_\gamma \lesssim \log n$  and some constants  $\gamma, \rho > 0$ , provided that step size  $\eta \asymp 1$  and sample size  $m \gtrsim n \text{polylog } m$

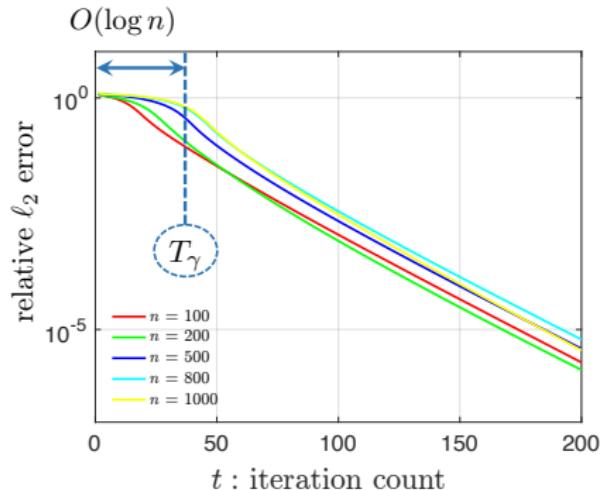
# Our theory

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\star) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\star\|_2, \quad t \geq T_\gamma \asymp \log n$$



# Our theory

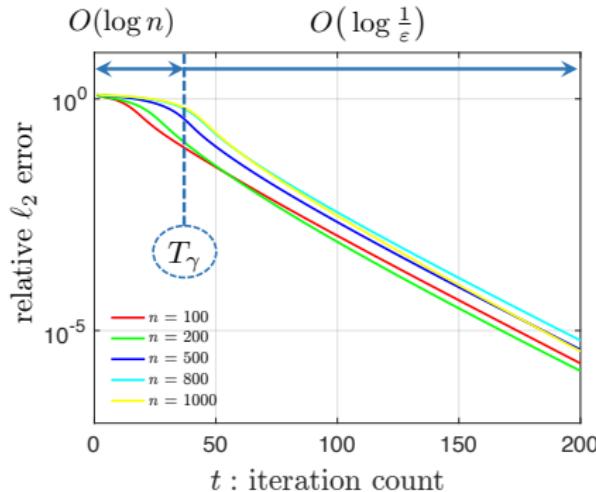
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- Stage 1: takes  $O(\log n)$  iterations to reach  $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$  (e.g.  $\gamma = 0.1$ )

# Our theory

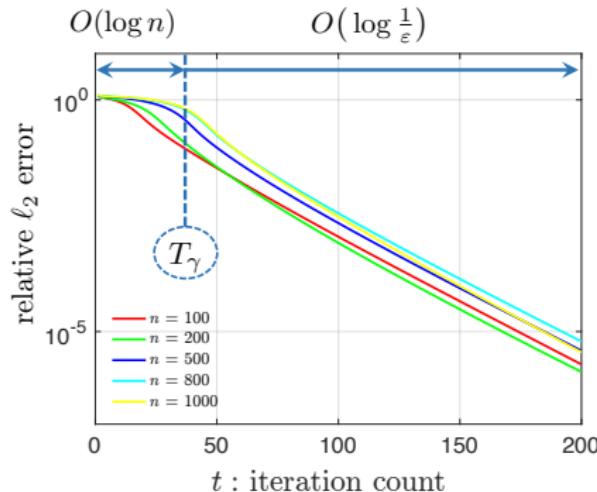
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t-T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *Stage 1:* takes  $O(\log n)$  iterations to reach  $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$  (e.g.  $\gamma = 0.1$ )
- *Stage 2:* linear (geometric) convergence

# Our theory

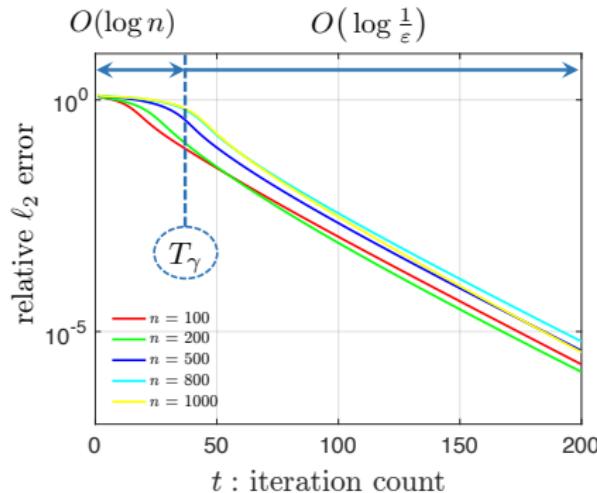
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\star) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\star\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *near-optimal computational cost:*
  - $O(\log n + \log \frac{1}{\varepsilon})$  iterations to yield  $\varepsilon$  accuracy

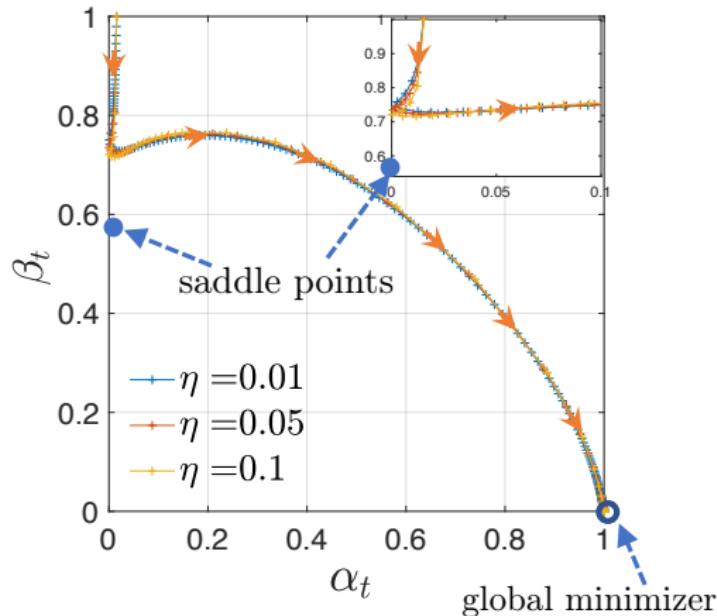
# Our theory

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\star) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\star\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *near-optimal computational cost:*
  - $O(\log n + \log \frac{1}{\varepsilon})$  iterations to yield  $\varepsilon$  accuracy
- *near-optimal sample size:*  $m \gtrsim n \text{poly} \log m$

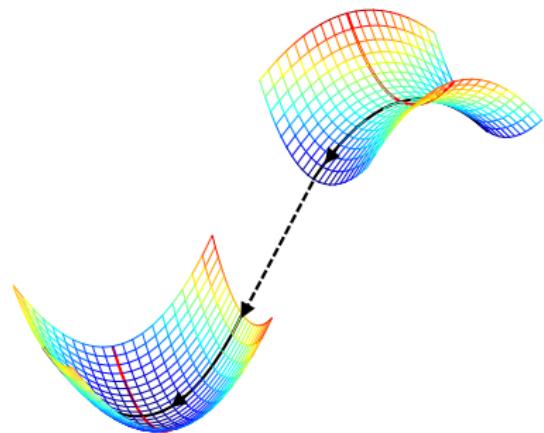
# Automatic saddle avoidance



Randomly initialized GD never hits saddle points!

# Other saddle-escaping schemes based on generic landscape analysis

	iteration complexity
<b>trust-region</b> (Sun et al. '16)	$n^7 + \log \log \frac{1}{\varepsilon}$
<b>perturbed GD</b> (Jin et al. '17)	$n^3 + n \log \frac{1}{\varepsilon}$
<b>perturbed accelerated GD</b> (Jin et al. '17)	$n^{2.5} + \sqrt{n} \log \frac{1}{\varepsilon}$
<b>GD (ours)</b> (Chen et al. '18)	$\log n + \log \frac{1}{\varepsilon}$



Generic optimization theory yields highly suboptimal convergence guarantees

Even **simplest** nonconvex methods  
are remarkably **efficient** under suitable statistical models

smart initialization	extra regularization	sample splitting	saddle escaping
			

1. "Gradient Descent with Random Initialization: ...", Y. Chen, Y. Chi, J. Fan, C. Ma, *Mathematical Programming*, vol. 176, no. 1-2, pp. 5-37, 2019
2. "Implicit regularization in nonconvex statistical estimation: ...", C. Ma, K. Wang, Y. Chi, Y. Chen, accepted to *Foundations of Computational Mathematics*, 2019
3. "Nonconvex optimization meets low-rank matrix factorization: An overview", Y. Chi, Y. Lu, Y. Chen, *IEEE Trans. Signal Processing*, vol. 67, no. 20, pp. 5239-5269, 2019

*Inference and uncertainty quantification for  
noisy matrix completion*

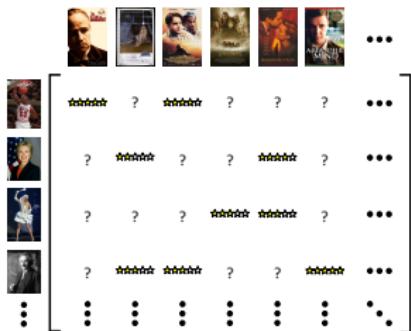
# Low-rank matrix completion

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

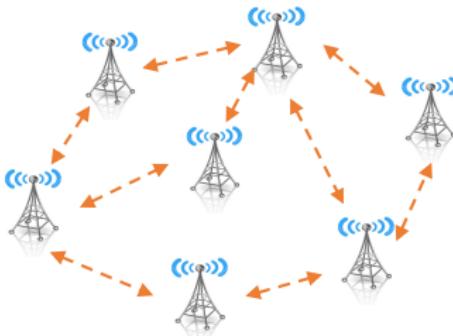


figure credit: E. J. Candès

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



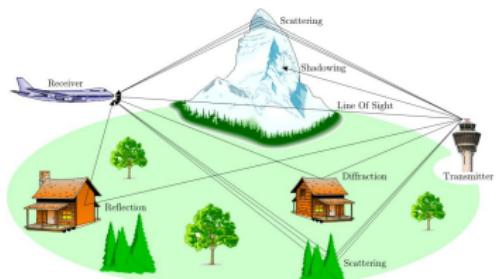
recommendation systems



localization



shape matching



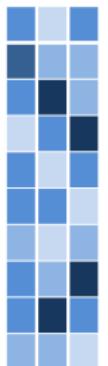
channel estimation

# Noisy low-rank matrix completion

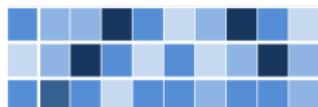
---

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

goal: estimate  $M^*$



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

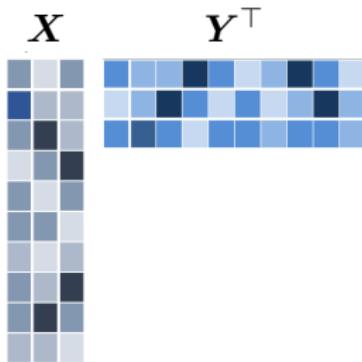


✓	?	?	?	✓	?
?	?	✓	✓	?	?
✓	?	?	✓	?	?
?	?	✓	?	?	✓
✓	?	?	?	?	?
?	✓	?	?	✓	?
?	?	✓	✓	?	?

sampling set  $\Omega$

# Nonconvex matrix completion

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



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

# Nonconvex matrix completion

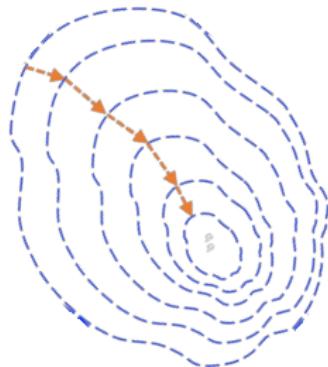
---

- 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
- Charisopoulos, Chen, Davis, Diaz, Ding, Drusvyatskiy '19
- ...

# Nonconvex matrix completion

---

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



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

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

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

# 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

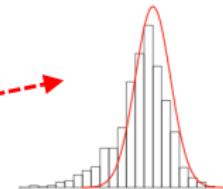
# One step further: reasoning about uncertainty?

	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



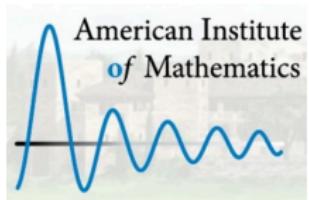
How to assess uncertainty, or “confidence”, of obtained estimates due to imperfect data acquisition?

- noise
- incomplete measurements
- ...

## INFERENCE IN HIGH DIMENSIONAL REGRESSION

organized by

Peter Bühlmann, 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

---

$$\boldsymbol{M}^{\text{ncvx}} \leftarrow \arg \min_{\boldsymbol{X}, \boldsymbol{Y}} \underbrace{f(\boldsymbol{X}, \boldsymbol{Y}; \text{data})}_{\text{empirical loss}} + \text{reg}(\boldsymbol{X}, \boldsymbol{Y})$$

- very challenging to pin down distributions of obtained estimates  
→ due to nonconvexity

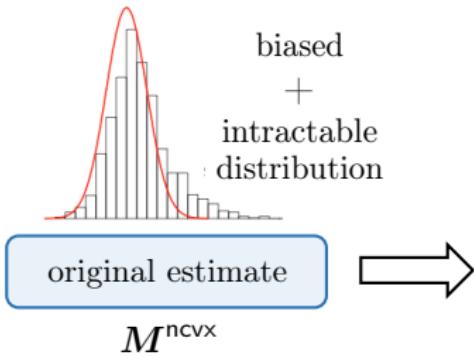
# Challenges

---

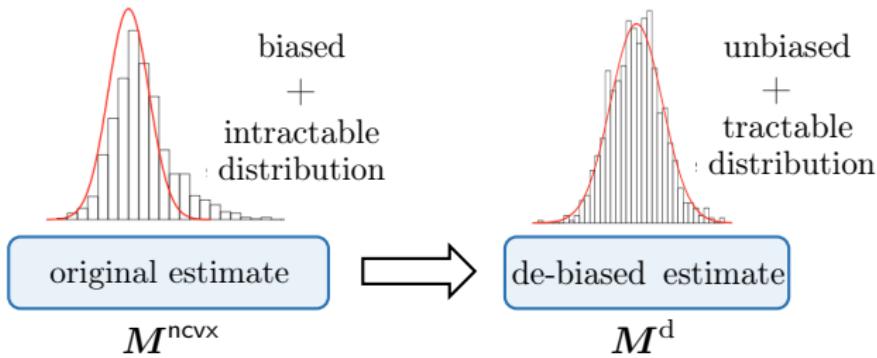
$$\boldsymbol{M}^{\text{ncvx}} \leftarrow \arg \min_{\boldsymbol{X}, \boldsymbol{Y}} \underbrace{f(\boldsymbol{X}, \boldsymbol{Y}; \text{data})}_{\text{empirical loss}} + \text{reg}(\boldsymbol{X}, \boldsymbol{Y})$$

- very challenging to pin down distributions of obtained estimates
  - due to nonconvexity
- existing estimation error bounds are highly sub-optimal
  - overly wide confidence intervals

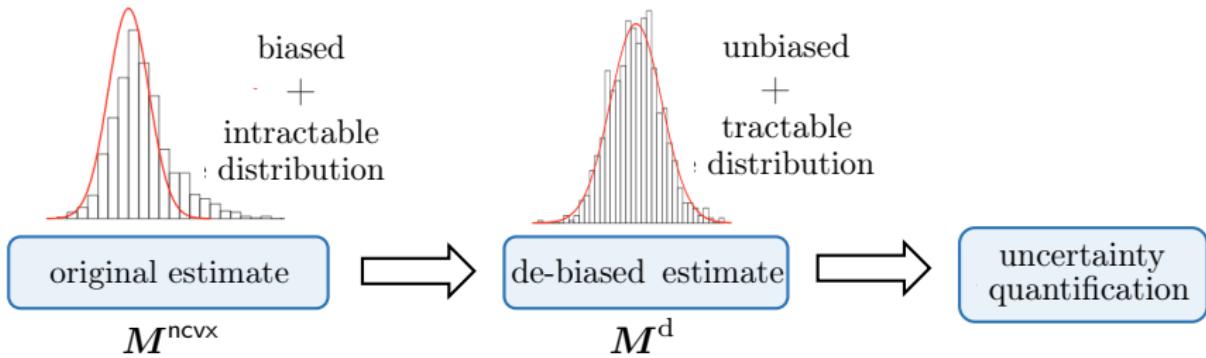
— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



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



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



# Model

---

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

goal: estimate  $\mathbf{M}^*$

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

# De-biasing nonconvex estimate

---

$$M^{\text{ncvx}} \xrightarrow{\text{de-biasing}} M^{\text{ncvx}} + \underbrace{\underbrace{\frac{1}{p} \mathcal{P}_\Omega(M^* + \text{noise} - M^{\text{ncvx}})}_{\text{mean: } \mathcal{I}}}_{(\text{heuristically}) \text{ unbiased estimate of } M^*}$$

# De-biasing nonconvex estimate

---

$$M^{\text{ncvx}} \xrightarrow{\text{de-biasing}} M^{\text{ncvx}} + \underbrace{\frac{1}{p} \mathcal{P}_\Omega(M^* + \text{noise} - M^{\text{ncvx}})}_{\text{mean: } \mathcal{I}} \underbrace{\quad}_{(\text{heuristically}) \text{ unbiased estimate of } M^*}$$

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

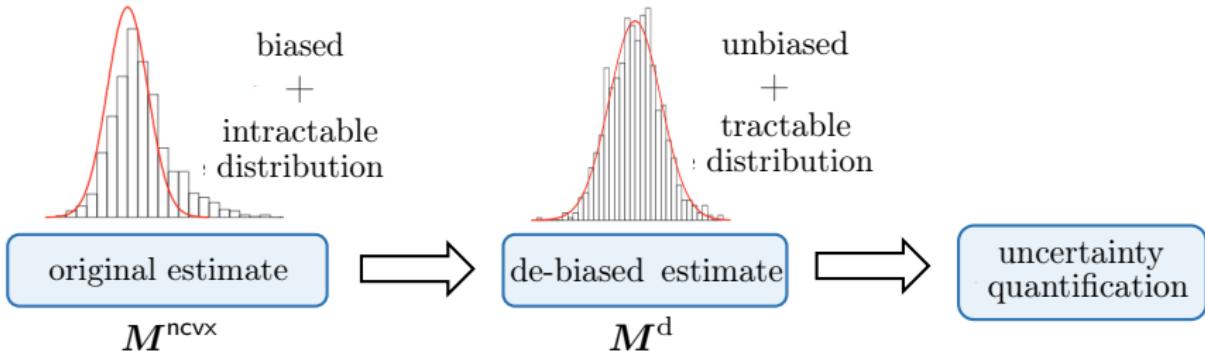
# De-biasing nonconvex estimate

---

$$\mathbf{M}^{\text{ncvx}} \xrightarrow{\text{de-biasing}} \text{proj}_{\text{rank-}r} \left( \mathbf{M}^{\text{ncvx}} + \underbrace{\frac{1}{p} \mathcal{P}_\Omega (\mathbf{M}^* + \text{noise} - \mathbf{M}^{\text{ncvx}})}_{\text{mean: } \mathcal{I}} \right) =: \mathbf{M}^d$$

1 iteration of singular value projection (Jain, Meka, Dhillon '10)

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



# 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

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

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

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

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

$$\mathbf{X}^d \mathbf{R} - \mathbf{X}^* \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,\cdot}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

$$\mathbf{Y}^d \mathbf{R} - \mathbf{Y}^* \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,\cdot}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

# Implications

---

$$\mathbf{X}^d \mathbf{R} - \mathbf{X}^* \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,:}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

$$\mathbf{Y}^d \mathbf{R} - \mathbf{Y}^* \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,:}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

- accurate uncertainty quantification for low-rank factors
  - *asymptotically optimal*

# Implications

---

$$\mathbf{X}^d \mathbf{R} - \mathbf{X}^* \approx \mathbf{Z}^X, \quad \mathbf{Z}_{i,\cdot}^X \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

$$\mathbf{Y}^d \mathbf{R} - \mathbf{Y}^* \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,\cdot}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \text{Cramer-Rao})$$

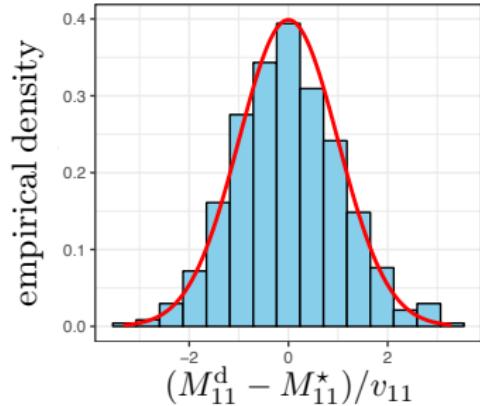
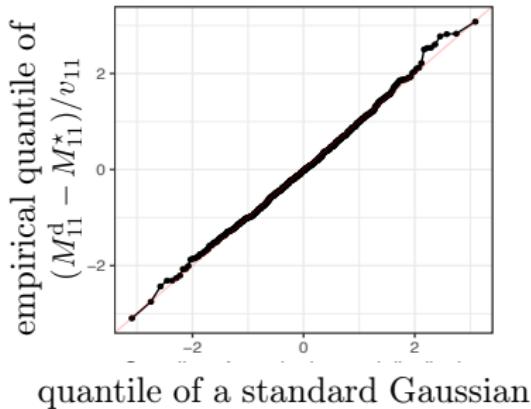
- accurate uncertainty quantification for matrix entries: if  $\|\mathbf{X}_{i,\cdot}^*\|_2 + \|\mathbf{Y}_{j,\cdot}^*\|_2$  is not too small, then

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

— *asymptotically optimal*

# Numerical experiments

---



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

## 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 3 (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)

*Bridging convex and nonconvex optimization in  
noisy matrix completion*

# Convex relaxation for low-rank structure

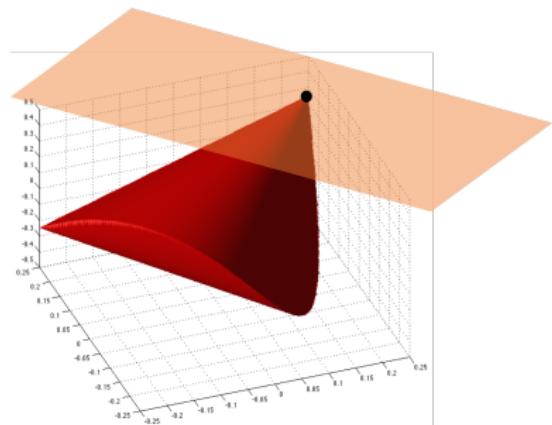
$$\underset{\mathbf{Z}}{\text{minimize}} \quad \|\mathbf{Z}\|_* := \sum_i \sigma_i(\mathbf{Z})$$

subj. to      **noiseless** data constraints



low-rank matrix

*figure credit: Piet Mondrian*



semidefinite relaxation

# Convex relaxation for low-rank structure

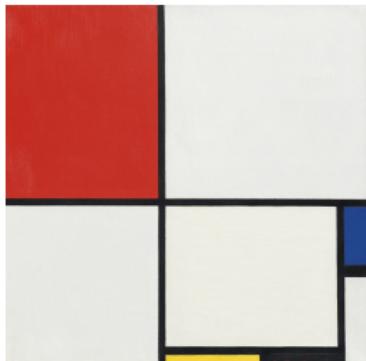
$$\begin{array}{ll} \text{minimize}_{\mathbf{Z}} & \|\mathbf{Z}\|_* := \sum_i \sigma_i(\mathbf{Z}) \\ \text{subj. to} & \text{noiseless data constraints} \end{array}$$

- ✓ 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)

•

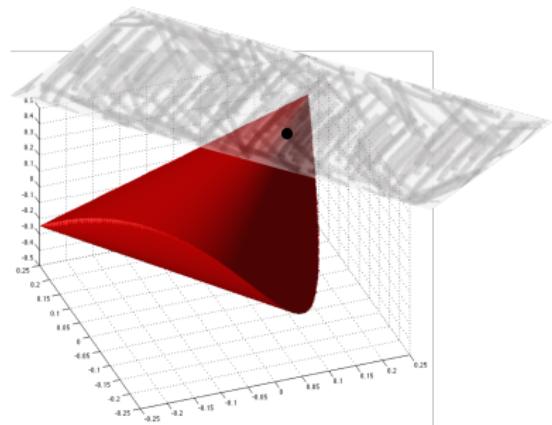
# Stability of convex relaxation against noise

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



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
- ...

# Noisy low-rank matrix completion

---

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

goal: estimate  $M^*$

**convex relaxation:**

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

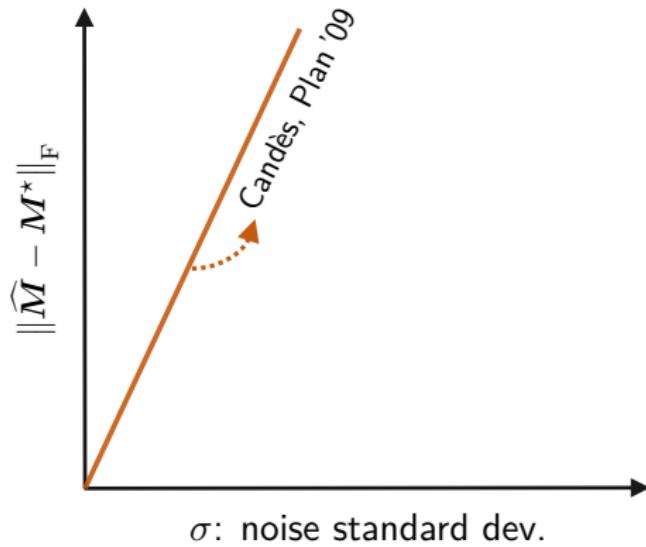
# Prior statistical guarantees for convex relaxation

---

- **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

$\sigma n^{1.5}$

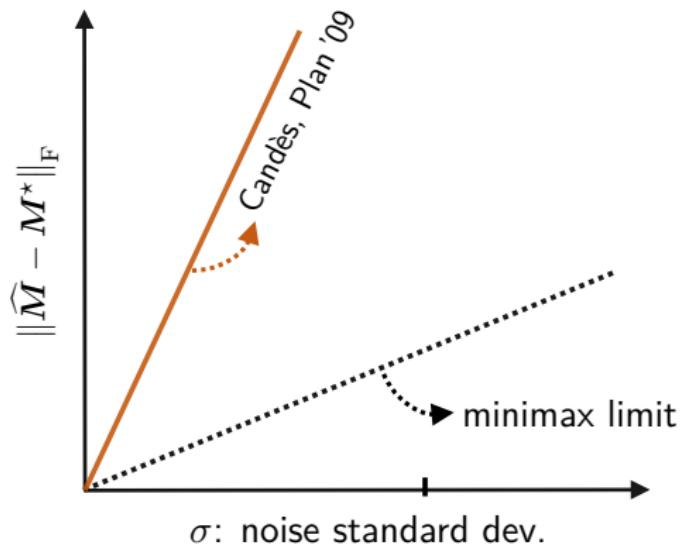


minimax limit

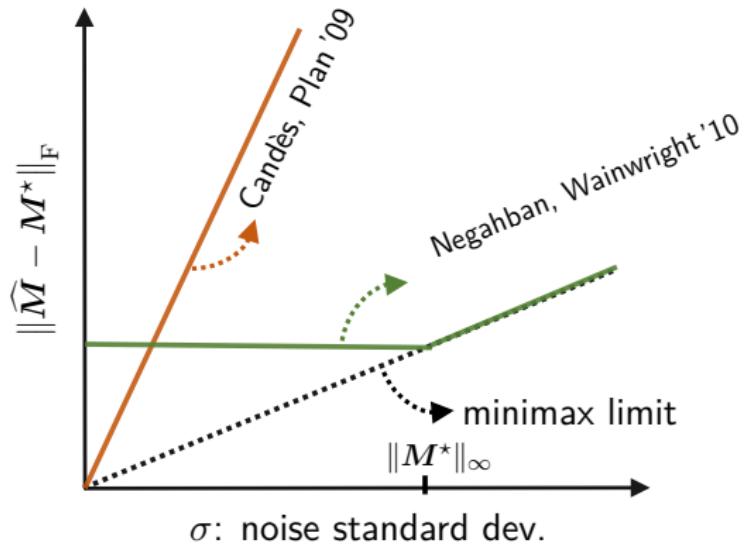
$$\sigma\sqrt{n/p}$$

Candès, Plan '09

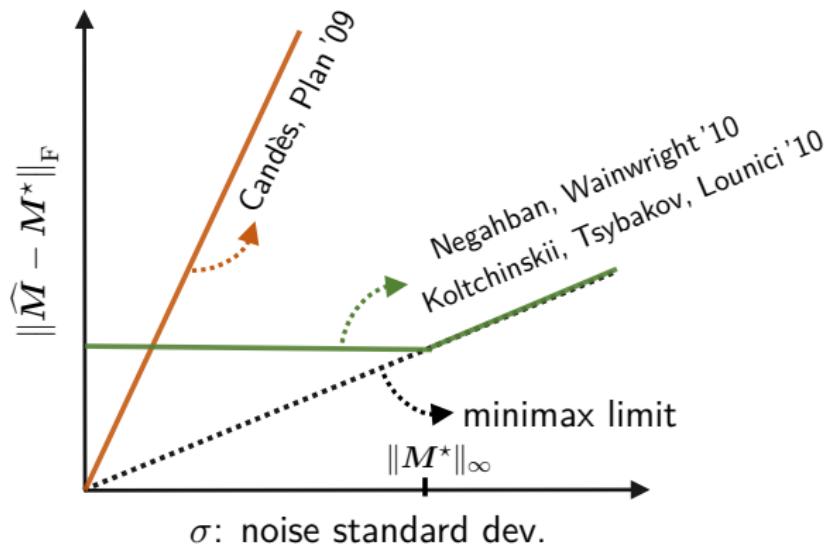
$$\sigma n^{1.5}$$



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

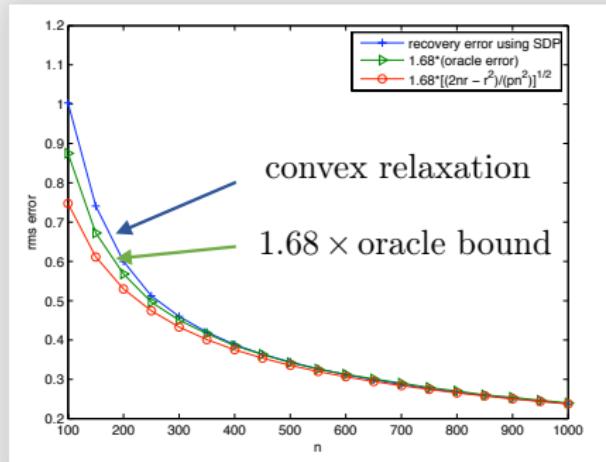


minimax limit	$\sigma\sqrt{n/p}$
Candès, Plan '09	$\sigma n^{1.5}$
Negahban, Wainwright '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$
Koltchinskii, Tsybakov, Lounici '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$



# 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:  $\widehat{M}_{\text{cvx}}$  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 motivating experiment

---

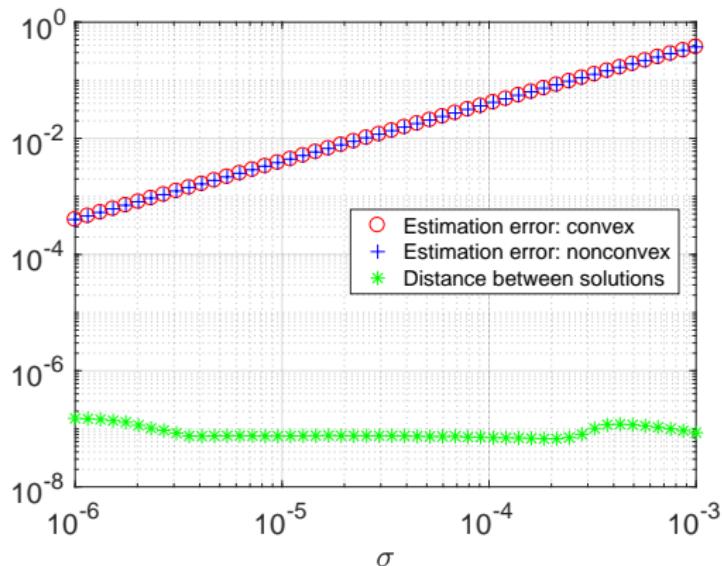
**convex:**  $\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$

**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

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



Convex and nonconvex solutions are exceedingly close!

convex

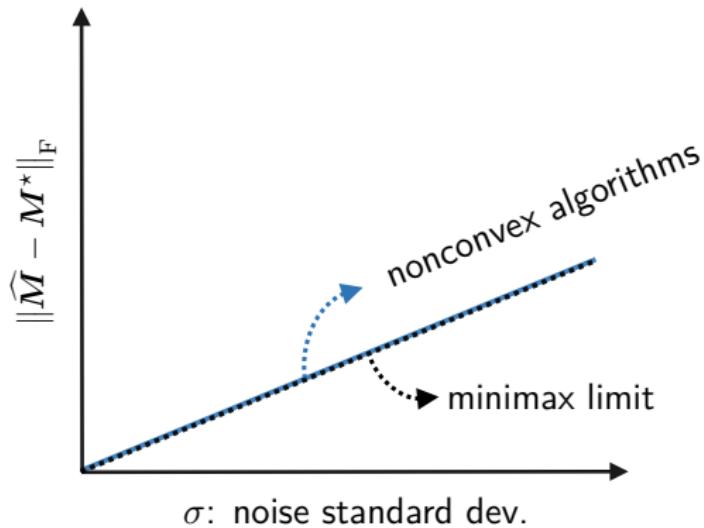


nonconvex



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

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



— Ma, Wang, Chi, Chen '17

## 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  $\sigma^2$
- 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  $\sigma^2$
- 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 4 (Chen, Chi, Fan, Ma, Yan '19)

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

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

## 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  $\sigma^2$
- 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 4 (Chen, Chi, Fan, Ma, Yan '19)

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

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

2.  $\|\widehat{M}_{\text{cvx}} \text{proj}_{\mathcal{M}^*} \widehat{M}_{\text{cvx}}\|_F \leq \frac{1}{n^5} \cdot \sigma \sqrt{\frac{n}{p}}$

## 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  $\sigma^2$
- 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 4 (Chen, Chi, Fan, Ma, Yan '19)

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

1.  $\widehat{M}_{\text{cvx}}$  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  $\sigma^2$
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

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

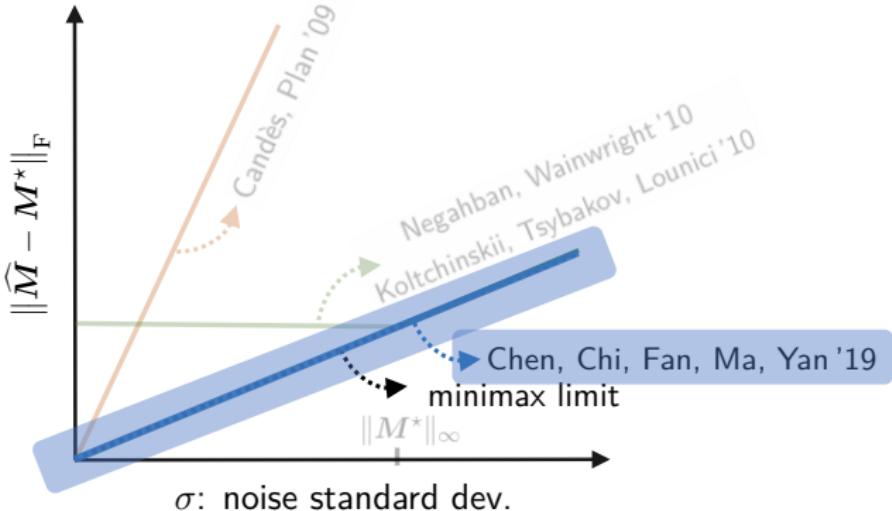
### Theorem 4 (Chen, Chi, Fan, Ma, Yan '19)

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

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

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

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



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

convex



nonconvex

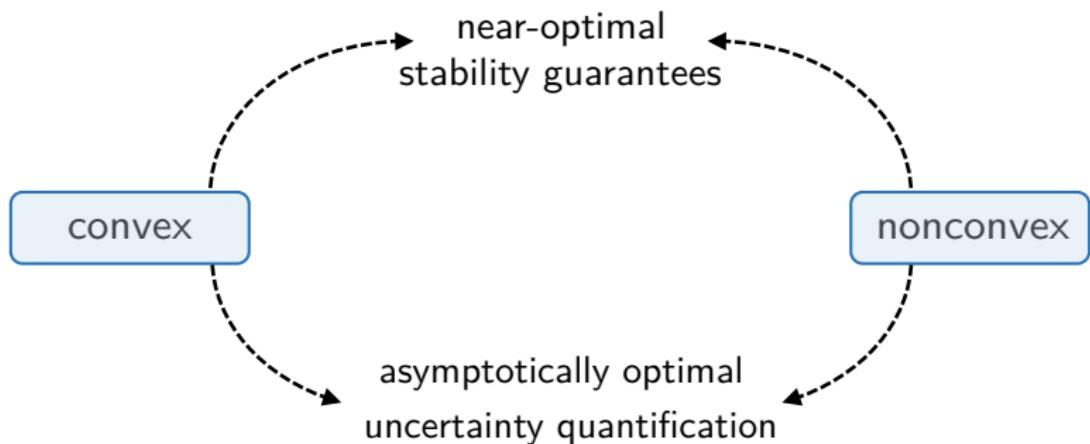


inference  $(\text{convex})$



inference  $(\text{nonconvex})$

Same inference procedures work for both cvx & noncvx estimates!



1. "Inference and uncertainty quantification for noisy matrix completion", accepted to *Proceedings of the National Academy of Sciences (PNAS)*, Y. Chen, J. Fan, C. Ma, Y. Yan, 2019
2. "Noisy matrix completion: understanding statistical guarantees for convex relaxation via nonconvex optimization", Y. Chen, Y. Chi, J. Fan, C. Ma, Y. Yan, 2019