

# Inference and Uncertainty Quantification for Low-Rank Models



Yuling Yan

Princeton University

# A ubiquitous low-complexity model

---



*Composition C by Piet Mondrian*

reconstructing **low-rank structure**  
from imperfect measurements

# A ubiquitous low-complexity model

---



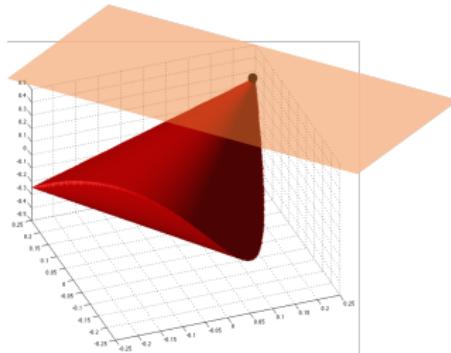
*Composition C* by Piet Mondrian

reconstructing **low-rank structure**  
from imperfect measurements

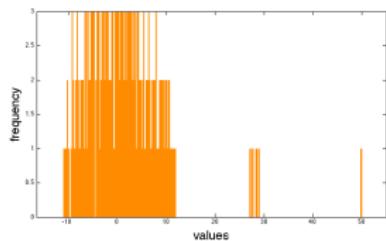
- matrix completion
- phase retrieval
- blind deconvolution
- tensor completion
- localization
- PCA / factor models
- community recovery
- joint shape mapping
- linear neural networks
- ...

# Various estimation schemes have been proposed

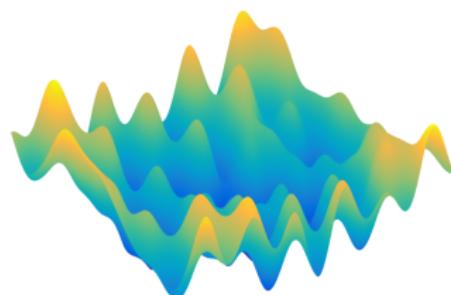
---



convex relaxation



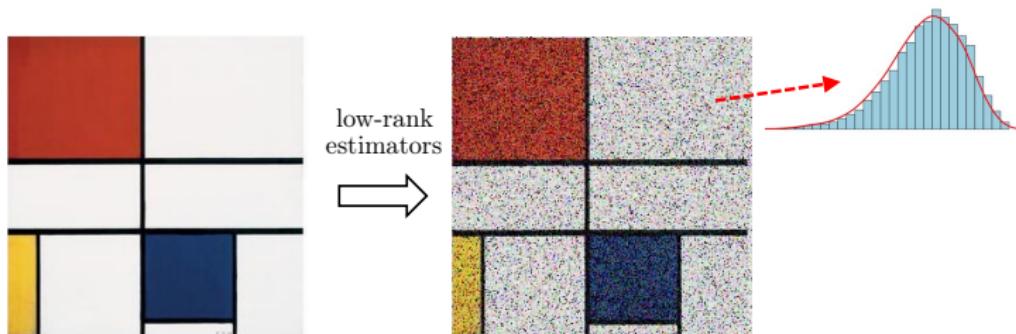
spectral methods



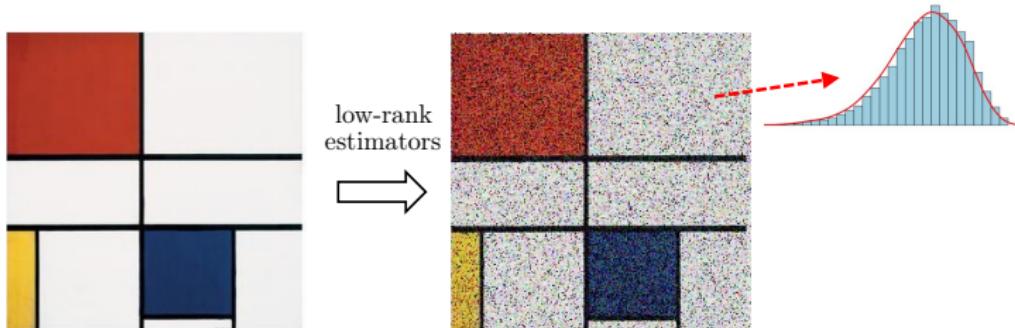
nonconvex optimization

# One step further: reasoning about uncertainty?

---

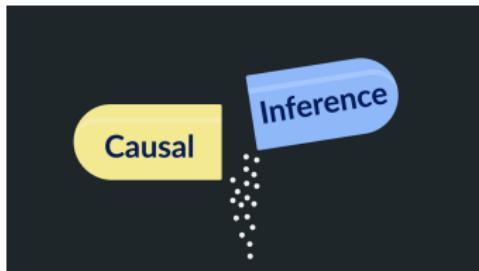


# One step further: reasoning about uncertainty?



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

- noise
- missing data
- outliers
- ...



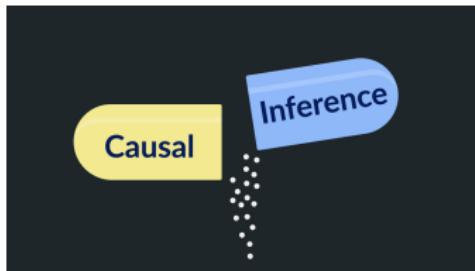
causal inference with panel data



structure from motion



anomaly detection



causal inference with panel data

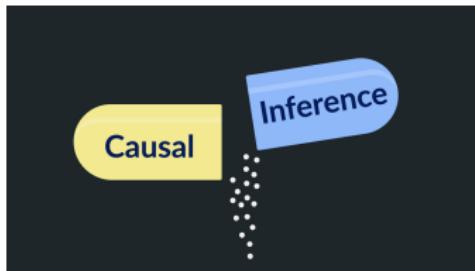


structure from motion



anomaly detection

"Matrix completion methods for causal panel data models," S. Athey, M. Bayati, N. Doudchenko, G. Imbens, K. Khosravi *Journal of the American Statistical Association*



causal inference with panel data

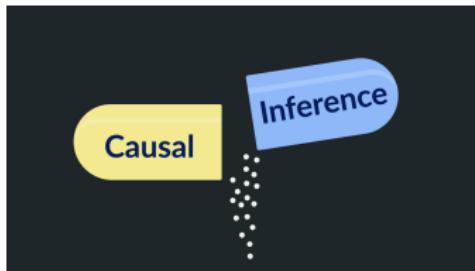


structure from motion



anomaly detection

"A Closed-Form Uncertainty Propagation in Non-Rigid Structure From Motion," J. Song,  
M. Patel, A. Jasour, M. Ghaffari, *IEEE Robotics and Automation Letters*



causal inference with panel data



structure from motion



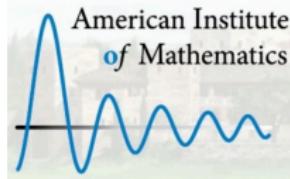
anomaly detection

"Near-optimal entrywise anomaly detection for low-rank matrices with sub-exponential noise," V. Farias, A. Li, A. Jasour, T. Peng, *ICML*

## INFERENCE IN HIGH DIMENSIONAL REGRESSION

organized by

Peter Buehlmann, Andrea Montanari, and Jonathan Taylor



The open problems discussion was also very productive, and led to formulating a selection of special topics addressed in the working groups. These were

•  
•  
•

- (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.

## This talk: two recent examples

---

1. Inference for noisy matrix completion
2. Inference for heteroskedastic PCA with missing data

## This talk: two recent examples

---

1. Inference for noisy matrix completion
  - convex optimization, nonconvex optimization
2. Inference for heteroskedastic PCA with missing data
  - spectral methods

## *Part 1: Inference for noisy matrix completion*



Yuxin Chen  
UPenn Wharton



Jianqing Fan  
Princeton ORFE



Cong Ma  
U Chicago Stats

# Noisy low-rank matrix completion

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

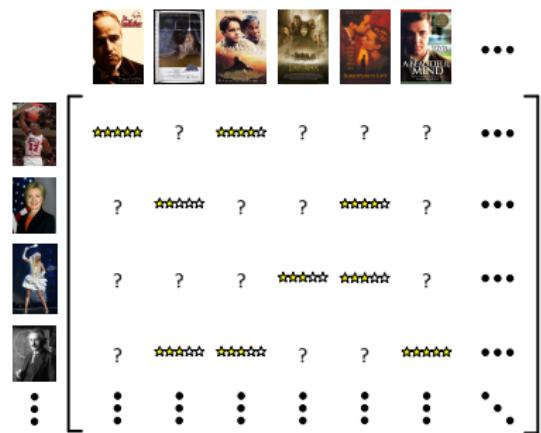


figure credit: E. Candès, M. Soltanolkotabi

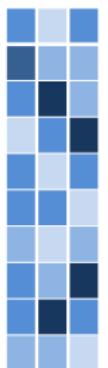
Given partial noisy entries of a low-rank matrix  $M^*$ , fill in missing entries

# 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  $M^*$

A natural least squares formulation:

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

# Noisy low-rank matrix completion

---

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

goal: estimate  $M^*$

A widely used convex relaxation method:

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

# Challenges

---

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{\sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

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

# Challenges

---

$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{\sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

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

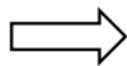
# Challenges

---

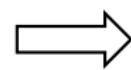
$$\boldsymbol{M}^{\text{cvx}} \triangleq \arg \min_{\boldsymbol{Z}} \underbrace{\sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2}_{\text{empirical loss}} + \lambda \|\boldsymbol{Z}\|_*$$

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

sharpened  
estimation bounds



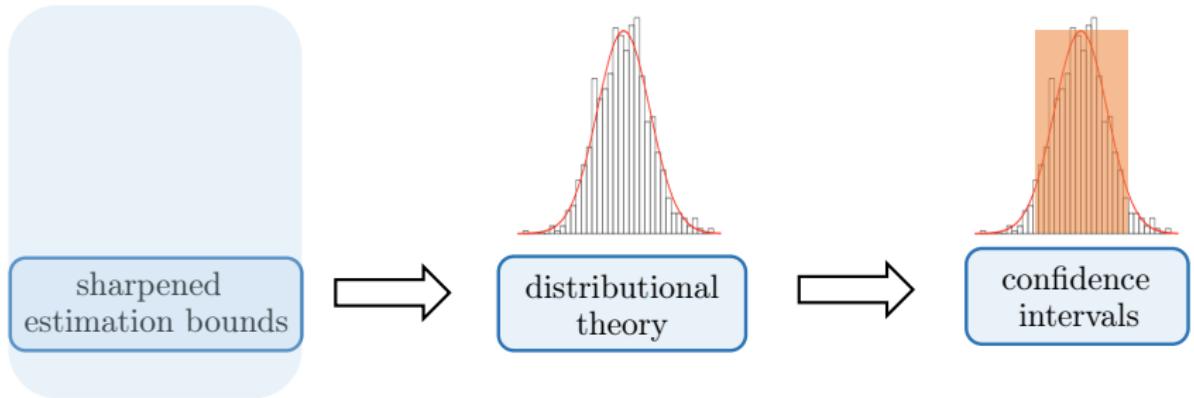
distributional  
theory



confidence  
intervals



## *Step 1: sharpening estimation guarantees*



# Prior statistical guarantees for convex relaxation

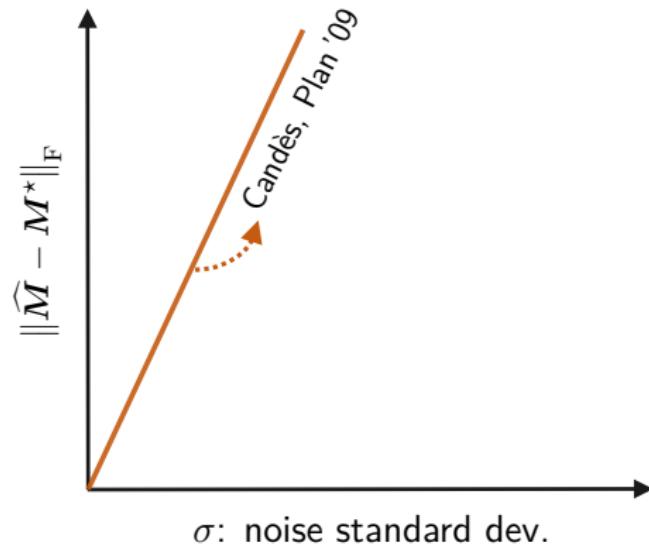
---

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

---

Candès, Plan '09

$\sigma n^{1.5}$

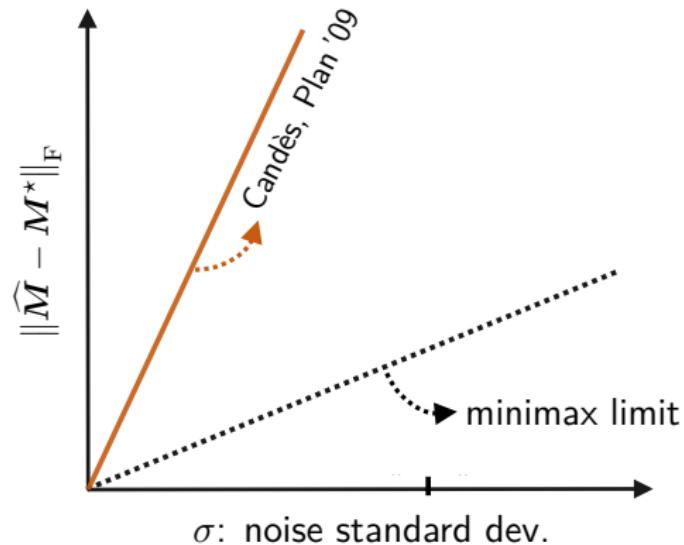


minimax limit

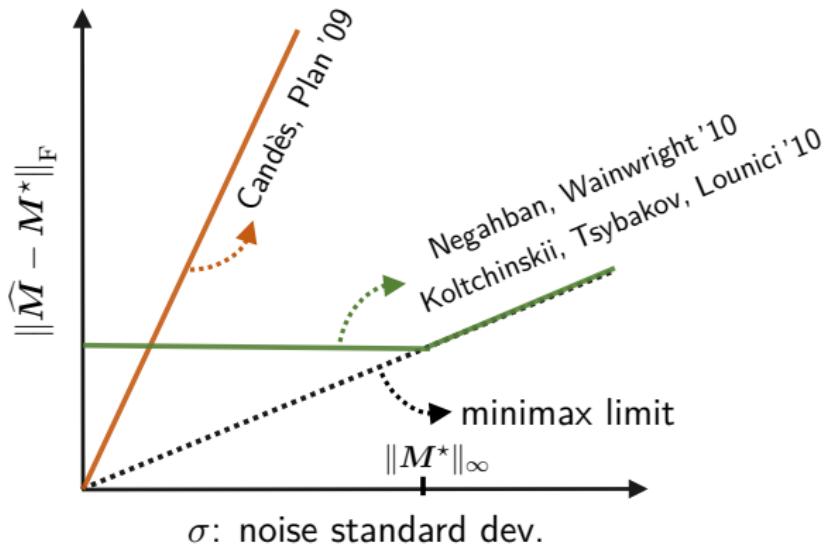
$\sigma\sqrt{n/p}$

Candès, Plan '09

$\sigma n^{1.5}$

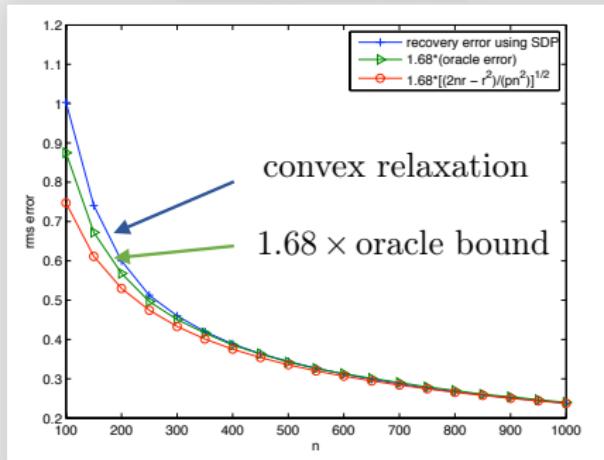


minimax limit	$\sigma\sqrt{n/p}$
Candès, Plan '09	$\sigma n^{1.5}$
Negahban, Wainwright '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$
Koltchinskii, Tsybakov, Lounici '10	$\max\{\sigma, \ \mathbf{M}^*\ _\infty\} \sqrt{n/p}$



# Matrix Completion with Noise

Emmanuel J. Candès and Yaniv Plan



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

dual certification (golfing scheme)



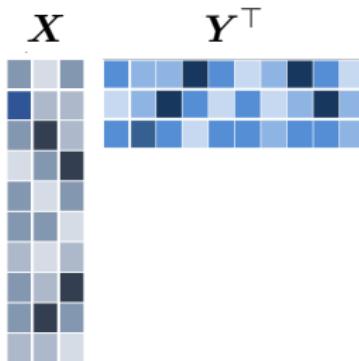
dual certification (golfing scheme)



nonconvex optimization

# A detour: nonconvex optimization

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

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


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

# A detour: nonconvex optimization

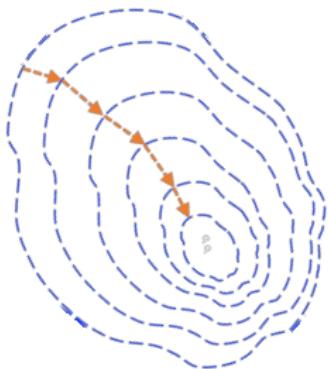
---

- 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

---

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



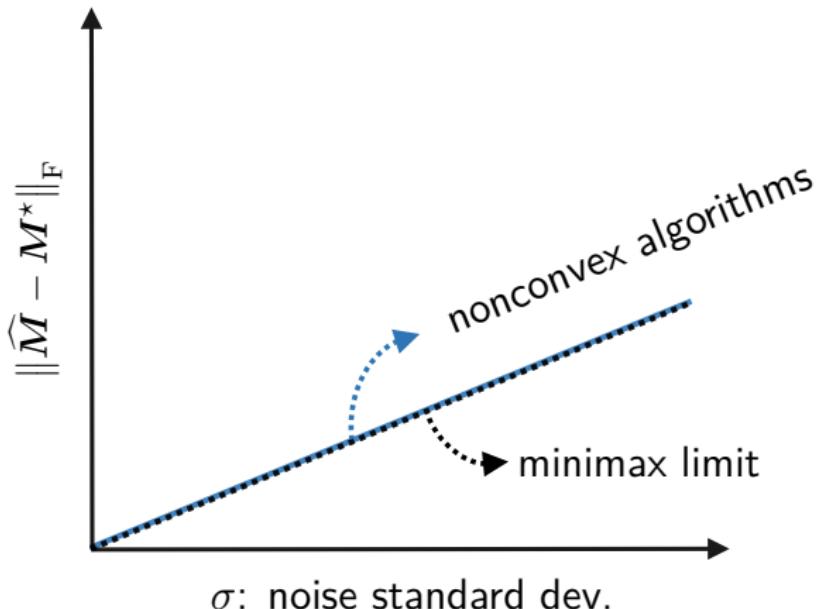
- **suitable initialization:**  $(\mathbf{X}^0, \mathbf{Y}^0)$
- **gradient descent:** for  $t = 0, 1, \dots$

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

# A detour: nonconvex optimization

—Ma, Wang, Chi, Chen '20

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



# A motivating experiment

---

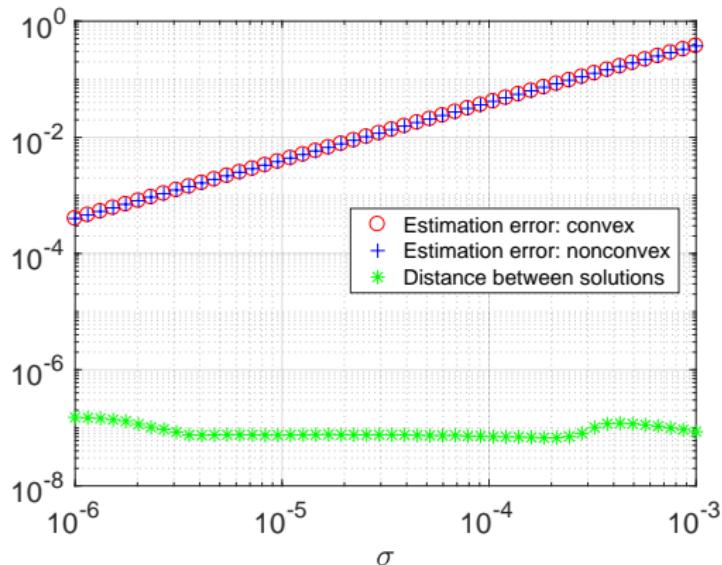
**convex:**  $\underset{\mathbf{Z} \in \mathbb{R}^{n \times n}}{\text{minimize}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$

**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

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



Convex and nonconvex solutions are exceedingly close!

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

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

convex



nonconvex



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

# Sharpened estimation guarantees

---

**Assumptions** (omitting log factors)

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{1}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  satisfying  $\sigma \lesssim \sqrt{np} \|M^*\|_\infty$
- true matrix  $M^* \in \mathbb{R}^{n \times n}$ :  $r = O(1)$ , incoherent, well-conditioned

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

# Sharpened estimation guarantees

**Assumptions** (omitting log factors)

- **random sampling:** each  $(i, j) \in \Omega$  with prob.  $p \gtrsim \frac{1}{n}$
- **random noise:** i.i.d. sub-Gaussian with variance  $\sigma^2$  satisfying  $\sigma \lesssim \sqrt{np} \|M^*\|_\infty$
- 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

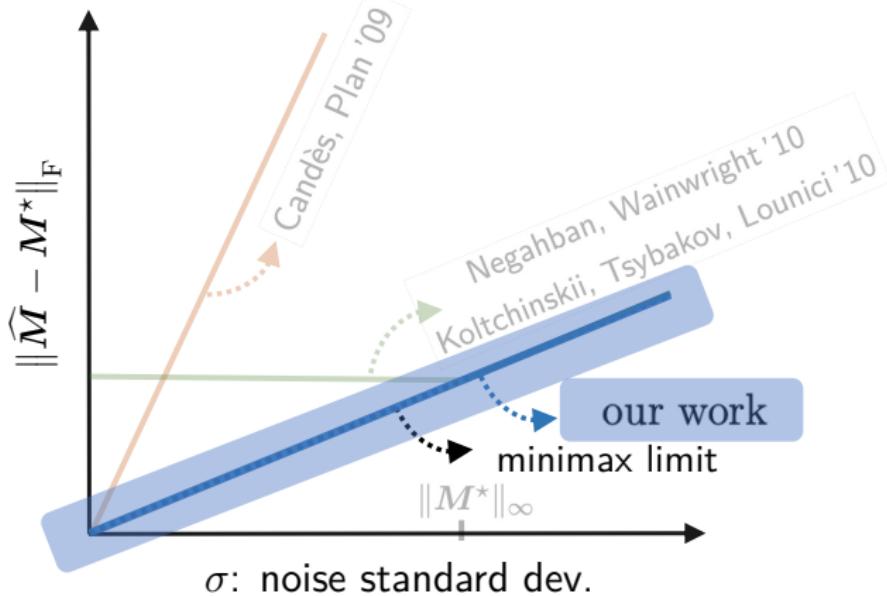
With high prob., any minimizer  $M^{\text{cvx}}$  of convex program obeys

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

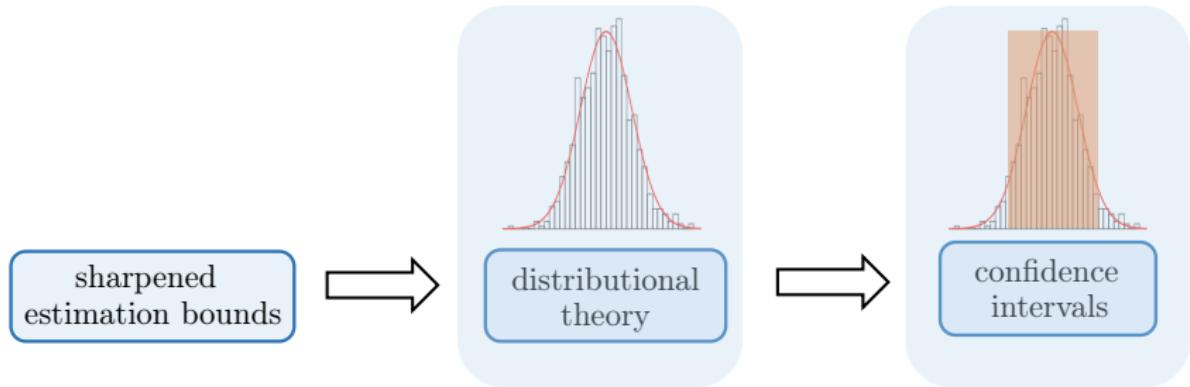
2.

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

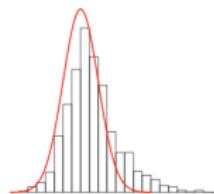
$$\|\mathbf{M}^{\text{cvx}} - \mathbf{M}^*\|_{\text{F}} \lesssim \sigma \sqrt{\frac{n}{p}} : \text{ minimax optimal}$$



## *Step 2: from estimation to inference ...*



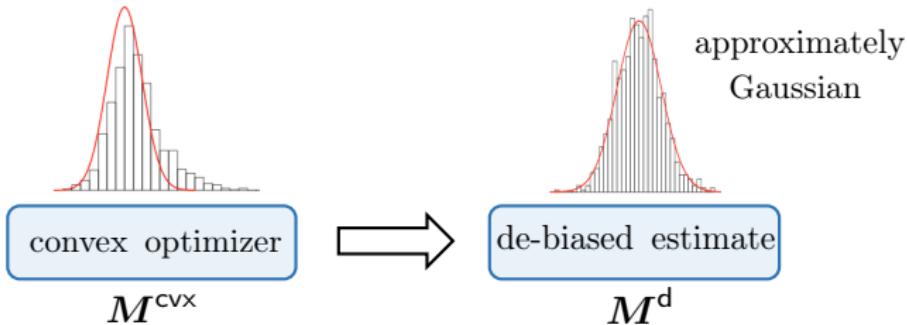
— inspired by Zhang, Zhang '11, van de Geer et al. '13, Javanmard, Montanari '13



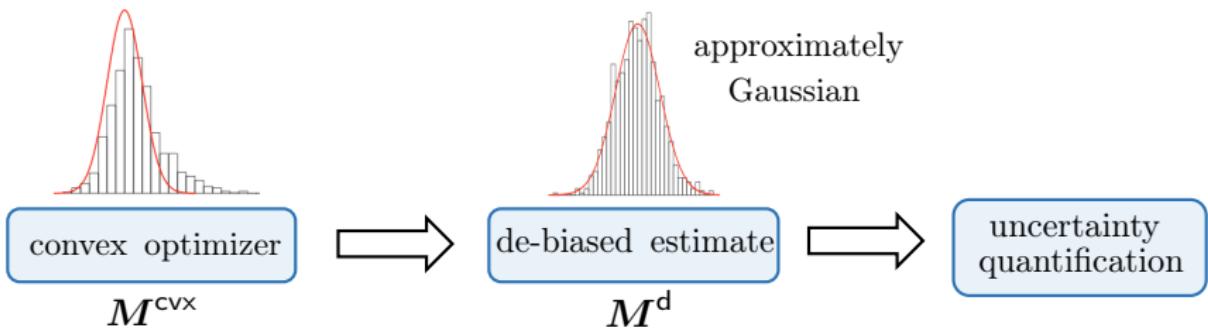
convex optimizer

$M^{\text{cvx}}$

— 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

---

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

# De-biasing convex estimate

---

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

# De-biasing convex estimate

---

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

# De-biasing convex estimate

---

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

## De-biasing convex estimate

---

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

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

## De-biasing convex estimate

---

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

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

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \boldsymbol{\Sigma}^* \mathbf{V}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}\left(0, \frac{\sigma^2}{p}(\|\mathbf{U}_{i,\cdot}^*\|_2^2 + \|\mathbf{V}_{j,\cdot}^*\|_2^2)\right) + \text{negligible term}$$

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \boldsymbol{\Sigma}^* \mathbf{V}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}\left(0, \frac{\sigma^2}{p}(\|\mathbf{U}_{i,\cdot}^*\|_2^2 + \|\mathbf{V}_{j,\cdot}^*\|_2^2)\right) + \text{negligible term}$$

- key decomposition:

$$M_{i,j}^d - M_{i,j}^* \approx \underbrace{\frac{\sigma}{\sqrt{p}}(\mathbf{V}_{j,\cdot}^* \mathbf{z} + \mathbf{U}_{i,\cdot}^* \mathbf{w})}_{\text{first-order Gaussian term}} + \underbrace{\frac{\sigma^2}{p} \mathbf{z}^\top \boldsymbol{\Sigma}^* \mathbf{w}}_{\text{second-order term}}$$

where  $\mathbf{z}, \mathbf{w} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_r)$

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \Sigma^* \mathbf{V}^{*\top} \xrightarrow{\text{bal. decom.}} \mathbf{X}^* \mathbf{Y}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}\left(0, \frac{\sigma^2}{p}(\|\mathbf{U}_{i,\cdot}^*\|_2^2 + \|\mathbf{V}_{j,\cdot}^*\|_2^2)\right) + \text{negligible term}$$

- oracle information:  $\{\mathbf{X}_{k,\cdot}^*\}_{k: k \neq i}$ ,  $\{\mathbf{Y}_{k,\cdot}^*\}_{k: k \neq j}$

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \Sigma^* \mathbf{V}^{*\top} \xrightarrow{\text{bal. decom.}} \mathbf{X}^* \mathbf{Y}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}\left(0, \frac{\sigma^2}{p}(\|\mathbf{U}_{i,\cdot}^*\|_2^2 + \|\mathbf{V}_{j,\cdot}^*\|_2^2)\right) + \text{negligible term}$$

- oracle information:  $\{\mathbf{X}_{k,\cdot}^*\}_{k: k \neq i}$ ,  $\{\mathbf{Y}_{k,\cdot}^*\}_{k: k \neq j}$
- parameters:  $\mathbf{X}_{i,\cdot}^*$ ,  $\mathbf{Y}_{j,\cdot}^*$

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \Sigma^* \mathbf{V}^{*\top} \xrightarrow{\text{bal. decom.}} \mathbf{X}^* \mathbf{Y}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

$$M_{i,j}^d - M_{i,j}^* \sim \mathcal{N}\left(0, \frac{\sigma^2}{p}(\|\mathbf{U}_{i,\cdot}^*\|_2^2 + \|\mathbf{V}_{j,\cdot}^*\|_2^2)\right) + \text{negligible term}$$

- oracle information:  $\{\mathbf{X}_{k,\cdot}^*\}_{k: k \neq i}$ ,  $\{\mathbf{Y}_{k,\cdot}^*\}_{k: k \neq j}$
- parameters:  $\mathbf{X}_{i,\cdot}^*$ ,  $\mathbf{Y}_{j,\cdot}^*$
- estimate  $M_{i,j}^* = \mathbf{X}_{i,\cdot}^* (\mathbf{Y}_{j,\cdot}^*)^\top$  based on a linear model

$$M_{i,k} = \mathbf{X}_{i,\cdot}^* (\mathbf{Y}_{k,\cdot}^*)^\top + E_{i,k}, \quad (i, k) \in \Omega$$

$$M_{k,j} = \mathbf{X}_{k,\cdot}^* (\mathbf{Y}_{j,\cdot}^*)^\top + E_{k,j}, \quad (k, j) \in \Omega$$

# Distributional theory

$$\mathbf{M}^* \xrightarrow{\text{rank-}r \text{ svd}} \mathbf{U}^* \boldsymbol{\Sigma}^* \mathbf{V}^{*\top} \xrightarrow{\text{bal. decom.}} \mathbf{X}^* \mathbf{Y}^{*\top}$$

## Theorem 2

Consider any  $(i, j)$  s.t.  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{V}_{j,\cdot}^*\|_2$  is not too small. Then

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

- oracle information:  $\{\mathbf{X}_{k,\cdot}^*\}_{k: k \neq i}$ ,  $\{\mathbf{Y}_{k,\cdot}^*\}_{k: k \neq j}$
- parameters:  $\mathbf{X}_{i,\cdot}^*$ ,  $\mathbf{Y}_{j,\cdot}^*$
- estimate  $M_{i,j}^* = \mathbf{X}_{i,\cdot}^* (\mathbf{Y}_{j,\cdot}^*)^\top$  based on a linear model

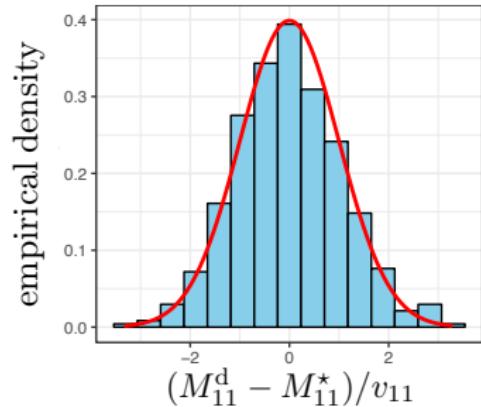
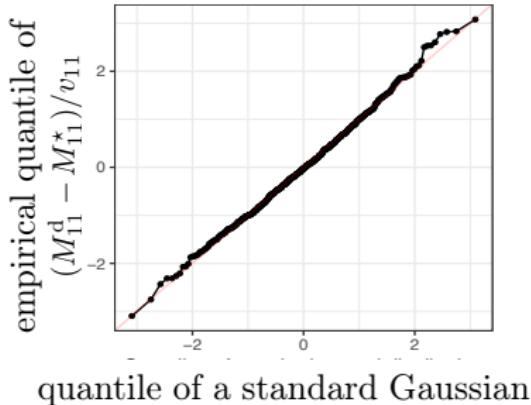
$$M_{i,k} = \mathbf{X}_{i,\cdot}^* (\mathbf{Y}_{k,\cdot}^*)^\top + E_{i,k}, \quad (i, k) \in \Omega$$

$$M_{k,j} = \mathbf{X}_{k,\cdot}^* (\mathbf{Y}_{j,\cdot}^*)^\top + E_{k,j}, \quad (k, j) \in \Omega$$

— *asymptotically optimal!*

# Numerical experiments

---



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

## A bit of analysis

---

$$\boldsymbol{M}^d = \text{proj}_{\text{rank-}r} \left( \boldsymbol{M}^{\text{cvx}} + \frac{1}{p} \mathcal{P}_\Omega(\boldsymbol{M}^\star + \boldsymbol{E} - \boldsymbol{M}^{\text{cvx}}) \right)$$

- challenging to characterize the distribution of  $\boldsymbol{M}^d$

## A bit of analysis

---

$$\boldsymbol{M}^d = \text{proj}_{\text{rank-}r} \left( \boldsymbol{M}^{\text{cvx}} + \frac{1}{p} \mathcal{P}_\Omega(\boldsymbol{M}^\star + \boldsymbol{E} - \boldsymbol{M}^{\text{cvx}}) \right)$$

- challenging to characterize the distribution of  $\boldsymbol{M}^d$
- solution: resort to nonconvex optimization



nonconvex optimization

## A bit of analysis

---

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

## A bit of analysis

---

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

⇓

nonconvex solution  $(\mathbf{X}, \mathbf{Y}) \implies \mathbf{M}^{\text{cvx}} \approx \mathbf{X} \mathbf{Y}^\top$

## A bit of analysis

---

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

⇓

$$\text{nonconvex solution } (\mathbf{X}, \mathbf{Y}) \implies \mathbf{M}^{\text{cvx}} \approx \mathbf{X} \mathbf{Y}^\top$$

⇓

$$\mathbf{X}^d = \mathbf{X} \left[ \mathbf{I}_r + \frac{\lambda}{p} (\mathbf{X} \mathbf{X}^\top)^{-1} \right]^{1/2}, \quad \mathbf{Y}^d = \mathbf{Y} \left[ \mathbf{I}_r + \frac{\lambda}{p} (\mathbf{Y} \mathbf{Y}^\top)^{-1} \right]^{1/2}$$

# A bit of analysis

---

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

↓

$$\text{nonconvex solution } (\mathbf{X}, \mathbf{Y}) \implies \mathbf{M}^{\text{cvx}} \approx \mathbf{X} \mathbf{Y}^\top$$

↓

$$\mathbf{X}^d = \mathbf{X} \left[ \mathbf{I}_r + \frac{\lambda}{p} (\mathbf{X} \mathbf{X}^\top)^{-1} \right]^{1/2}, \quad \mathbf{Y}^d = \mathbf{Y} \left[ \mathbf{I}_r + \frac{\lambda}{p} (\mathbf{Y} \mathbf{Y}^\top)^{-1} \right]^{1/2}$$

↓

$$\mathbf{M}^d \approx \mathbf{X}^d \mathbf{Y}^{d\top}$$

## A bit of analysis

---

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

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

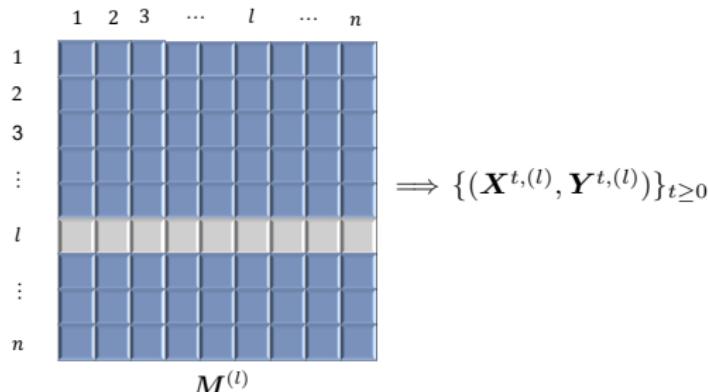
$$\mathbf{Y}^d \mathbf{H} - \mathbf{Y}^* \approx \mathbf{Z}^Y, \quad \mathbf{Z}_{i,:}^Y \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \frac{\sigma^2}{p} (\mathbf{X}^{*\top} \mathbf{X}^*)^{-1})$$

# A bit of analysis

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

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

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



Leave-one-out analysis: leave out a small amount of information from data and run GD...

convex



nonconvex



inference  $(\text{convex})$

$\approx$

inference  $(\text{nonconvex})$

## 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

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

# Back to estimation: de-biased estimator is optimal

Distributional theory in turn allows us to track estimation accuracy

## Theorem 3

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

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

*Part 2: Inference for heteroskedastic PCA with missing data*

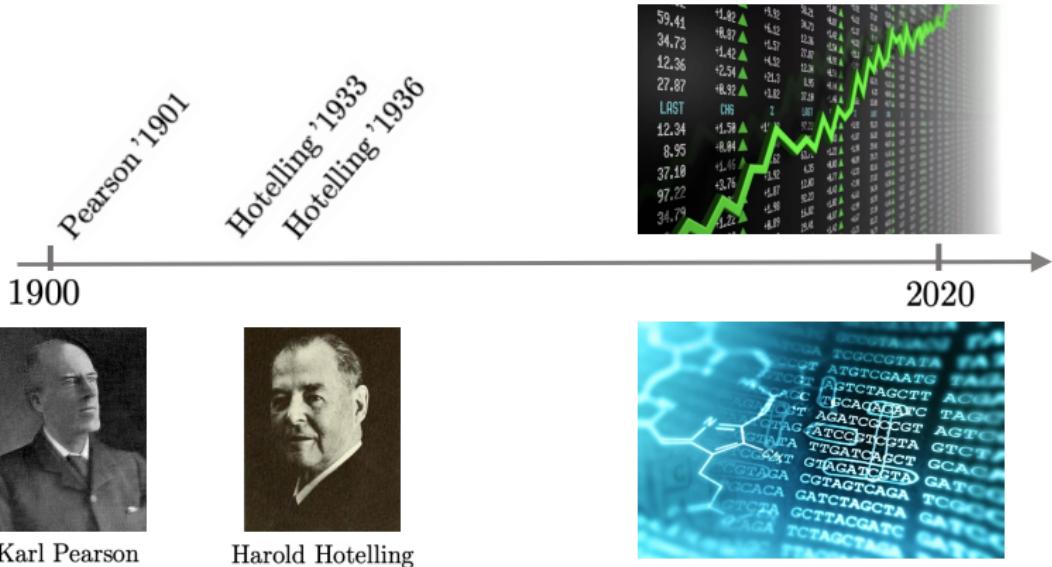


Yuxin Chen  
UPenn Wharton



Jianqing Fan  
Princeton ORFE

# Principal component analysis



# Principal component analysis

---

A spiked covariance model:

$$\text{data} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \underbrace{\mathbf{S}^*}_{\text{low-rank}} + \sigma^2 \mathbf{I}_d)$$

# Principal component analysis

---

A spiked covariance model:

$$\text{data} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \underbrace{\boldsymbol{S}^*}_{\text{low-rank}} + \sigma^2 \mathbf{I}_d)$$

Goal: estimate

- spiked covariance matrix  $\boldsymbol{S}^*$
- principal subspace  $\boldsymbol{U}^*$  — the eigenspace of  $\boldsymbol{S}^*$

# Principal component analysis

---

A spiked covariance model:

$$\text{data} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \underbrace{\boldsymbol{S}^*}_{\text{low-rank}} + \sigma^2 \mathbf{I}_d)$$

Goal: estimate

- spiked covariance matrix  $\boldsymbol{S}^*$
- principal subspace  $\boldsymbol{U}^*$  — the eigenspace of  $\boldsymbol{S}^*$

- Johnstone '03
- Baik, Arous, Péché '05
- Paul '07
- Nadler '08
- Donoho, Gavish, Johnstone '18
- Bao, Ding, Wang, Wang '22
- Zhang, Cai, Wu '22
- ...

# Principal component analysis

---

A spiked covariance model:

$$\text{data} \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \underbrace{\boldsymbol{S}^*}_{\text{low-rank}} + \sigma^2 \mathbf{I}_d)$$

Goal: estimate

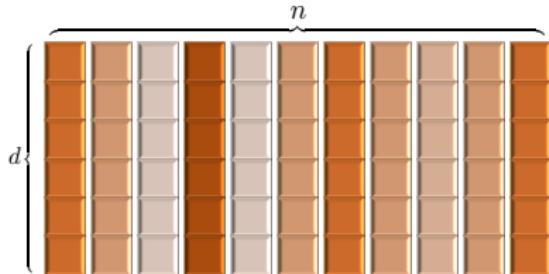
- spiked covariance matrix  $\boldsymbol{S}^*$
- principal subspace  $\boldsymbol{U}^*$  — the eigenspace of  $\boldsymbol{S}^*$

- Johnstone '03
- Baik, Arous, Péché '05
- Paul '07
- Nadler '08
- Donoho, Gavish, Johnstone '18
- Bao, Ding, Wang, Wang '22
- Zhang, Cai, Wu '22
- ...

Today's talk: inference for PCA in high dimension

# Principal component analysis

---



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$

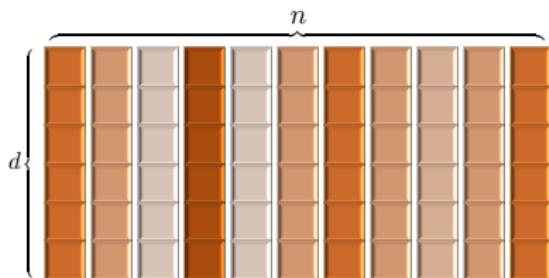
- Ground-truth data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

$$\text{where } \mathbf{S}^* = \mathbf{U}^* \boldsymbol{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$$

# Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbf{U}^*$  ( $r$ -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$

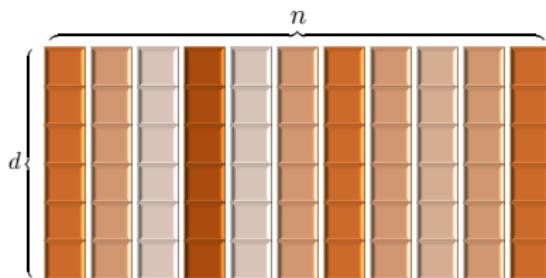
- Ground-truth data

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

where  $\mathbf{S}^* = \mathbf{U}^* \boldsymbol{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$  has rank  $r \ll d$

# Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbf{U}^*$  ( $r$ -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$

noise matrix:  $\mathbf{E}$

- Ground-truth data

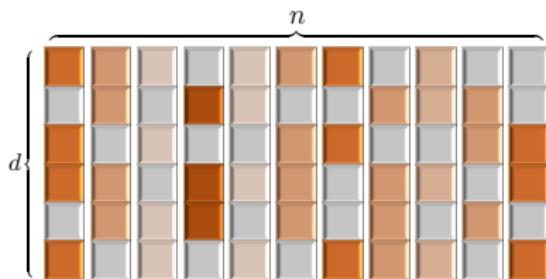
$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}, \quad \mathbf{x}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*)$$

where  $\mathbf{S}^* = \mathbf{U}^* \boldsymbol{\Lambda}^* \mathbf{U}^{*\top} \in \mathbb{R}^{d \times d}$  has rank  $r \ll d$

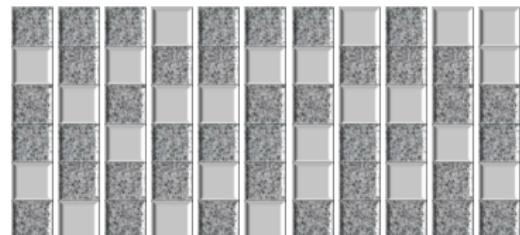
- Noisy observations:  $\mathbf{X} + \mathbf{E}$  (a.k.a. spiked covariance model)

# Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq U^*$  ( $r$ -dimensional)



$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$$



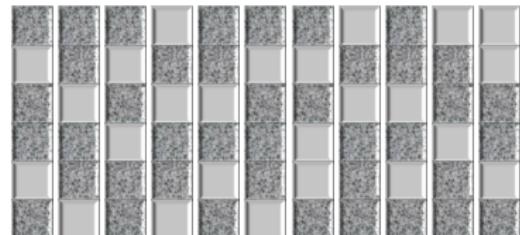
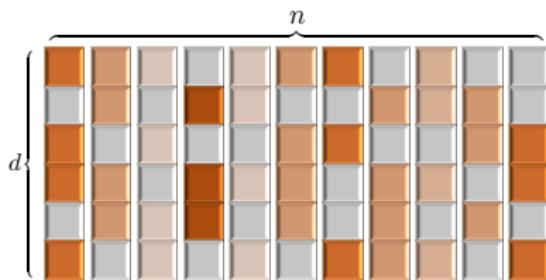
noise matrix:  $\mathbf{E}$

- Incomplete observations  $\longrightarrow$  sampling set  $\Omega$ :

$$Y_{i,j} = \begin{cases} X_{i,j} + E_{i,j}, & (i, j) \in \Omega \\ 0, & \text{else} \end{cases} \quad \text{or} \quad \mathbf{Y} = \mathcal{P}_{\Omega}(\mathbf{X} + \mathbf{E})$$

# Principal component analysis

$\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbf{U}^*$  ( $r$ -dimensional)



- **Goal:**

- Construct confidence regions for principal subspace  $\mathbf{U}^*$
- Construct entrywise confidence intervals for covariance matrix  $\mathbf{S}^*$

# What we consider here . . .

---

- **Heteroskedastic noise:**  $\{E_{i,j}\}$  are ind. sub-Gaussian obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] = \omega_i^{*2} \leq \omega_{\max}^2, \quad \underbrace{\|E_{i,j}\|_{\psi_2}}_{\text{sub-Gaussian norm}} = O(\omega_i^*)$$

- noise variance  $\{\omega_i^{*2}\}$ : **unknown**, location-varying

# What we consider here . . .

---

- **Heteroskedastic noise:**  $\{E_{i,j}\}$  are ind. sub-Gaussian obeying

$$\mathbb{E}[E_{i,j}] = 0, \quad \mathbb{E}[E_{i,j}^2] = \omega_i^{*2} \leq \omega_{\max}^2, \quad \underbrace{\|E_{i,j}\|_{\psi_2}}_{\text{sub-Gaussian norm}} = O(\omega_i^*)$$

- **Random sampling:**  $(i, j) \in \Omega$  independently with prob.  $p$

## What we consider here . . .

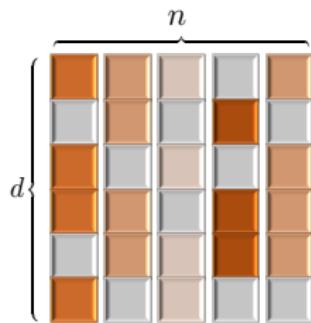
---

**Our focus:** estimating/inferring column subspace of  $\mathbf{X}$  from noisy and incomplete observations  $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$

# What we consider here ...

---

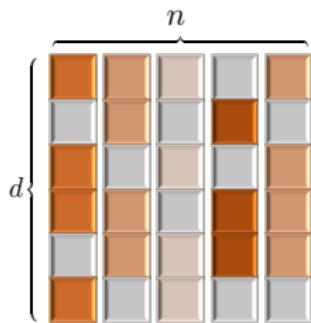
**Our focus:** estimating/inferring column subspace of  $\mathbf{X}$  from noisy and incomplete observations  $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$



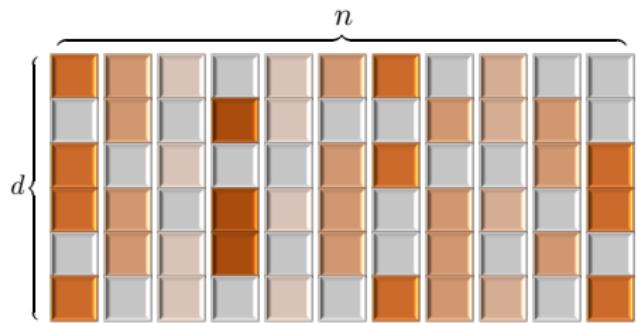
$n \lesssim d$ : solvable via *matrix completion* methods  
(e.g., Chen, Fan, Ma, Yan '19)

# What we consider here ...

**Our focus:** estimating/inferring column subspace of  $\mathbf{X}$  from noisy and incomplete observations  $\mathbf{Y} = \mathcal{P}_{\Omega}(\mathbf{X} + \mathbf{E})$  when  $\underbrace{n \gg d}_{\text{more challenging regime}}$



$n \lesssim d$ : solvable via *matrix completion* methods  
(e.g., Chen, Fan, Ma, Yan '19)



$n \gg d$ : sometimes it's only feasible to estimate col-space instead of whole matrix

## A natural SVD-based algorithm

---

- **Compute:** rank- $r$  SVD  $\mathbf{U}\Sigma\mathbf{V}^\top$  of  $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$
- **Output:**  $\mathbf{U}$   $\longrightarrow$  estimate of  $\mathbf{U}^*$

# A natural SVD-based algorithm

---

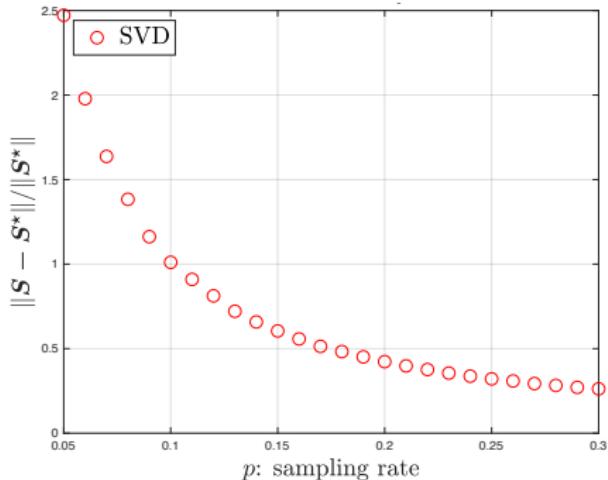
- **Compute:** rank- $r$  SVD  $\mathbf{U}\Sigma\mathbf{V}^\top$  of  $\mathbf{Y} = \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E})$
- **Output:**  $\mathbf{U}$   $\longrightarrow$  estimate of  $\mathbf{U}^*$

**Rationale:** under zero-mean noise and random sampling, we have

$$\text{col-space}(\mathbb{E}[\mathbf{Y}]) = \text{col-space}(\mathbf{X}) = \mathbf{U}^*$$

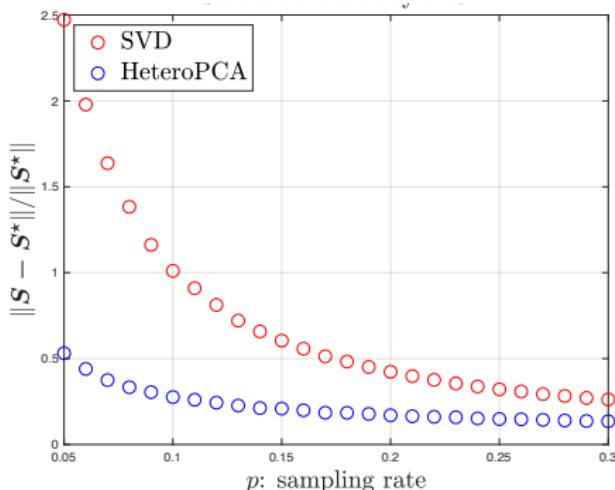
# Numerical suboptimality of SVD-based approach

---



$n = 2000, \ d = 100, \ r = 3, \ \omega_1^*, \dots, \omega_d^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0.025, 0.1]$

# Numerical suboptimality of SVD-based approach



$$n = 2000, \quad d = 100, \quad r = 3, \quad \omega_1^*, \dots, \omega_d^* \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0.025, 0.1]$$

Plain SVD is suboptimal in the presence of missing data if  $n \gg d$

## Diagnosis: diagonal entries need special treatment

---

$$\text{col-space}(\mathbf{Y}) = \text{eig-space}(\mathbf{Y}\mathbf{Y}^T)$$

## Diagnosis: diagonal entries need special treatment

---

$$\text{col-space}(\mathbf{Y}) = \text{eig-space}(\mathbf{Y}\mathbf{Y}^\top)$$

**Large bias in diagonal entries:**

$$\frac{1}{p^2} \mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \underbrace{\mathbf{X}\mathbf{X}^\top}_{\checkmark} + \underbrace{\left(\frac{1}{p} - 1\right) \mathcal{P}_{\text{diag}}(\mathbf{X}\mathbf{X}^\top)}_{\text{potentially large diagonal matrix!}} + \frac{n}{p} \text{diag}\{[\omega_i^{*2}]\}$$

## Diagnosis: diagonal entries need special treatment

---

$$\text{col-space}(\mathbf{Y}) = \text{eig-space}(\mathbf{Y}\mathbf{Y}^\top)$$

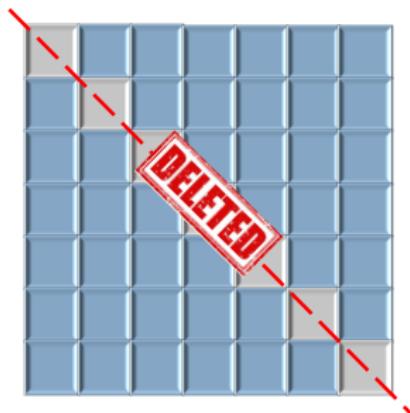
**Large bias in diagonal entries:**

$$\frac{1}{p^2} \mathbb{E}[\mathbf{Y}\mathbf{Y}^\top] = \underbrace{\mathbf{X}\mathbf{X}^\top}_{\checkmark} + \underbrace{\left(\frac{1}{p} - 1\right) \mathcal{P}_{\text{diag}}(\mathbf{X}\mathbf{X}^\top)}_{\text{potentially large diagonal matrix!}} + \frac{n}{p} \text{diag}\left\{[\omega_i^{*2}]\right\}$$

- a common issue under missing data or heteroskedastic noise

## Two spectral algorithms that take care of diagonals

---

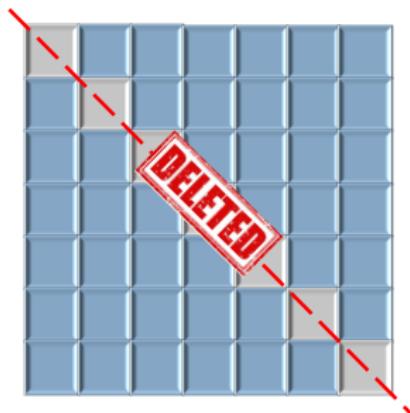


diagonal-deleted/reweighted PCA

- remove/reweight  $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

## Two spectral algorithms that take care of diagonals

---



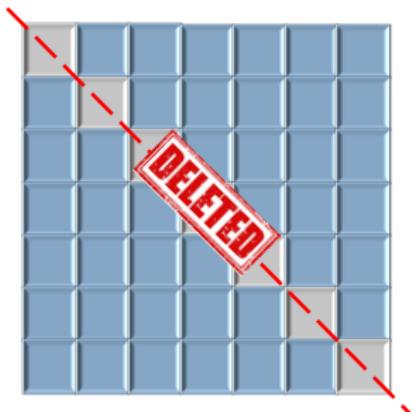
diagonal-deleted/reweighted PCA

- remove/reweight  $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

- Loh, Wainwright '12
- Lounici '13 '14
- Florescu and Perkins '16
- Montanari and Sun '18
- Zhu, Wang, Samworth '19
- Cai, Li, Chi, Poor, Chen '19
- ...

# Two spectral algorithms that take care of diagonals

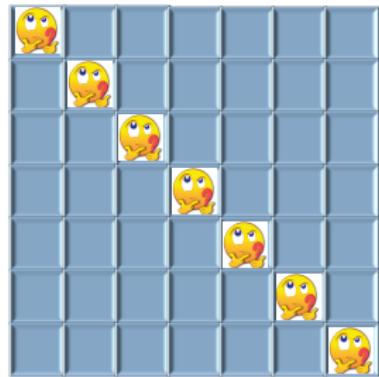
---



diagonal-deleted/reweighted PCA

- remove/reweight  $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

— Loh, Wainwright '12  
— Lounici '13 '14  
— Florescu and Perkins '16  
— Montanari and Sun '18  
— Zhu, Wang, Samworth '19  
— Cai, Li, Chi, Poor, Chen '19  
— ...



HeteroPCA (Zhang et al '18)

- iteratively estimate  $\text{diag}(\mathbf{Y}\mathbf{Y}^\top)$

# HeteroPCA (Zhang, Cai, Wu '18)

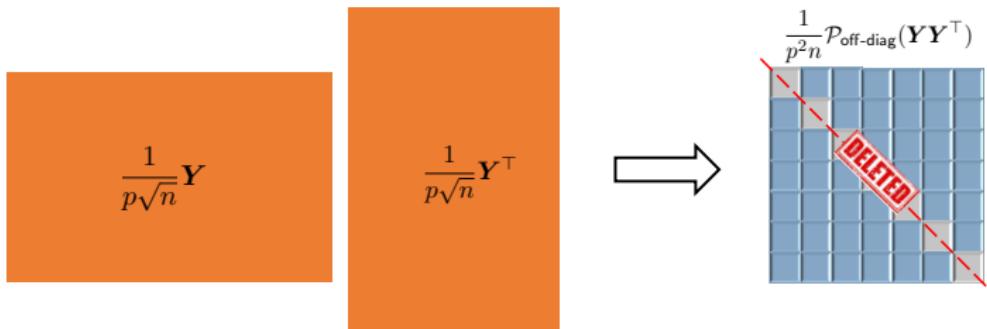
---

$$\frac{1}{p\sqrt{n}} \mathbf{Y}$$

$$\frac{1}{p\sqrt{n}} \mathbf{Y}^\top$$

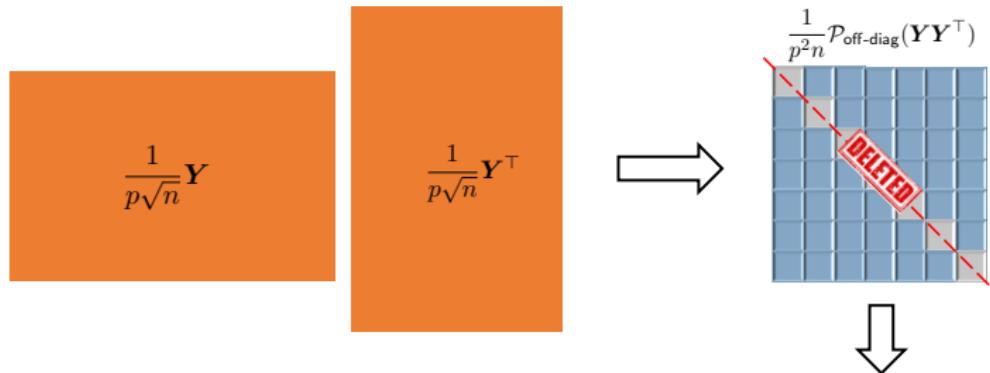
# HeteroPCA (Zhang, Cai, Wu '18)

---



# HeteroPCA (Zhang, Cai, Wu '18)

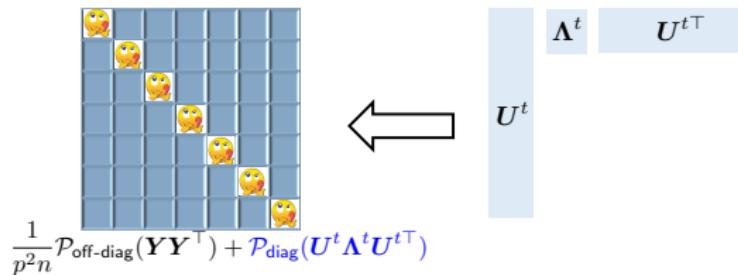
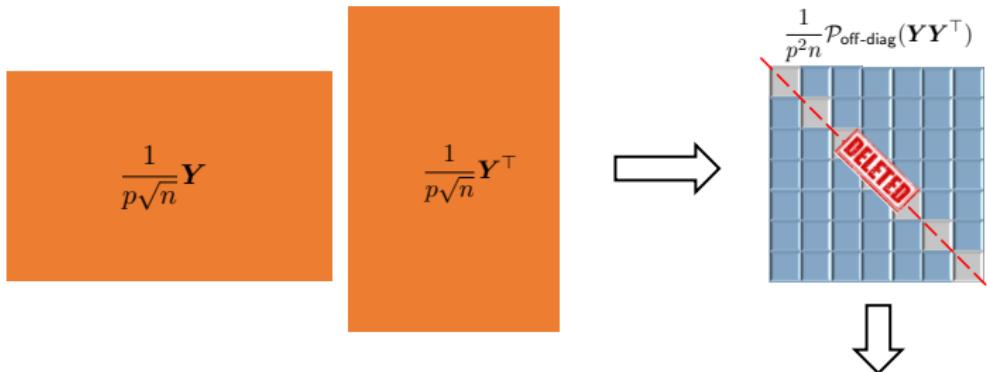
---



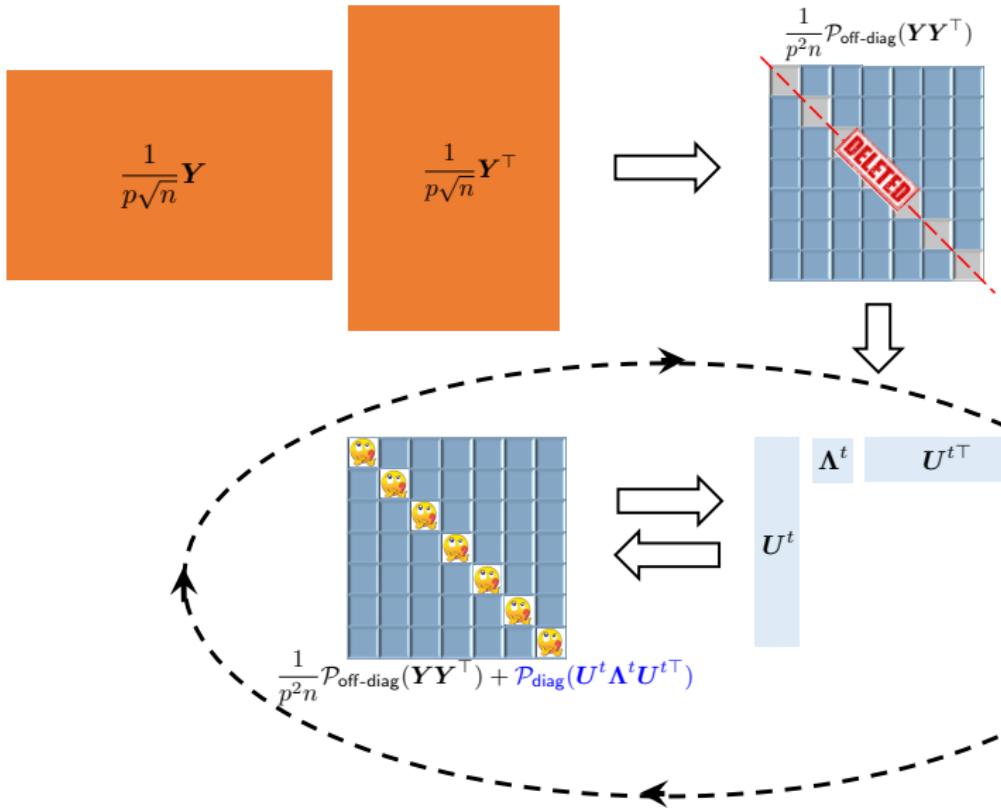
$$\begin{matrix} & \Lambda^t & U^{t\top} \\ U^t & & \end{matrix}$$

# HeteroPCA (Zhang, Cai, Wu '18)

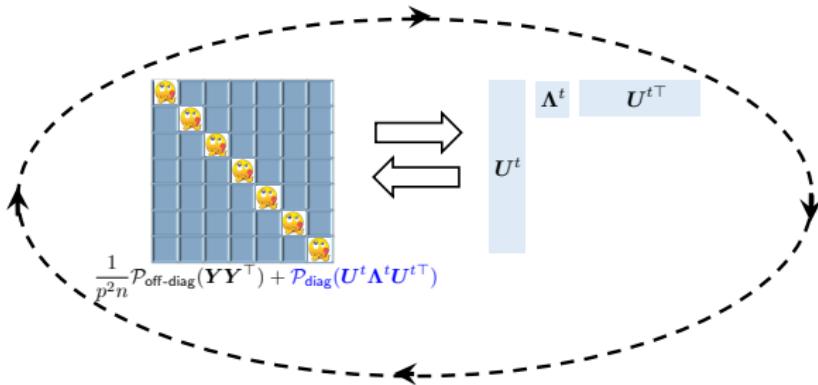
---



# HeteroPCA (Zhang, Cai, Wu '18)



# HeteroPCA (Zhang, Cai, Wu '18)

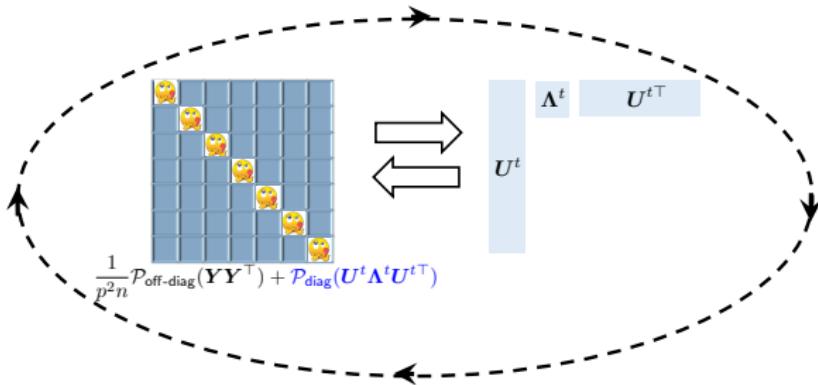


- **Initialize:**  $\mathbf{G}^0 = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$
- **Iterative update:** for  $t = 0, 1, \dots, t_0$

$$(\mathbf{U}^t, \boldsymbol{\Lambda}^t) = \text{eigs}(\mathbf{G}^t, r)$$

$$\mathbf{G}^{t+1} = \mathbf{G}^0 + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \boldsymbol{\Lambda}^t \mathbf{U}^{t\top})$$

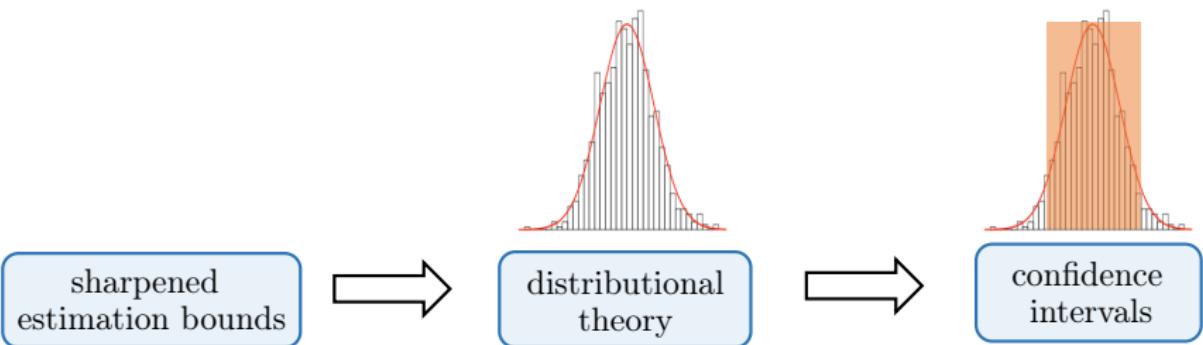
# HeteroPCA (Zhang, Cai, Wu '18)



- **Initialize:**  $\mathbf{G}^0 = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(\mathbf{Y}\mathbf{Y}^\top)$
- **Iterative update:** for  $t = 0, 1, \dots, t_0$   
$$(\mathbf{U}^t, \boldsymbol{\Lambda}^t) = \text{eigs}(\mathbf{G}^t, r)$$
$$\mathbf{G}^{t+1} = \mathbf{G}^0 + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \boldsymbol{\Lambda}^t \mathbf{U}^{t\top})$$

- **Output:**  $\mathbf{U} := \mathbf{U}^{t_0} \longrightarrow \text{estimate of } \mathbf{U}^*$

## Our contributions: estimation and inference based on HeteroPCA



# Sharpened estimation guarantees for HeteroPCA

---

## Assumptions (omitting log factors)

- $S^*$  has rank  $r = O(1)$ , incoherent, well-conditioned
- sampling rate exceeds certain threshold

$$p \gtrsim \max \left\{ \frac{1}{\sqrt{nd}}, \frac{1}{n} \right\}$$

- per-entry signal-to-noise ratio (SNR) cannot be too low:

$$\frac{\omega_{\max}^2}{\lambda_r(S^*)/d} \lesssim \min \left\{ pn, p\sqrt{nd} \right\}$$

# Sharpened estimation guarantees for HeteroPCA

---

## Assumptions (omitting log factors)

- $S^*$  has rank  $r = O(1)$ , incoherent, well-conditioned
- sampling rate exceeds certain threshold

$$p \gtrsim \underbrace{\frac{1}{\sqrt{nd}}}_{\text{HeteroPCA}} \quad \text{vs.} \quad \underbrace{\frac{1}{d}}_{\text{matrix completion}} \quad (n \gg d)$$

- per-entry signal-to-noise ratio (SNR) cannot be too low:

$$\frac{\omega_{\max}^2}{\lambda_r(S^*)/d} \lesssim \underbrace{p\sqrt{nd}}_{\text{HeteroPCA}} \quad \text{vs.} \quad \underbrace{pd}_{\text{matrix completion}} \quad (n \gg d)$$

# Sharpened estimation guarantees for HeteroPCA

## Theorem 4

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

$$\|\mathbf{UH} - \mathbf{U}^*\| \lesssim \zeta, \quad \|\mathbf{UH} - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{1}{\sqrt{d}} \zeta$$

where

$$\zeta := \frac{1}{\sqrt{nd} p} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}} = o(1).$$

# Sharpened estimation guarantees for HeteroPCA

## Theorem 4

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

$$\|\mathbf{UH} - \mathbf{U}^*\| \lesssim \zeta, \quad \|\mathbf{UH} - \mathbf{U}^*\|_{2,\infty} \lesssim \frac{1}{\sqrt{d}} \zeta$$

where

$$\zeta := \frac{1}{\sqrt{nd} p} + \frac{\omega_{\max}^2}{p \lambda_r^*} \sqrt{\frac{d}{n}} + \sqrt{\frac{1}{np}} + \frac{\omega_{\max}}{\sqrt{\lambda_r^*}} \sqrt{\frac{d}{np}} = o(1).$$

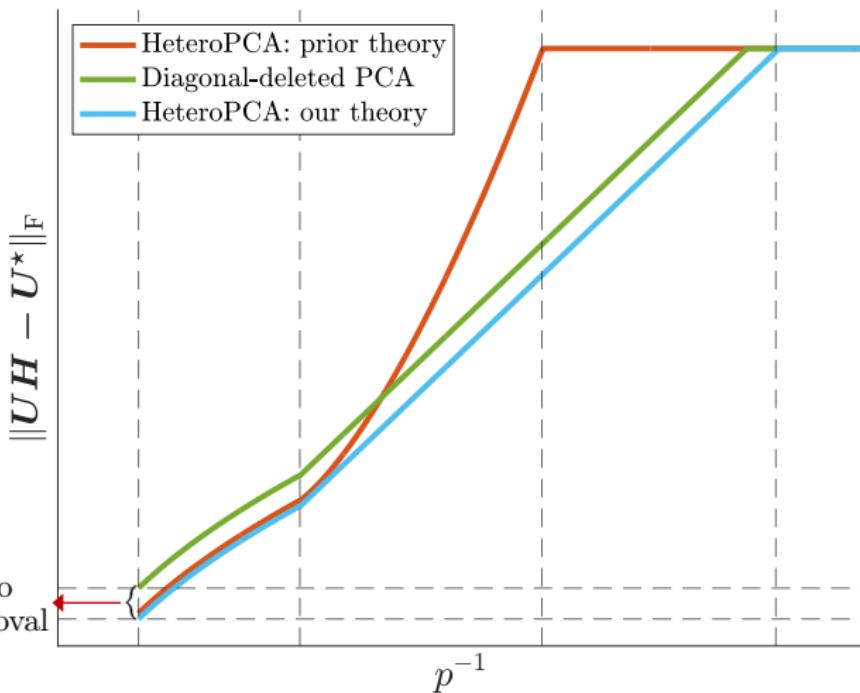
- rotational ambiguity: for any orthonormal matrix  $\mathbf{O} \in \mathcal{O}^{r \times r}$ ,  $\mathbf{U}$  and  $\mathbf{UO}$  represents the same subspace
- global rotation matrix  $\mathbf{H} = \arg \min_{\mathbf{O} \in \mathcal{O}^{r \times r}} \|\mathbf{UO} - \mathbf{U}^*\|_{\text{F}}$

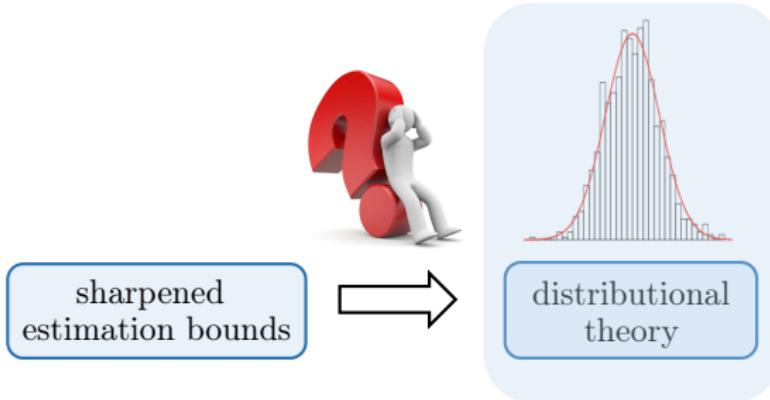
# Sharpened estimation guarantees for HeteroPCA

$$p = \frac{1}{d}$$

$$p = \frac{1}{n^{1/3} d^{2/3}}$$

$$p = \frac{1}{\sqrt{nd}}$$





*Given HeteroPCA is an appealing estimator, can we take one step further to obtain distributional characterizations?*

# Distributional theory for $\mathbf{U}$

---

## Theorem 5

Consider any  $1 \leq l \leq d$ . Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P}((\mathbf{UH} - \mathbf{U}^*)_{l,\cdot} \in \mathcal{C}) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

# Distributional theory for $\mathbf{U}$

---

## Theorem 5

Consider any  $1 \leq l \leq d$ . Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P}((\mathbf{UH} - \mathbf{U}^*)_{l,\cdot} \in \mathcal{C}) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Each row of  $\mathbf{U}$  is approximately Gaussian
  - nearly unbiased + tractable covariance

# Distributional theory for $\mathbf{U}$

## Theorem 5

Consider any  $1 \leq l \leq d$ . Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P}((\mathbf{UH} - \mathbf{U}^*)_{l,\cdot} \in \mathcal{C}) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Each row of  $\mathbf{U}$  is approximately Gaussian
  - nearly unbiased + tractable covariance

$$\boldsymbol{\Sigma}_{U,l}^* := \left( \frac{1-p}{np} S_{l,l}^* + \frac{\omega_l^{*2}}{np} \right) (\boldsymbol{\Lambda}^*)^{-1} + \frac{2(1-p)}{np} \mathbf{U}_{l,\cdot}^{*\top} \mathbf{U}_{l,\cdot}^*$$

$$+ (\boldsymbol{\Lambda}^*)^{-1} \mathbf{U}^{*\top} \text{diag} \left\{ [d_{l,i}^*]_{1 \leq i \leq d} \right\} \mathbf{U}^* (\boldsymbol{\Lambda}^*)^{-1}$$

$$d_{l,i}^* := \frac{1}{np^2} \left[ \omega_l^{*2} + (1-p) S_{l,l}^{*2} \right] \left[ \omega_i^{*2} + (1-p) S_{i,i}^{*2} \right] + \frac{2(1-p)^2}{np^2} S_{l,i}^{*2}$$

# Distributional theory for $U$

## Theorem 5

Consider any  $1 \leq l \leq d$ . Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P}\left((\mathbf{UH} - \mathbf{U}^*)_{l,\cdot} \in \mathcal{C}\right) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- Key decomposition:

$$\mathbf{UH} - \mathbf{U}^* \approx \left[ \underbrace{\Delta \mathbf{X}^\top}_{\text{linear term}} + \underbrace{\mathcal{P}_{\text{off-diag}}(\Delta \Delta^\top)}_{\text{quadratic term}} \right] \mathbf{U}^* (\Lambda^*)^{-1}$$

where  $\Delta$  is the “equivalent noise matrix”

$$\Delta := \frac{1}{p} \mathcal{P}_\Omega(\mathbf{X} + \mathbf{E}) - \mathbf{X}$$

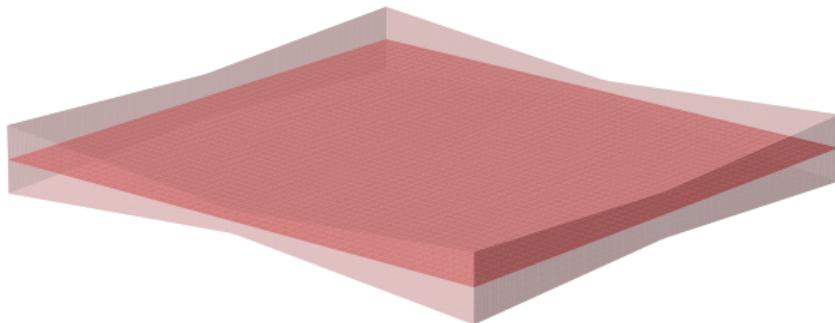
# Distributional theory for $U$

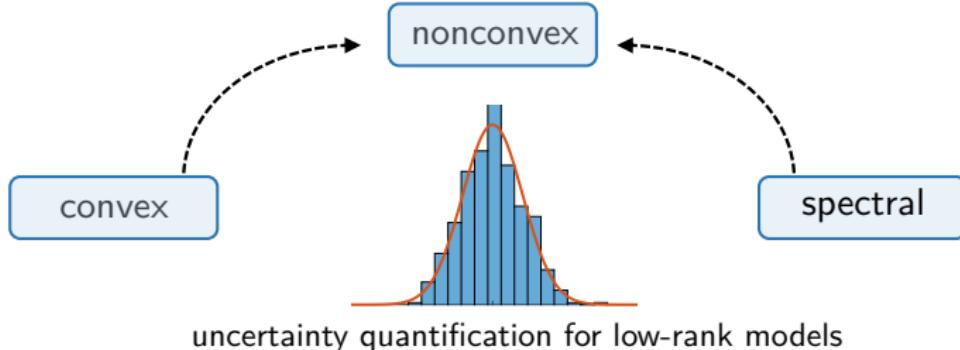
## Theorem 5

Consider any  $1 \leq l \leq d$ . Under previous assumptions, we have

$$\sup_{\text{cvx set } \mathcal{C}} \left| \mathbb{P}((\mathbf{UH} - \mathbf{U}^*)_{l,\cdot} \in \mathcal{C}) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*) \{ \mathcal{C} \} \right| = o(1)$$

- confidence regions for the principal subspace



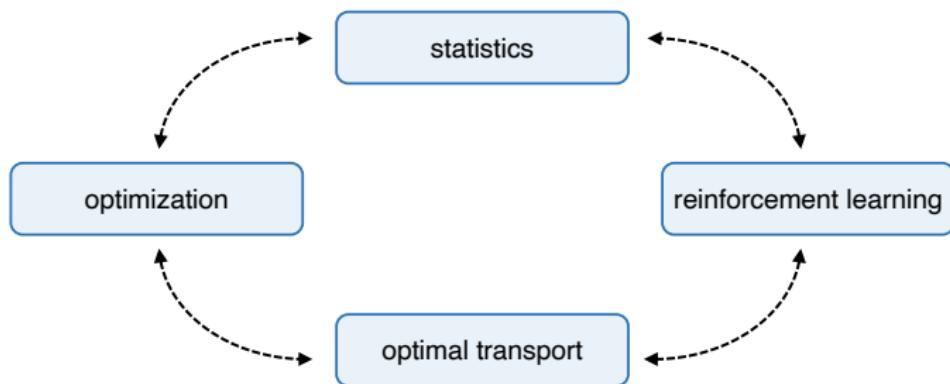


"Inference and uncertainty quantification for noisy matrix completion," Y. Chen, J. Fan, C. Ma, Y. Yan, *Proceedings of National Academy of Sciences (PNAS)*

"Bridging convex and nonconvex optimization in robust PCA: noise, outliers, and missing data," Y. Chen, J. Fan, C. Ma, Y. Yan, *Annals of Statistics*

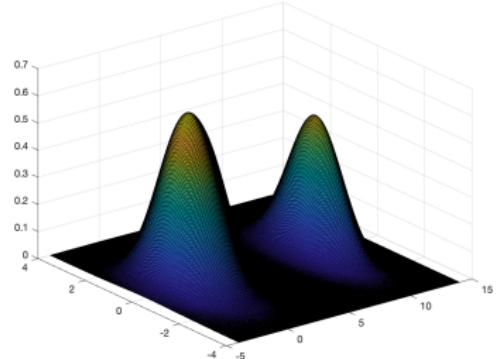
"Inference for heteroskedastic PCA with missing data," Y. Yan, Y. Chen, J. Fan

## *Other highlights of my research*



# Overparameterization for Gaussian mixture model

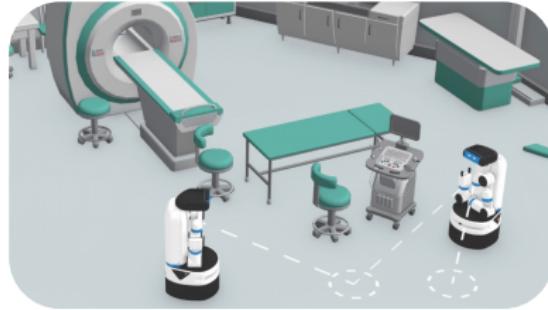
- **observations:**  $X_1, \dots, X_N \stackrel{\text{ind}}{\sim} \rho^* * \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$
- **algorithm:** an interacting particle system where the **weight** and **location** of particles evolve according to the **Wasserstein-Fisher-Rao** gradient flow
- **theory:** convergence to nonparametric MLE in mean-field limit



"Learning Gaussian mixtures using the Wasserstein-Fisher-Rao gradient flow," Y. Yan,  
K. Wang, P. Rigollet

# Offline multi-agent reinforcement learning

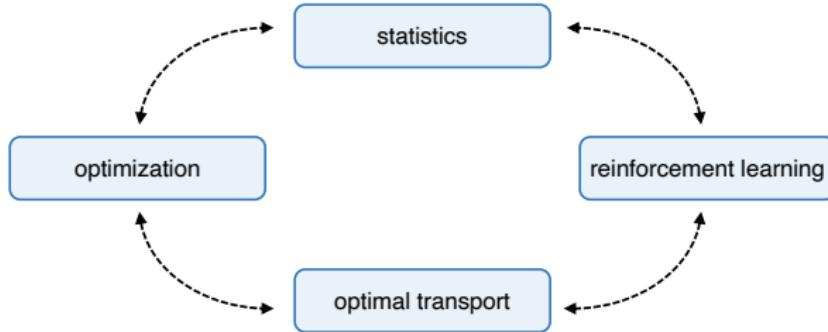
---



"Model-based reinforcement learning for offline zero-sum Markov games," Y. Yan, G. Li, Y. Chen, J. Fan

# Many future opportunities...

---



- robust estimation/inference
- uncertainty quantification for deep learning
- multi-agent reinforcement learning
- applications in economics and social sciences ...

*Thank you!*