

# Geometry and (Implicit) Regularization in Nonconvex Low-Rank Estimation

Yuejie Chi

**Carnegie Mellon University**

SPARS 2019

Jul. 2019

# Acknowledgements



Yuxin Chen  
Princeton



Jianqing Fan  
Princeton



Yingbin Liang  
Ohio State



Cong Ma  
Princeton



Kaizheng Wang  
Princeton



Yuanxin Li  
CMU



Huishuai Zhang  
MSRA

## Empirical risk minimization

Given data  $z$ , estimate parameters  $\mathbf{x} \in \mathbb{R}^n$ :

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) := \frac{1}{m} \sum_{i=1}^m \ell(z_i; \mathbf{x})$$

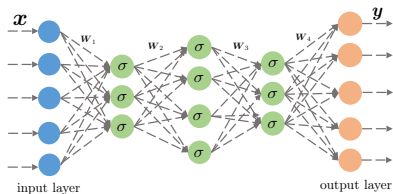
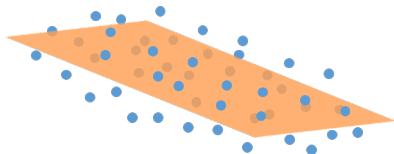
where  $\ell(z_i; \mathbf{x})$  is the sample loss.

# Empirical risk minimization

Given data  $z$ , estimate parameters  $\mathbf{x} \in \mathbb{R}^n$ :

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) := \frac{1}{m} \sum_{i=1}^m \ell(z_i; \mathbf{x})$$

where  $\ell(z_i; \mathbf{x})$  is the sample loss.

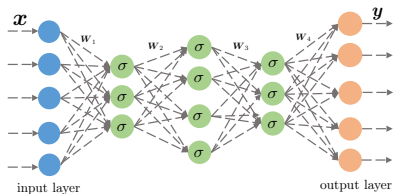
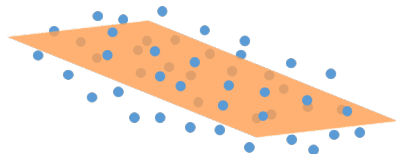


# Empirical risk minimization

Given data  $z$ , estimate parameters  $\mathbf{x} \in \mathbb{R}^n$ :

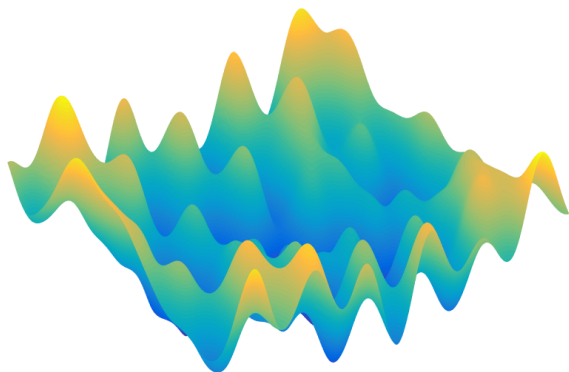
$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) := \frac{1}{m} \sum_{i=1}^m \ell(z_i; \mathbf{x})$$

where  $\ell(z_i; \mathbf{x})$  is the sample loss.



Often lead to nonconvex problems that are deemed intractable!

## Nonconvex problems are hard!



*“...in fact, the great watershed in optimization isn't between linearity and nonlinearity, but convexity and nonconvexity.*”

R. T. Rockafellar, in SIAM Review, 1993

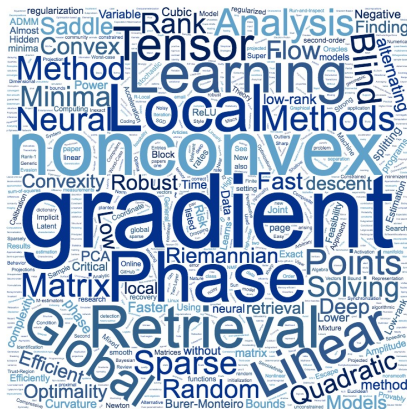
## Nonconvex problems are hard!



*“...in fact, the great watershed in optimization isn't between linearity and nonlinearity, but convexity and nonconvexity.*”

R. T. Rockafellar, in SIAM Review, 1993

# Recent developments: provable nonconvex optimization



*Only an incomplete list...*

**Phase retrieval:** Gerchberg, Saxton '72, Netrapalli et al. '13, Candès, Li, Soltanolkotabi '14, Chen, Candès '15, Cai, Li, Ma '15, Zhang et al. '16, Wang et al. '16, Sun, Qu, Wright '16, Ma et al. '17, Chen et al. '18, Soltani, Hegde '18, Ruan and Duchi, '18, ...

**Matrix sensing/completion:** Keshavan et al. '09, Jain et al. '09, Hardt '13, Jain et al. '13, Sun, Luo '15, Chen, Wainwright '15, Tu et al. '15, Zheng, Lafferty '15, Bhojanapalli et al. '16, Ge, Lee, Ma '16, Jin et al. '16, Ma et al. '17, Chen and Li '17, Cai et al. '18, Li, Zhu, Tang, Wakin '18, Charisopoulos et al. '19, ...

**Blind deconvolution/demixing:** Li et al. '16, Lee et al. '16, Cambareri, Jacques '16, Ling, Strohmer '16, Huang, Hand '16, Ma et al. '17, Zhang et al. '18, Li, Bresler '18, Dong, Shi '18, ...

**Dictionary learning:** Arora et al. '14, Sun et al. '15, Chatterji, Bartlett '17, Bai et al. '18, Gilboa et al. '18, Rambhatla et al. '19, ...

**Robust principal component analysis:** Netrapalli et al. '14, Yi et al. '16, Gu et al. '16, Ge et al. '17, Cherapanamjeri et al. '17, Vaswani et al. '18, Maunu et al. '19, ...

**Deep learning:** Zhong et al. '17, Bai, Mei, Montanari '17, Du et al. '17, Ge, Lee, Ma '17, Du et al. '18, Soltanolkotabi and Oymak, '18...



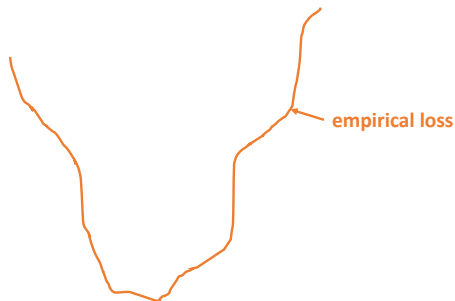
# Statistical thinking in nonconvex optimization

Data/measurements follow certain **statistical models** and hence are not worst-case instances.

# Statistical thinking in nonconvex optimization

Data/measurements follow certain **statistical models** and hence are not worst-case instances.

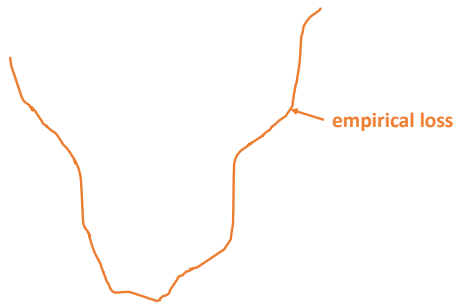
$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{x})$$



# Statistical thinking in nonconvex optimization

Data/measurements follow certain **statistical models** and hence are not worst-case instances.

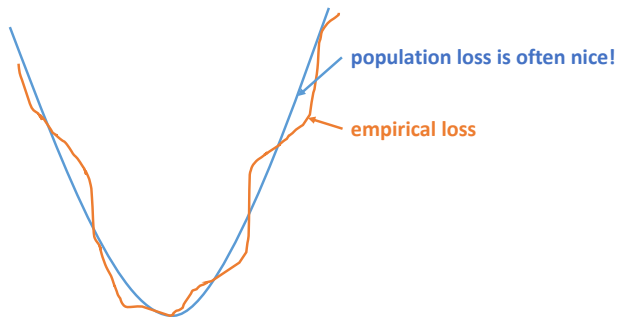
$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{x}) \quad \xrightarrow{m \rightarrow \infty} \quad \mathbb{E}[\ell(y; \mathbf{x})]$$



# Statistical thinking in nonconvex optimization

Data/measurements follow certain **statistical models** and hence are not worst-case instances.

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{x}) \quad \xrightarrow{m \rightarrow \infty} \quad \mathbb{E}[\ell(y; \mathbf{x})]$$

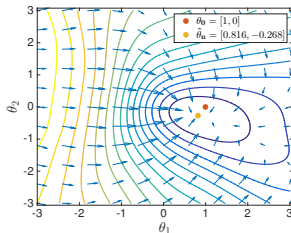


*We will detail an example of “nice” population landscape later.*

# From population to empirical risk: a geometric perspective

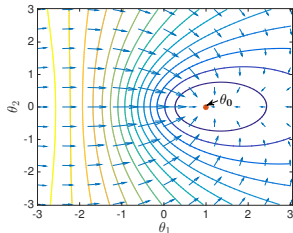
**Geometric analysis:** uniform concentration of Hessians and gradients, along some descent directions;

- one-to-one correspondence between critical points;
- preservation of geometric curvatures.



empirical risk

$\approx$



population risk

---

Bai et al. '16, Sun et al. '15, Sun et al. '16, Ge et al. '16; Figure credit: Bai, Mei, and Montanari

# This talk: sample-starved regime

## Sample-starved regime:

*sample size  $\gtrsim O(\text{number of unknowns})$*

# This talk: sample-starved regime

## Sample-starved regime:

*sample size  $\gtrsim O(\text{number of unknowns})$*

Even when  $\mathbb{E}[f(\mathbf{x})]$  is locally strongly convex and smooth,

- *$f(\mathbf{x})$  may be much more ill-conditioned than  $\mathbb{E}[f(\mathbf{x})]$ ;  
smaller step size and more computation*

# This talk: sample-starved regime

## Sample-starved regime:

*sample size  $\gtrsim O(\text{number of unknowns})$*

Even when  $\mathbb{E}[f(\mathbf{x})]$  is locally strongly convex and smooth,

- $f(\mathbf{x})$  may be much more ill-conditioned than  $\mathbb{E}[f(\mathbf{x})]$ ;  
*smaller step size and more computation*
- $f(\mathbf{x})$  may lack curvatures in certain regions.  
*complicated regularization*



# This talk: sample-starved regime

## Sample-starved regime:

$$\text{sample size} \gtrsim O(\text{number of unknowns})$$

Even when  $\mathbb{E}[f(\mathbf{x})]$  is locally strongly convex and smooth,

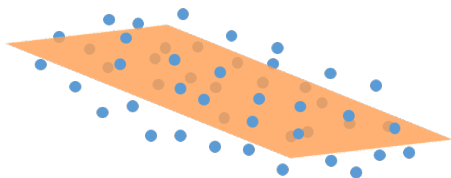
- $f(\mathbf{x})$  may be much more ill-conditioned than  $\mathbb{E}[f(\mathbf{x})]$ ;  
*smaller step size and more computation*
- $f(\mathbf{x})$  may lack curvatures in certain regions.  
*complicated regularization*



**Does the geometric gap between  $f(\mathbf{x})$  vs  $\mathbb{E}[f(\mathbf{x})]$  hurt optimization efficacy?**

*a case study with low-rank matrix completion*

## Revisiting PCA: in search of low-rank representation



Given  $\mathbf{M} \succeq 0 \in \mathbb{R}^{n \times n}$  (e.g. sample covariance matrix), find its best rank- $r$  approximation:

$$\underbrace{\widehat{\mathbf{M}} = \operatorname{argmin}_{\mathbf{Z}} \|\mathbf{Z} - \mathbf{M}\|_{\text{F}}^2 \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{Z}) \leq r}_{\text{nonconvex optimization!}}$$

## An optimization viewpoint

**Low-rank factorization:** if we factorize  $Z = \mathbf{X}\mathbf{X}^\top$  with  $\mathbf{X} \in \mathbb{R}^{n \times r}$ , then it leads to a nonconvex problem:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

# An optimization viewpoint

**Low-rank factorization:** if we factorize  $\mathbf{Z} = \mathbf{X}\mathbf{X}^\top$  with  $\mathbf{X} \in \mathbb{R}^{n \times r}$ , then it leads to a nonconvex problem:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

## Theorem (Baldi and Hornik, 1989)

Suppose  $\mathbf{M}$  has a strict eigen-gap between  $\lambda_r$  and  $\lambda_{r+1}$ , the critical points of  $f(\mathbf{X})$  can be categorized into

- global minima;
- strict saddle points, from which there exist directions to strictly decrease  $f(\mathbf{X})$ .

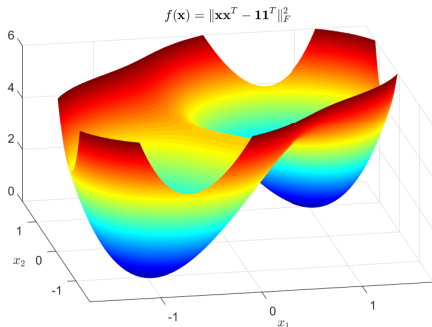
In other words, *all local minima are global minima!*

---

Baldi and Hornik. "Neural networks and principal component analysis: Learning from examples without local minima." Neural networks 2.1 (1989): 53-58.

# Benign landscape of PCA

For example, for 2-dimensional case  $f(\mathbf{x}) = \left\| \mathbf{x}\mathbf{x}^\top - \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right\|_F^2$

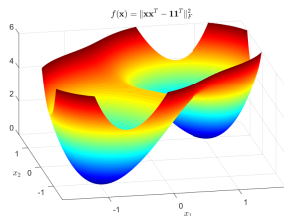


global minima:  $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ; strict saddles:  $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , and  $\pm \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

— No “spurious” local minima!

# Parameter recovery via gradient descent

a two-step recovery strategy:



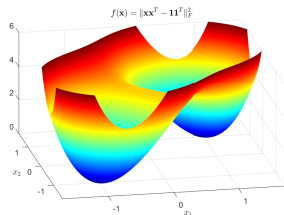
- Find an initial point that falls into a “basin of attraction”
- Gradient iterations:

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \nabla f(\mathbf{X}^t)$$

for  $t = 0, 1, \dots$

# Parameter recovery via gradient descent

a two-step recovery strategy:



- Find an initial point that falls into a “basin of attraction”
- Gradient iterations:

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \nabla f(\mathbf{X}^t)$$

for  $t = 0, 1, \dots$

- The spectral method can be used for initialization;
- Low-complexity local refinements via gradient descent.



## Low-rank matrix completion: dealing with missing data

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

Given partial samples of a *low-rank* matrix  $\mathbf{M}$  in an index set  $\Omega$ , fill in missing entries.

$$\text{find low-rank } \widehat{\mathbf{M}} \quad \text{s.t.} \quad \mathcal{P}_{\Omega}(\widehat{\mathbf{M}}) = \mathcal{P}_{\Omega}(\mathbf{M})$$

*Applications: recommendation systems, ...*

# Incoherence

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}$$

vs.

$$\underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}$$

# Incoherence

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}_{\text{hard}} \quad \text{vs.} \quad \underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}_{\text{easy}}$$

## Definition (Incoherence for matrix completion)

A rank- $r$  matrix  $M^{\natural}$  with eigendecomposition  $M^{\natural} = U^{\natural} \Sigma^{\natural} U^{\natural \top}$  is said to be  $\mu$ -incoherent if

$$\|U^{\natural}\|_{2,\infty} \leq \sqrt{\frac{\mu}{n}} \|U^{\natural}\|_{\text{F}} = \sqrt{\frac{\mu r}{n}}.$$

Note:  $\|U\|_{2,\infty} = \max_i \|e_i^{\top} U\|_2$ .

Lower bound [Candès and Tao]:  $p \gtrsim \mu r \log n/n$ .

# Incoherence

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}_{\text{hard } \mu=n} \quad \text{vs.} \quad \underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}_{\text{easy } \mu=1}$$

## Definition (Incoherence for matrix completion)

A rank- $r$  matrix  $M^{\natural}$  with eigendecomposition  $M^{\natural} = U^{\natural} \Sigma^{\natural} U^{\natural \top}$  is said to be  $\mu$ -incoherent if

$$\|U^{\natural}\|_{2,\infty} \leq \sqrt{\frac{\mu}{n}} \|U^{\natural}\|_{\text{F}} = \sqrt{\frac{\mu r}{n}}.$$

Note:  $\|U\|_{2,\infty} = \max_i \|e_i^{\top} U\|_2$ .

Lower bound [Candès and Tao]:  $p \gtrsim \mu r \log n/n$ .

## A natural least-squares formulation

given:  $\mathcal{P}_\Omega(\mathbf{M})$

↓

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \left\| \mathcal{P}_\Omega(\mathbf{X}\mathbf{X}^\top - \mathbf{M}) \right\|_F^2$$

## A natural least-squares formulation

given:  $\mathcal{P}_\Omega(\mathbf{M})$

↓

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \left\| \mathcal{P}_\Omega(\mathbf{X}\mathbf{X}^\top - \mathbf{M}) \right\|_F^2$$

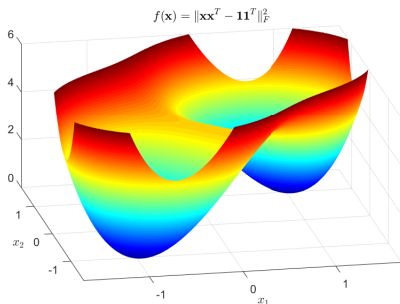
- **Bernoulli sampling:** Assume every entry is observed i.i.d. with  $0 < p \leq 1$ :

$$\mathbb{E}[f(\mathbf{X})] = p \left\| \mathbf{X}\mathbf{X}^\top - \mathbf{M} \right\|_F^2.$$

# What does the population level look like?

**Population level ( $p = 1$ ): this is PCA.**

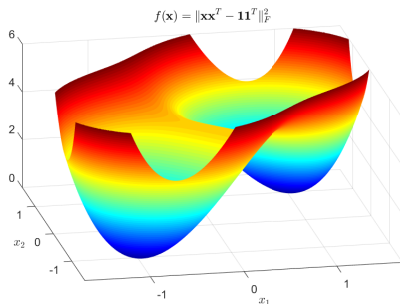
$f(\mathbf{X})$  restricted strongly convex and smooth along descent direction  $\mathbf{V}$  when  $\mathbf{X}$  is close to  $\mathbf{X}^\natural$ .



## What does the population level look like?

**Population level ( $p = 1$ ): this is PCA.**

$f(\mathbf{X})$  restricted strongly convex and smooth along descent direction  $\mathbf{V}$  when  $\mathbf{X}$  is close to  $\mathbf{X}^\natural$ .



**Consequence:** GD converges within  $O(\log \frac{1}{\epsilon})$  iterations if  $p = 1$ .



## What does the finite-sample level look like?

Assume every entry is observed i.i.d. with probability  $0 < p \leq 1$ .

$$f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left( \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

## What does the finite-sample level look like?

Assume every entry is observed i.i.d. with probability  $0 < p \leq 1$ .

$$f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left( \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

**Finite-sample level** ( $p \asymp \frac{\text{polylog}n}{n}$ )

$f(\mathbf{X})$  restricted strongly convex and smooth

along descent direction  $\mathbf{V}$  **only when  $\mathbf{X}$  is incoherent:**

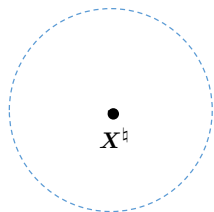
$$\|\mathbf{X} - \mathbf{X}^\natural\|_{2,\infty} \ll \|\mathbf{X}^\natural\|_{2,\infty}$$

## Incoherence region

Which region enjoys both restricted strong convexity and smoothness?

## Incoherence region

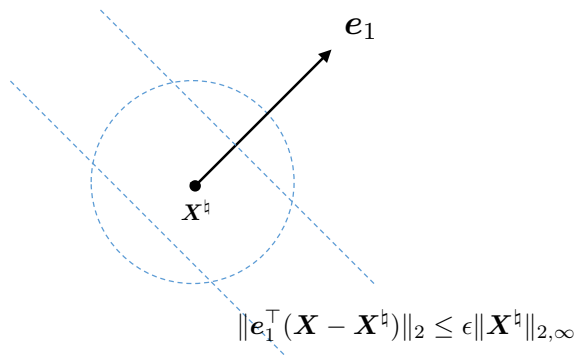
Which region enjoys both restricted strong convexity and smoothness?



- $\mathbf{X}$  is not far away from  $\mathbf{X}^h$

## Incoherence region

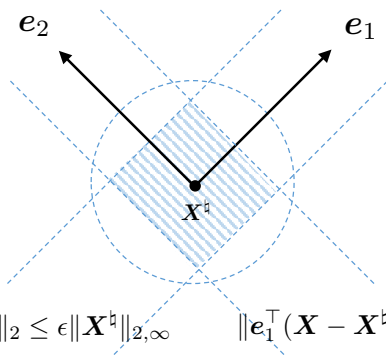
Which region enjoys both restricted strong convexity and smoothness?



- $\mathbf{X}$  is not far away from  $\mathbf{X}^b$
- $\mathbf{X}$  is incoherent w.r.t. coordinates (incoherence region)

## Incoherence region


Which region enjoys both restricted strong convexity and smoothness?

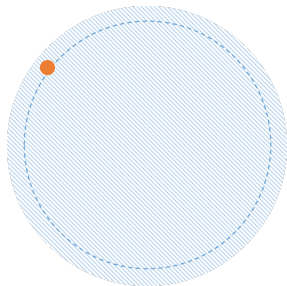


$$\|\mathbf{e}_2^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty} \quad \|\mathbf{e}_1^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty}$$

- $\mathbf{X}$  is not far away from  $\mathbf{X}^\dagger$
- $\mathbf{X}$  is incoherent w.r.t. coordinates (incoherence region)


# Vanilla gradient descent is at risk

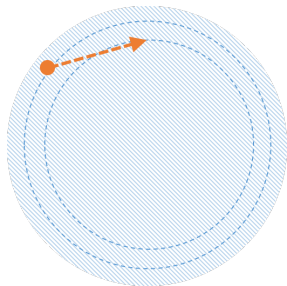
 region of local strong convexity + smoothness



*GD on the pop. loss*

# Vanilla gradient descent is at risk

 region of local strong convexity + smoothness

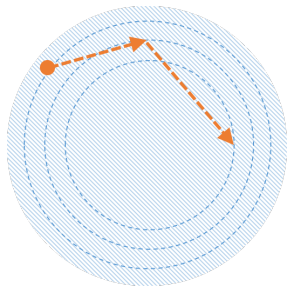


*GD on the pop. loss*



# Vanilla gradient descent is at risk

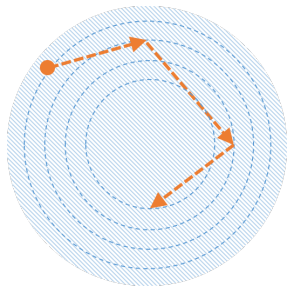
● region of local strong convexity + smoothness



*GD on the pop. loss*

# Vanilla gradient descent is at risk

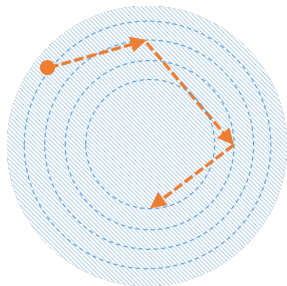
● region of local strong convexity + smoothness



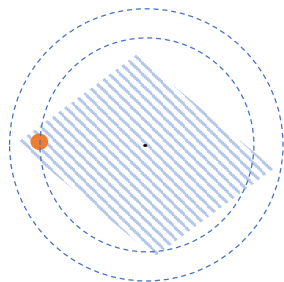
*GD on the pop. loss*

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

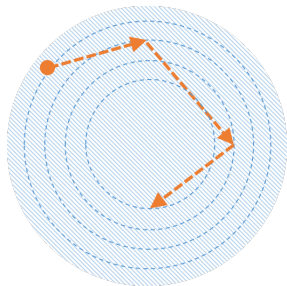


*GD on the emp. loss*

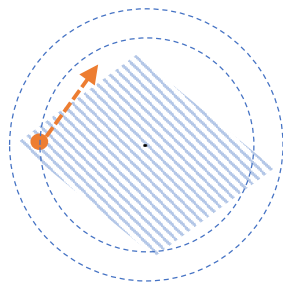
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

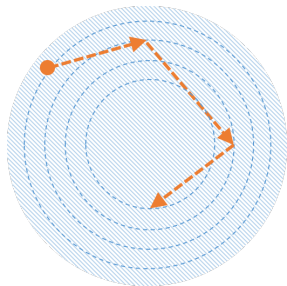


*GD on the emp. loss*

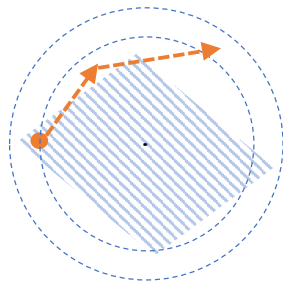
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

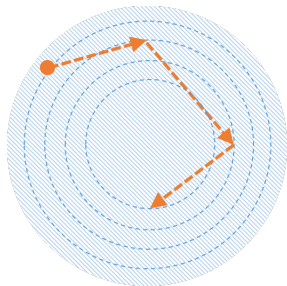


*GD on the emp. loss*

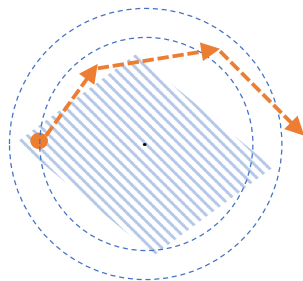
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

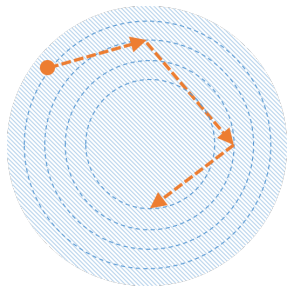


*GD on the emp. loss*

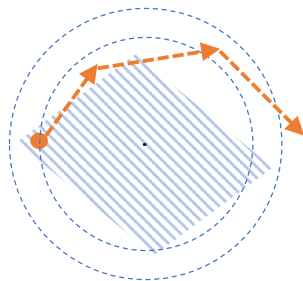
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

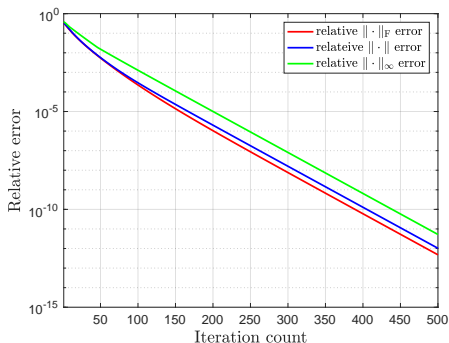


*GD on the emp. loss*

- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region
- Existing algorithms enforce regularization, or apply sample splitting to promote incoherence

# Matrix completion via vanilla GD

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left( \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$

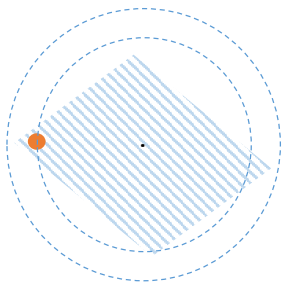


Vanilla GD converges fast without regularization!



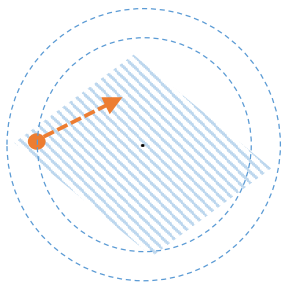
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



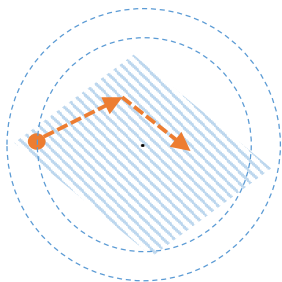
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



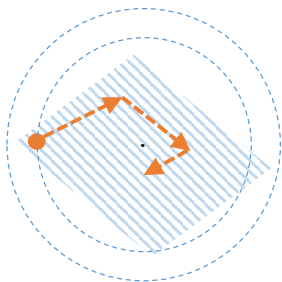
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



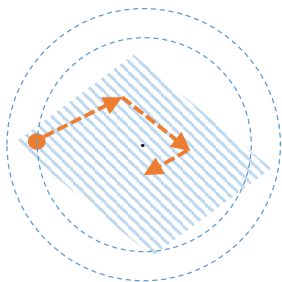
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**  
even without regularization

## Theoretical guarantees - noise-free case

### Theorem (Ma, Wang, Chi, Chen, FoCM 2019+)

Suppose  $M = \mathbf{X}^{\natural} \mathbf{X}^{\natural\top}$  is rank- $r$ , incoherent and well-conditioned. Vanilla GD (with spectral initialization) achieves

- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{\text{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^{\natural}\|_{\text{F}},$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^{\natural}\|, \quad (\text{spectral})$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{2,\infty} \lesssim \rho^t \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^{\natural}\|_{2,\infty}, \quad (\text{incoherence})$

where  $\rho = 1 - \frac{\sigma_{\min} \eta}{5} < 1$ , if step size  $\eta \asymp 1/\sigma_{\max}$  and sample complexity  $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$ .

## Theoretical guarantees - noise-free case

### Theorem (Ma, Wang, Chi, Chen, FoCM 2019+)

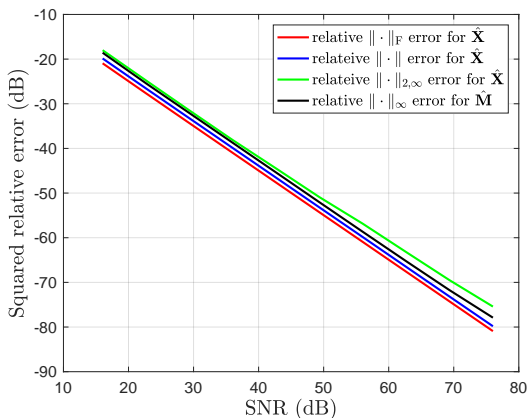
Suppose  $M = \mathbf{X}^{\natural} \mathbf{X}^{\natural\top}$  is rank- $r$ , incoherent and well-conditioned. Vanilla GD (with spectral initialization) achieves

- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{\text{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^{\natural}\|_{\text{F}},$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^{\natural}\|, \quad (\text{spectral})$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^{\natural}\|_{2,\infty} \lesssim \rho^t \mu r \sqrt{\frac{\log n}{np}} \|\mathbf{X}^{\natural}\|_{2,\infty}, \quad (\text{incoherence})$

where  $\rho = 1 - \frac{\sigma_{\min} \eta}{5} < 1$ , if step size  $\eta \asymp 1/\sigma_{\max}$  and sample complexity  $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$ .

- A recent follow-up by Xiaodong Li studied the rectangular case and improved the sample complexity to  $O(\mu^2 n r^2 \log n)$ .

# Noisy matrix completion via vanilla GD



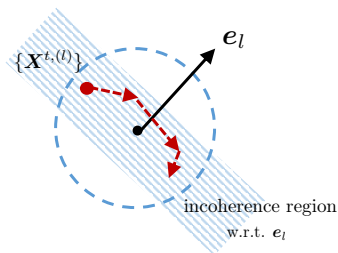
**Near-optimal entry-wise error control:**

$$\left\| \mathbf{X}^t \mathbf{X}^{t\top} - \mathbf{M}^\natural \right\|_\infty \lesssim \left( \rho^t \mu r \sqrt{\frac{\log n}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \left\| \mathbf{M}^\natural \right\|_\infty$$



## Key ingredient: leave-one-out analysis

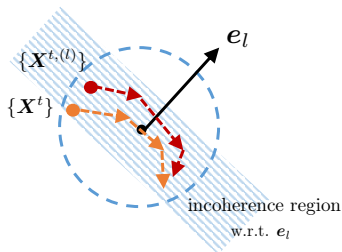
How to establish  $\|e_l^\top (\mathbf{X}^t - \mathbf{X}^\dagger)\|_2 \ll \|\mathbf{X}^\dagger\|_{2,\infty}$ ?



- Create auxiliary leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  are incoherent in the  $l$ th row;

## Key ingredient: leave-one-out analysis

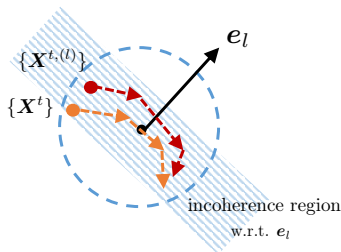
How to establish  $\|e_l^\top (\mathbf{X}^t - \mathbf{X}^\natural)\|_2 \ll \|\mathbf{X}^\natural\|_{2,\infty}$ ?



- Create auxiliary leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  are incoherent in the  $l$ th row;
- Leave-one-out iterates  $\mathbf{X}^{t,(l)} \approx$  true iterates  $\mathbf{X}^t$

## Key ingredient: leave-one-out analysis

How to establish  $\|e_l^\top (\mathbf{X}^t - \mathbf{X}^\natural)\|_2 \ll \|\mathbf{X}^\natural\|_{2,\infty}$ ?



- Create auxiliary leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  are incoherent in the  $l$ th row;
- Leave-one-out iterates  $\mathbf{X}^{t,(l)} \approx$  true iterates  $\mathbf{X}^t$
- $\|e_l^\top (\mathbf{X}^t - \mathbf{X}^\natural)\|_2 \leq \|e_l^\top (\mathbf{X}^{t,(l)} - \mathbf{X}^\natural)\|_2 + \|e_l^\top (\mathbf{X}^t - \mathbf{X}^{t,(l)})\|_2$

# An aside: stability of nuclear norm minimization

convex



nonconvex

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

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

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

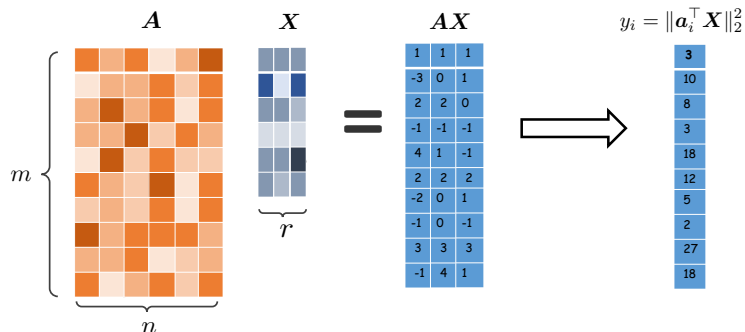
With high prob., any minimizer  $\widehat{\mathbf{M}}_{\text{cvx}}$  of convex program is nearly rank- $r$  and is minimax near-optimal:

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

Noisy Matrix Completion: Understanding Statistical Guarantees for Convex Relaxation via Nonconvex Optimization. [arXiv:1902.07698](https://arxiv.org/abs/1902.07698).

*The phenomenon is quite general*

# Generalized phase retrieval



Recover  $X^\natural \in \mathbb{R}^{n \times r}$  from  $m$  “random” quadratic measurements

$$y_i = \left\| a_i^\top X^\natural \right\|_2^2 = \langle a_i a_i^\top, X^\natural X^{\natural \top} \rangle, \quad i = 1, \dots, m$$

where  $a_i$ 's are i.i.d. Gaussian entries.

*Applications: optical imaging, phase space tomography ...*

## Implicit regularization for generalized phase retrieval

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \left\| \mathbf{a}_k^\top \mathbf{X} \right\|^2 - y_k \right)^2$$

# Implicit regularization for generalized phase retrieval

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \left\| \mathbf{a}_k^\top \mathbf{X} \right\|^2 - y_k \right)^2$$

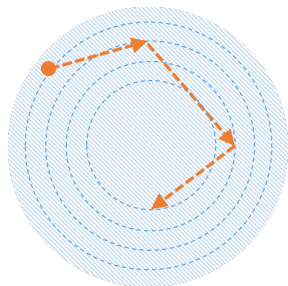
- region of local strong convexity + smoothness



# Implicit regularization for generalized phase retrieval

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \|\mathbf{a}_k^\top \mathbf{X}\|^2 - y_k \right)^2$$

● region of local strong convexity + smoothness

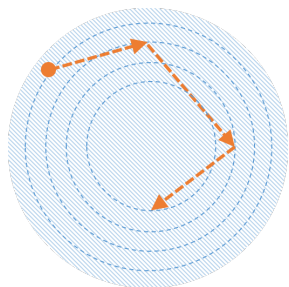


$$O(1) \preceq \nabla^2 f(\mathbf{x}) \preceq O(n)$$

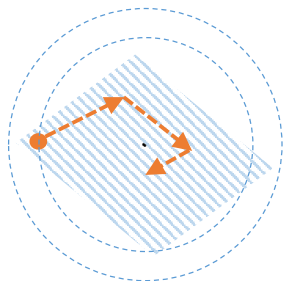
# Implicit regularization for generalized phase retrieval

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \|\mathbf{a}_k^\top \mathbf{X}\|^2 - y_k \right)^2$$

● region of local strong convexity + smoothness



$$O(1) \preceq \nabla^2 f(\mathbf{x}) \preceq O(n)$$



$$O(1) \preceq \nabla^2 f(\mathbf{x}) \preceq O(\log n)$$

# Theoretical guarantees

## Theorem (Li, Ma, Chen, Chi, AISTATS 2019)

*Under i.i.d. Gaussian design, GD achieves linear convergence*

- $\max_k \|\mathbf{a}_k^\top (\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural)\| \lesssim \sqrt{\log n} \frac{\sigma_r^2(\mathbf{X}^\natural)}{\|\mathbf{X}^\natural\|_F}$  (incoherence)

# Theoretical guarantees

## Theorem (Li, Ma, Chen, Chi, AISTATS 2019)

Under i.i.d. Gaussian design, GD achieves linear convergence

- $\max_k \|\mathbf{a}_k^\top (\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural)\| \lesssim \sqrt{\log n} \frac{\sigma_r^2(\mathbf{X}^\natural)}{\|\mathbf{X}^\natural\|_F}$  (incoherence)
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural\|_F \lesssim \left(1 - \frac{\sigma_r^2(\mathbf{X}^\natural)\eta}{2}\right)^t \|\mathbf{X}^\natural\|_F$  (linear convergence)

provided that  $\eta \asymp \frac{1}{(\log n \vee r)^2 \sigma_r^2(\mathbf{X}^\natural)}$  and  $m \gtrsim nr^4 \log n$ .

# Theoretical guarantees

## Theorem (Li, Ma, Chen, Chi, AISTATS 2019)

Under i.i.d. Gaussian design, GD achieves linear convergence

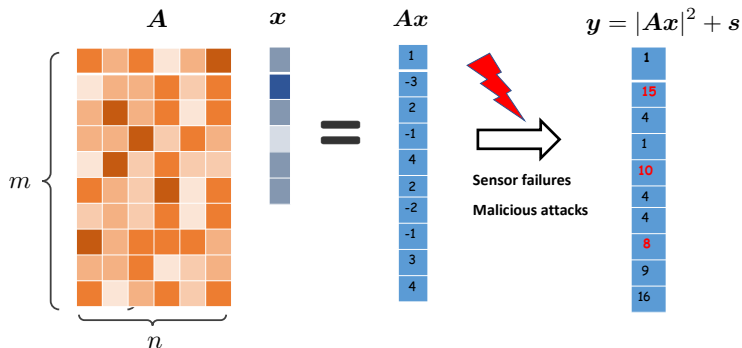
- $\max_k \|\mathbf{a}_k^\top (\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural)\| \lesssim \sqrt{\log n} \frac{\sigma_r^2(\mathbf{X}^\natural)}{\|\mathbf{X}^\natural\|_F}$  (incoherence)
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^\natural\|_F \lesssim \left(1 - \frac{\sigma_r^2(\mathbf{X}^\natural)\eta}{2}\right)^t \|\mathbf{X}^\natural\|_F$  (linear convergence)

provided that  $\eta \asymp \frac{1}{(\log n \vee r)^2 \sigma_r^2(\mathbf{X}^\natural)}$  and  $m \gtrsim nr^4 \log n$ .

**Big computational saving:** GD attains  $\varepsilon$ -accuracy within  $O((\log n \vee r)^2 \log \frac{1}{\varepsilon})$  iterations if  $m \asymp nr^4 \log n$ .

*Towards robust nonconvex statistical estimation*

# Outlier-corrupted phase retrieval

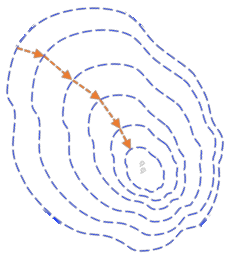


Recover  $\mathbf{x}^\dagger \in \mathbb{R}^n$  from  $m$  corrupted measurements

$$y_i = \left| \mathbf{a}_i^\top \mathbf{x}^\dagger \right|_2^2 + s_i, \quad i = 1, \dots, m$$

where  $\|s\|_0 \leq \alpha \cdot m$ ,  $0 \leq \alpha < 1$  is fraction of outliers.

## Existing approaches fail



- **Initialization would fail:**  $\mathbf{x}^0 \leftarrow$  leading eigenvector of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top$$

- **Gradient iterations would fail:**

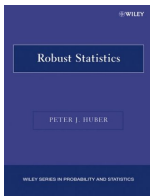
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_{i=1}^m \nabla l_i(y_i; \mathbf{x}^t)$$

for  $t = 0, 1, \dots$

Even a single outlier can fail the algorithm!

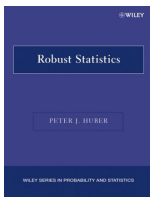


# Median-truncated gradient descent



**Key idea: “median-truncation”** —  
discard samples *adaptively* based on  
how large sample gradients / values  
deviate from median

# Median-truncated gradient descent



**Key idea: “median-truncation”** — discard samples *adaptively* based on how large sample gradients / values deviate from median

- **Robustifying spectral initialization:**  $\mathbf{x}^0 \leftarrow$  leading eigenvector of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top \mathbb{1}_{\{|y_i| \lesssim \text{median}\{y_i\}\}}$$

- **Robustifying gradient descent:**

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_{i \in \mathcal{T}_t} \nabla \ell_i(y_i; \mathbf{x}^t), \quad t = 0, 1, \dots$$

where  $\mathcal{T}_t = \{i : |y_i - |\mathbf{a}_i^\top \mathbf{x}^t|| \lesssim \text{median} \{|y_i - |\mathbf{a}_i^\top \mathbf{x}^t||\}\}$ .

# Theoretical guarantees

## Theorem (Zhang, Chi and Liang, TIT 2019)

*Under i.i.d. Gaussian design, median-truncated GD achieves linear convergence*

- $\|\mathbf{x}^t - \mathbf{x}^\natural\|_2 \lesssim \left(1 - \frac{\eta}{2}\right)^t \|\mathbf{x}^\natural\|_2$  (linear convergence)

*for  $\eta \asymp 1$ , provided that  $m \gtrsim n \log n$  and  $\alpha \lesssim \alpha_0$  for some constant  $\alpha_0$ .*

# Theoretical guarantees

## Theorem (Zhang, Chi and Liang, TIT 2019)

*Under i.i.d. Gaussian design, median-truncated GD achieves linear convergence*

- $\|\mathbf{x}^t - \mathbf{x}^\natural\|_2 \lesssim \left(1 - \frac{\eta}{2}\right)^t \|\mathbf{x}^\natural\|_2$  (linear convergence)

*for  $\eta \asymp 1$ , provided that  $m \gtrsim n \log n$  and  $\alpha \lesssim \alpha_0$  for some constant  $\alpha_0$ .*

**Add-on robustness:** GD attains  $\varepsilon$ -accuracy within  $O\left(\log \frac{1}{\varepsilon}\right)$  iterations if  $m \gtrsim n \log n$  even with a constant fraction of arbitrary outliers.

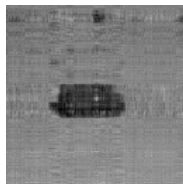
## Extension to low-rank matrix recovery

Similar idea for compressive low-rank matrix recovery:

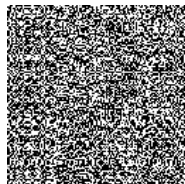
$$y_i = \langle \mathbf{A}_i, \mathbf{X}^\natural \rangle + s_i, \quad i = 1, \dots, m$$



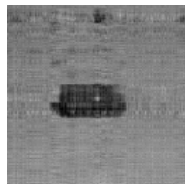
Ground truth



GD  
no outliers



GD  
1% outliers



median-TGD  
1% outliers

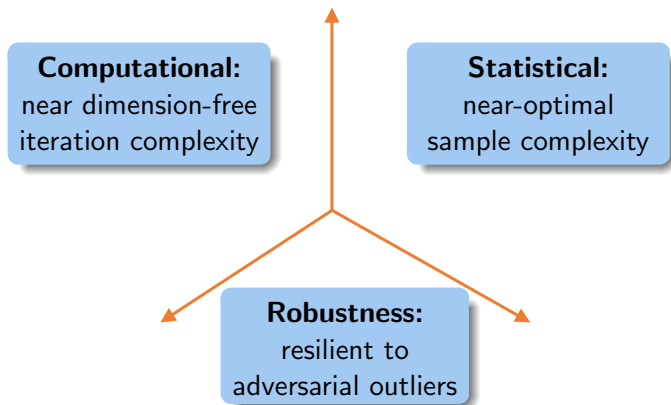
**Figure:** Recovery performance comparisons for compressive recovery of a  $128 \times 128$  image from  $m = 4600$  measurements and assumed rank  $r = 8$ .

---

Li, Chi, Zhang and Liang, "Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent", *Information and Inference: A Journal of the IMA*, 2019+.

*Final remarks*

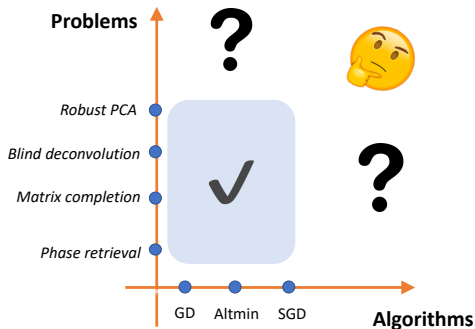
# Bridging the theory-practice gap



## Fusing statistical thinkings into nonconvex optimization:

- identification and exploitation of benign geometric properties;
- analyzing iterate trajectories beyond black-box optimization.

# Limitations



- current analysis is largely case-by-case: lengthy proofs, somewhat similar recipes;
- Is there a unified framework? E.g., RIP for sparsity.
- Can we relax strong randomness assumptions, e.g. Gaussian (phase retrieval), and uniform sampling (matrix completion)?



## Survey, tutorial articles:

1. Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview, **Y. Chi**, Y. M. Lu and Y. Chen, overview article, *IEEE Trans. on Signal Processing*, accepted, arXiv:1809.09573.
2. Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation, Y. Chen and **Y. Chi**, *IEEE Signal Processing Magazine*, 2018.

## Selected articles:

1. Implicit Regularization for Nonconvex Statistical Estimation, C. Ma, K. Wang, **Y. Chi** and Y. Chen, *Foundations of Computational Mathematics*, accepted.
2. Nonconvex Matrix Factorization from Rank-One Measurements, Y. Li, C. Ma, Y. Chen, and **Y. Chi**, AISTATS 2019.
3. Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent, Y. Li, Y. Chi, H. Zhang and Y. Liang, *Information and Inference: A Journal of the IMA*, 2019+.
4. Median-Truncated Nonconvex Approach for Phase Retrieval with Outliers, H. Zhang, **Y. Chi** and Y. Liang, *IEEE Trans. on Information Theory*, 2019.
5. Gradient Descent with Random Initialization: Fast Global Convergence for Nonconvex Phase Retrieval, Y. Chen, **Y. Chi**, J. Fan and C. Ma, *Mathematical Programming*, 2019.

# Thanks!

Our research is supported by NSF, ONR and ARO.



<https://users.ece.cmu.edu/~yuejiec/>