

# Geometry and Regularization in Nonconvex Low-Rank Estimation

Yuejie Chi

**Carnegie Mellon University**

University of Texas, Austin  
Feb. 2019

# Acknowledgements

Thanks to my collaborators:



Y. Chen  
Princeton



C. Ma  
Princeton



K. Wang  
Princeton



Y. Li  
CMU

This research is supported by NSF, ONR, AFOSR and ARO.



## Empirical risk minimization

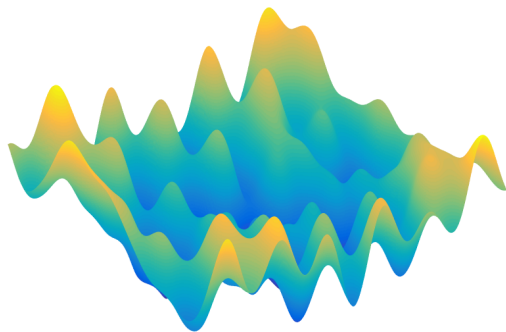
Given data  $z$ , estimate parameters  $\mathbf{x} \in \mathbb{R}^n$ :

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

# Empirical risk minimization

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

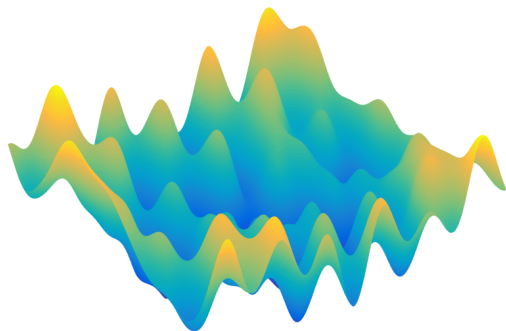
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) := \frac{1}{m} \sum_{i=1}^m \ell(z_i; \mathbf{x})$$



## Empirical risk minimization

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

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

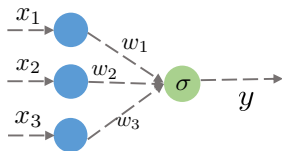


There may be exponentially many local optima

# Exponentially many local minima

Given training data  $\{\mathbf{x}_i, y_i\}_{i=1}^m$ ,

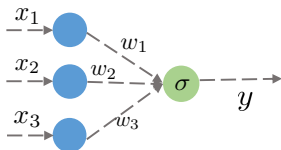
$$\text{minimize}_{\mathbf{w} \in \mathbb{R}^n} \ell_m(\mathbf{w}) := \frac{1}{2m} \sum_{i=1}^m \left( y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i) \right)^2$$



# Exponentially many local minima

Given training data  $\{\mathbf{x}_i, y_i\}_{i=1}^m$ ,

$$\text{minimize}_{\mathbf{w} \in \mathbb{R}^n} \ell_m(\mathbf{w}) := \frac{1}{2m} \sum_{i=1}^m \left( y_i - \sigma(\mathbf{w}^\top \mathbf{x}_i) \right)^2$$



## Theorem (Auer et al., 1996)

Let  $\sigma(\cdot)$  be sigmoid and  $\ell(\cdot)$  be the quadratic loss function. There *exists* a sequence of training samples  $\{\mathbf{x}_i, y_i\}_{i=1}^m$  such that  $\ell_m(\mathbf{w})$  has  $\lfloor \frac{m}{n} \rfloor^n$  distinct local minima.

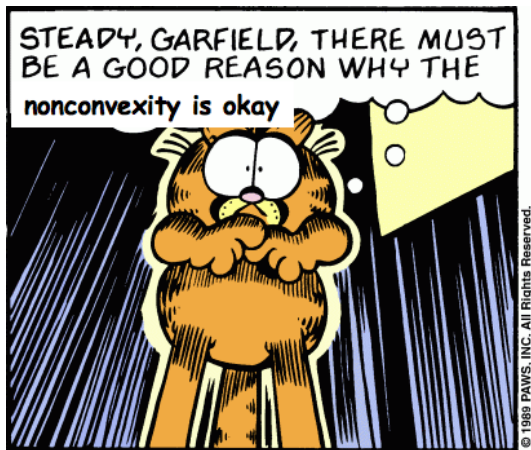
Nonconvex problems are hard!





## But they're solved on a daily basis in practice

Using simple algorithms such as gradient descent, e.g., “back propagation” for training deep neural networks...

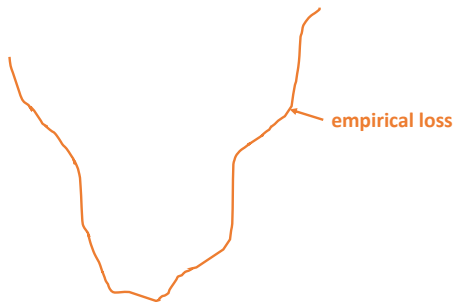




# Statistical thinking in nonconvex optimization

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.

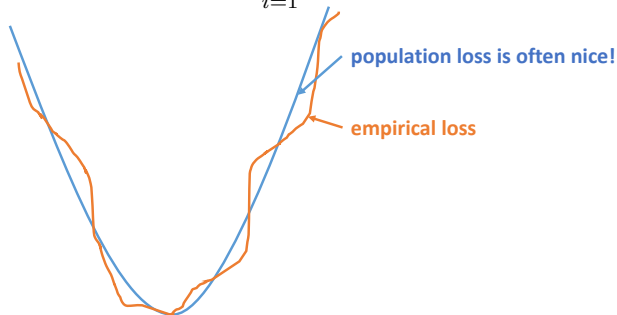
$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{x}) \quad \xrightarrow{m \rightarrow \infty} \quad \mathbb{E}[\ell(y; \mathbf{x})]$$



# Statistical thinking in nonconvex optimization

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

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



*When  $m \rightarrow \infty$ , empirical risk  $\approx$  population risk!*

## From population risk to empirical risk

**Sample-starved / finite-sample regime:**

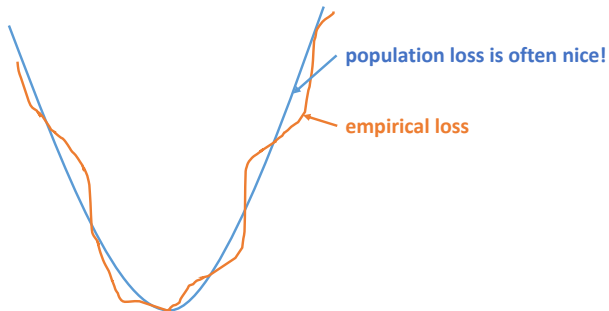
sample size  $m \approx \tilde{O}(n)$ , the number of parameters.

# From population risk to empirical risk

## Sample-starved / finite-sample regime:

sample size  $m \approx \tilde{O}(n)$ , the number of parameters.

*Even when  $\mathbb{E}[f(x)]$  is (locally) strongly convex and smooth, the empirical loss  $f(x)$  may **not** be when  $m \approx \tilde{O}(n)$ .*

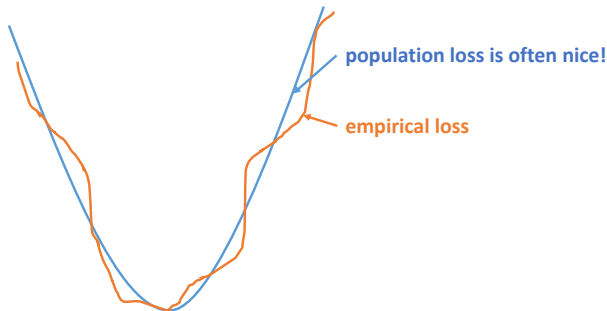


# From population risk to empirical risk

**Sample-starved / finite-sample regime:**

sample size  $m \approx \tilde{O}(n)$ , the number of parameters.

*Even when  $\mathbb{E}[f(\mathbf{x})]$  is (locally) strongly convex and smooth, the empirical loss  $f(\mathbf{x})$  may **not** be when  $m \approx \tilde{O}(n)$ .*



$f(\mathbf{x})$  may lack curvatures in certain regions/directions!

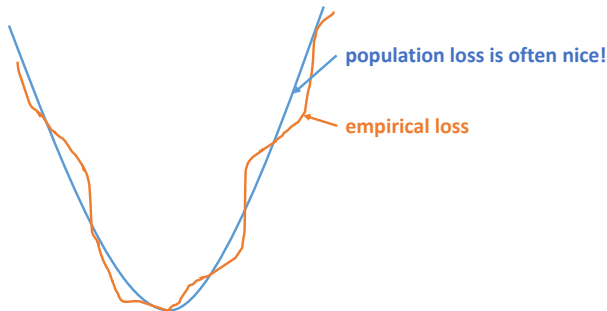


# From population risk to empirical risk

**Sample-starved / finite-sample regime:**

sample size  $m \approx \tilde{O}(n)$ , the number of parameters.

*Even when  $\mathbb{E}[f(\mathbf{x})]$  is (locally) strongly convex and smooth, the empirical loss  $f(\mathbf{x})$  may **not** be when  $m \approx \tilde{O}(n)$ .*



$f(\mathbf{x})$  may lack curvatures in certain regions/directions!

*Will this be problematic?*

*a case study with low-rank matrix estimation*

# Rethinking PCA for modern datasets

August 2018 | Volume 106 | Number 8

## Proceedings OF THE IEEE

SPECIAL ISSUE

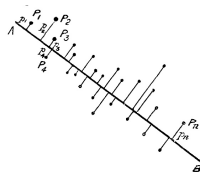
### Rethinking Principal Component Analysis (PCA) for Modern Data Sets

Point of View: The Twin Arts of Writing and Revising Technical Articles

Scanning Our Past: Imperial Science: Victorian Cable Telegraphy and the Making of "Maxwell's Equations"



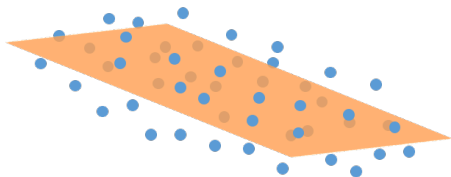
Classical PCA (Pearson 1901):



Modern PCA (Big, Messy Data):



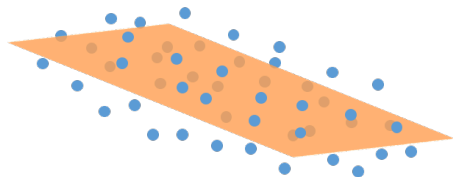
## Revisiting PCA: in search of low-rank representation



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

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

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



This problem admits a closed-form solution:

- let  $\mathbf{M} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$  be eigen-decomposition of  $\mathbf{M}$  ( $\lambda_1 \geq \dots \geq \lambda_n$ ), then

$$\widehat{\mathbf{M}} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

— *nonconvex, but tractable*

## An optimization viewpoint

**Burer-Monteiro 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:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

## An optimization viewpoint

**Burer-Monteiro 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:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

- **Pro:** reduce parameter space from  $O(n^2)$  to  $O(n)$ ;
- **Con:** nonconvex and susceptible to local optima.

## An optimization viewpoint

**Burer-Monteiro 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:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_F^2$$

- **Pro:** reduce parameter space from  $O(n^2)$  to  $O(nr)$ ;
- **Con:** nonconvex and susceptible to local optima.

**Theorem (PCA doesn't have spurious local minima, Baldi and Hornik, 1989)**

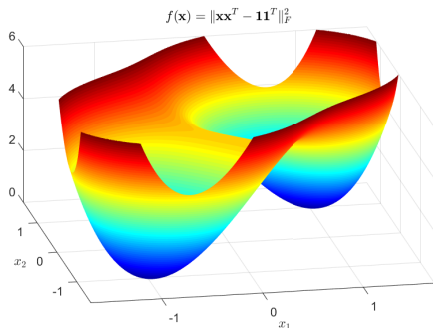
*Suppose  $\mathbf{M}$  has a strict eigen-gap between  $\lambda_r$  and  $\lambda_{r+1}$ , the critical points of  $f(\mathbf{X})$  can be categorized into*

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



# Benign landscape of PCA

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

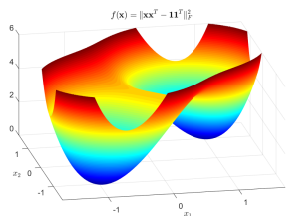


global minima:  $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ; strict saddles:  $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , and  $\pm \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

— No “spurious” local minima!

# Parameter recovery via gradient descent

a two-step recovery strategy:



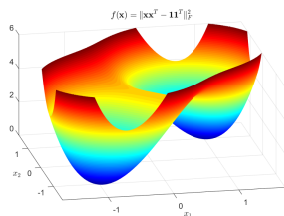
- Initialization by spectral method
- Gradient iterations:

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

for  $t = 0, 1, \dots$

# Parameter recovery via gradient descent

a two-step recovery strategy:



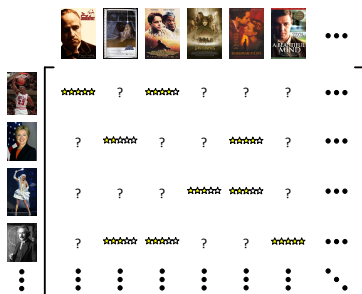
- **Initialization by spectral method**
- **Gradient iterations:**

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

for  $t = 0, 1, \dots$

- The initial point falls into a “basin of attraction”;
- Low-complexity local refinements via gradient descent.

# Low-rank matrix completion: dealing with missing data



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

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

*Applications: recommendation systems, ...*

## A natural least-squares formulation

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

$\Downarrow$

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

## A natural least-squares formulation

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

$\Downarrow$

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

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

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

*does not imply optimization efficiency!*

# 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

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

## 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]:  $p \gtrsim \mu r \log n/n$ .

## What does the population 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 population level look like?

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

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

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

$$\begin{aligned} & \text{vec}(\mathbf{V})^\top \mathbb{E} [\nabla^2 f(\mathbf{X})] \text{vec}(\mathbf{V}) \\ &= \underbrace{\frac{1}{2} \left\| \mathbf{V} \mathbf{X}^\top + \mathbf{X} \mathbf{V}^\top \right\|_F^2 + \left\langle \mathbf{X} \mathbf{X}^\top - \mathbf{X}^\natural \mathbf{X}^{\natural\top}, \mathbf{V} \mathbf{V}^\top \right\rangle}_{\text{locally restricted strongly convex and smooth}} \end{aligned}$$

along descent direction  $\mathbf{V}$  when  $\mathbf{X}$  is close to  $\mathbf{X}^\natural$ .

**Consequence:** GD converges within  $O\left(\log \frac{1}{\varepsilon}\right)$  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** ( $p \asymp \frac{\text{polylog} n}{n}$ )

$\nabla^2 f(\mathbf{X})$  strongly convex and smooth

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

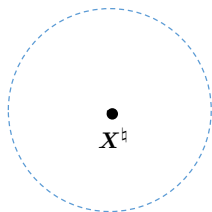
$$\|\mathbf{X} - \mathbf{X}^\natural\|_{2,\infty} \ll \|\mathbf{X}^\natural\|_{2,\infty}$$

## Incoherence region

Which region enjoys both restricted strong convexity and smoothness?

## Incoherence region

Which region enjoys both restricted strong convexity and smoothness?

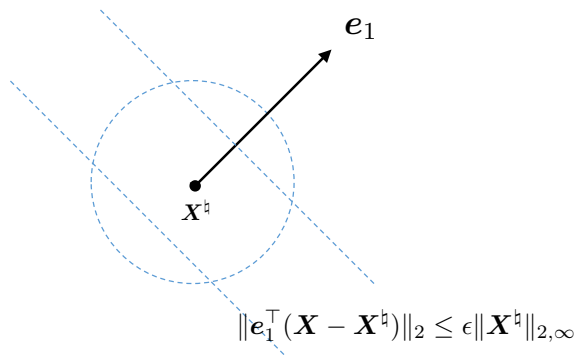


- $\mathbf{X}$  is not far away from  $\mathbf{X}^h$



## Incoherence region

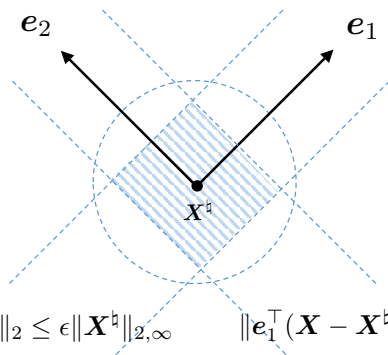
Which region enjoys both restricted strong convexity and smoothness?



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

## Incoherence region


Which region enjoys both restricted strong convexity and smoothness?

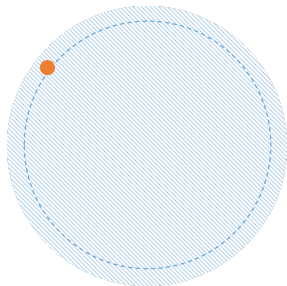


$$\|\mathbf{e}_2^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty} \quad \|\mathbf{e}_1^\top (\mathbf{X} - \mathbf{X}^\dagger)\|_2 \leq \epsilon \|\mathbf{X}^\dagger\|_{2,\infty}$$

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

# Vanilla gradient descent is at risk

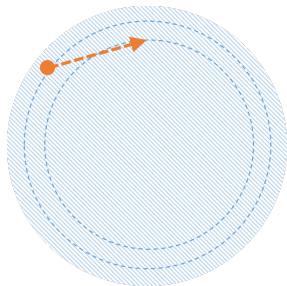
 region of local strong convexity + smoothness



*GD on the pop. loss*

# Vanilla gradient descent is at risk

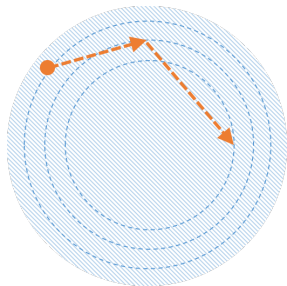
● region of local strong convexity + smoothness



*GD on the pop. loss*

# Vanilla gradient descent is at risk

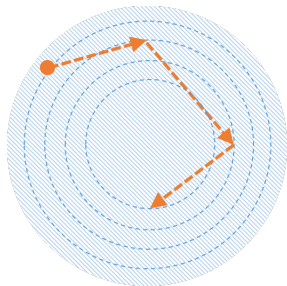
● region of local strong convexity + smoothness



*GD on the pop. loss*

# Vanilla gradient descent is at risk

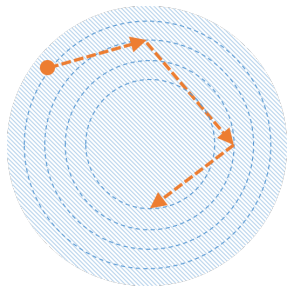
● region of local strong convexity + smoothness



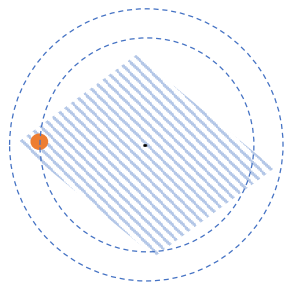
*GD on the pop. loss*

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

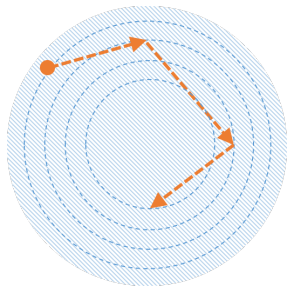


*GD on the emp. loss*

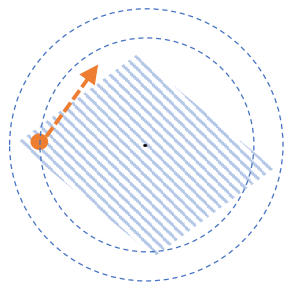
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*



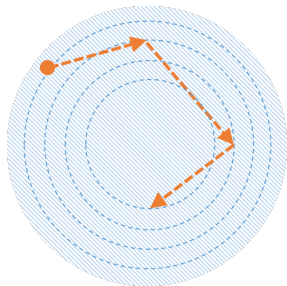
*GD on the emp. loss*

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

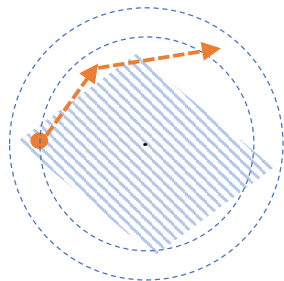


# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

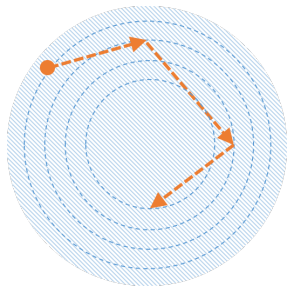


*GD on the emp. loss*

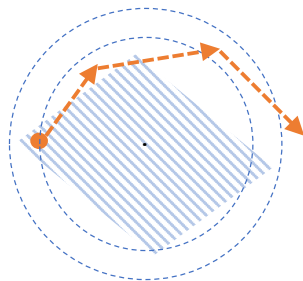
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

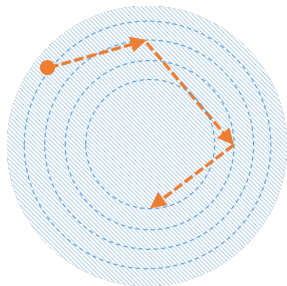


*GD on the emp. loss*

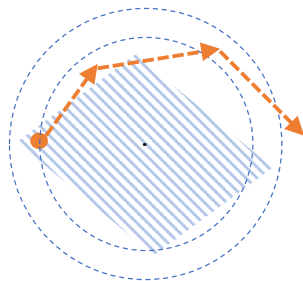
- Generic optimization theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



*GD on the pop. loss*

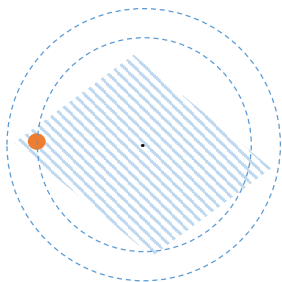


*GD on the emp. loss*

- 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

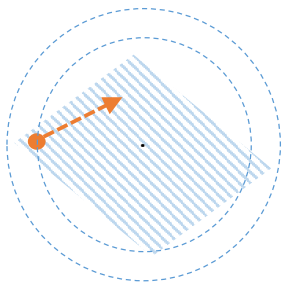
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



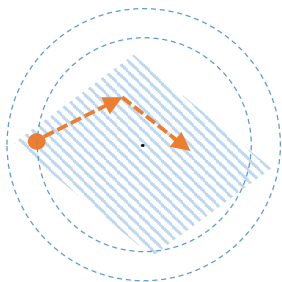
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



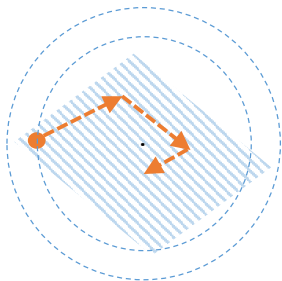
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



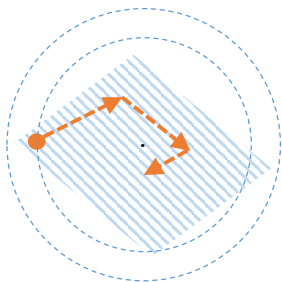
## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



## Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness

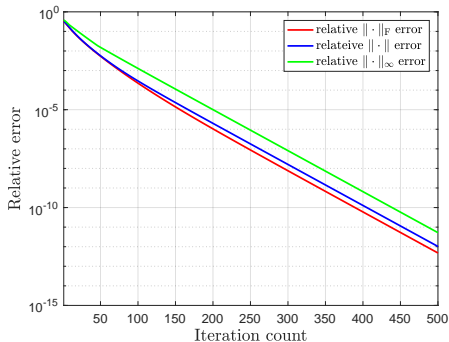


GD implicitly forces iterates to remain **incoherent**  
even without regularization



# Matrix completion via vanilla GD

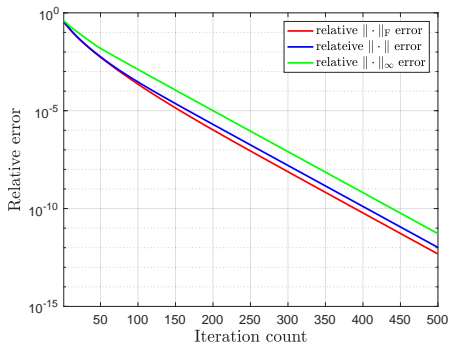
$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \sum_{(j,k) \in \Omega} \left( \mathbf{e}_j^\top \mathbf{X} \mathbf{X}^\top \mathbf{e}_k - M_{j,k} \right)^2$$



Vanilla GD converges fast without regularization!

# Matrix completion via vanilla GD

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



Our finding: GD is implicitly regularized to stay incoherent!

## Theoretical guarantees - noise-free case

### Theorem (Ma, Wang, Chi, Chen)

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

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

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)

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 improves the sample complexity to  $O(\mu^2 n r^2 \log n)$ .

# Noisy matrix completion via vanilla GD

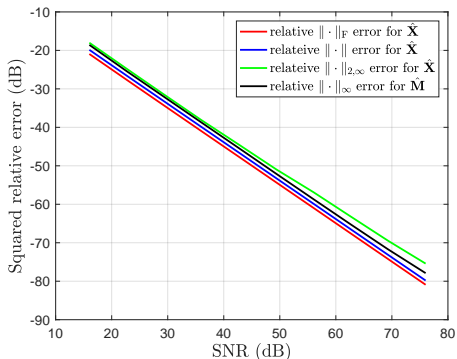


Figure: The relative error of  $\hat{\mathbf{X}}$  in  $\|\cdot\|_F$ ,  $\|\cdot\|$ ,  $\|\cdot\|_{2,\infty}$  vs.  $\text{SNR} := \frac{\|\mathbf{M}^\natural\|_F^2}{n^2\sigma^2}$ .

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

## Key ingredient: leave-one-out analysis

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

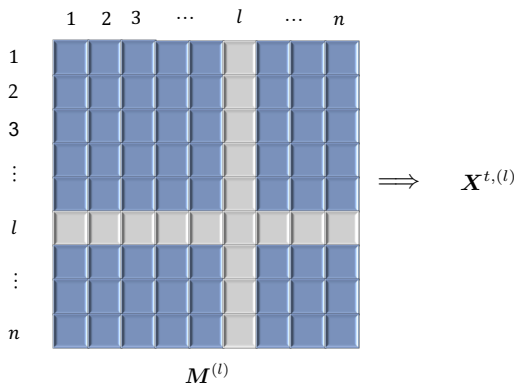
Technical difficulty:  $\mathbf{X}^t$  is statistically dependent with  $\Omega$ ;

## Key ingredient: leave-one-out analysis

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

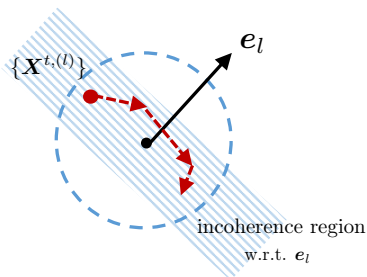
**Technical difficulty:**  $\mathbf{X}^t$  is statistically dependent with  $\Omega$ ;

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



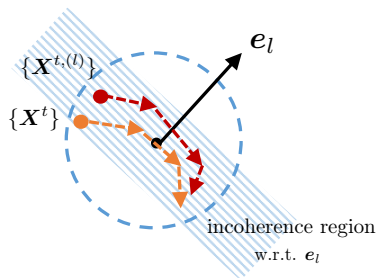


## Key ingredient: leave-one-out analysis



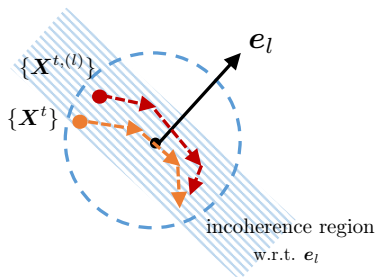
- Leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  contains more information of  $l$ th row of  $\mathbf{X}^{\dagger}$ ; indep. of randomness in  $l$ th row

## Key ingredient: leave-one-out analysis



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

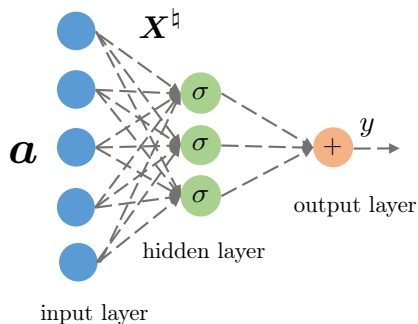
## Key ingredient: leave-one-out analysis



- Leave-one-out iterates  $\{\mathbf{X}^{t,(l)}\}$  contains more information of  $l$ th row of  $\mathbf{X}^\natural$ ; indep. of randomness in  $l$ th row
- Leave-one-out iterates  $\mathbf{X}^{t,(l)} \approx$  true iterates  $\mathbf{X}^t$
- $\|e_l^\top (\mathbf{X}^t - \mathbf{X}^\natural)\|_2 \leq \|e_l^\top (\mathbf{X}^{t,(l)} - \mathbf{X}^\natural)\|_2 + \|e_l^\top (\mathbf{X}^t - \mathbf{X}^{t,(l)})\|_2$

*The phenomenon is quite general*

## Shallow neural network with quadratic activation

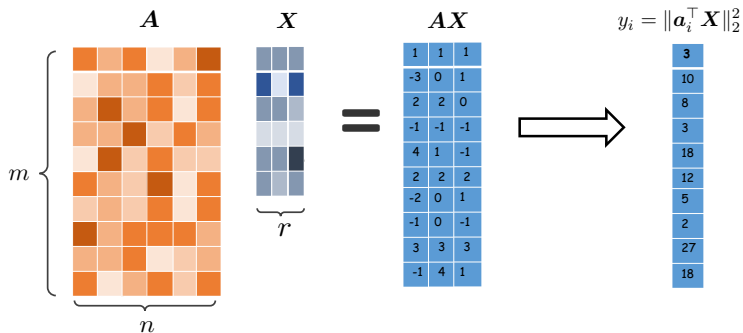


Set  $\mathbf{X}^h = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r]$ , then

$$y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i) \stackrel{\sigma(z) := z^2}{=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i)^2 = \left\| \mathbf{a}^\top \mathbf{X}^h \right\|_2^2.$$

*Identifiability up to orthonormal transform of  $\mathbf{X}^h$ .*

# Generalized phase retrieval



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

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

*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

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

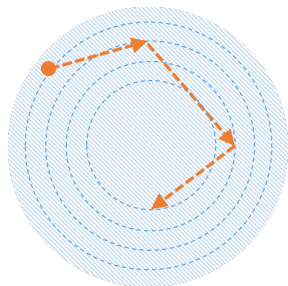
- region of local strong convexity + smoothness



# 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

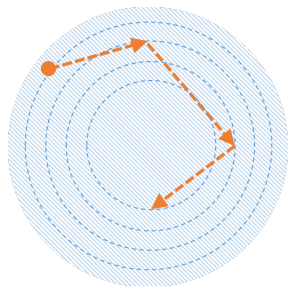


$$O(1) \preceq \nabla^2 f(\mathbf{x}) \preceq O(n)$$

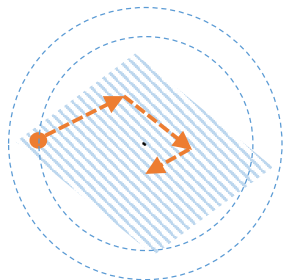
# Implicit regularization for generalized phase retrieval

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

● region of local strong convexity + smoothness



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



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

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

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

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

# Theoretical guarantees

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

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)

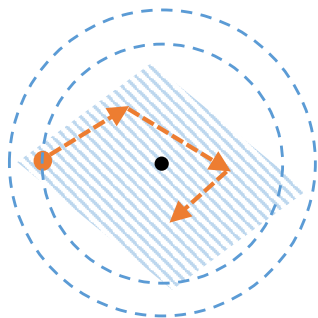
Under i.i.d. Gaussian design, GD achieves linear convergence

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

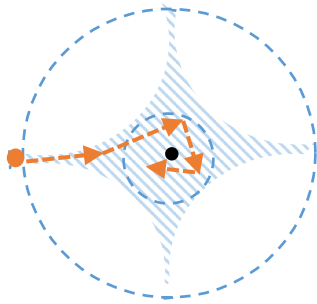
provided that  $\eta \asymp \frac{1}{(\log n \vee r)^2 \sigma_r^2(\mathbf{X}^\natural)}$  and  $m \gtrsim nr^4 \log n$ .

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

# Incoherence region in high dimensions



2-dimensional



high-dimensional

incoherence region is vanishingly small

# Conclusions

**From population loss from empirical loss:** vanilla gradient descent exploits local hidden convexity as if it almost runs on the population loss!

**Computational:**  
near dimension-free  
iteration complexity

**Statistical:**  
near-optimal  
sample complexity

**Analytical toolkits:** leave-one-out perturbation argument to establish near-independence. Useful in other problems!

## References

1. Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview, Y. Chi, Y. M. Lu and Y. Chen, arXiv:1809.09573.
2. Implicit Regularization for Nonconvex Statistical Estimation, C. Ma, K. Wang, Y. Chi and Y. Chen, arXiv:1711.10467.
3. Nonconvex Matrix Factorization from Rank-One Measurements, Y. Li, C. Ma, Y. Chen, and Y. Chi, arXiv:1802.06286.
4. Gradient Descent with *Random Initialization*: Fast Global Convergence for Nonconvex Phase Retrieval, Y. Chen, Y. Chi, J. Fan and C. Ma, arXiv:1803.07726.
5. Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation, Y. Chen and Y. Chi, survey article on IEEE Signal Processing Magazine, Jul. 2018.

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

Thank you!