

# Nonconvex Matrix Completion without Regularization

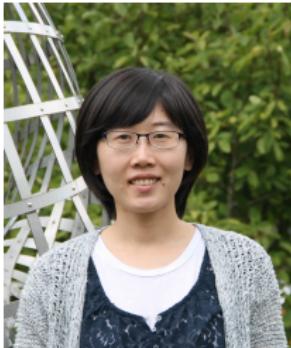


Cong Ma

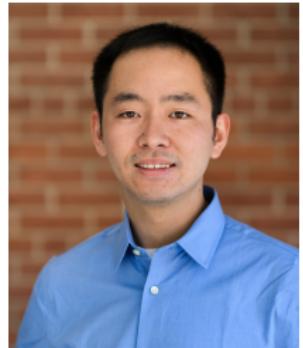
ORFE, Princeton University



Kaizheng Wang  
Princeton ORFE



Yuejie Chi  
CMU ECE



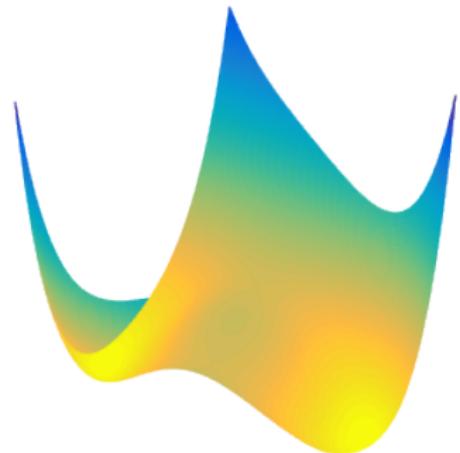
Yuxin Chen  
Princeton EE

# Nonconvex problems are everywhere

---

Empirical risk minimization is usually nonconvex

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



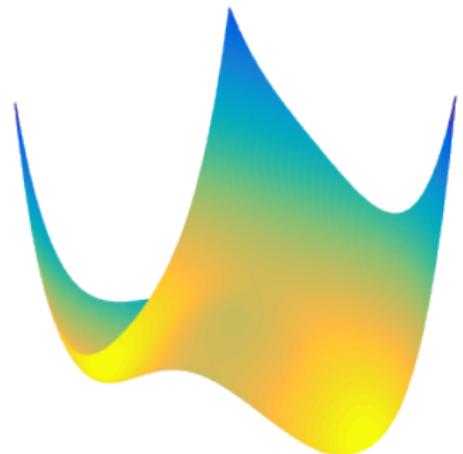
# Nonconvex problems are everywhere

---

Empirical risk minimization is usually nonconvex

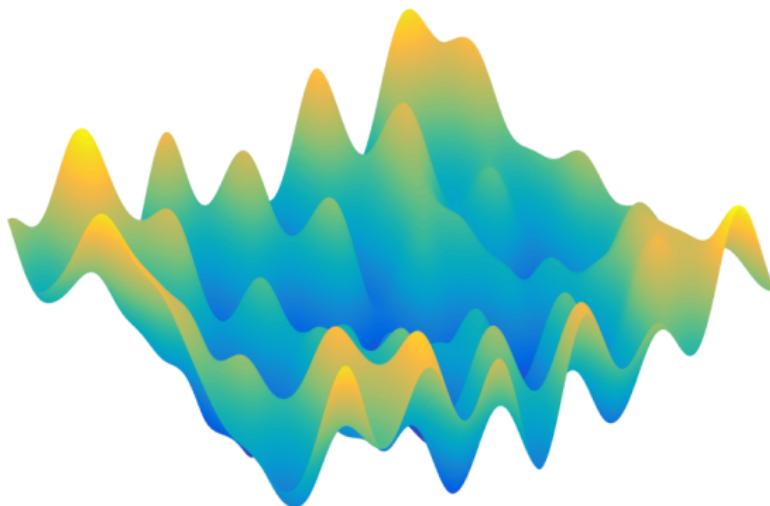
$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}; \text{data})$$

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



# Nonconvex optimization may be super scary

---

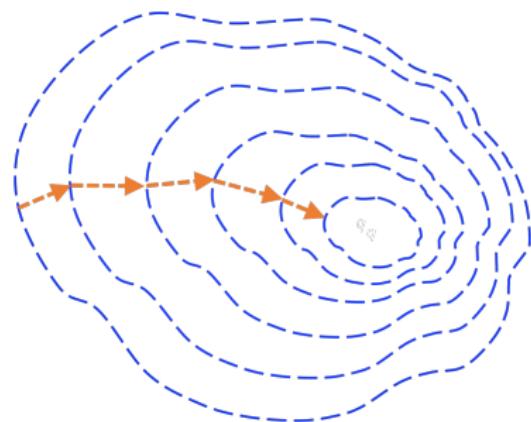
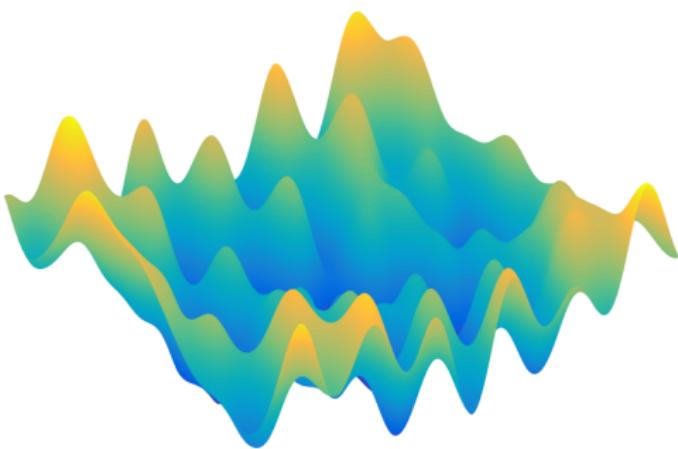


There may be bumps everywhere and exponentially many local optima

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

# Nonconvex optimization may be super scary

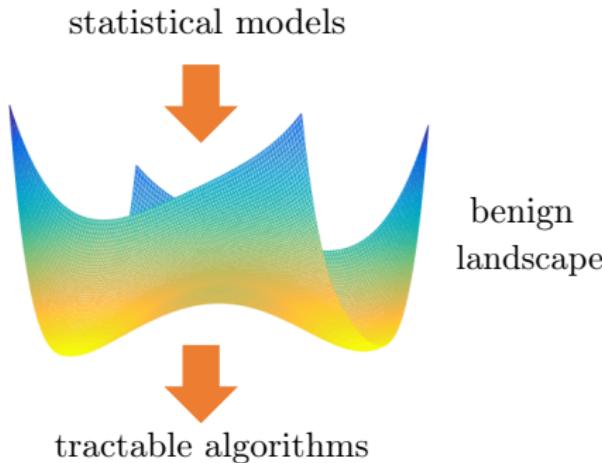
---



But they are solved on a daily basis via simple algorithms like  
*(stochastic) gradient descent*

# Statistical models come to rescue

---



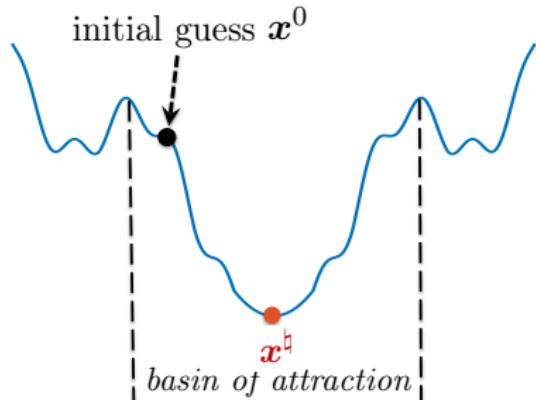
When data are generated by certain statistical models, problems are often much nicer than worst-case instances

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

Chi, Lu, Chen '18

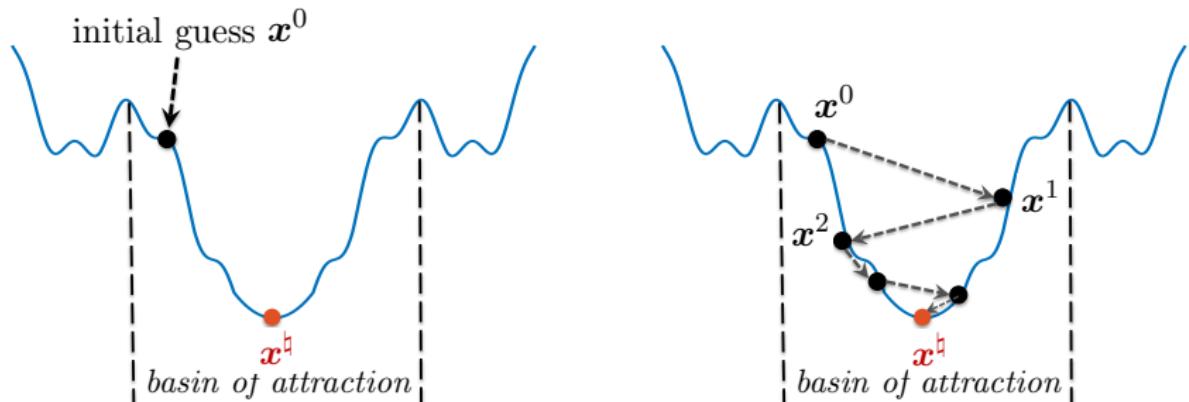
## A popular two-stage approach

---



1. initialize within local basin sufficiently close to  $x^*$   
(restricted) strongly convex and smooth

# A popular two-stage approach



1. initialize within local basin sufficiently close to  $x^*$   
(restricted) strongly convex and smooth
2. iterative refinement

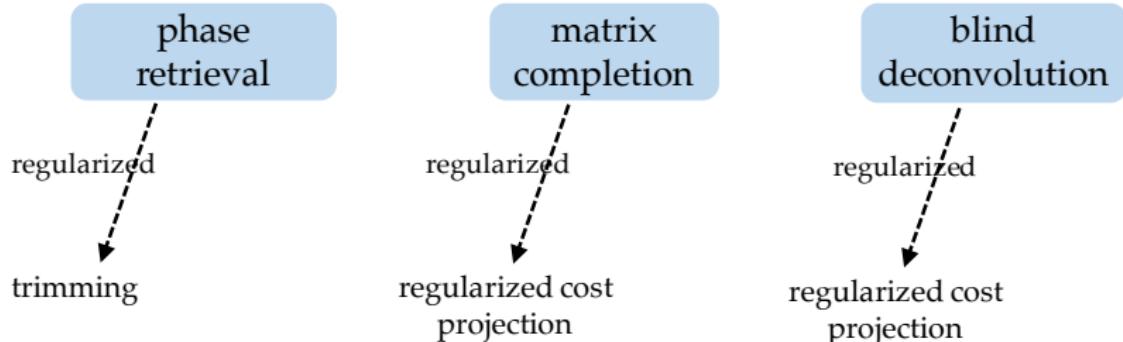
# Proper regularization is *often* recommended

---

Improves computation by stabilizing search directions

# Proper regularization is *often* recommended

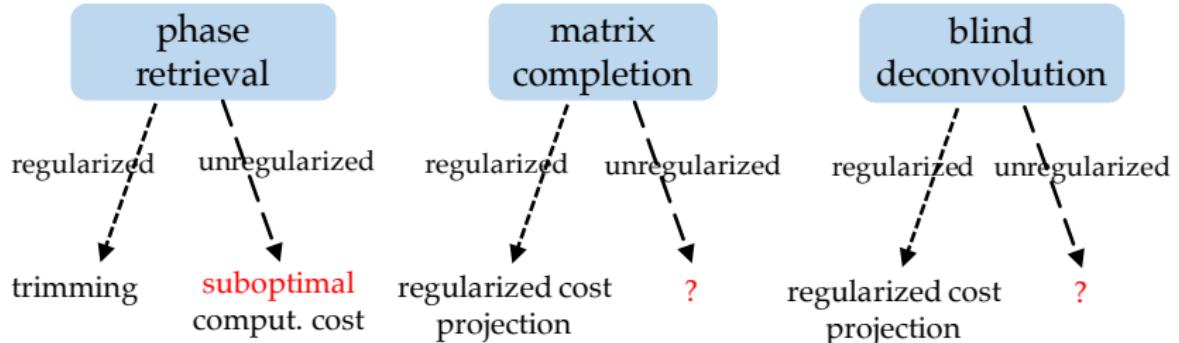
Improves computation by stabilizing search directions



- trimming:  $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathcal{T}(\nabla f(\mathbf{x}^t))$
- regularized cost:  $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t (\nabla f(\mathbf{x}^t) + \nabla \mathcal{R}(\mathbf{x}^t))$
- projection:  $\mathbf{x}^{t+1} = \mathcal{P}(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t))$

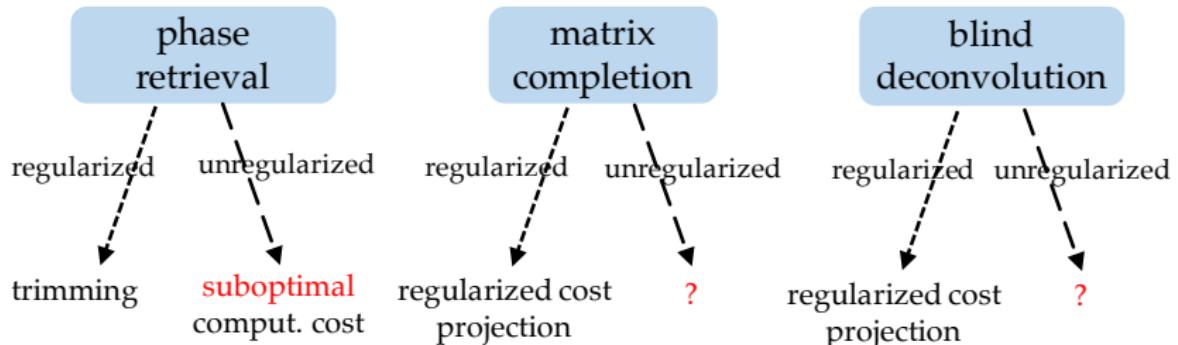
# How about unregularized gradient methods?

Improves computation by stabilizing search directions



# How about unregularized gradient methods?

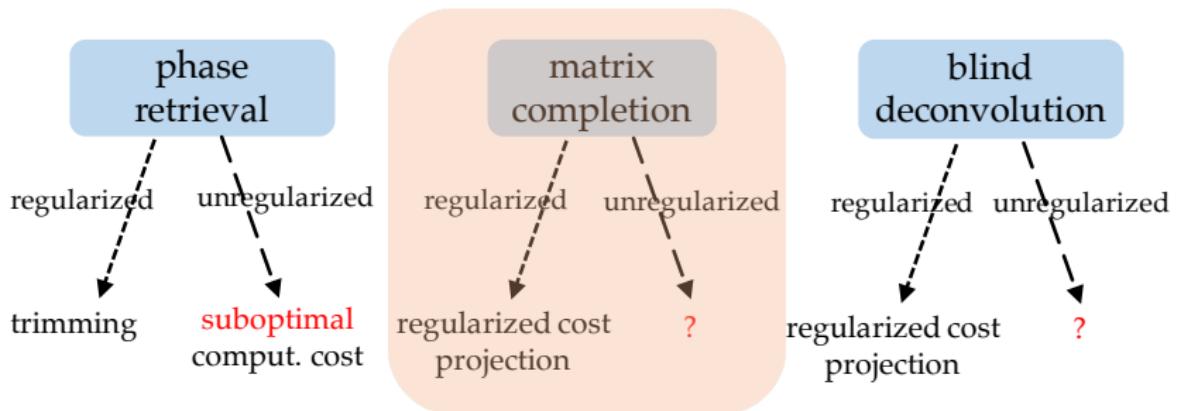
Improves computation by stabilizing search directions



*Are unregularized methods suboptimal for nonconvex estimation?*

# How about unregularized gradient methods?

Improves computation by stabilizing search directions



*Are unregularized methods suboptimal for nonconvex estimation?*

# Low-rank matrix completion

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

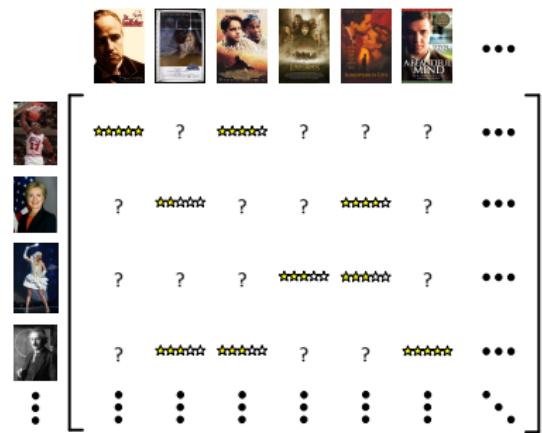


figure credit: Emmanuel Candès

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

## A natural least squares formulation

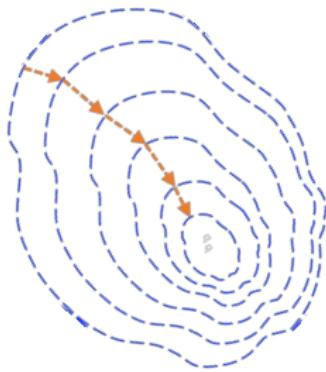
---

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

# A natural least squares formulation

---

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

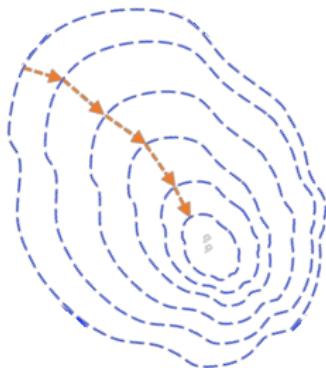


- **spectral initialization:**  $\mathbf{X}^0 \leftarrow$  leading eigenvectors of data matrix

# A natural least squares formulation

---

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



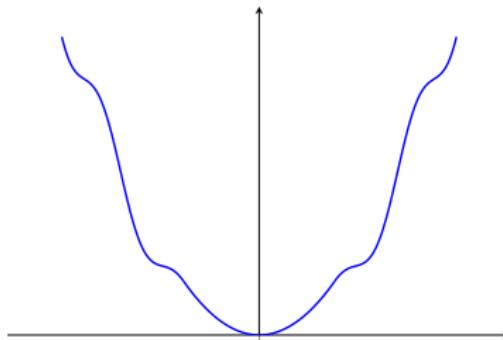
- **spectral initialization:**  $\mathbf{X}^0 \leftarrow$  leading eigenvectors of data matrix

- **gradient descent:**

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

# Gradient descent theory revisited

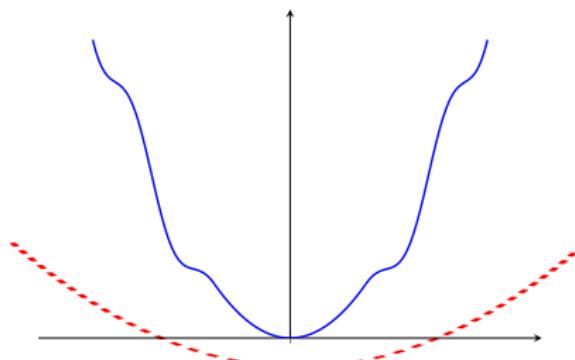
---



Two standard conditions that enable geometric convergence of GD

# Gradient descent theory revisited

---

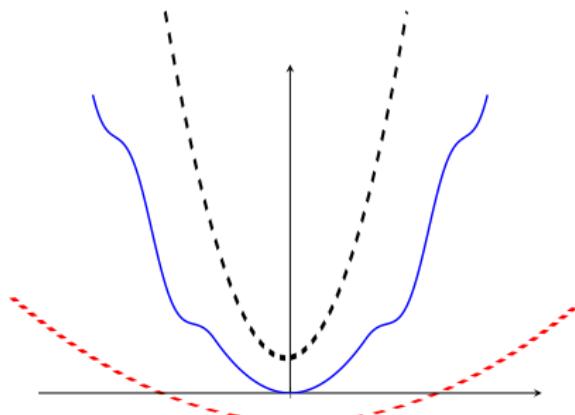


Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity

# Gradient descent theory revisited

---



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity
- (local) smoothness

# Gradient descent theory revisited

---

$f$  is said to be  $\alpha$ -strongly convex and  $\beta$ -smooth if

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

**$\ell_2$  error contraction:** GD with  $\eta = 1/\beta$  obeys

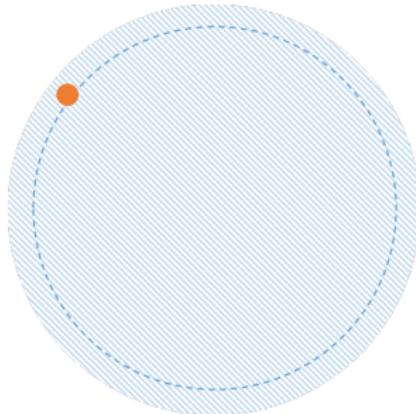
$$\|\mathbf{X}^{t+1} - \mathbf{X}^\natural\|_{\text{F}} \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{X}^t - \mathbf{X}^\natural\|_{\text{F}}$$

# Gradient descent theory revisited

---

$$\|\mathbf{X}^{t+1} - \mathbf{X}^\natural\|_F \leq (1 - \alpha/\beta) \|\mathbf{X}^t - \mathbf{X}^\natural\|_F$$

- region of local strong convexity + smoothness

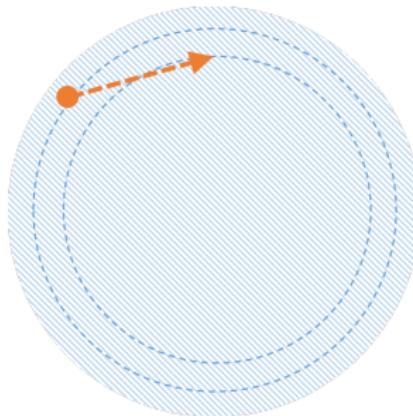


# Gradient descent theory revisited

---

$$\|\mathbf{X}^{t+1} - \mathbf{X}^\natural\|_F \leq (1 - \alpha/\beta) \|\mathbf{X}^t - \mathbf{X}^\natural\|_F$$

- region of local strong convexity + smoothness

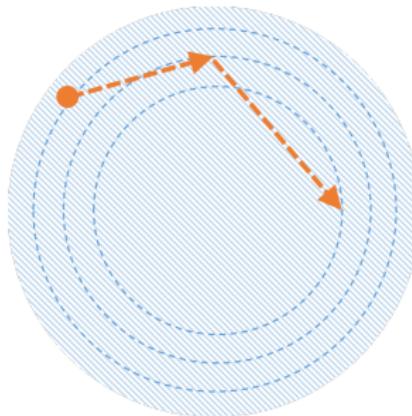


# Gradient descent theory revisited

---

$$\|\mathbf{X}^{t+1} - \mathbf{X}^\natural\|_F \leq (1 - \alpha/\beta) \|\mathbf{X}^t - \mathbf{X}^\natural\|_F$$

- region of local strong convexity + smoothness

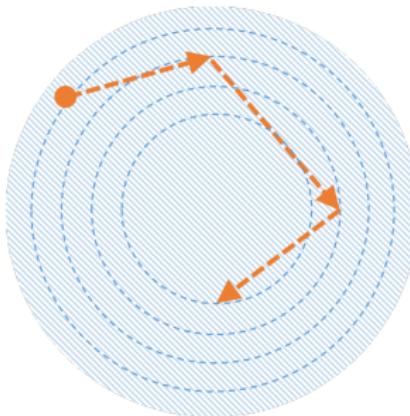


# Gradient descent theory revisited

---

$$\|\mathbf{X}^{t+1} - \mathbf{X}^\natural\|_F \leq (1 - \alpha/\beta) \|\mathbf{X}^t - \mathbf{X}^\natural\|_F$$

- region of local strong convexity + smoothness



# What does this optimization theory say about GD?

---

Independent Bernoulli sampling:  $M_{j,k}^\natural$  is observed with probability  $p$

**Population level** ( $p = 1$ )

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

$$\|\mathbf{V}\|_{\text{F}}^2 \lesssim \text{vec}(\mathbf{V})^\top \mathbb{E} \left[ \nabla^2 f(\mathbf{X}) \right] \text{vec}(\mathbf{V}) \lesssim \|\mathbf{V}\|_{\text{F}}^2$$

for all  $\mathbf{X}$  and  $\mathbf{V}$  such that

# What does this optimization theory say about GD?

---

Independent Bernoulli sampling:  $M_{j,k}^\natural$  is observed with probability  $p$

**Population level** ( $p = 1$ )

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

$$\|\mathbf{V}\|_{\text{F}}^2 \lesssim \text{vec}(\mathbf{V})^\top \mathbb{E} \left[ \nabla^2 f(\mathbf{X}) \right] \text{vec}(\mathbf{V}) \lesssim \|\mathbf{V}\|_{\text{F}}^2$$

for all  $\mathbf{X}$  and  $\mathbf{V}$  such that

- $\mathbf{X}$  is not far away from  $\mathbf{X}^\natural$  in  $\ell_{\text{F}}$  sense.
- $\mathbf{V}$  points towards  $\mathbf{X}^\natural$

# What does this optimization theory say about GD?

---

Independent Bernoulli sampling:  $M_{j,k}^\natural$  is observed with probability  $p$

**Finite-sample level** ( $p \asymp \text{poly log } n/n$ )

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

$$\|\mathbf{V}\|_{\text{F}}^2 \lesssim \text{vec}(\mathbf{V})^\top \mathbb{E} \left[ \nabla^2 f(\mathbf{X}) \right] \text{vec}(\mathbf{V}) \lesssim \|\mathbf{V}\|_{\text{F}}^2$$

for all  $\mathbf{X}$  and  $\mathbf{V}$  such that

- $\mathbf{X}$  is not far away from  $\mathbf{X}^\natural$  in  $\ell_{\text{F}}$  sense.
- $\mathbf{V}$  points towards  $\mathbf{X}^\natural$

# What does this optimization theory say about GD?

---

Independent Bernoulli sampling:  $M_{j,k}^\natural$  is observed with probability  $p$

**Finite-sample level** ( $p \asymp \text{poly log } n/n$ )

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

$$\|\mathbf{V}\|_{\text{F}}^2 \lesssim \text{vec}(\mathbf{V})^\top \mathbb{E} \left[ \nabla^2 f(\mathbf{X}) \right] \text{vec}(\mathbf{V}) \lesssim \|\mathbf{V}\|_{\text{F}}^2$$

for all  $\mathbf{X}$  and  $\mathbf{V}$  such that

- $\mathbf{X}$  is not far away from  $\mathbf{X}^\natural$  in  $\ell_{\text{F}}$  sense.
- $\mathbf{V}$  points towards  $\mathbf{X}^\natural$
- $\mathbf{X}$  is **incoherent** w.r.t. sampling vectors  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$

# Incoherence

## Definition 1 (Incoherence for matrix completion)

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

$$\max_i \|\mathbf{e}_i^\top \mathbf{U}\|_2 = \|\mathbf{U}\|_{2,\infty} \leq \sqrt{\frac{\mu}{n}} \|\mathbf{U}\|_{\text{F}} = \sqrt{\frac{\mu r}{n}}.$$

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

## Incoherence region

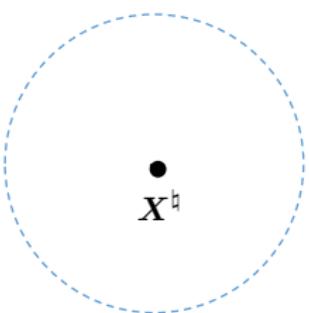
---

Which region enjoys both restricted strong convexity and smoothness?

## Incoherence region

---

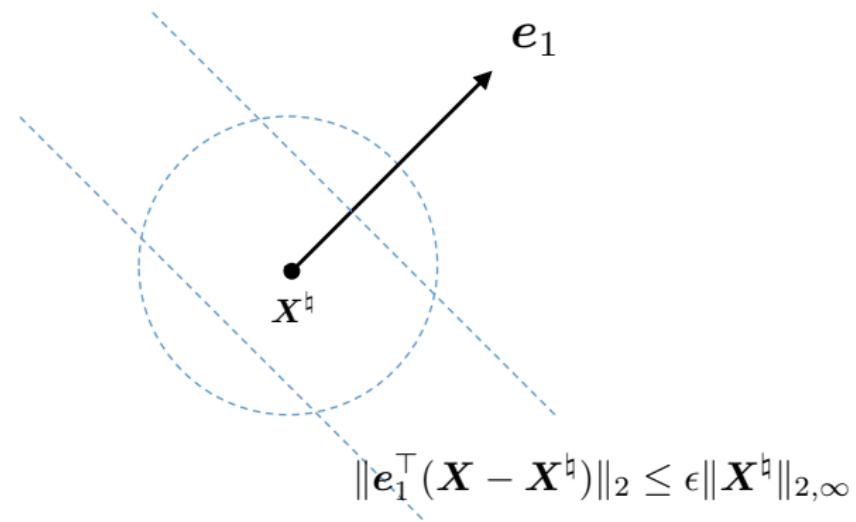
Which region enjoys both restricted strong convexity and smoothness?



- $X$  is not far away from  $X^\natural$

## Incoherence region

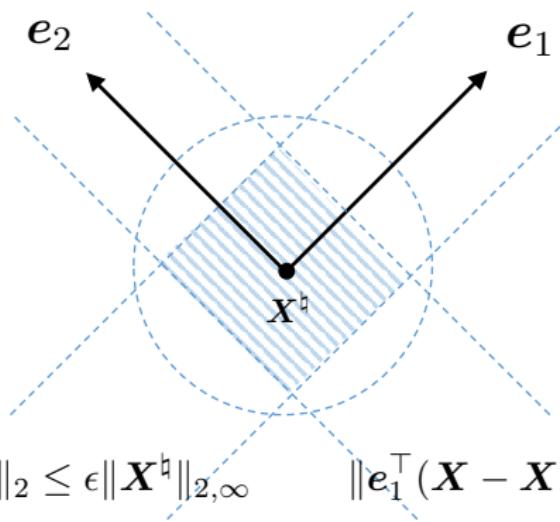
Which region enjoys both restricted strong convexity and smoothness?



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

# Incoherence region

Which region enjoys both restricted strong convexity and smoothness?



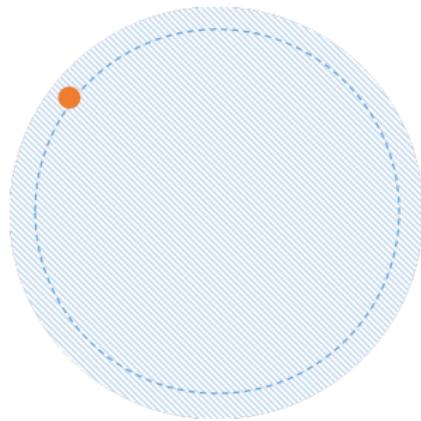
- $X$  is not far away from  $X^\natural$
- $X$  is incoherent w.r.t. sampling vectors (incoherence region)

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



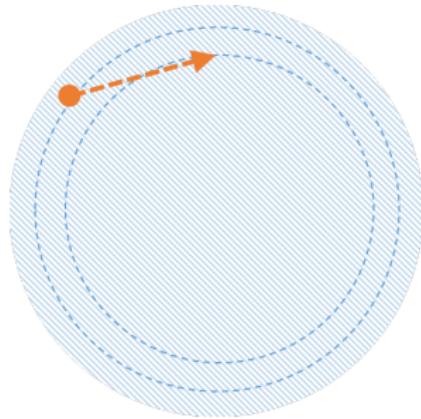
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

## “Failure” of gradient descent?

---



region of local strong convexity + smoothness



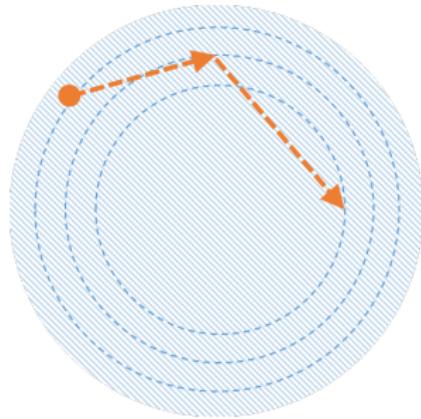
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



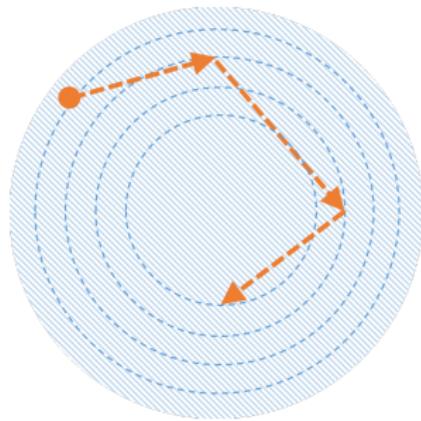
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



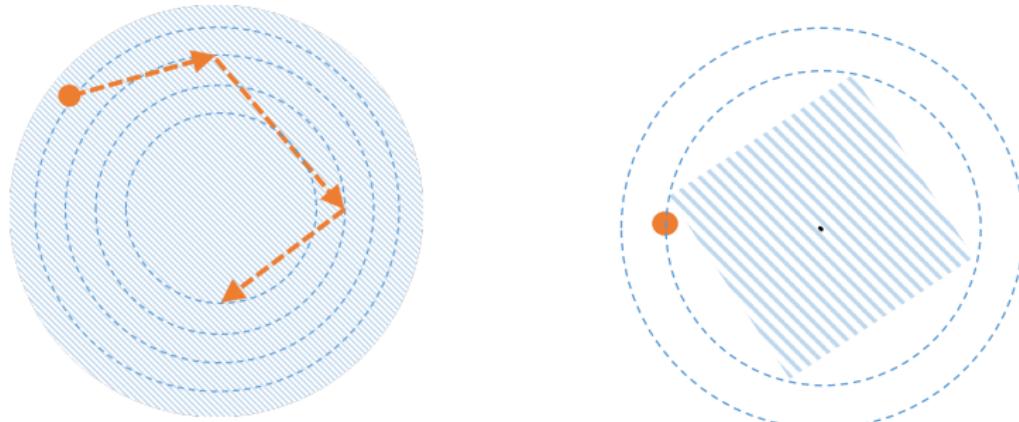
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



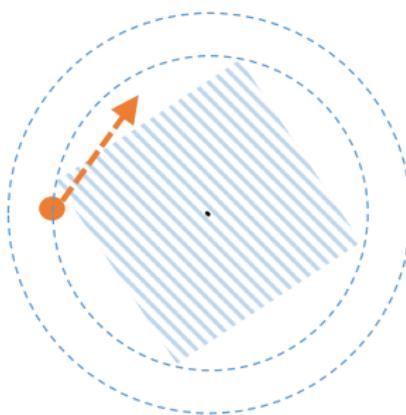
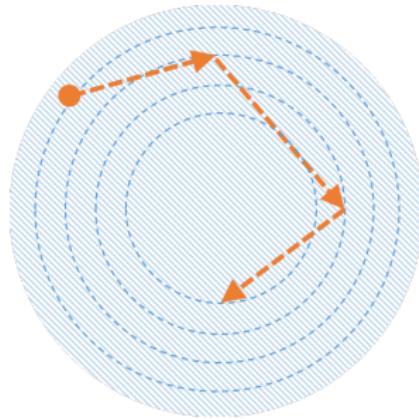
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



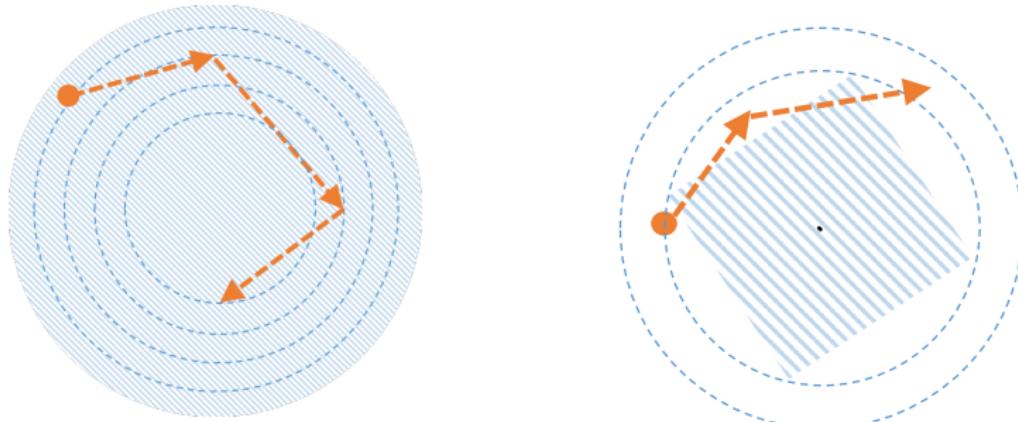
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



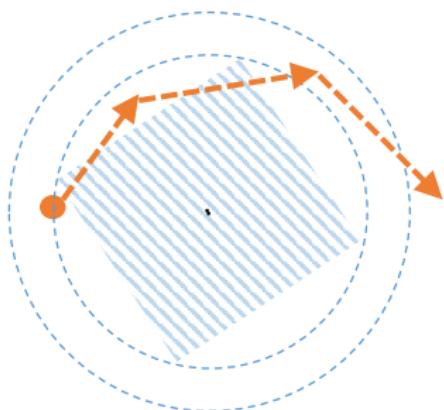
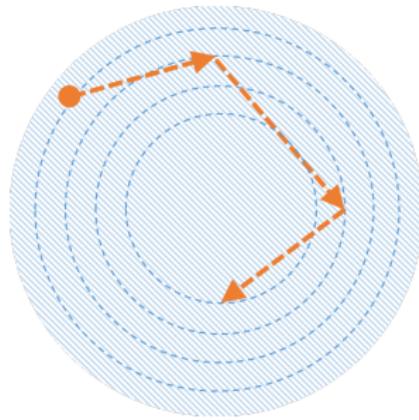
- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# “Failure” of gradient descent?

---



region of local strong convexity + smoothness



- Prior theory only ensures that iterates remain in  $\ell_F$  ball but not incoherence region

# Existing solutions

---

- regularized loss (solve  $\text{minimize}_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$  instead)
  - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16, Chen, Li '17

# Existing solutions

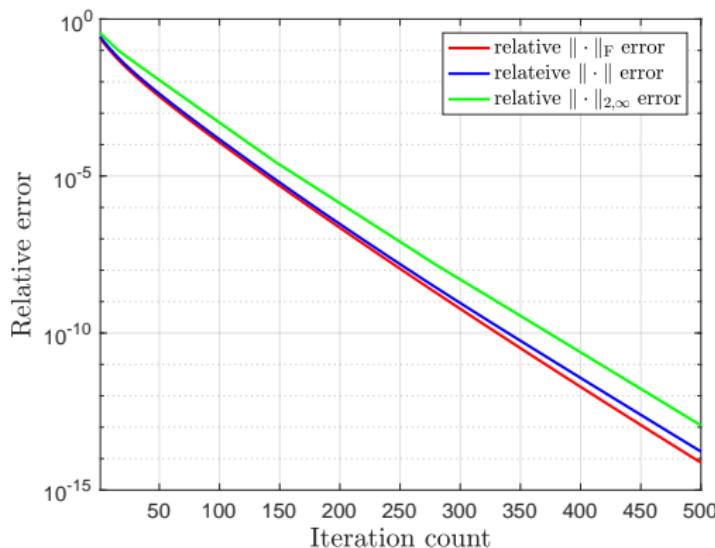
---

- regularized loss (solve  $\text{minimize}_{\mathbf{X}} f(\mathbf{X}) + R(\mathbf{X})$  instead)
  - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16, Chen, Li '17
- projection onto set of incoherent matrices
  - e.g. Chen, Wainwright '15, Zheng, Lafferty '16

*Is regularization necessary for nonconvex matrix completion?*

# Numerical surprise with unregularized GD

$$n = 1000, r = 10, p = 0.1, \eta = 0.2$$



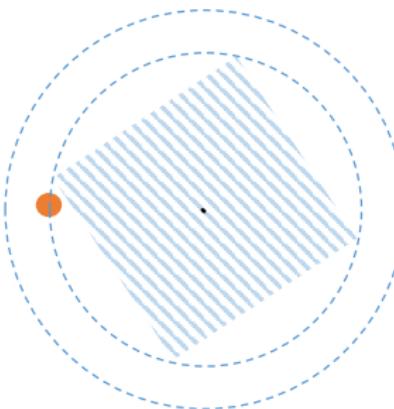
Vanilla GD without regularization converges fast for MC!

# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

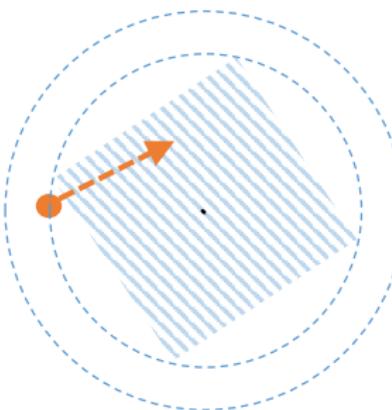


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

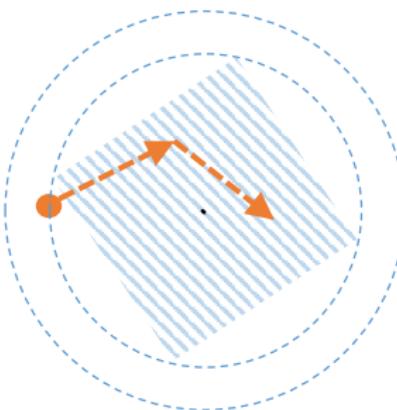


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

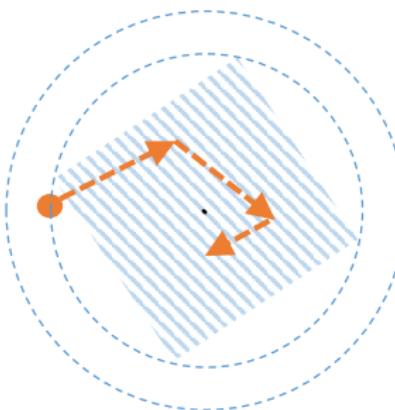


# Our findings: GD is implicitly regularized

---



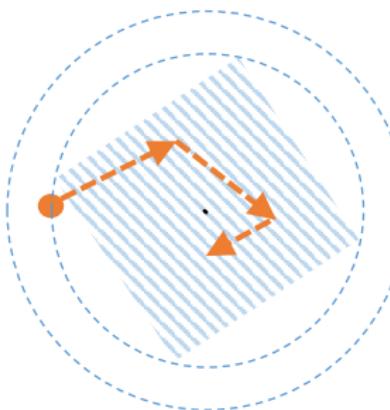
region of local strong convexity + smoothness



# Our findings: GD is implicitly regularized



region of local strong convexity + smoothness

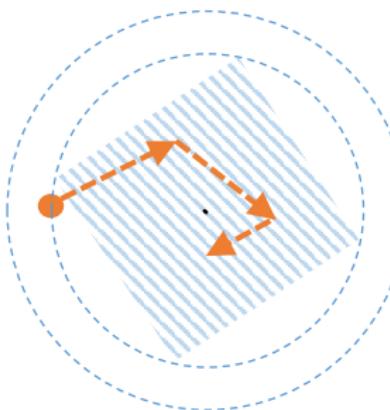


GD implicitly forces iterates to remain **incoherent**

# Our findings: GD is implicitly regularized



region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**

- cannot be derived from generic optimization theory; relies on finer statistical analysis for entire trajectory of GD

# Theoretical guarantees

---

## Theorem 2 (Matrix completion)

Suppose  $M$  is rank- $r$ , incoherent and well-conditioned. **Vanilla gradient descent** (with spectral initialization) achieves  $\varepsilon$  accuracy

- in  $O(\log \frac{1}{\varepsilon})$  iterations

if step size  $\eta \lesssim 1/\sigma_{\max}(M)$  and sample size  $\gtrsim nr^3 \log^3 n$

# Theoretical guarantees

## Theorem 2 (Matrix completion)

Suppose  $M$  is rank- $r$ , incoherent and well-conditioned. *Vanilla gradient descent* (with spectral initialization) achieves  $\varepsilon$  accuracy

- in  $O(\log \frac{1}{\varepsilon})$  iterations w.r.t.  $\|\cdot\|_F$ ,  $\|\cdot\|$ , and  $\underbrace{\|\cdot\|_{2,\infty}}_{\text{incoherence}}$

if step size  $\eta \lesssim 1/\sigma_{\max}(M)$  and sample size  $\gtrsim nr^3 \log^3 n$

# Theoretical guarantees

## Theorem 2 (Matrix completion)

Suppose  $M$  is rank- $r$ , incoherent and well-conditioned. **Vanilla gradient descent** (with spectral initialization) achieves  $\varepsilon$  accuracy

- in  $O(\log \frac{1}{\varepsilon})$  iterations w.r.t.  $\|\cdot\|_F$ ,  $\|\cdot\|$ , and  $\underbrace{\|\cdot\|_{2,\infty}}_{\text{incoherence}}$

if step size  $\eta \lesssim 1/\sigma_{\max}(M)$  and sample size  $\gtrsim nr^3 \log^3 n$

- Byproduct: vanilla GD controls **entrywise error**
  - errors are spread out across all entries

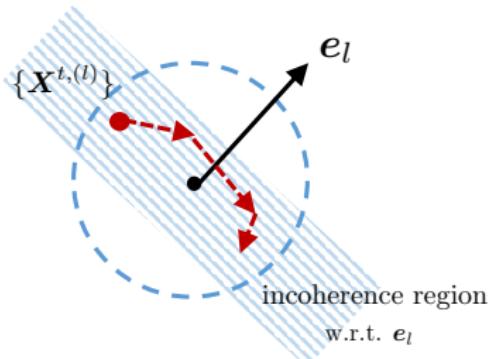
## Key ingredient: leave-one-out analysis

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

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

## Key ingredient: leave-one-out analysis

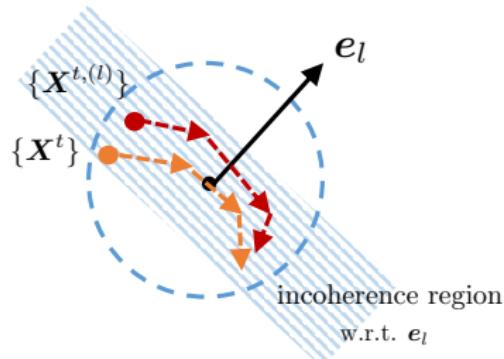
---



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

## Key ingredient: leave-one-out analysis

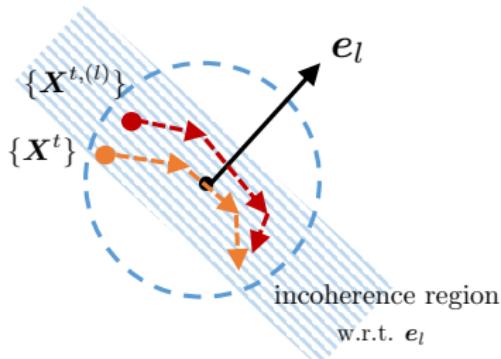
---



- Leave-one-out iterates  $\{X^{t,(l)}\}$  contains more information of  $l$ th row of  $X^t$ ; 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

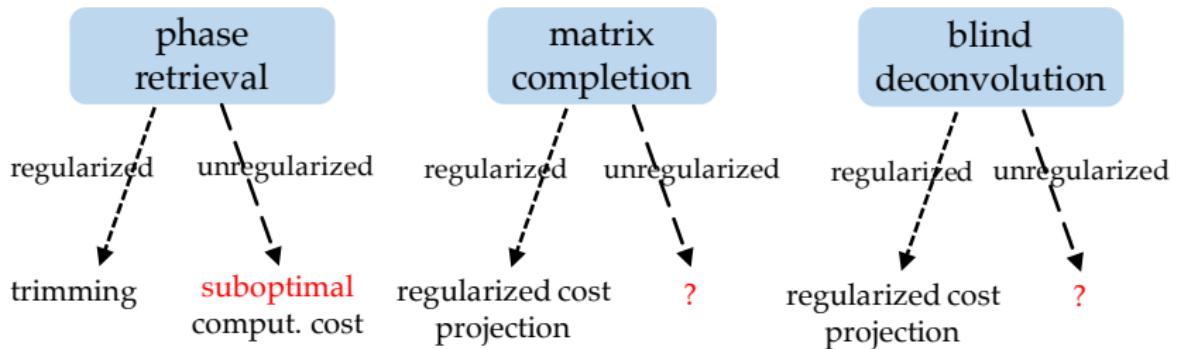
---



- 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  $\mathbf{X}^t$
- $\|\mathbf{e}_l^\top (\mathbf{X}^t - \mathbf{X}^\natural)\|_2 \leq \|\mathbf{e}_l^\top (\mathbf{X}^{t,(l)} - \mathbf{X}^\natural)\|_2 + \|\mathbf{e}_l^\top (\mathbf{X}^t - \mathbf{X}^{t,(l)})\|_2$

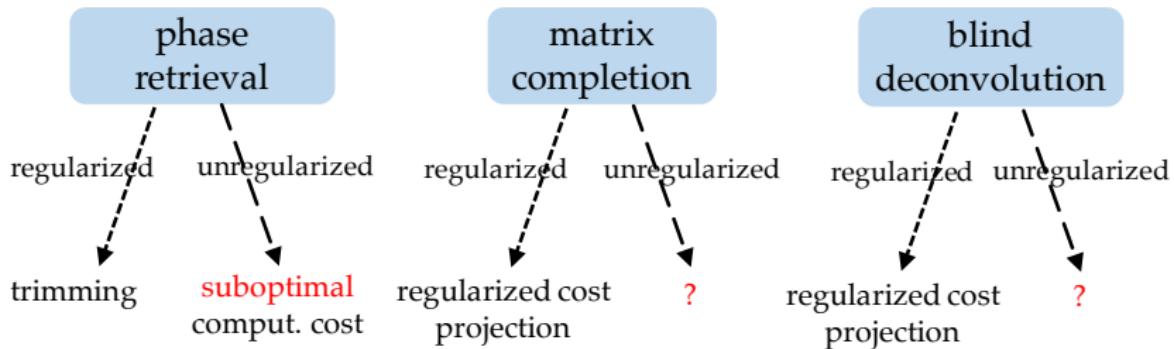
*This phenomenon is quite general*

# How about unregularized gradient methods?



*Are unregularized methods suboptimal for nonconvex estimation?*

# How about unregularized gradient methods?



*Are unregularized methods suboptimal for nonconvex estimation?*

No, for a variety of nonconvex statistical estimation problems.

# Summary

---

- **Implicit regularization:** vanilla gradient descent automatically forces iterates to stay *incoherent*

# Summary

---

- **Implicit regularization:** vanilla gradient descent automatically forces iterates to stay *incoherent*
- Enable error controls in a much stronger sense (e.g. *entrywise error control*)

## Paper:

“Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution”,  
Cong Ma, Kaizheng Wang, Yuejie Chi, Yuxin Chen, arXiv:1711.10467