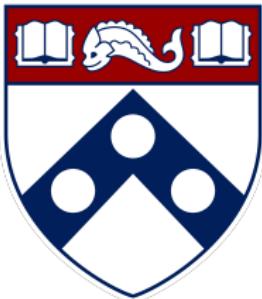


Nonconvex Optimization for High-Dimensional Estimation (Part 2)



Yuxin Chen

Wharton Statistics & Data Science, Spring 2022

A case study: solving quadratic systems of equations

Solving quadratic systems of equations

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

A diagram illustrating the computation of quadratic measurements. On the left, a matrix A of size $m \times n$ is shown as a grid of colored squares. To its right is a vector x^* represented by a vertical stack of blue squares. An equals sign follows. To the right of that is the product Ax^* , shown as a vertical stack of blue squares with numerical values: 1, -3, 2, -1, 4, 2, -2, -1, 3, 4. A large arrow points from this row to the final column, which contains the values 1, 9, 4, 1, 16, 4, 4, 1, 9, 16.

$y = Ax^* ^2$
1
9
4
1
16
4
4
1
9
16

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

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

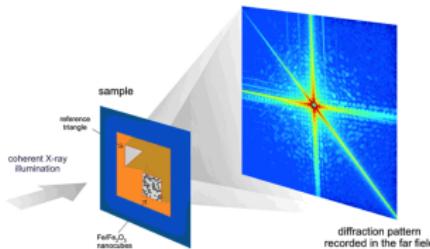
assume w.l.o.g. $\|\mathbf{x}^*\|_2 = 1$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

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

figure credit: Stanford SLAC



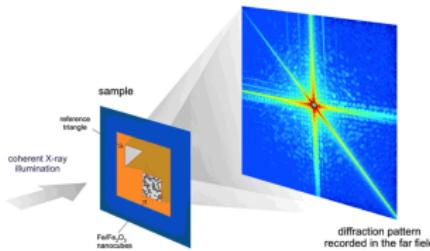
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

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

figure credit: Stanford SLAC

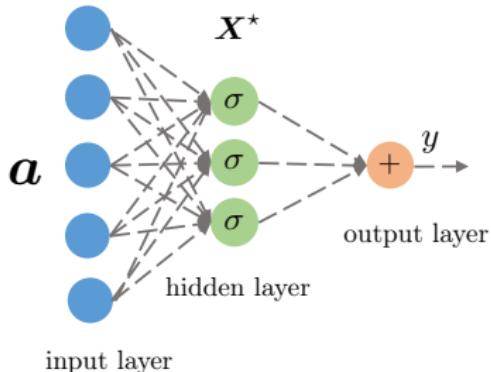


$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

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

Motivation: learning neural nets with quadratic activation

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

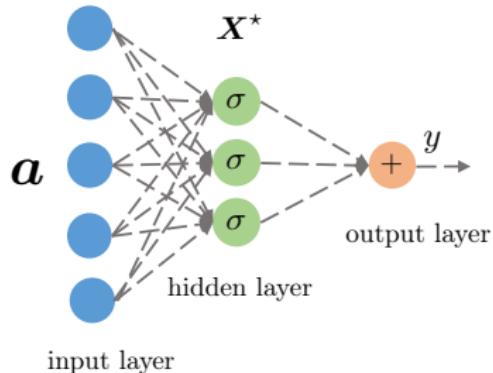


input features: a ; weights: $\mathbf{X}^* = [x_1^*, \dots, x_r^*]$

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

Motivation: learning neural nets with quadratic activation

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

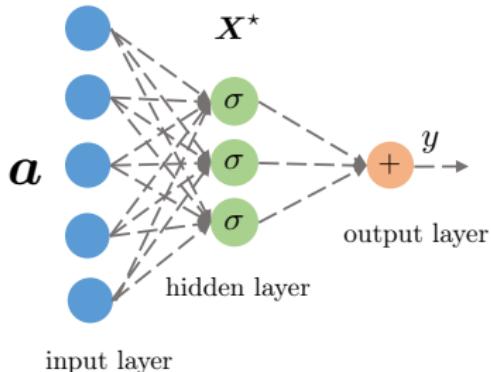


input features: a ; weights: $\mathbf{X}^* = [x_1^*, \dots, x_r^*]$

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

Motivation: learning neural nets with quadratic activation

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



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

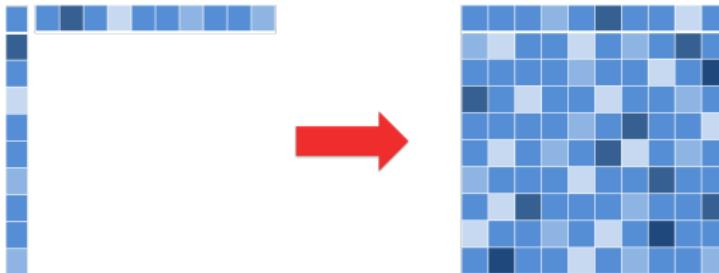
$$\text{output: } y = \sum_{i=1}^r \sigma(a^\top x_i^*) \stackrel{\sigma(z)=z^2}{=} \sum_{i=1}^r (a^\top x_i^*)^2$$

We consider simplest model when $r = 1$ (higher r is similar)

An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

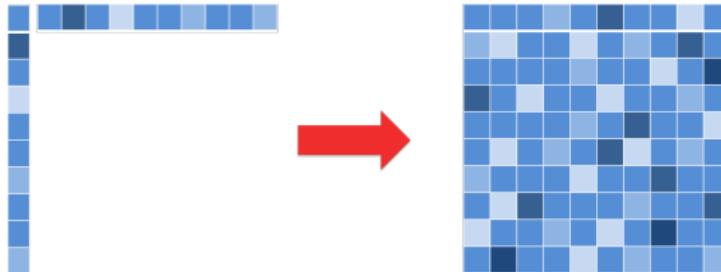
$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a} \implies y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a} \implies y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



find \mathbf{X}

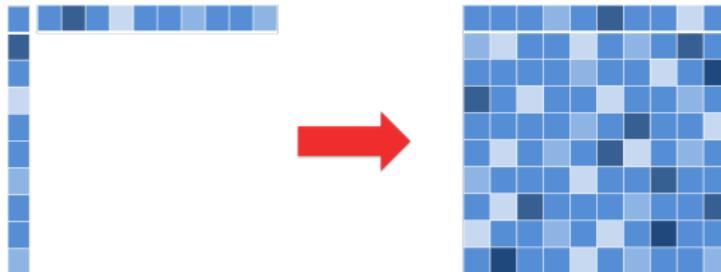
$$\text{s.t. } y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m$$

$$\text{rank}(\mathbf{X}) = 1$$

An equivalent view: low-rank factorization

Introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^\top$ to linearize constraints

$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a} \implies y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



find \mathbf{X}

$$\text{s.t. } y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m$$

$$\text{rank}(\mathbf{X}) = 1$$

Solving quadratic systems is essentially low-rank matrix completion

A natural least-squares formulation

given: $y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$



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

A natural least-squares formulation

$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

⇓

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$

- **pros:** often exact as long as sample size is sufficiently large

A natural least-squares formulation

$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$



$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$

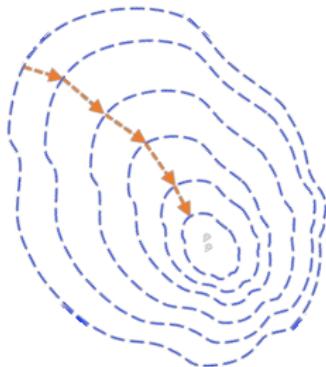
- **pros:** often exact as long as sample size is sufficiently large
- **cons:** $f(\cdot)$ is highly nonconvex
→ *computationally challenging!*

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

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

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

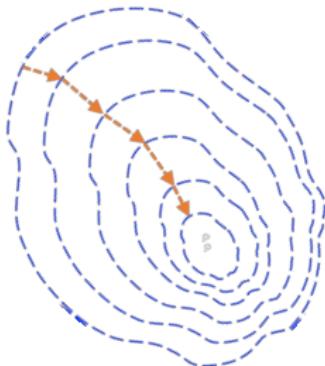
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



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

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

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



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

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

Spectral initialization

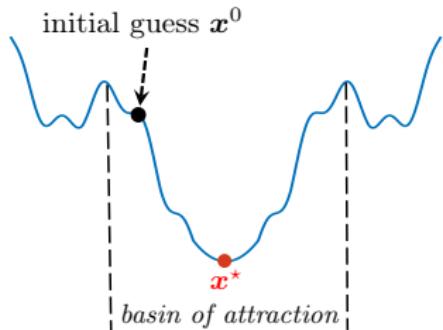
$\boldsymbol{x}^0 \leftarrow$ leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top$$

Rationale: under random Gaussian design $\mathbf{a}_i \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$,

$$\mathbb{E}[\mathbf{Y}] := \mathbb{E} \left[\frac{1}{m} \sum_{k=1}^m \mathbf{y}_k \mathbf{a}_k \mathbf{a}_k^\top \right] = \underbrace{\|\mathbf{x}^*\|_2^2 \mathbf{I} + 2\mathbf{x}^* \mathbf{x}^{*\top}}_{\text{leading eigenvector: } \pm \mathbf{x}^*}$$

Rationale of two-stage approach



1. initialize within local basin sufficiently close to x^*
(restricted) strongly convex; no saddles / spurious local mins

Rationale of two-stage approach



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

A highly incomplete list of two-stage methods

phase retrieval:

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

other problems:

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

First theory of WF

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

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size $\eta \lesssim 1/n$ and sample size:
 $m \gtrsim n \log n$.

First theory of WF

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

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size $\eta \lesssim 1/n$ and sample size:
 $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$

First theory of WF

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

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size $\eta \lesssim 1/n$ and sample size:
 $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$

First theory of WF

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

Theorem 1 (Candès, Li, Soltanolkotabi '14)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{4}\right)^{t/2} \|\mathbf{x}^*\|_2,$$

with high prob., provided that step size and sample size: .

- Iteration complexity: $O(n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on (worst-case) local geometry

Improved theory of WF

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

Theorem 2 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \left(1 - \frac{\eta}{2}\right)^t \|\mathbf{x}^*\|_2$$

with high prob., provided that step size $\eta \asymp 1/\log n$ and sample size $m \gtrsim n \log n$.

- Iteration complexity: $O(n \log \frac{1}{\epsilon}) \searrow O(\log n \log \frac{1}{\epsilon})$
- Sample complexity: $O(n \log n)$
- Derived based on finer analysis of GD trajectory

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$$

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ $\underbrace{\text{but ill-conditioned}}_{\text{condition number } \asymp n}$ (even locally)

What does optimization theory say about WF?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

Finite-sample level ($m \asymp n \log n$)

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ $\underbrace{\text{but ill-conditioned}}_{\text{condition number } \asymp n}$ (even locally)

Consequence (Candès et al '14): WF attains ε -accuracy within $O(n \log \frac{1}{\varepsilon})$ iterations if $m \asymp n \log n$

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$



This choice is suggested by **worst-case** optimization theory

Generic optimization theory gives pessimistic bounds

WF converges in $O(n)$ iterations



Step size taken to be $\eta = O(1/n)$

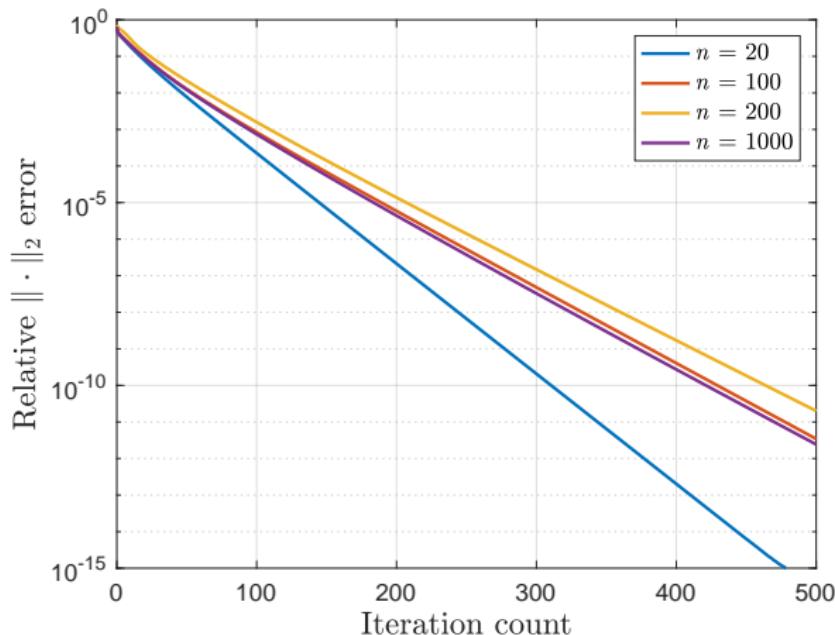


This choice is suggested by worst-case optimization theory



Does it capture what really happens?

Numerical efficiency with $\eta_t = 0.1$



Vanilla GD (WF) converges fast for a constant step size!

A second look at gradient descent theory

Which local region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x}^*)^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

A second look at gradient descent theory

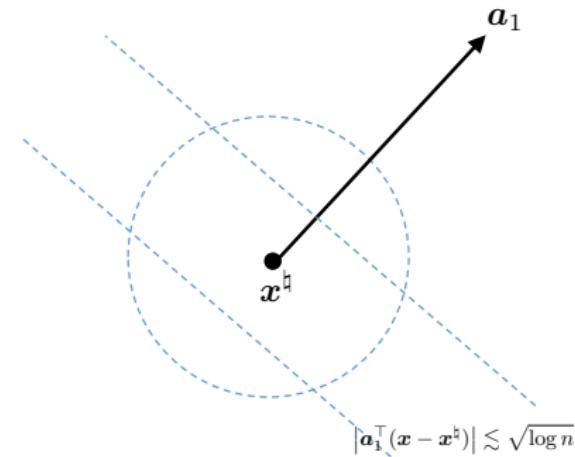
Which local region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x}^*)^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

- Not sufficiently smooth if \mathbf{x} and \mathbf{a}_k are too close (coherent)

A second look at gradient descent theory

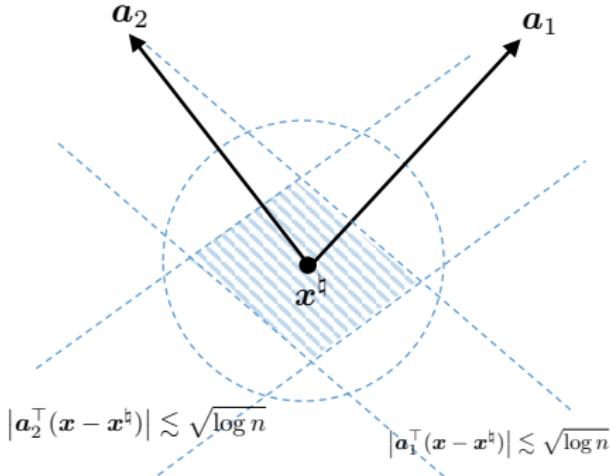
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

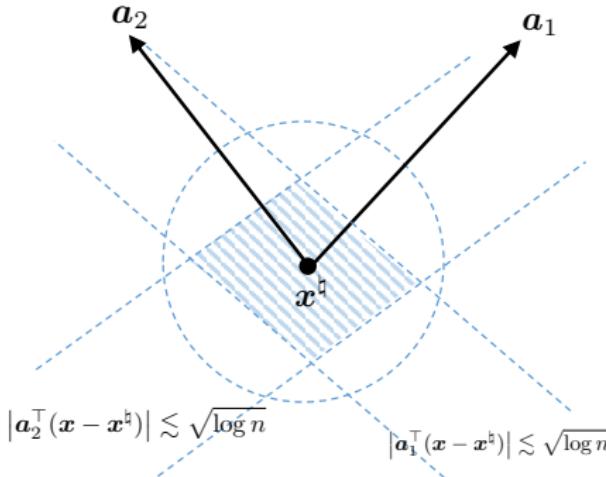
Which local region enjoys both strong convexity and smoothness?



- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

A second look at gradient descent theory

Which local region enjoys both strong convexity and smoothness?



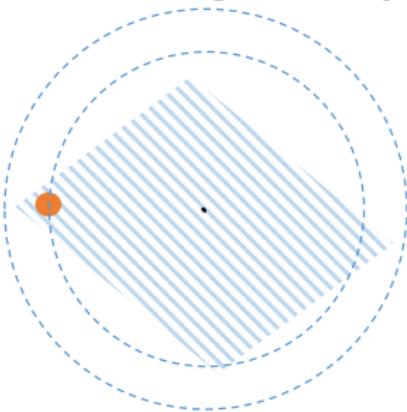
- x is incoherent w.r.t. sampling vectors $\{a_k\}$ (incoherence region)

Prior works suggest enforcing **regularization** (e.g. truncation, projection, regularized loss) to promote incoherence

Encouraging message: GD is implicitly regularized



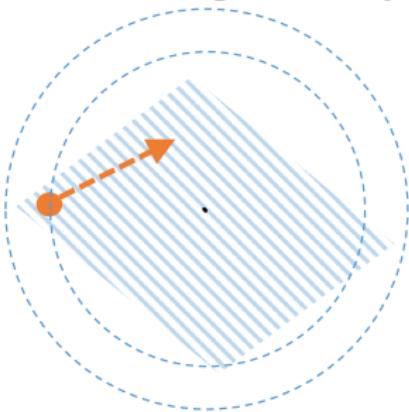
region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized



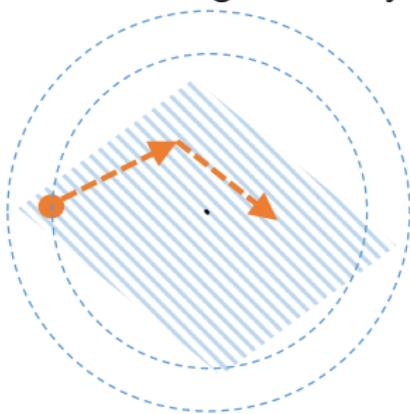
region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized



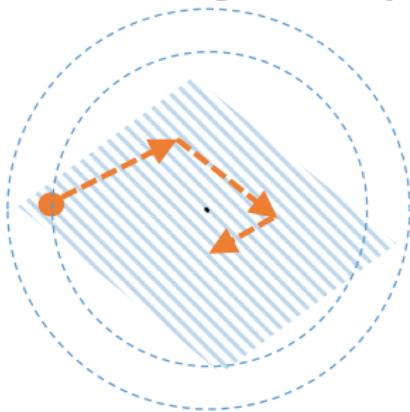
region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized



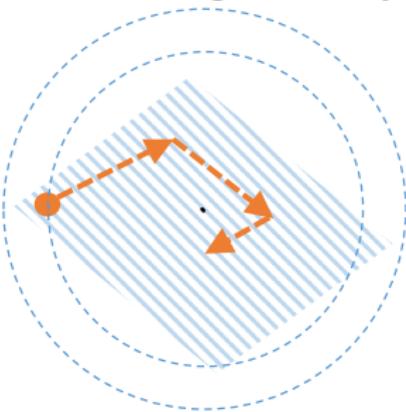
region of local strong convexity + smoothness



Encouraging message: GD is implicitly regularized



region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent with $\{a_k\}$**

$$\max_k |a_k^\top (x^t - x^*)| \lesssim \sqrt{\log n} \|x^*\|_2, \quad \forall t$$

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

Theoretical guarantees for local refinement stage

Theorem 3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

- $\max_k |\mathbf{a}_k^\top \mathbf{x}^t| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2$ (incoherence)

Theoretical guarantees for local refinement stage

Theorem 3 (Ma, Wang, Chi, Chen '17)

Under i.i.d. Gaussian design, WF with spectral initialization achieves

- $\max_k |\mathbf{a}_k^\top \mathbf{x}^t| \lesssim \sqrt{\log n} \|\mathbf{x}^*\|_2$ (incoherence)
- $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim (1 - \frac{\eta}{2})^t \|\mathbf{x}^*\|_2$ (linear convergence)

provided that step size $\eta \asymp 1/\log n$ and sample size $m \gtrsim n \log n$.

- Attains ε accuracy within $O(\log n \log \frac{1}{\varepsilon})$ iterations

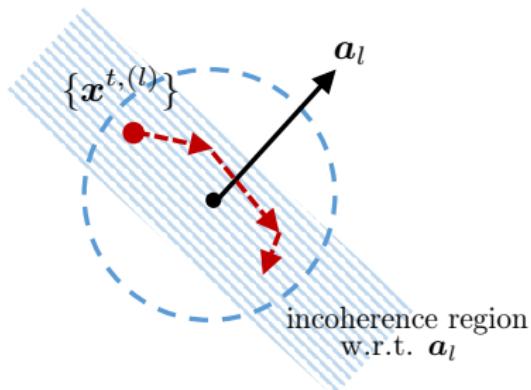
Key proof idea: leave-one-out analysis

For each $1 \leq l \leq m$, introduce leave-one-out iterates $\mathbf{x}^{t,(l)}$ by dropping l th measurement

$$\begin{array}{c} \mathbf{A}^{(l)} \\ \hline \mathbf{a}_l^\top \end{array} \quad \mathbf{x}^* \quad = \quad \begin{array}{c} \mathbf{A}^{(l)} \mathbf{x}^* \\ \hline \end{array} \quad \mathbf{y}^{(l)} = |\mathbf{A}^{(l)} \mathbf{x}^*|^2$$

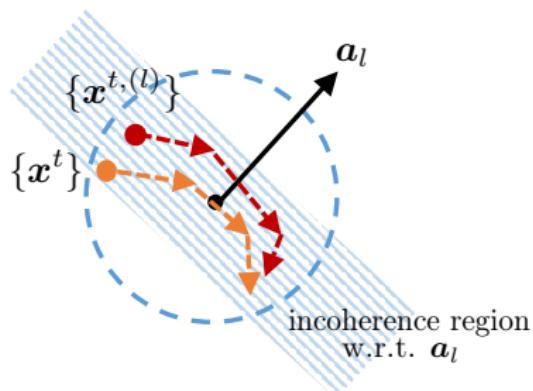
The diagram illustrates the computation of $\mathbf{A}^{(l)} \mathbf{x}^*$. On the left, a matrix $\mathbf{A}^{(l)}$ is shown with a red row \mathbf{a}_l^\top highlighted. This row is explicitly subtracted from the matrix to form the result. The result is a vector with entries 1, -3, 2, -1, 4. An arrow points to another vector with entries 1, 9, 4, 1, 16, which is labeled $\mathbf{y}^{(l)} = |\mathbf{A}^{(l)} \mathbf{x}^*|^2$.

Key proof idea: leave-one-out analysis



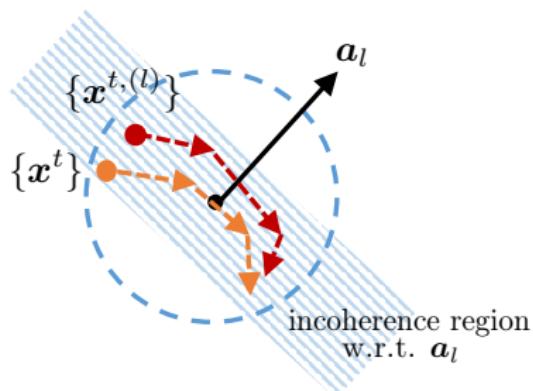
- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l

Key proof idea: leave-one-out analysis



- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l
- Leave-one-out iterate $x^{t,(l)}$ \approx true iterate x^t

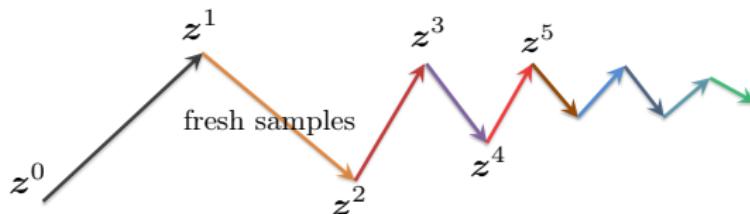
Key proof idea: leave-one-out analysis



- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t
 $\implies x^t$ is nearly independent of a_l
 nearly orthogonal to

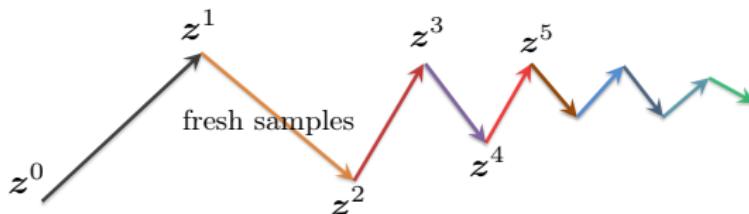
No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis

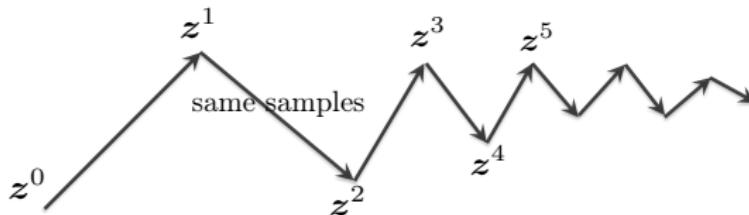


No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis



- This tutorial:** reuses all samples in all iterations



Other examples: low-rank matrix estimation

Low-rank matrix completion

Problem: complete a rank- r matrix M from partial entries: $M_{i,j}$, $(i, j) \in \Omega$

- *random sampling*: (i, j) is included in Ω independently with prob. p

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

Low-rank matrix completion

Problem: complete a rank- r matrix M from partial entries: $M_{i,j}$, $(i, j) \in \Omega$

- *random sampling*: (i, j) is included in Ω independently with prob. p

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

Strong convexity and smoothness do not hold in general

→ need to regularize loss function by promoting **incoherent** solutions

Incoherence for matrix completion

Definition 4 (Incoherence for matrix completion)

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

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

e.g.,

$\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}$	vs.	$\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}$
--	-----	--

$\underbrace{\hspace{10em}}$ hard $\mu=n$

$\underbrace{\hspace{10em}}$ easy $\mu=1$

Gradient descent for matrix completion

Let $M = X^* X^{*\top}$. Observe

$$Y_{i,j} = M_{i,j} + E_{i,j}, \quad (i, j) \in \Omega$$

where $(i, j) \in \Omega$ independently with prob. p , and $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ ¹

$$\text{minimize } \left\| \mathcal{P}_\Omega(\widehat{M} - Y) \right\|_F^2 \quad \text{s.t.} \quad \text{rank}(\widehat{M}) \leq r$$

¹can be relaxed to sub-Gaussian noise and the asymmetric case

Gradient descent for matrix completion

Let $M = X^* X^{*\top}$. Observe

$$Y_{i,j} = M_{i,j} + E_{i,j}, \quad (i, j) \in \Omega$$

where $(i, j) \in \Omega$ independently with prob. p , and $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ ¹

$$\text{minimize } \left\| \mathcal{P}_\Omega(\widehat{M} - Y) \right\|_F^2 \quad \text{s.t. } \text{rank}(\widehat{M}) \leq r$$

$$\text{minimize}_{X \in \mathbb{R}^{n \times r}} \quad f(X) = \underbrace{\sum_{(j,k) \in \Omega} (e_j^\top X X^\top e_k - Y_{j,k})^2}_{\text{unregularized least-squares loss}}$$

¹can be relaxed to sub-Gaussian noise and the asymmetric case

Gradient descent for matrix completion

1. **spectral initialization:** let $\mathbf{U}^0 \boldsymbol{\Sigma}^0 \mathbf{U}^{0\top}$ be rank- r eigendecomposition of

$$\frac{1}{p} \mathcal{P}_\Omega(\mathbf{Y}).$$

and set $\mathbf{X}^0 = \mathbf{U}^0 (\boldsymbol{\Sigma}^0)^{1/2}$

2. **gradient descent updates:**

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

Gradient descent for matrix completion

Define the optimal rotation from the t th iterate \mathbf{X}^t to \mathbf{X}^* as

$$\mathbf{Q}^t := \operatorname{argmin}_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{X}^t \mathbf{R} - \mathbf{X}^*\|_{\text{F}}$$

where $\mathcal{O}^{r \times r}$ is the set of $r \times r$ orthonormal matrices

- orthogonal Procrustes problem

Gradient descent for matrix completion

Theorem 5 (Noiseless MC, Ma, Wang, Chi, Chen '17)

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

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

where $0 < \rho < 1$, if the step size $\eta \asymp 1/\sigma_{\max}$ and the sample complexity $n^2 p \gtrsim \mu^3 nr^3 \log^3 n$

Gradient descent for matrix completion

Theorem 5 (Noiseless MC, Ma, Wang, Chi, Chen '17)

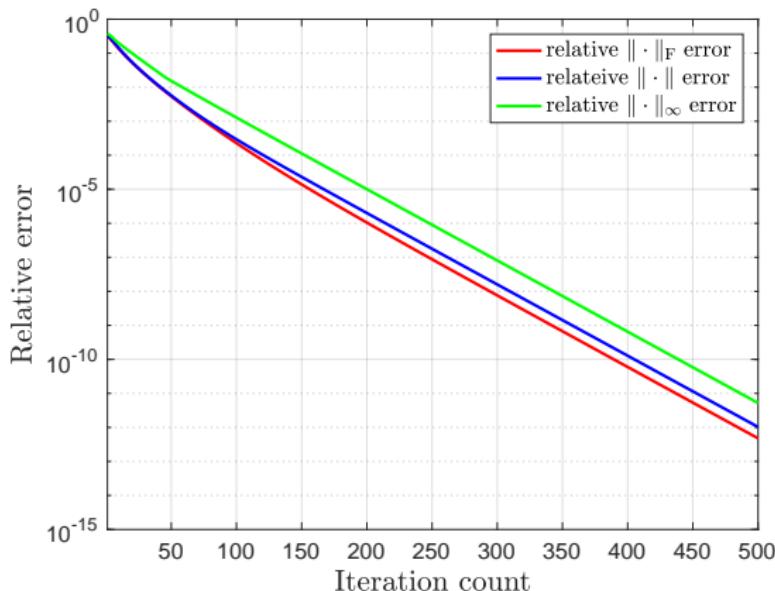
Suppose $\mathbf{M} = \mathbf{X}^* \mathbf{X}^{*\top}$ is rank- r , incoherent and well-conditioned.
Vanilla GD (with spectral initialization) achieves

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

where $0 < \rho < 1$, if the step size $\eta \asymp 1/\sigma_{\max}$ and the sample complexity $n^2 p \gtrsim \mu^3 n r^3 \log^3 n$

- vanilla GD converges linearly for matrix completion!

Numerical evidence for noiseless data



Relative error of $\mathbf{X}^t \mathbf{X}^{t\top}$ (measured by $\|\cdot\|_F$, $\|\cdot\|$, $\|\cdot\|_\infty$) vs. iteration count for matrix completion, where $n = 1000$, $r = 10$, $p = 0.1$, and $\eta_t = 0.2$

Related theory

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

Related theory promotes incoherence explicitly:

Related theory

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

Related theory promotes incoherence explicitly:

- regularized loss (solve $\min_{\mathbf{X}} f(\mathbf{X}) + Q(\mathbf{X})$ instead)
 - e.g. Keshavan, Montanari, Oh '10, Sun, Luo '14, Ge, Lee, Ma '16

Related theory

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

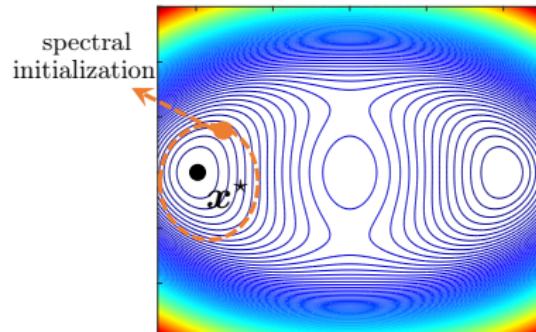
Related theory promotes incoherence explicitly:

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

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

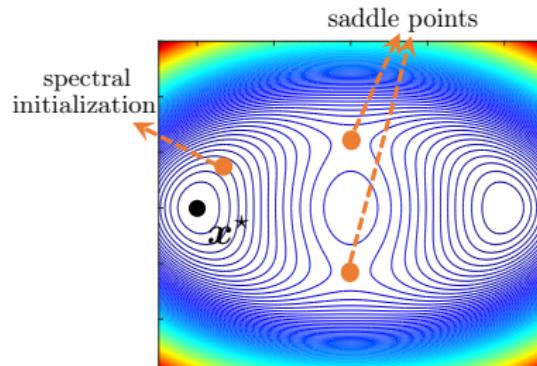
Are carefully-designed initialization or saddle-point escaping schemes necessary for fast convergence?

Initialization



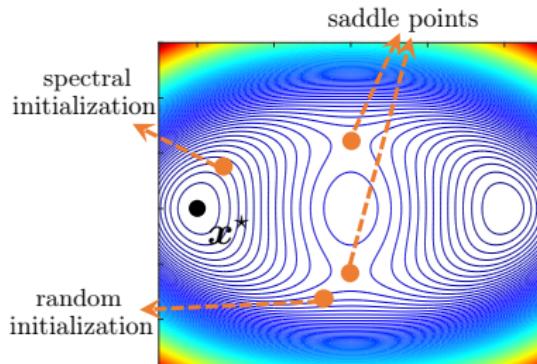
- Spectral initialization gets us reasonably close to truth

Initialization



- Spectral initialization gets us reasonably close to truth
- Cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

Initialization

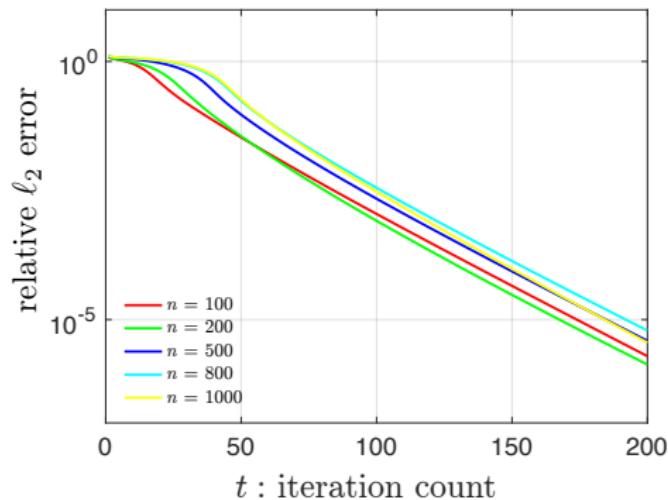


- Spectral initialization gets us reasonably close to truth
- Cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

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

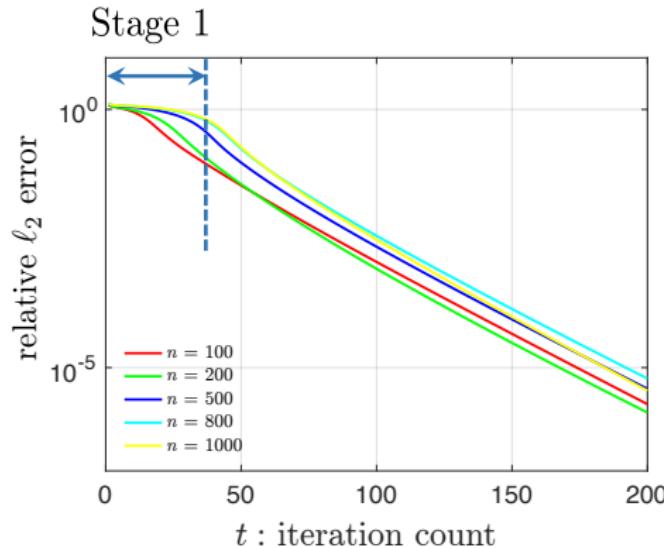
Numerical efficiency of randomly initialized GD

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



Numerical efficiency of randomly initialized GD

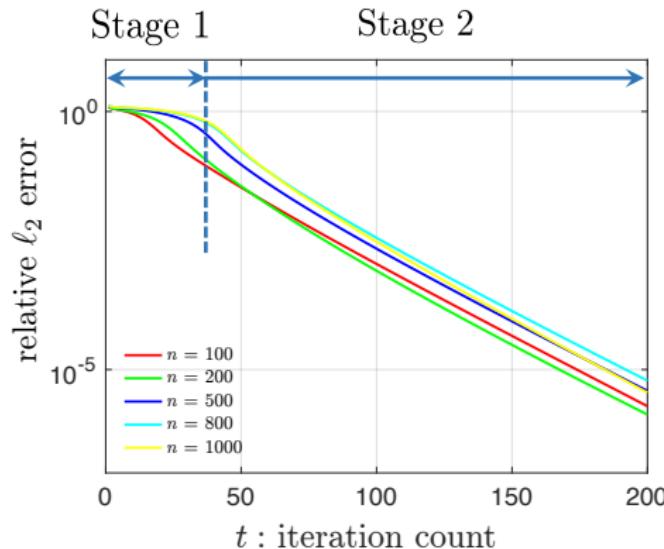
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within a few iterations

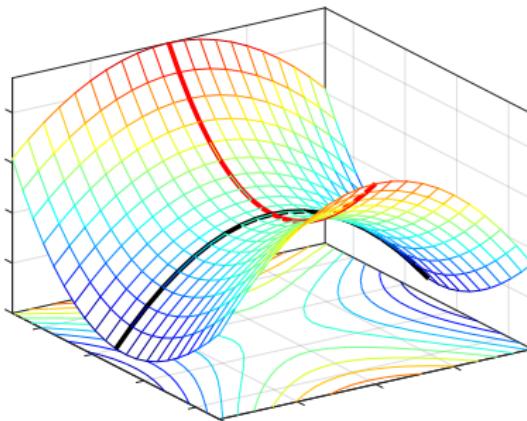
Numerical efficiency of randomly initialized GD

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



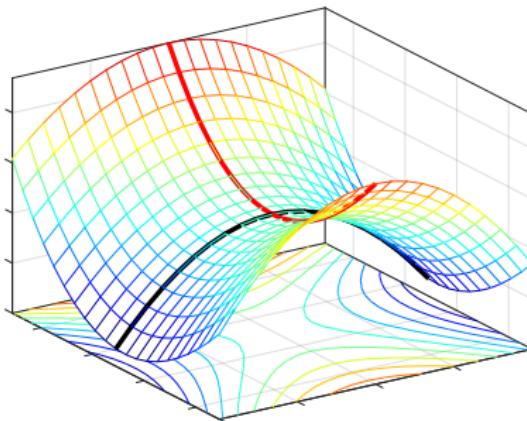
Randomly initialized GD enters local basin within a few iterations

A geometric analysis



- if $m \gtrsim n \log^3 n$, then (Sun et al. '16)
 - there is no spurious local mins
 - all saddle points are strict (i.e. associated Hessian matrices have at least one sufficiently negative eigenvalue)

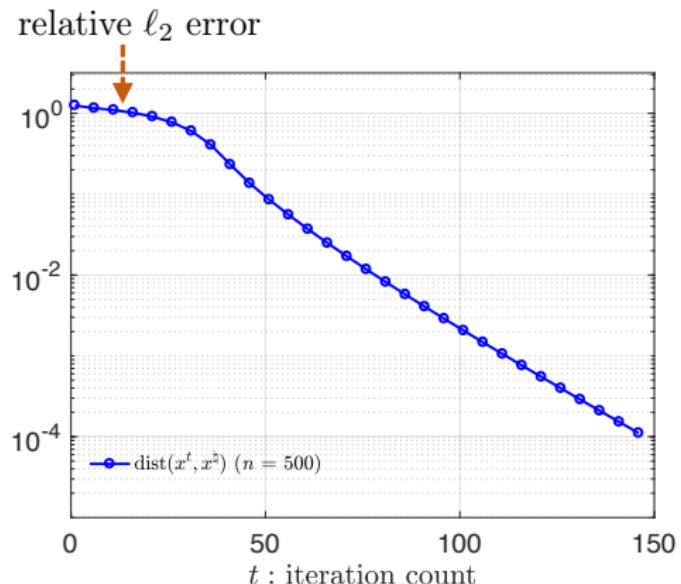
A geometric analysis



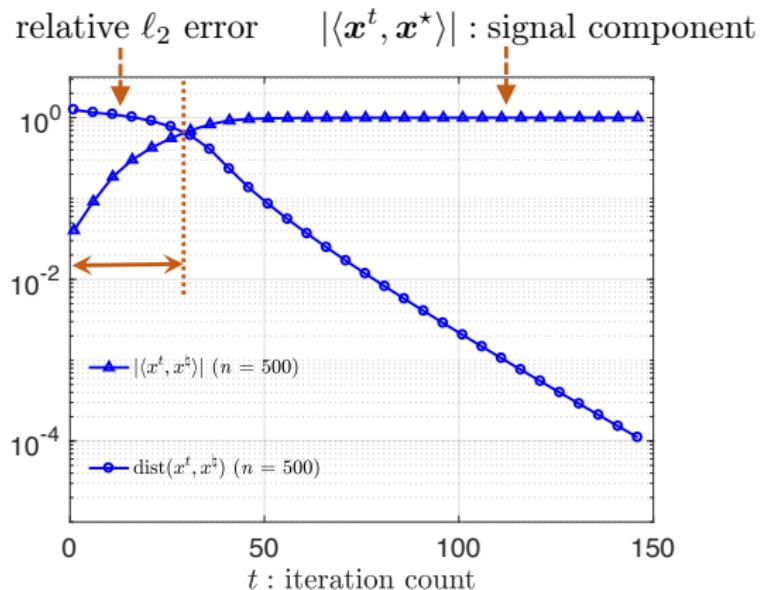
- With such benign landscape, GD with random initialization converges to global min **almost surely** (Lee et al. '16)

No convergence rate guarantees for vanilla GD!

Exponential growth of signal strength in Stage 1

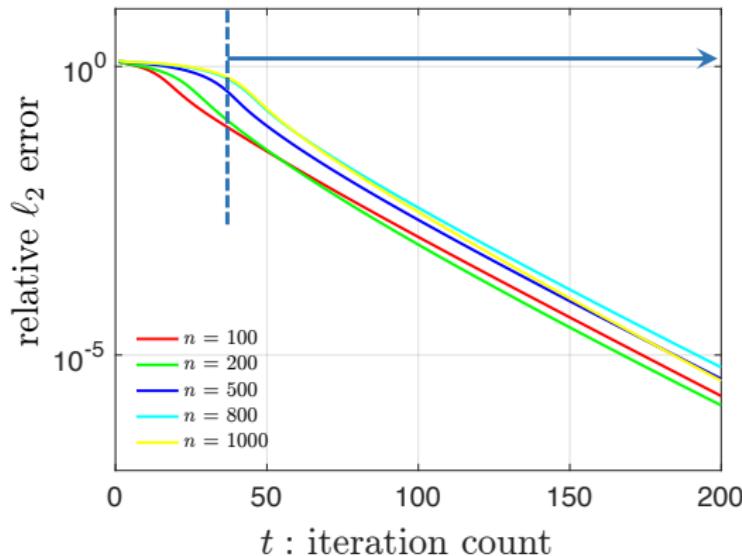


Exponential growth of signal strength in Stage 1

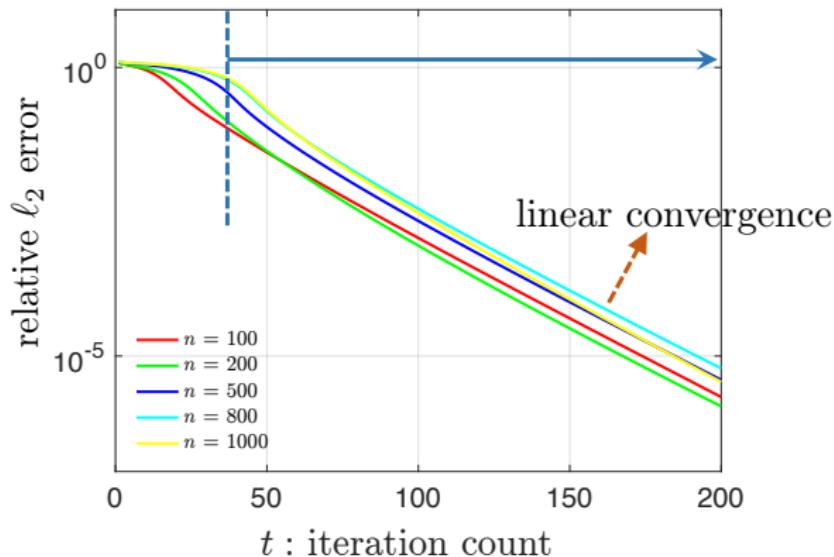


Numerically, $O(\log n)$ iterations are enough to enter local region

Linear / geometric convergence in Stage 2



Linear / geometric convergence in Stage 2



Numerically, GD converges linearly within local region

Theoretical guarantees for randomly initialized GD

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

Theorem 6 (Chen, Chi, Fan, Ma '18)

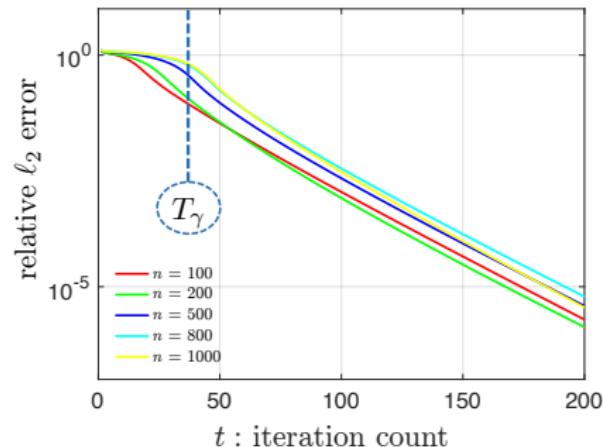
Under i.i.d. Gaussian design, GD with $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$ achieves

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

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

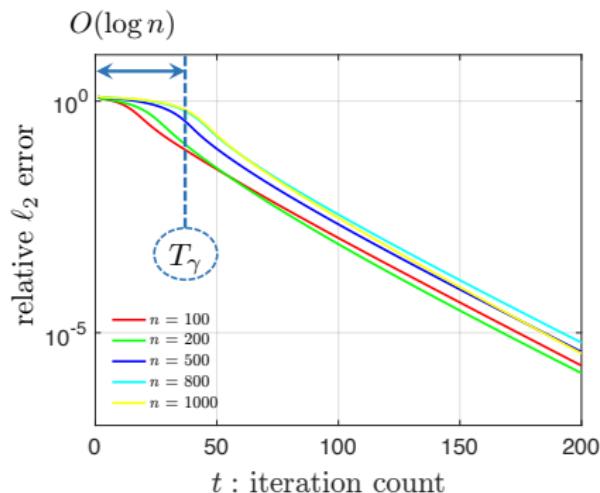
Theoretical guarantees for randomly initialized GD

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



Theoretical guarantees for randomly initialized GD

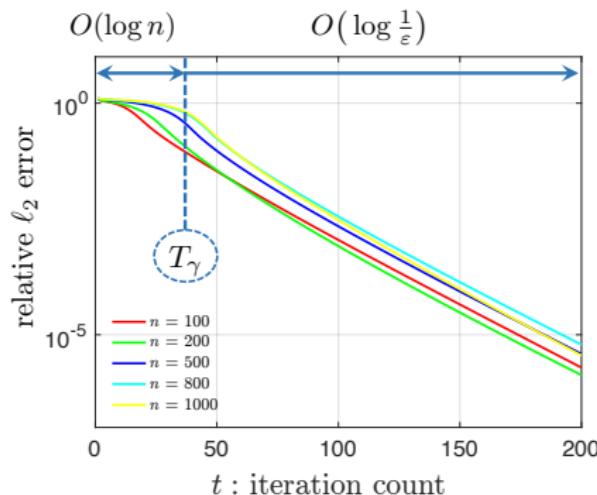
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

Theoretical guarantees for randomly initialized GD

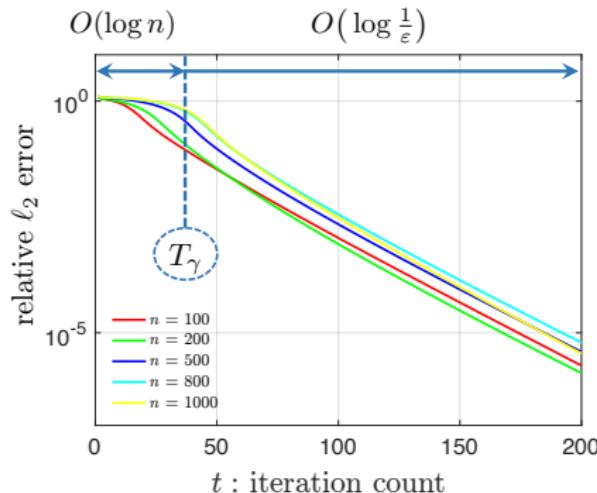
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

Theoretical guarantees for randomly initialized GD

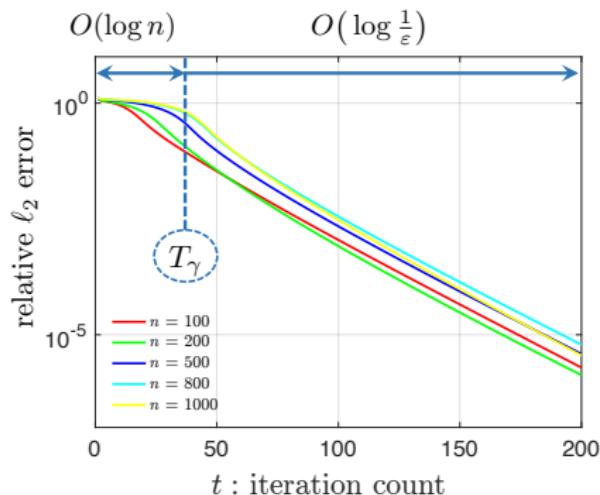
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

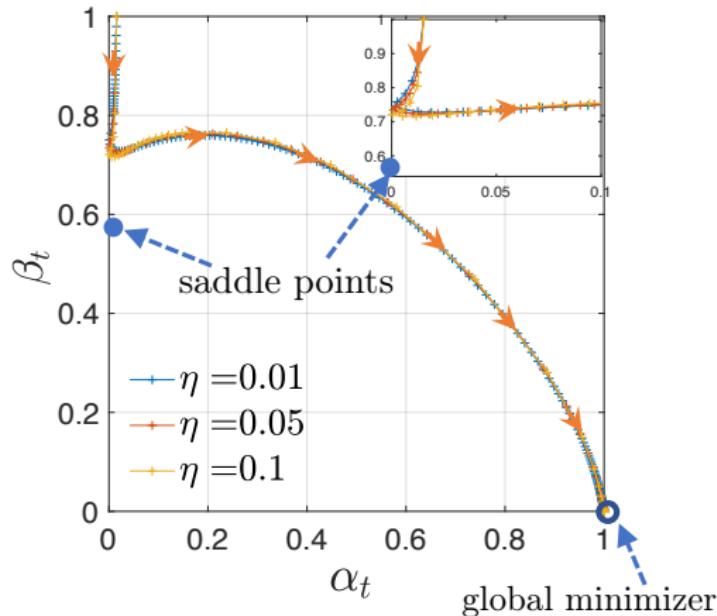
Theoretical guarantees for randomly initialized GD

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



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

Saddle-escaping schemes?



Randomly initialized GD never hits saddle points in phase retrieval!

Other saddle-escaping schemes

	iteration complexity	num of iterations needed to escape saddles	local iteration complexity
Trust-region (Sun et al. '16)	$n^7 + \log \log \frac{1}{\varepsilon}$	n^7	$\log \log \frac{1}{\varepsilon}$
Perturbed GD (Jin et al. '17)	$n^3 + n \log \frac{1}{\varepsilon}$	n^3	$n \log \frac{1}{\varepsilon}$
Perturbed accelerated GD (Jin et al. '17)	$n^{2.5} + \sqrt{n} \log \frac{1}{\varepsilon}$	$n^{2.5}$	$\sqrt{n} \log \frac{1}{\varepsilon}$
GD (Chen et al. '18)	$\log n + \log \frac{1}{\varepsilon}$	$\log n$	$\log \frac{1}{\varepsilon}$

Generic optimization theory yields highly suboptimal convergence guarantees