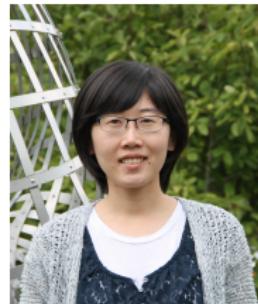


# **Taming Nonconvexity in Information Science**



**Yuxin Chen**  
Princeton



**Yuejie Chi**  
CMU

**ITW 2018 Tutorial**  
**Guangzhou, China**

# Acknowledgement

---

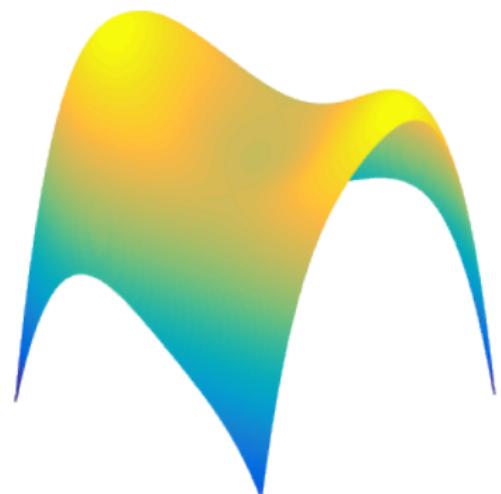
- Our collaborators: Emmanuel J. Candès, Jianqing Fan, Yuanxin Li, Yingbin Liang, Yue M. Lu, Cong Ma, Kaizheng Wang, Huishuai Zhang
- This work is supported in part by ARO W911NF-18-1-0303, AFOSR FA9550-19-1-0030 and FA9550-15-1-0205, ONR N00014-18-1-2142, and NSF ECCS-1818571 and CCF-1806154.

# Nonconvex estimation problems are everywhere

---

Empirical risk minimization is usually nonconvex

$\text{minimize}_x \quad f(x; \text{data}) \quad \rightarrow \quad \text{loss function may be nonconvex}$



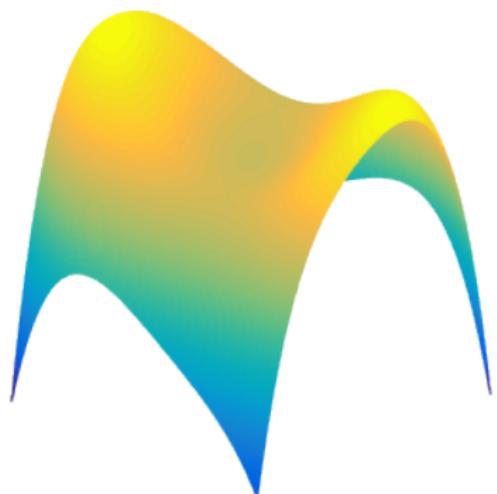
# Nonconvex estimation problems are everywhere

---

Empirical risk minimization is usually nonconvex

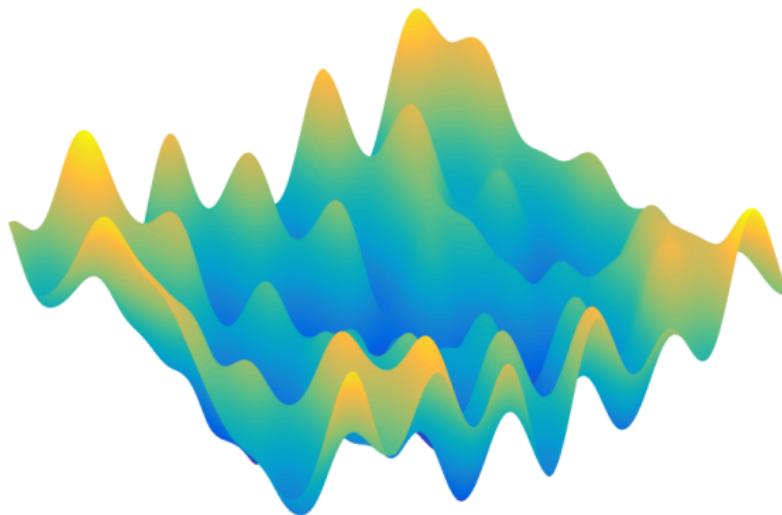
$\text{minimize}_x \quad f(x; \text{data}) \quad \rightarrow \quad \text{loss function may be nonconvex}$

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



# Nonconvex optimization may be super scary

---



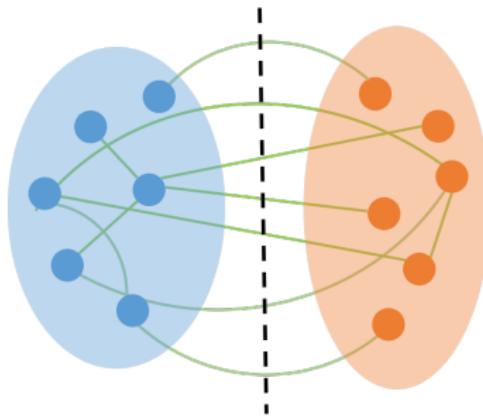
There may be bumps everywhere and exponentially many local optima

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

## Example: solving quadratic programs is hard

Finding maximum cut in a graph is about solving a quadratic program

$$\begin{array}{ll}\text{maximize}_x & x^\top W x \\ \text{subj. to} & x_i^2 = 1, \quad i = 1, \dots, n\end{array}$$



## Example: solving quadratic programs is hard

---



"I can't find an efficient algorithm, but neither can all these people."

*figure credit: coding horror*

**\$1,000,000 question**

## One strategy: convex relaxation

---

Can relax into convex problems by

- finding convex surrogates (e.g. matrix completion)
- lifting into higher dimensions (e.g. Max-Cut)

# Example of convex surrogate: matrix completion

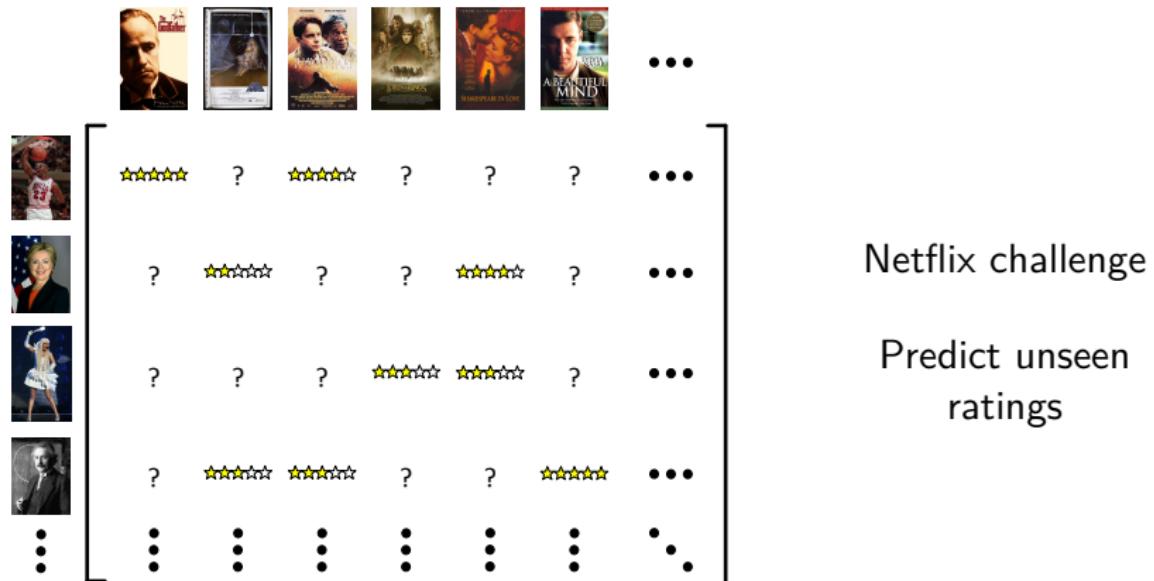


figure credit: Candès et al.

# Example of convex surrogate: matrix completion

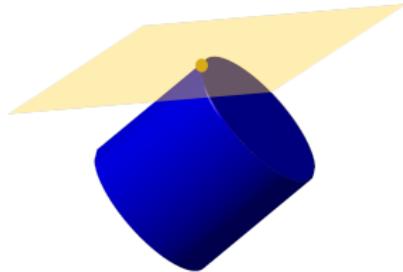
---

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$\text{minimize}_M \text{rank}(M)$  subj. to data constraints

↓ cvx surrogate

$\text{minimize}_M \text{nuc-norm}(M)$  subj. to data constraints



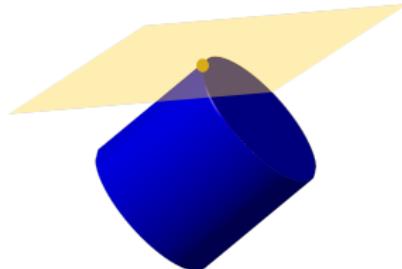
# Example of convex surrogate: matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$\text{minimize}_M \text{rank}(M)$  subj. to data constraints

↓ cvx surrogate

$\text{minimize}_M \text{nuc-norm}(M)$  subj. to data constraints



*Robust variation used everyday by  
Netflix* Candès, Li, Ma, Wright '10

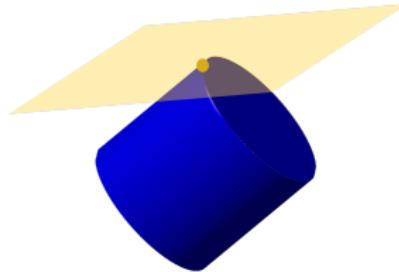
# Example of convex surrogate: matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$\text{minimize}_M \text{rank}(M)$  subj. to data constraints

↓ cvx surrogate

$\text{minimize}_M \text{nuc-norm}(M)$  subj. to data constraints



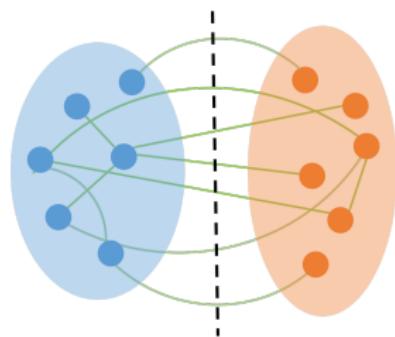
*Robust variation used everyday by  
Netflix* Candès, Li, Ma, Wright '10

**Problem:** operate in *full* matrix space even though  $X$  is low-rank

# Example of lifting: Max-Cut

---

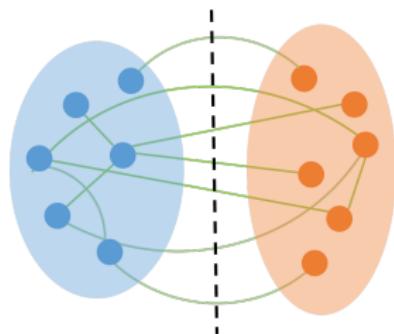
— Goemans, Williamson '95



$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

# Example of lifting: Max-Cut

— Goemans, Williamson '95



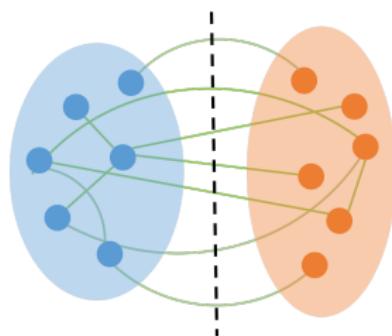
$$\begin{aligned} & \text{maximize}_{\mathbf{x}} && \mathbf{x}^\top \mathbf{W} \mathbf{x} \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓  
let  $\mathbf{X}$  be  $\mathbf{x}\mathbf{x}^\top$

$$\begin{aligned} & \text{maximize}_{\mathbf{X}} && \langle \mathbf{X}, \mathbf{W} \rangle \\ & \text{subj. to} && \mathbf{X}_{i,i} = 1, \quad i = 1, \dots, n \\ & && \mathbf{X} \succeq \mathbf{0} \\ & && \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

# Example of lifting: Max-Cut

— Goemans, Williamson '95



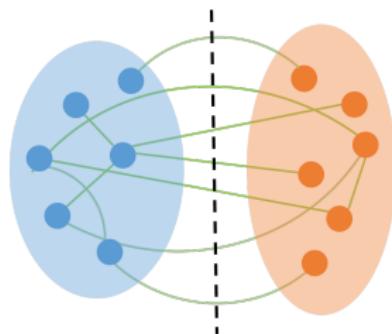
$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓  
let  $X$  be  $xx^\top$

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

# Example of lifting: Max-Cut

— Goemans, Williamson '95



$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓  
let  $X$  be  $xx^\top$

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

**Problem:** explosion in dimensions ( $\mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$ )

*How about optimizing nonconvex problems directly  
without lifting?*

## **Nonconvex optimization**

---

Complicated nonconvex problems are solved on a daily basis via simple algorithms such as stochastic gradient descent

# Nonconvex optimization

---

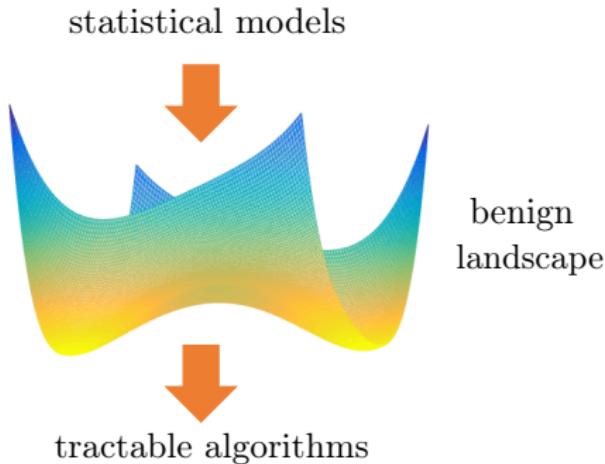
Complicated nonconvex problems are solved on a daily basis via simple algorithms such as stochastic gradient descent



- How come simple nonconvex algorithms work so well in practice?

# Statistical models come to rescue

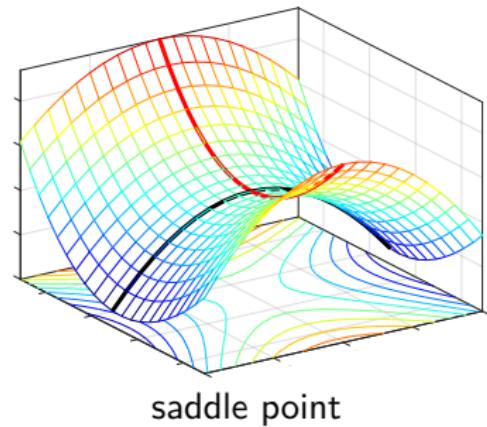
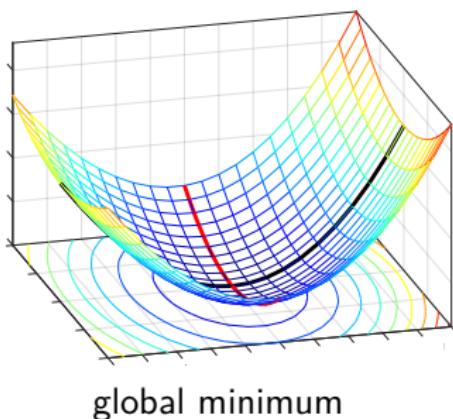
---



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

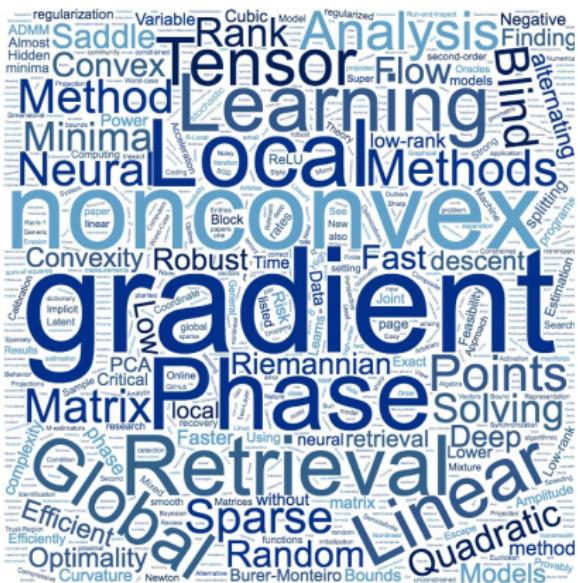
# Sometimes they are much nicer than we think

Under certain **statistical models**,  
we see benign global geometry: **no spurious local optima**



*Even the simplest possible nonconvex methods  
might be remarkably efficient under suitable statistical models*

# Nonconvex optimization with performance guarantees



- <http://sunju.org/research/nonconvex/>
  - “Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview,” Chi, Lu, Chen ’18

**Phase retrieval:** Gerchberg-Saxton '72, Netrapalli et al. '13, Candès, Li, Soltanolkotabi '14, Chen, Candès '15, Cai, Li, Ma '15, Zhang et al. '16, Wang et al. '16, Sun et al. '16, Ma et al. '17, Chen et al. '18, ...

**Matrix completion:** Keshavan et al. '09, Jain et al. '09, Hardt '13, Sun, Luo '15, Chen, Wainwright '15, Zheng, Lafferty '16, Ge et al. '16, Jin et al. '16, Ma et al. '17, ...

**Matrix sensing:** Jain et al. '13, Tu et al. '15, Zheng, Lafferty '15, Bhojanapalli et al. 16, Li, Zhu, Tang '18, ...

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

**Dictionary learning:** Arora et al. '14, Sun et al. '15, Chatterji, Bartlett '17, ...

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

# Outline

---

- Part I: Overview
- Part II: Phase retrieval: a case study
  - Spectral initialization
  - Local refinement: algorithm and analysis
- Part III: Low-rank matrix estimation
- Part IV: Closing remarks

# 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  $A$  is shown as an  $m \times n$  grid of orange squares. To its right is a vector  $x^*$  consisting of  $n$  blue squares. An equals sign follows. To the right of that is a vector  $Ax^*$  consisting of  $m$  blue squares, each containing a numerical value. A large arrow points from  $Ax^*$  to the final column  $y$ , which consists of  $m$  blue squares.

$Ax^*$	$y =  Ax^* ^2$
1	1
-3	9
2	4
-1	1
4	16
2	4
-2	4
-1	1
3	9
4	16

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

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

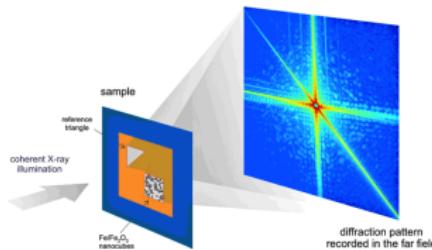
assume w.l.o.g.  $\|\boldsymbol{x}^*\|_2 = 1$

# 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*



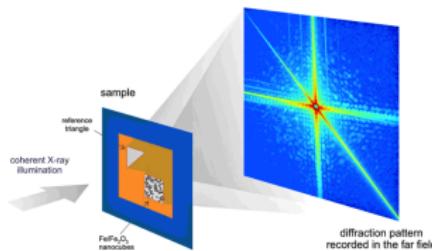
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

# 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: covariance estimation from quadratic sketches

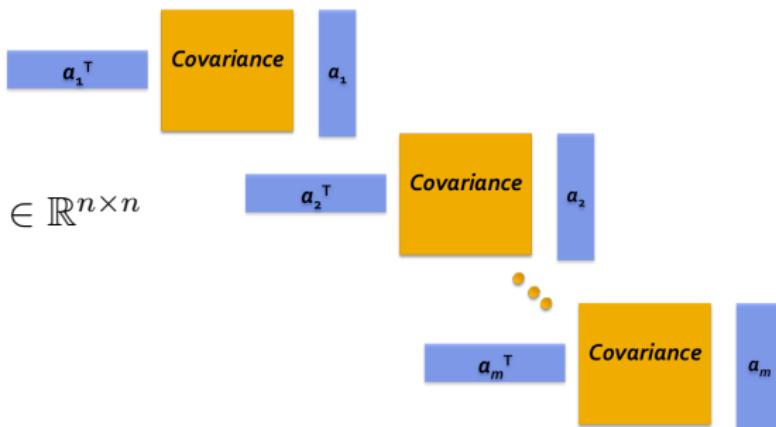
---

— Chen, Chi, Goldsmith '13, Cai, Zhang '13

- **Data:**  $m$  quadratic measurements about *low-rank* covariance matrix  $\Sigma$

$$y_i = \mathbf{a}_i^\top \Sigma \mathbf{a}_i + \text{noise}, \quad i = 1, \dots, m$$

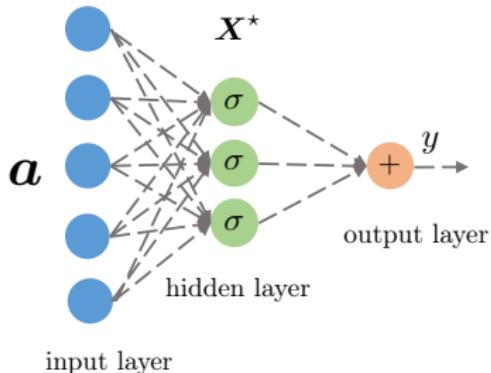
- **Goal:** recover  $\Sigma \in \mathbb{R}^{n \times n}$



# Motivation: learning neural nets with quadratic activation

---

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

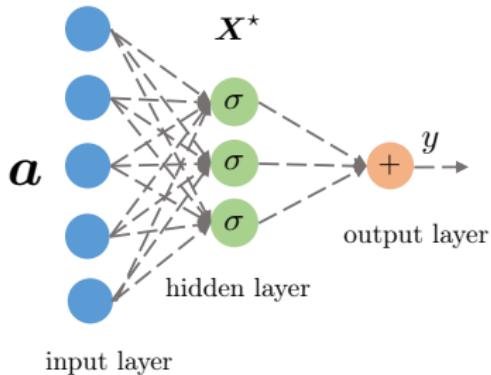


input features:  $\mathbf{a}$ ; weights:  $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

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

# Motivation: learning neural nets with quadratic activation

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

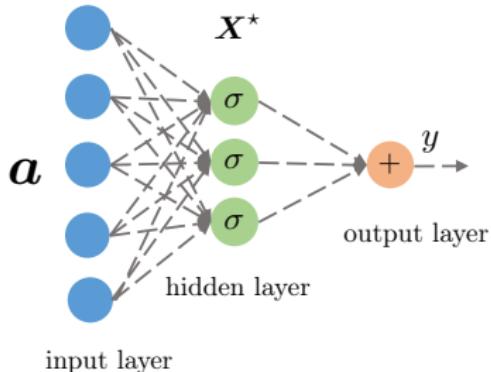


input features:  $a$ ; weights:  $\mathbf{X}^* = [x_1^*, \dots, x_r^*]$

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

# Motivation: learning neural nets with quadratic activation

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



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

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

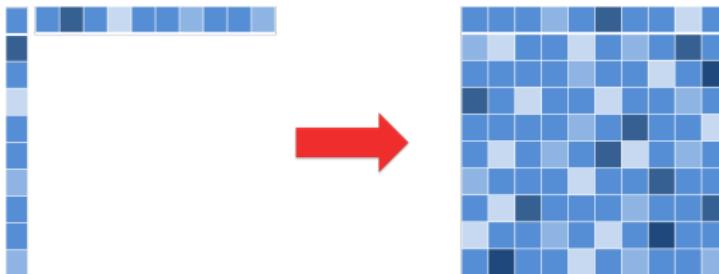
We consider simplest model when  $r = 1$

## An equivalent view: low-rank factorization

---

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

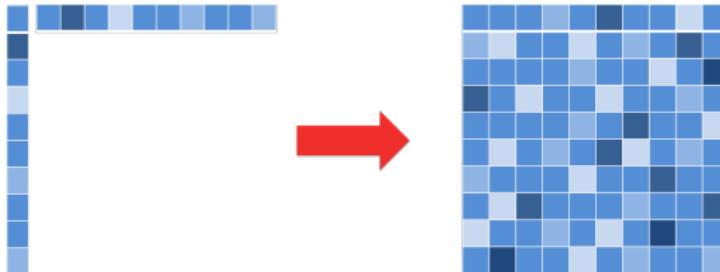
$$y_k = (\mathbf{a}_k^\top \mathbf{x})^2 = \mathbf{a}_k^\top (\mathbf{x}\mathbf{x}^\top) \mathbf{a} \implies y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k$$



## 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  $\mathbf{X}$

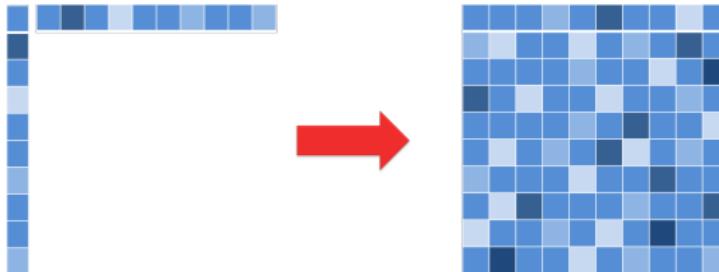
$$\text{s.t. } y_k = \mathbf{a}_k^\top \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m$$

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

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

$\Downarrow$

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

## A natural least squares formulation

---

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

$\Downarrow$

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

- **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)

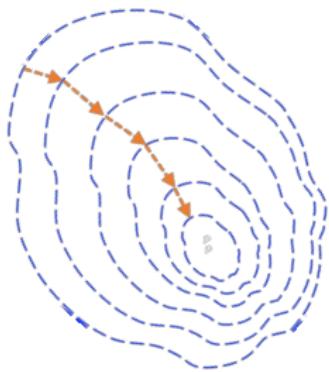
---

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

# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

---

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

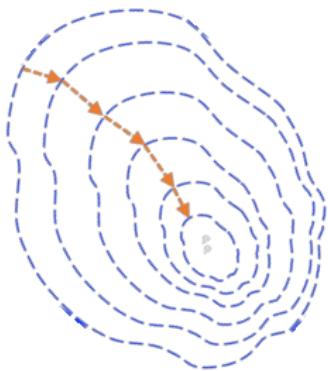


- **spectral initialization:**  $\boldsymbol{x}^0 \leftarrow$  leading eigenvector of certain data matrix

# 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

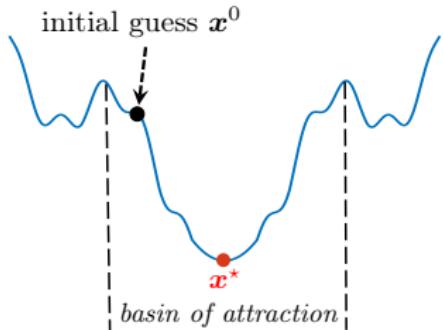
$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top$$

**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  $\underbrace{\text{local basin sufficiently close to } x^*}_{\text{(restricted) strongly convex; no saddles / spurious local mins}}$

# 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
- ...

# Computational cost

---

$$\mathbf{A}\mathbf{x} := [\mathbf{a}_k^\top \mathbf{x}]_{1 \leq k \leq m}$$

- **Spectral initialization:** leading eigenvector  $\rightarrow$  a few applications of  $\mathbf{A}$  and  $\mathbf{A}^\top$

$$\frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top = \frac{1}{m} \mathbf{A}^\top \operatorname{diag}\{y_k\} \mathbf{A}$$

# Computational cost

---

$$\mathbf{A}\mathbf{x} := [\mathbf{a}_k^\top \mathbf{x}]_{1 \leq k \leq m}$$

- **Spectral initialization:** leading eigenvector  $\rightarrow$  a few applications of  $\mathbf{A}$  and  $\mathbf{A}^\top$

$$\frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top = \frac{1}{m} \mathbf{A}^\top \operatorname{diag}\{y_k\} \mathbf{A}$$

- **Iterations:** one application of  $\mathbf{A}$  and  $\mathbf{A}^\top$  per iteration

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

# Asymptotic notation

---

- $f(n) \lesssim g(n)$  or  $f(n) = O(g(n))$  means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}$$

- $f(n) \gtrsim g(n)$  means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \geq \text{const}$$

- $f(n) \asymp g(n)$  means

$$\text{const}_1 \leq \lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}_2$$

# 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

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

## Theorem 1 (Candès, Li, Soltanolkotabi '14)

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

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

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  $m \gtrsim n \log n$ .

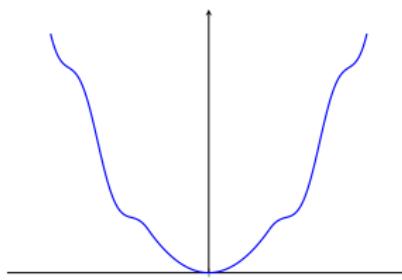
- Iteration complexity:  $O(n \log \frac{1}{\epsilon}) \searrow O(\log n \log \frac{1}{\epsilon})$
- Sample complexity:  $O(n \log n)$
- Derived based on finer analysis of GD trajectory

# Gradient descent theory revisited

---

Consider unconstrained optimization problem

$$\text{minimize}_x \quad f(x)$$



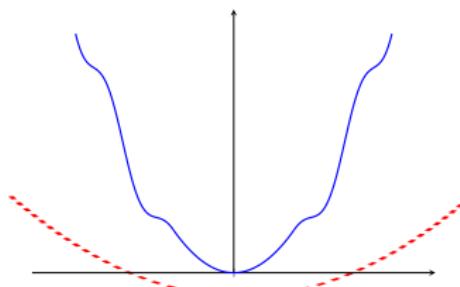
Two standard conditions that enable geometric convergence of GD

# Gradient descent theory revisited

---

Consider unconstrained optimization problem

$$\text{minimize}_x \quad f(x)$$



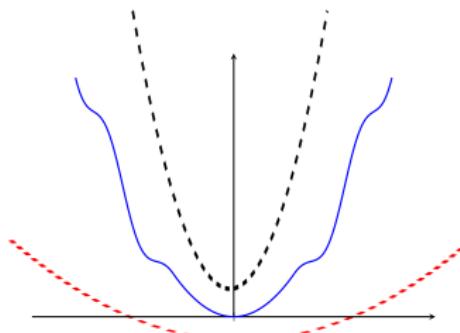
Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)

# Gradient descent theory revisited

Consider unconstrained optimization problem

$$\text{minimize}_x \quad f(x)$$



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)
- (local) smoothness

$$\nabla^2 f(x) \succ 0 \quad \text{and} \quad \text{is well-conditioned}$$

# 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

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{\star}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^{\star}\|_2$$

# 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

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^*\|_2$$

- Condition number  $\beta/\alpha$  determines rate of convergence

# 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

$$\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^*\|_2$$

- Condition number  $\beta/\alpha$  determines rate of convergence
- Attains  $\varepsilon$ -accuracy within  $O\left(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon}\right)$  iterations

# What does this 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 this 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 this 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** ( $m \asymp n \log n$ )

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

# What does this optimization theory say about WF?

---

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

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

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

**Consequence (Candès et al '14):** WF attains  $\varepsilon$ -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  $\eta = O(1/n)$



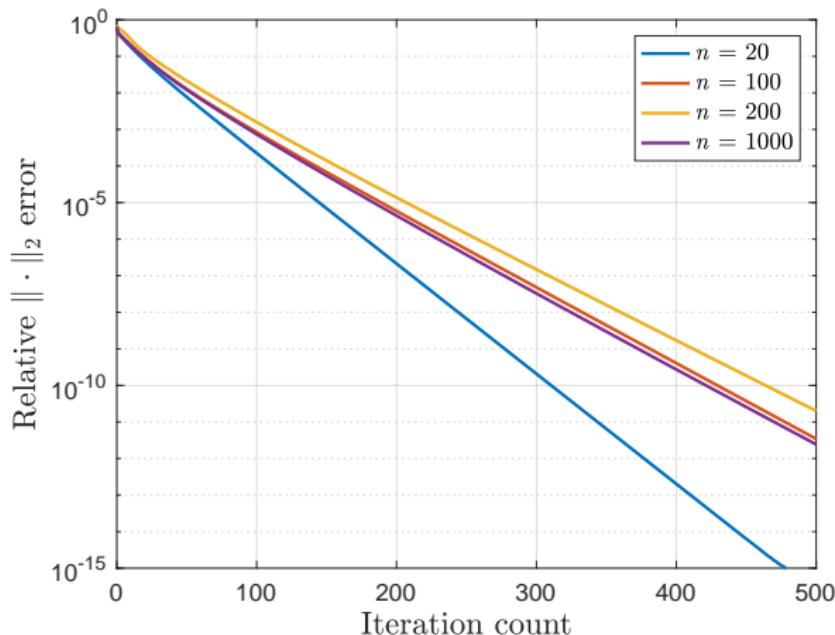
This choice is suggested by **worst-case** optimization theory



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

---

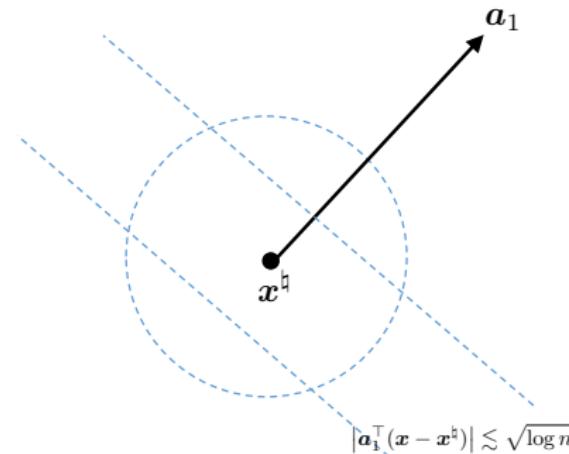
Which local region enjoys both strong convexity and smoothness?

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

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

# A second look at gradient descent theory

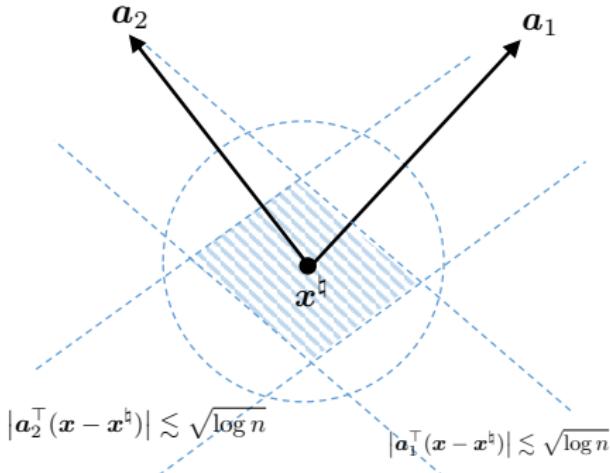
Which local region enjoys both strong convexity and smoothness?



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

# A second look at gradient descent theory

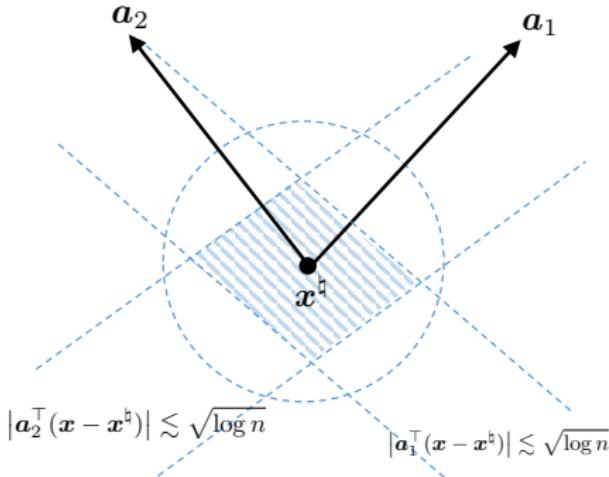
Which local region enjoys both strong convexity and smoothness?



- $x$  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**)

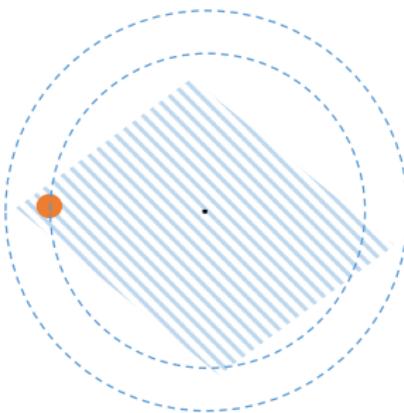
Prior works suggest enforcing **regularization** (e.g. truncation, projection, regularized loss) to promote incoherence

# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

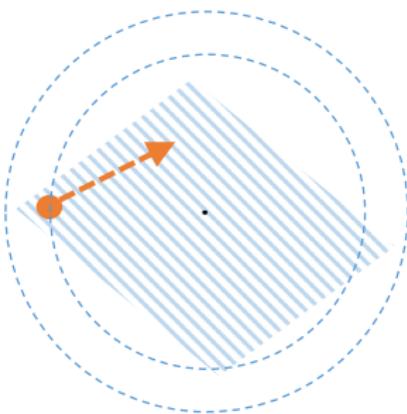


# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

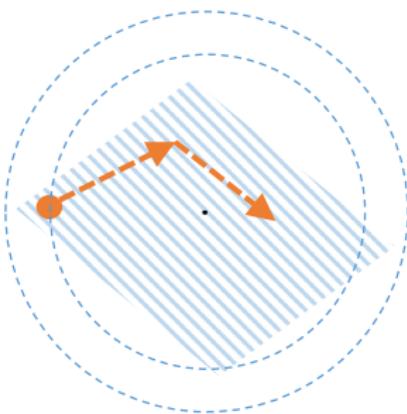


# Encouraging message: GD is implicitly regularized

---



region of local strong convexity + smoothness

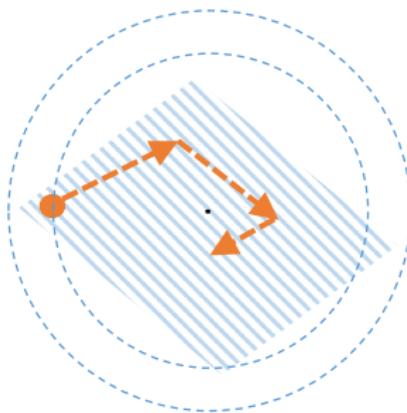


# Encouraging message: GD is implicitly regularized

---



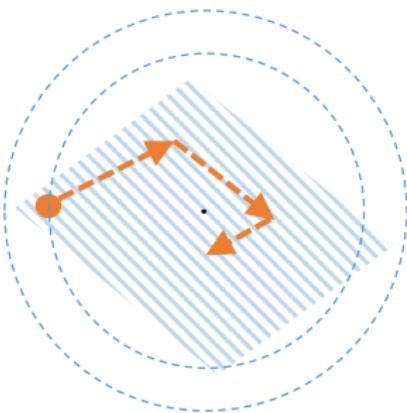
region of local strong convexity + smoothness



# 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  $\varepsilon$  accuracy within  $O(\log n \log \frac{1}{\varepsilon})$  iterations

# Key proof idea: leave-one-out analysis

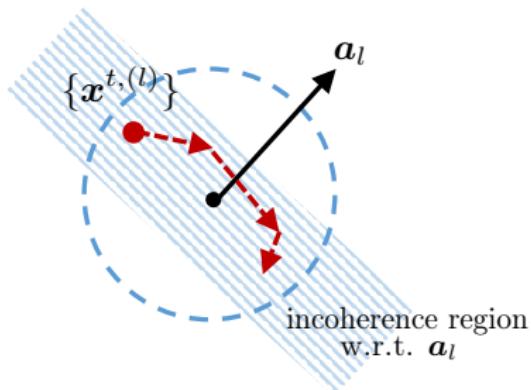
For each  $1 \leq l \leq m$ , introduce leave-one-out iterates  $\mathbf{x}^{t,(l)}$  by dropping  $l$ th measurement

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

The diagram illustrates the computation of a leave-one-out iterate. On the left, a matrix  $A^{(l)}$  is shown as a 4x4 grid of orange and brown squares. Below it, its transpose  $a_l^\top$  is shown as a 4x1 vector of orange and brown squares. An equals sign follows. To the right, the product  $A^{(l)} \mathbf{x}^*$  is shown as a 4x1 vector with entries 1, -3, 2, and -1. Below this, a horizontal line with a gap indicates the row corresponding to  $a_l^\top$  is omitted. An arrow points to the result  $y^{(l)} = |A^{(l)} \mathbf{x}^*|^2$ , which is a 4x1 vector with entries 1, 9, 4, and 1. Below this, another horizontal line with a gap indicates the row corresponding to  $a_l^\top$  is included again.

## Key proof idea: leave-one-out analysis

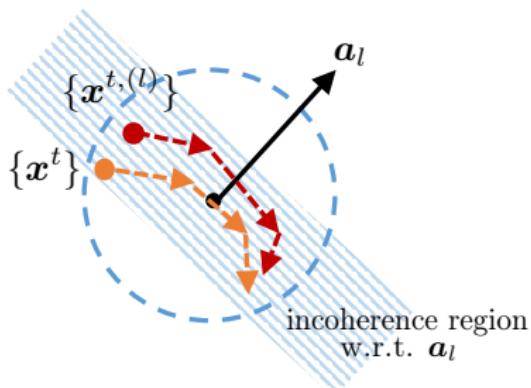
---



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

## Key proof idea: leave-one-out analysis

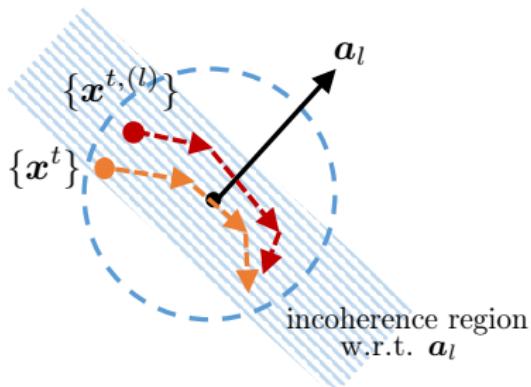
---



- 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

---

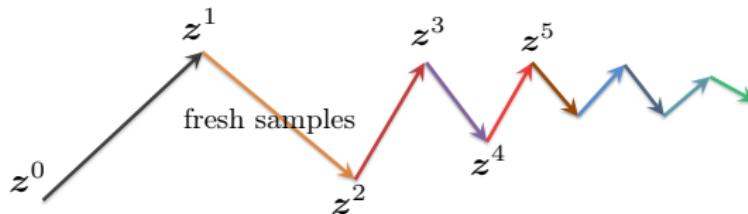


- Leave-one-out iterate  $x^{t,(l)}$  is independent of  $a_l$
- Leave-one-out iterate  $x^{t,(l)} \approx$  true iterate  $x^t$   
 $\implies x^t$  is nearly independent of  $a_l$   
nearly orthogonal to

# No need of sample splitting

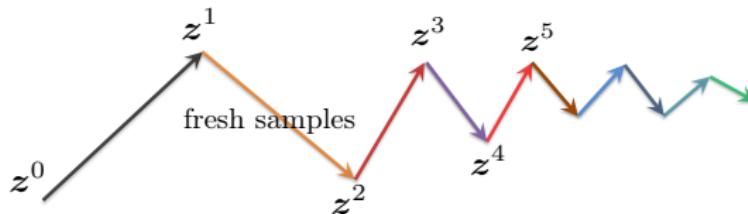
---

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

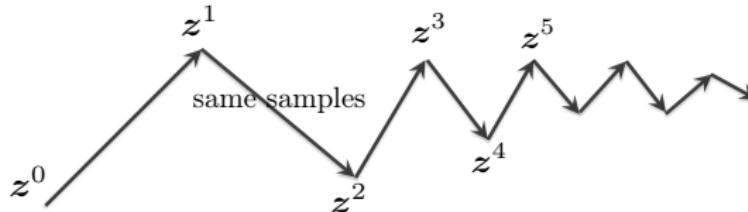


# No need of sample splitting

- 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



# Questions

---

So far we have presented theory for

spectral initialization + vanilla gradient descent (WF)

# Questions

---

So far we have presented theory for

spectral initialization + vanilla gradient descent (WF)

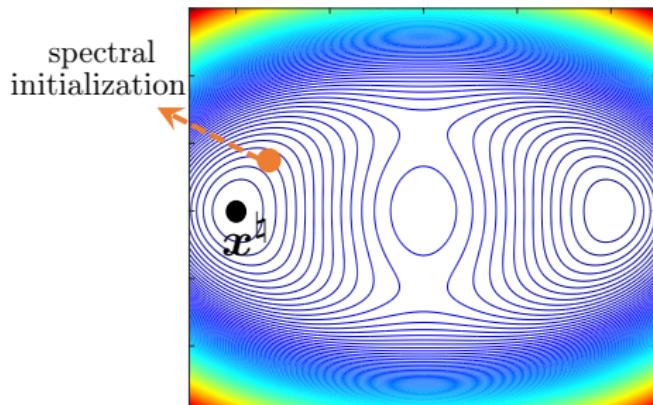
## Questions:

- Is carefully-designed initialization necessary for fast convergence?
- Can we further improve sample complexity?
- Robustness vis a vis noise and outliers?

*Is carefully-designed initialization necessary for fast convergence?*

# Initialization

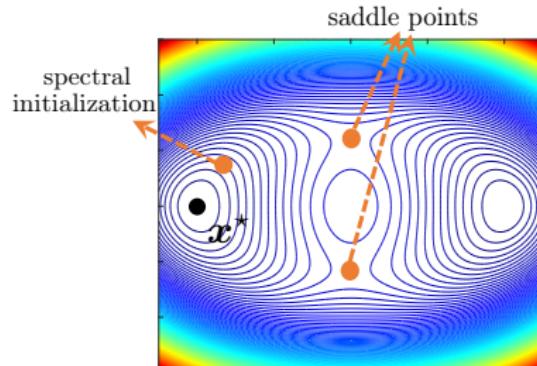
---



- Spectral initialization gets us reasonably close to truth

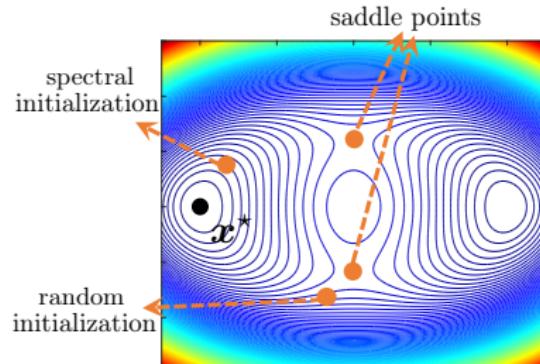
# Initialization

---



- 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



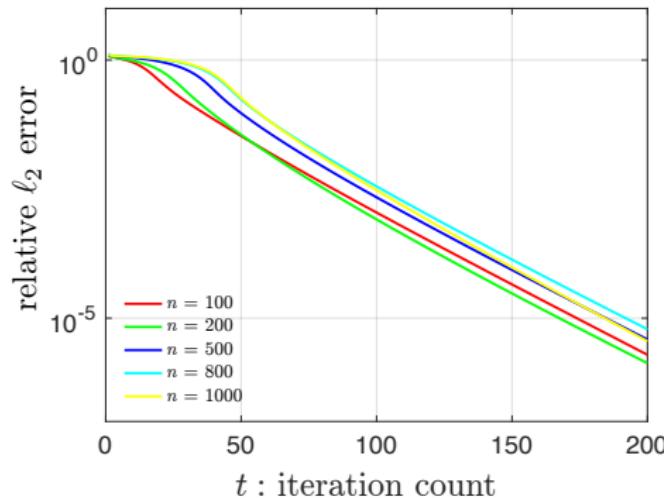
- Spectral initialization gets us reasonably close to truth
- Cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

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

# Numerical efficiency of randomly initialized GD

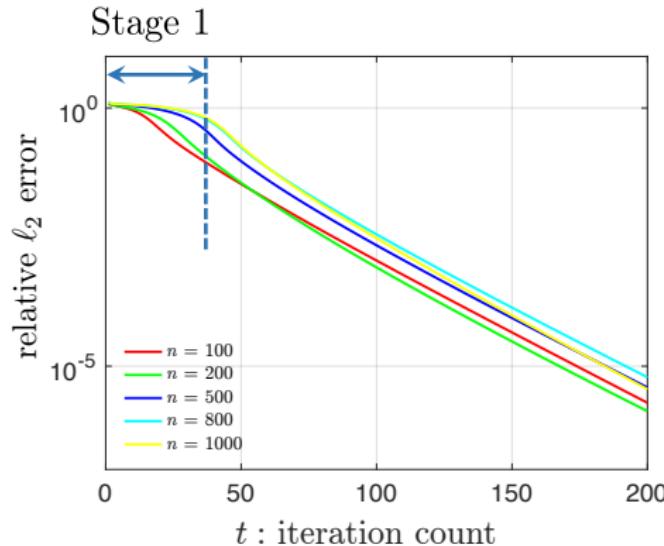
---

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



# Numerical efficiency of randomly initialized GD

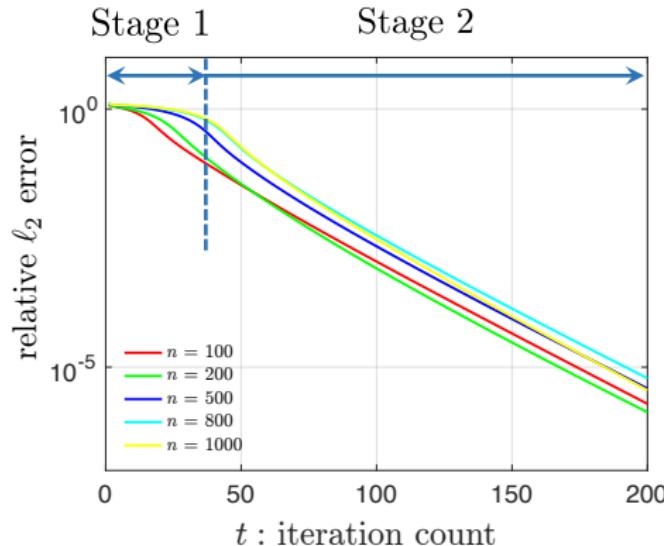
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **a few iterations**

# Numerical efficiency of randomly initialized GD

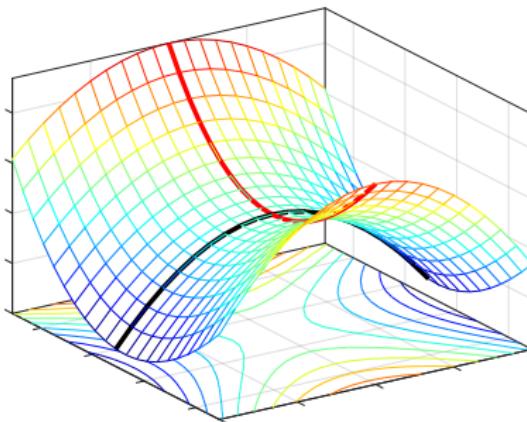
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **a few iterations**

# A geometric analysis

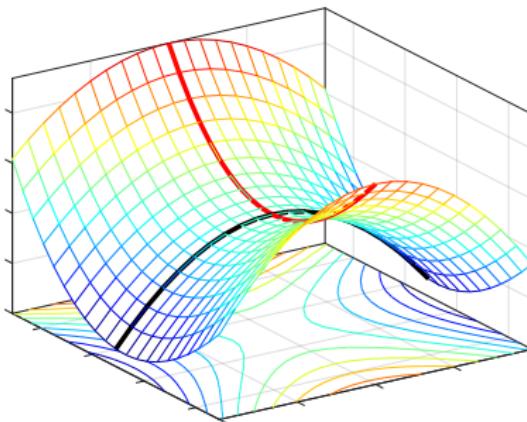
---



- 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

---

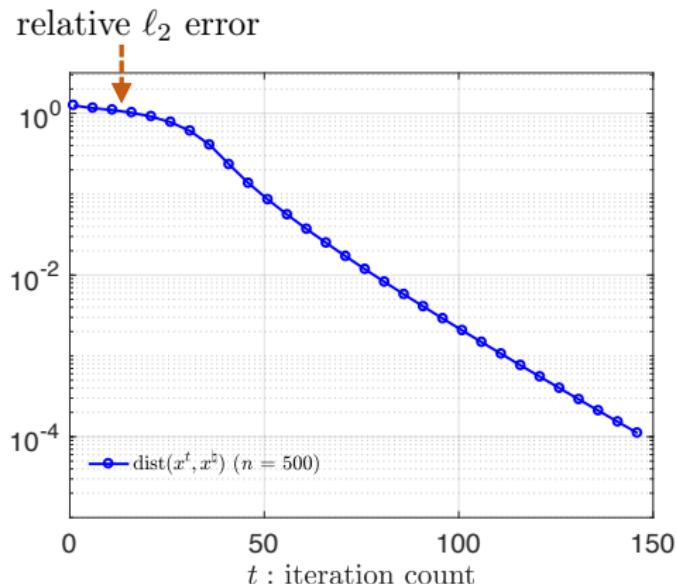


- With such benign landscape, GD with random initialization converges to global min **almost surely** (Lee et al. '16)

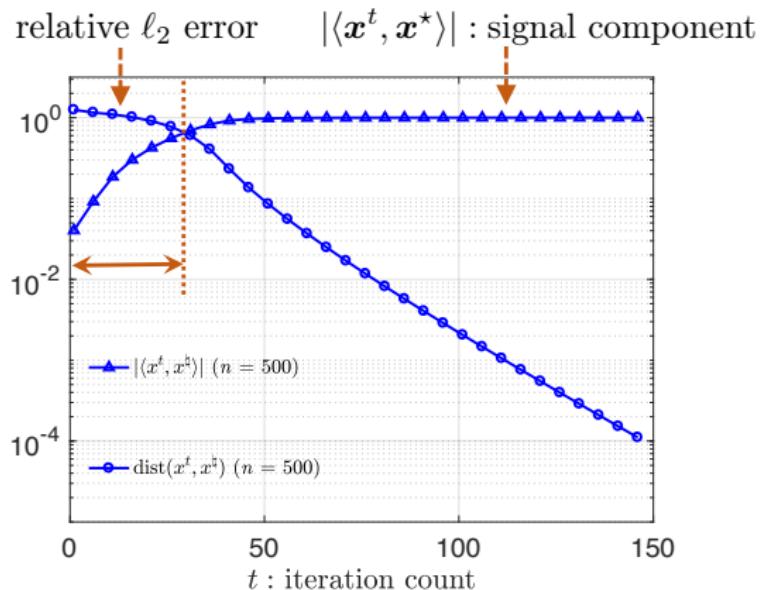
No convergence rate guarantees for vanilla GD!

# Exponential growth of signal strength in Stage 1

---



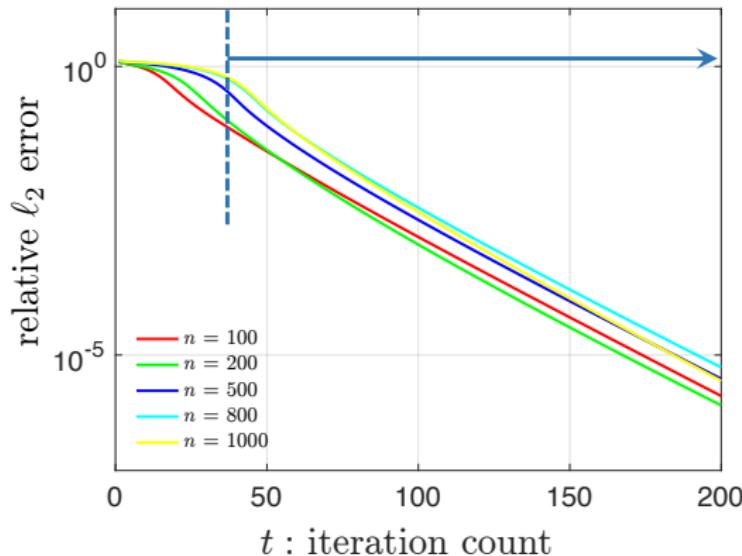
# Exponential growth of signal strength in Stage 1



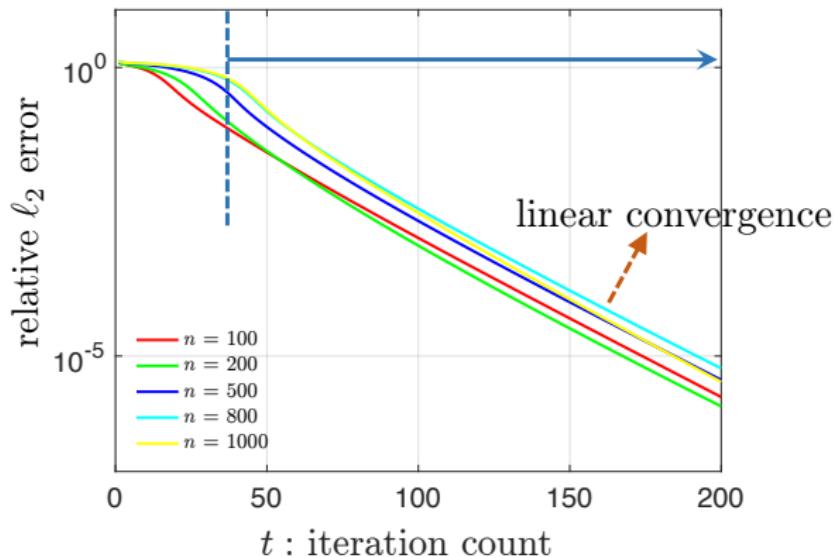
Numerically,  $O(\log n)$  iterations are enough to enter local region

# 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 4 (Chen, Chi, Fan, Ma '18)

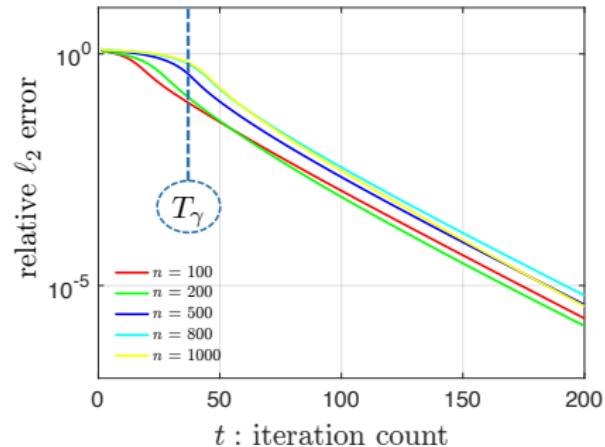
Under i.i.d. Gaussian design, GD with  $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$  achieves

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

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

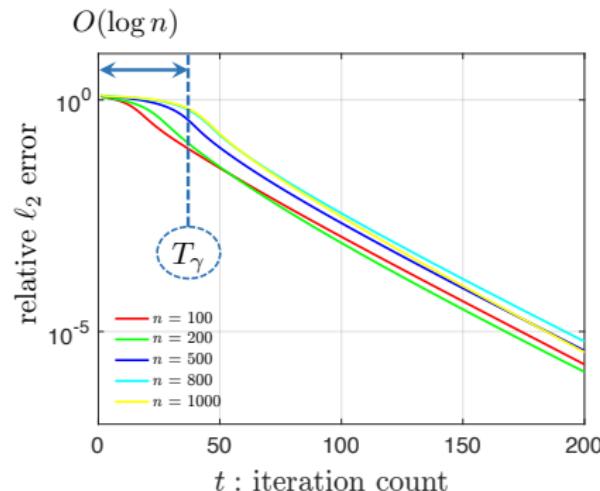
# Theoretical guarantees for randomly initialized GD

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



# Theoretical guarantees for randomly initialized GD

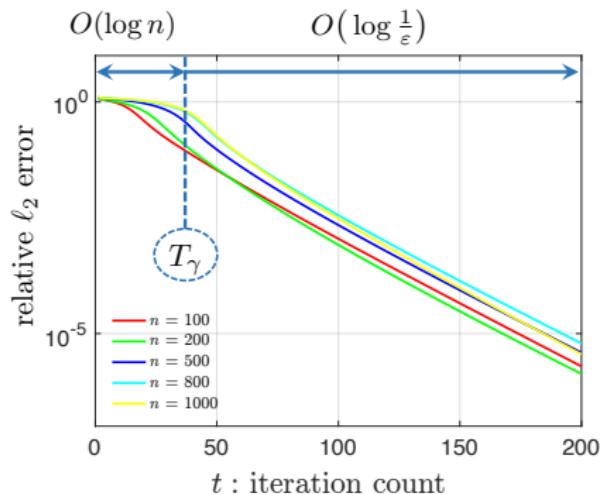
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Theoretical guarantees for randomly initialized GD

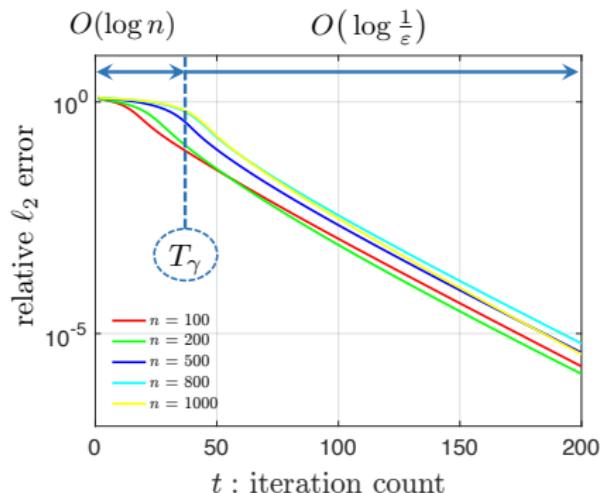
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- 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

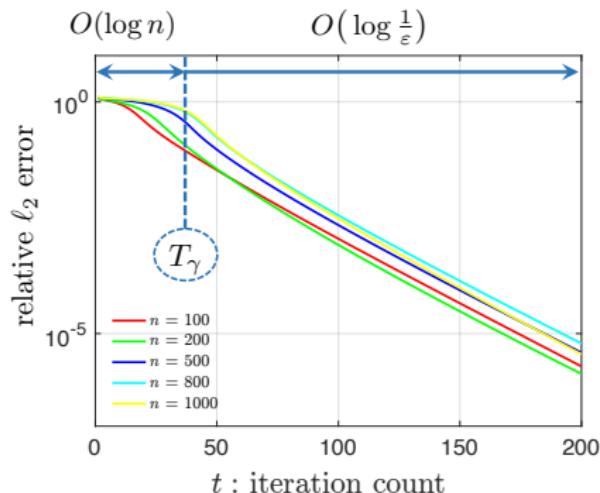
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\star) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\star\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Theoretical guarantees for randomly initialized GD

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



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

# Experiments on images

---



- coded diffraction patterns
- $x^* \in \mathbb{R}^{256 \times 256}$
- $m/n = 12$

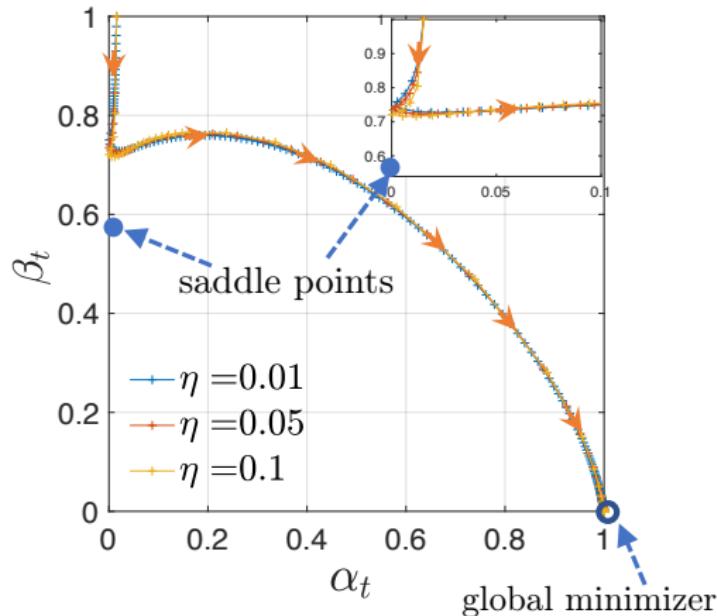
# GD with random initialization

---

$x^t$	$\langle x^t, x^* \rangle x^*$	$x^t - \langle x^t, x^* \rangle x^*$
GD iterate	signal component	perpendicular component

*use Adobe Acrobat to see animation*

# Saddle-escaping schemes?



Randomly initialized GD never hits saddle points in phase retrieval!

# Other saddle-escaping schemes

---

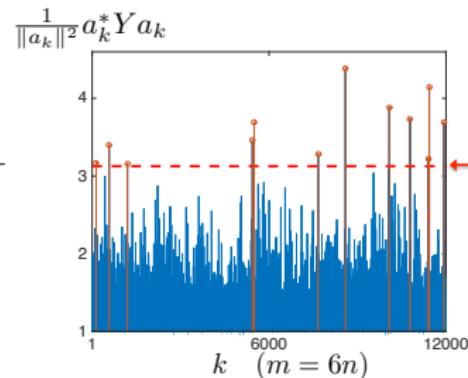
	iteration complexity	num of iterations needed to escape saddles	local iteration complexity
<b>Trust-region</b> (Sun et al. '16)	$n^7 + \log \log \frac{1}{\varepsilon}$	$n^7$	$\log \log \frac{1}{\varepsilon}$
<b>Perturbed GD</b> (Jin et al. '17)	$n^3 + n \log \frac{1}{\varepsilon}$	$n^3$	$n \log \frac{1}{\varepsilon}$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

Generic optimization theory yields highly suboptimal convergence guarantees

*Can we further improve sample complexity?*

# Truncated spectral initialization

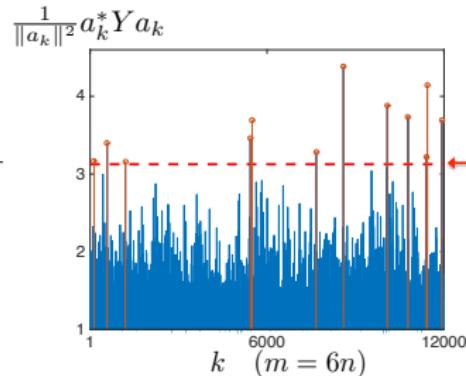
$$\mathbb{E}[\mathbf{Y}] := \mathbb{E} \left[ \frac{1}{m} \sum_{k=1}^m \mathbf{y}_k \mathbf{a}_k \mathbf{a}_k^\top \right] = \|\mathbf{x}^*\|_2^2 \mathbf{I} + 2\mathbf{x}^* \mathbf{x}^{*\top}$$



**problem:** unless  $m \gg n$ , dangerous to use empirical average because large observations  $y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2$  bear too much influence

# Truncated spectral initialization

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



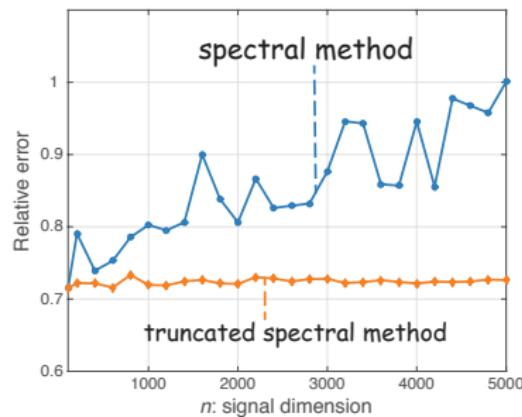
**problem:** unless  $m \gg n$ , dangerous to use empirical average because large observations  $y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2$  bear too much influence

**solution:** discard high leverage samples and compute leading eigenvector of truncated sum

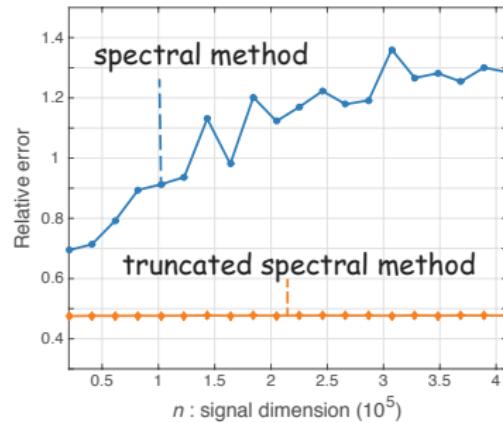
$$\frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^\top \cdot \mathbf{1}_{\{|y_k| \leq \alpha^2 \text{Avg}(|y_j|)\}}$$

# Importance of truncated spectral initialization

---



real Gaussian  $m = 6n$



complex CDP  $m = 12n$

# Importance of truncated spectral initialization

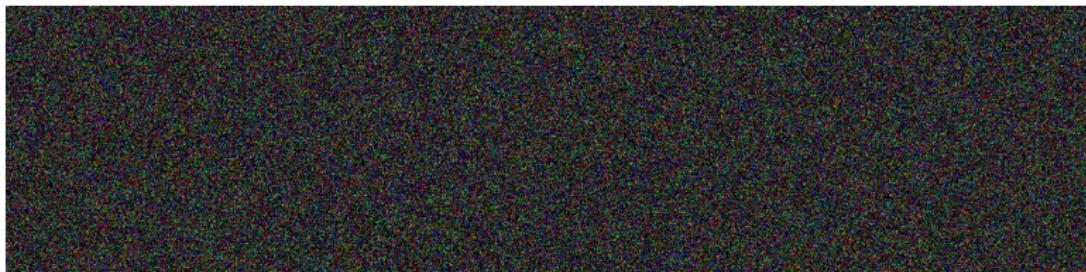
---



Original image

# Importance of truncated spectral initialization

---



Spectral initialization

# Importance of truncated spectral initialization

---



Spectral initialization



Truncated spectral initialization

# Precise asymptotic characterization (Lu, Li '17)

- $m/n \asymp 1$
- i.i.d. Gaussian design

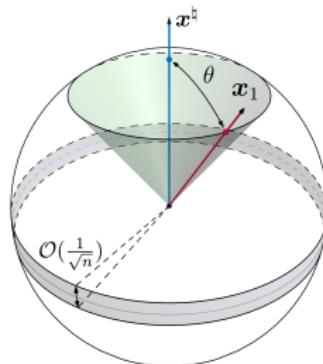


Fig. credit: Lu, Li '17

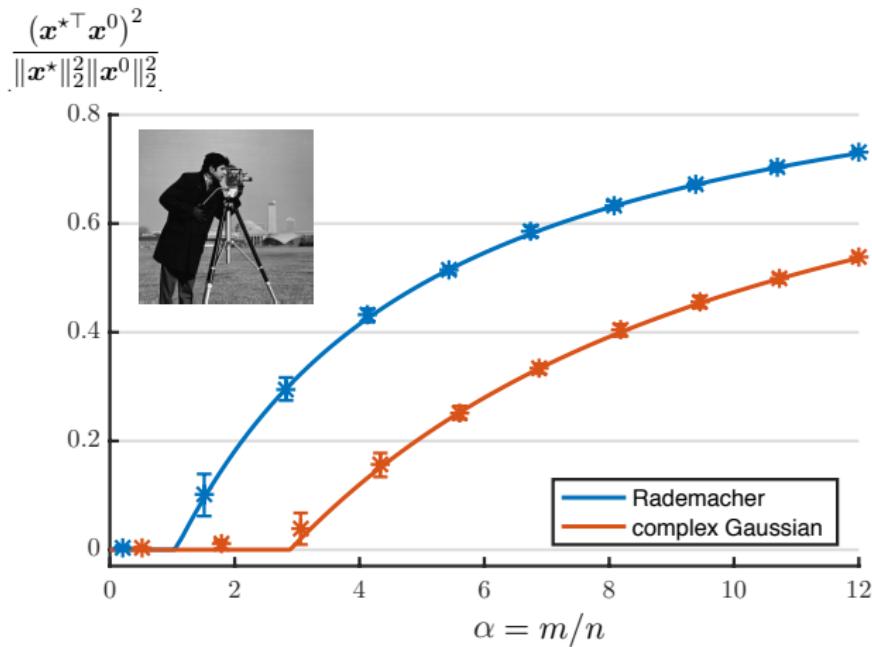
## Theorem 5 (Lu, Li '17, Mondelli, Montanari '17)

There exist analytical formulas  $\rho(\cdot)$  and constants  $\alpha_{\min}$  and  $\alpha_{\max}$  s.t.

$$\underbrace{\frac{(\mathbf{x}^{\star \top} \mathbf{x}^0)^2}{\|\mathbf{x}^{\star}\|_2^2 \|\mathbf{x}^0\|_2^2}}_{\text{cosine similarity}} \rightarrow \begin{cases} 0, & \text{if } m/n < \alpha_{\min} \\ \rho(m/n), & \text{if } m/n > \alpha_{\max} \end{cases}$$

# Theoretical prediction vs. simulations

image size:  $64 \times 64$



*Fig. credit: Lu, Li '17*

# Improving search directions

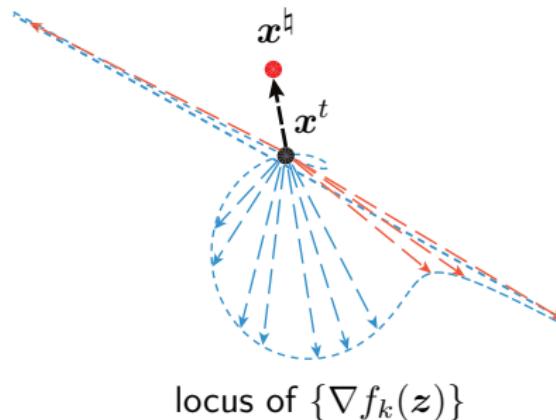
---

$$\text{WF (GD): } \boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \frac{\eta}{m} \sum_k \nabla f_k(\boldsymbol{x}^t)$$

# Improving search directions

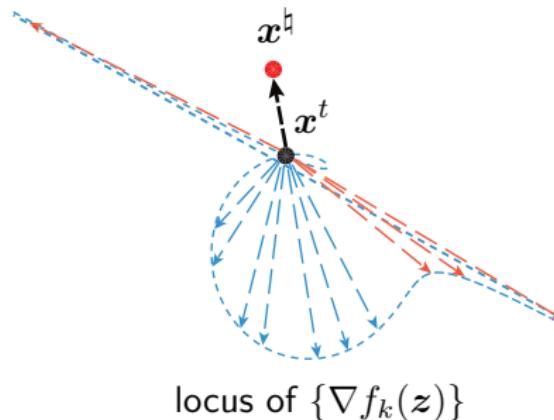
---

$$\text{WF (GD): } \mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_k \nabla f_k(\mathbf{x}^t)$$



# Improving search directions

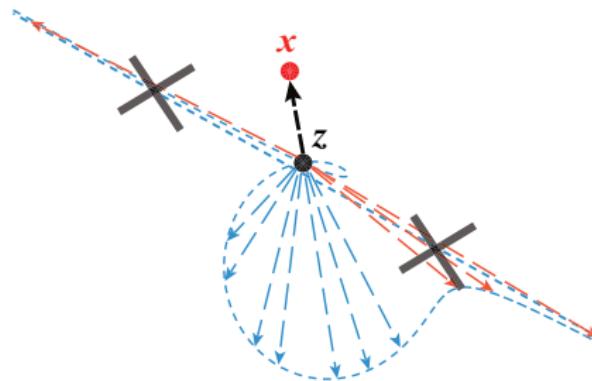
$$\text{WF (GD): } \boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \frac{\eta}{m} \sum_k \nabla f_k(\boldsymbol{x}^t)$$



**Problem:** descent direction might have large variability

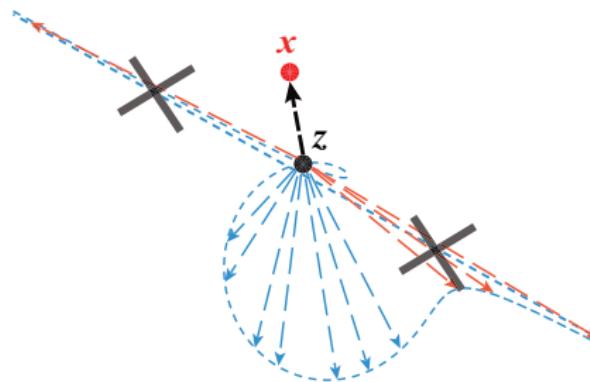
## Solution: variance reduction via trimming

More adaptive rule:  $\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_{k \in \mathcal{T}_t} \nabla f_k(\mathbf{x}^t)$



## Solution: variance reduction via trimming

More adaptive rule:  $\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_{k \in \mathcal{T}_t} \nabla f_k(\mathbf{x}^t)$



- $\mathcal{T}_t$  trims away excessively large grad components

$$\mathcal{T}_t := \left\{ k : \quad \|\nabla f_k(\mathbf{x}^t)\|_2 \lesssim \text{typical-size} \left\{ \|\nabla f_l(\mathbf{x}^t)\|_2 \right\}_{1 \leq l \leq m} \right\}$$

Slight bias + much reduced variance

## Summary: truncated Wirtinger flow

---

(1) **Regularized spectral initialization:**  $x^0 \leftarrow$  principal component of

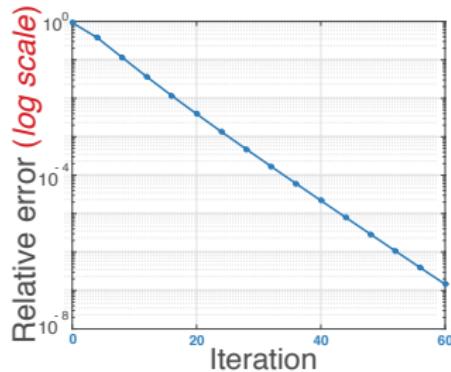
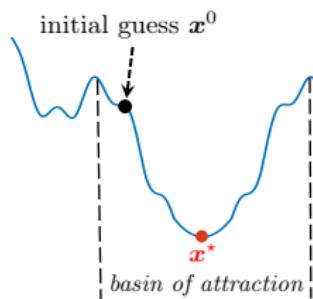
$$\frac{1}{m} \sum_{k \in \mathcal{T}_0} y_k \mathbf{a}_k \mathbf{a}_k^\top$$

(2) Follow **adaptive gradient descent**

$$\mathbf{x}^t = \mathbf{x}^t - \frac{\eta_t}{m} \sum_{k \in \mathcal{T}_t} \nabla f_k(\mathbf{x}^t)$$

**Adaptive and iteration-varying rules:** discard high-leverage data  
 $\{y_k : k \notin \mathcal{T}_t\}$

# Theoretical guarantees (noiseless data)



## Theorem 6 (Chen, Candès '15)

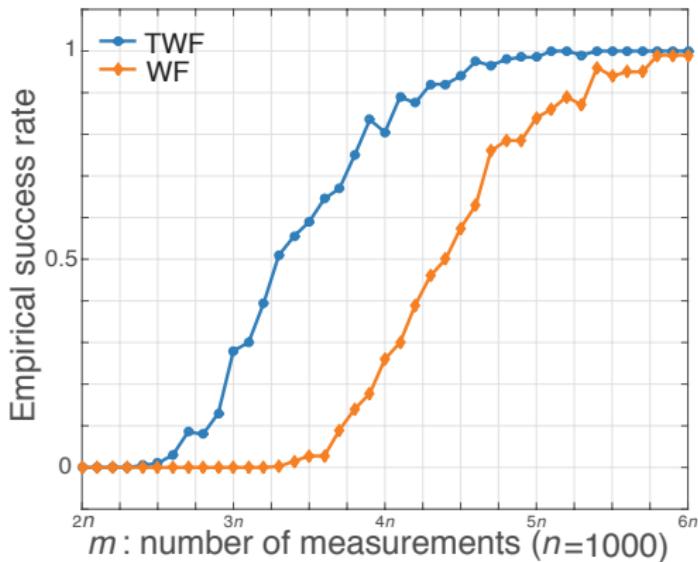
Suppose  $\mathbf{a}_k \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$  and sample size  $m \gtrsim n$ . With high prob.,

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) := \min \|\mathbf{x}^t \pm \mathbf{x}^*\|_2 \leq \nu (1 - \rho)^t \|\mathbf{x}^*\|_2$$

where  $0 < \nu, \rho < 1$  are universal constants

# Empirical success rate (noiseless data)

---



Empirical success rate vs. sample size

*Stability vis a vis noise and outliers?*

# Stability under noisy data

---

- Noisy data:  $y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2 + \eta_k$
- Signal-to-noise ratio:

$$\text{SNR} := \frac{\sum_k (\mathbf{a}_k^\top \mathbf{x}^*)^4}{\sum_k \eta_k^2} \approx \frac{3m \|\mathbf{x}^*\|_2^4}{\|\boldsymbol{\eta}\|_2^2}$$

- i.i.d. Gaussian design  $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$

# Stability under noisy data

---

- Noisy data:  $y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2 + \eta_k$
- Signal-to-noise ratio:

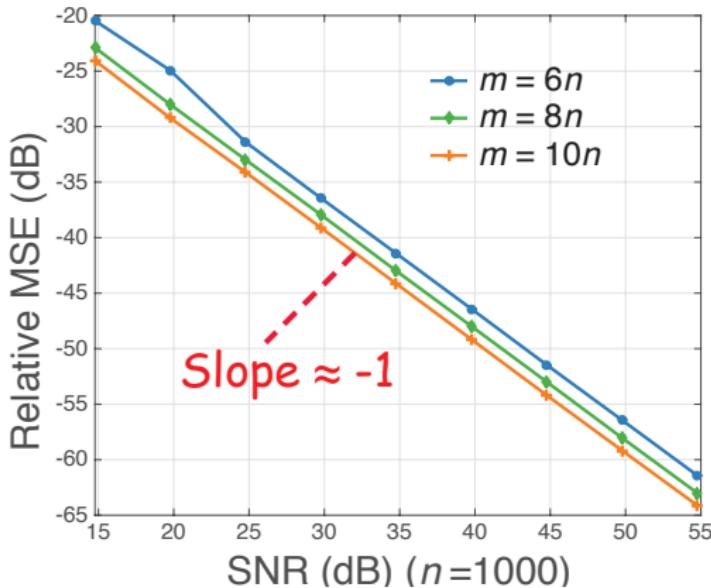
$$\text{SNR} := \frac{\sum_k (\mathbf{a}_k^\top \mathbf{x}^*)^4}{\sum_k \eta_k^2} \approx \frac{3m \|\mathbf{x}^*\|_2^4}{\|\boldsymbol{\eta}\|_2^2}$$

- i.i.d. Gaussian design  $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$

## Theorem 7 (Chen, Candès '15)

Relative error of TWF converges to  $O(\frac{1}{\sqrt{\text{SNR}}})$

# Relative MSE vs. SNR (Poisson data)

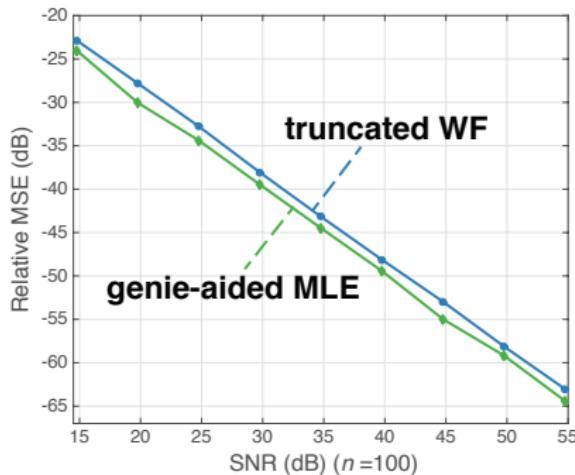


Empirical evidence: relative MSE scales inversely with SNR

# This accuracy is nearly un-improvable (empirically)

Comparison with ideal MLE (with phase info. revealed)

**ideal knowledge:**  $y_k \sim \text{Poisson}(|\mathbf{a}_k^\top \mathbf{x}^*|^2)$  and  $\varepsilon_k = \text{sign}(\mathbf{a}_k^\top \mathbf{x}^*)$



Little loss due to missing phases!

# This accuracy is nearly un-improvable (theoretically)

---

- Poisson data:  $y_k \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\mathbf{a}_k^\top \mathbf{x}^*|^2)$
- Signal-to-noise ratio:

$$\text{SNR} \approx \frac{\sum_k |\mathbf{a}_k^\top \mathbf{x}^*|^4}{\sum_k \text{Var}(y_k)} \approx 3\|\mathbf{x}^*\|_2^2$$

# This accuracy is nearly un-improvable (theoretically)

- Poisson data:  $y_k \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\mathbf{a}_k^\top \mathbf{x}^*|^2)$
- Signal-to-noise ratio:

$$\text{SNR} \approx \frac{\sum_k |\mathbf{a}_k^\top \mathbf{x}^*|^4}{\sum_k \text{Var}(y_k)} \approx 3\|\mathbf{x}^*\|_2^2$$

## Theorem 8 (Chen, Candès '15)

Under i.i.d. Gaussian design, for any estimator  $\hat{\mathbf{x}}$ ,

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x}^*: \|\mathbf{x}^*\|_2 \geq \log^{1.5} m} \frac{\mathbb{E} [\text{dist}(\hat{\mathbf{x}}, \mathbf{x}^*) \mid \{\mathbf{a}_k\}]}{\|\mathbf{x}^*\|_2} \gtrsim \frac{1}{\sqrt{\text{SNR}}},$$

provided that sample size  $m \asymp n$

# Robust recovery vis a vis outliers

---

Consider now two sources of corruption: *sparse outliers* and *bounded noise*

$$y_i = |\mathbf{a}_i^\top \mathbf{x}^*|^2 + \eta_i + w_i, \quad i = 1, \dots, m,$$

- $\|\boldsymbol{\eta}\|_0 \leq s \cdot m$ : sparse outlier, where  $0 \leq s < 1$  is fraction of outliers
- $\mathbf{w}$ : bounded noise

**Motivation:** outliers happen with sensor failures, malicious attacks ...

# Robust recovery vis a vis outliers

---

**Goal:** develop algorithms that are *oblivious* to outliers, and statistically and computationally efficient

- performs equally well regardless of existence of outliers
- small sample size: ideally  $m \asymp n$
- large fraction of outliers: ideally  $s \asymp 1$
- low computational complexity and easy to implement

# Existing approaches are not robust in the presence of arbitrary outliers

---

- **Spectral initialization would fail:** leading eigenvector of  $\mathbf{Y}$  can be arbitrarily perturbed

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m \textcolor{red}{y_i} \mathbf{a}_i \mathbf{a}_i^\top \quad (\text{WF})$$

or  $\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top \mathbb{1}_{\{|y_i| \lesssim \text{mean}(\{\textcolor{red}{y}_i\})\}} \quad (\text{TWF})$

# Existing approaches are not robust in the presence of arbitrary outliers

---

- **Spectral initialization would fail:** leading eigenvector of  $\mathbf{Y}$  can be arbitrarily perturbed

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m \mathbf{y}_i \mathbf{a}_i \mathbf{a}_i^\top \quad (\text{WF})$$

or  $\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^\top \mathbb{1}_{\{|y_i| \lesssim \text{mean}(\{\mathbf{y}_i\})\}} \quad (\text{TWF})$

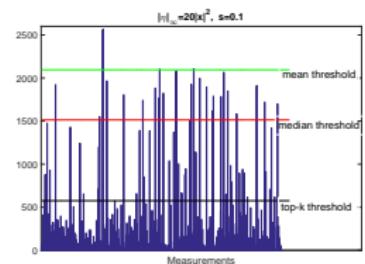
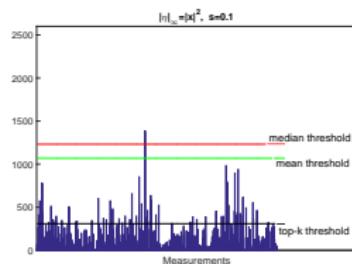
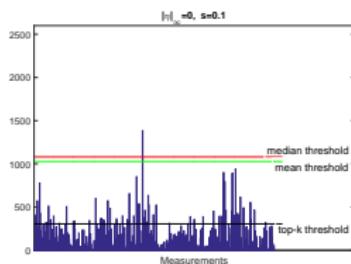
- **GD would fail:** search directions can be arbitrarily perturbed

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

# Solution: median truncation

Median is often more stable for various levels of outliers

- well-known in robust statistics to be outlier-resilient



no outliers

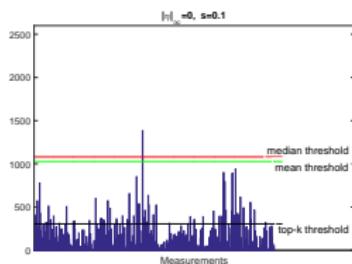
small outlier magnitudes

large outlier magnitudes

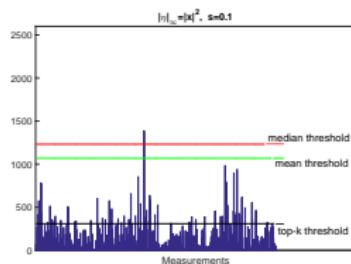
# Solution: median truncation

Median is often more stable for various levels of outliers

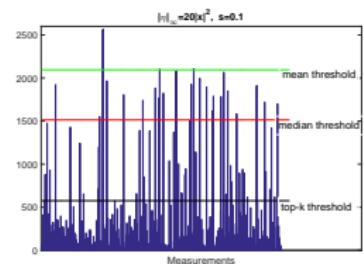
- well-known in robust statistics to be outlier-resilient



no outliers



small outlier magnitudes



large outlier magnitudes

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

# Median-truncated gradient descent

---

- (1) **Median-truncated spectral initialization:**  $\mathbf{x}^0 \leftarrow$  leading eigenvector of

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

- (2) **Median-truncated gradient descent:**

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \frac{\eta}{m} \sum_{k \in \mathcal{T}_t} \nabla f_k(\mathbf{x}^t),$$

where

$$\mathcal{T}_t = \{k : |y_k - |\mathbf{a}_k^\top \mathbf{x}^t|| \lesssim \text{median} (\{|y_k - |\mathbf{a}_k^\top \mathbf{x}^t|\|)\}$$

# Performance guarantees

## Theorem 9 (Zhang, Chi and Liang '16)

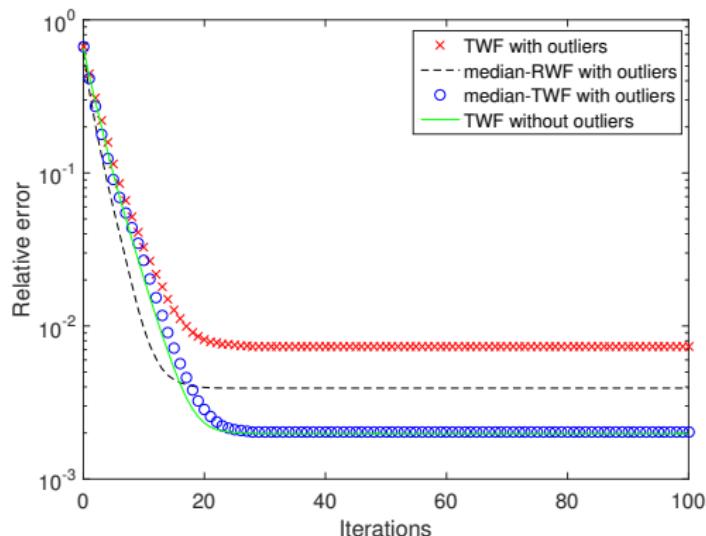
Assume  $\|\mathbf{w}\|_\infty \leq c_1 \|\mathbf{x}^*\|_2^2$ , and  $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . If  $m \gtrsim n \log n$  and  $s \lesssim s_0$ , then with high prob., median-TWF/RWF yields

$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \lesssim \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{x}^*\|_2} + (1 - \rho)^t \|\mathbf{x}^*\|_2, \quad t = 0, 1, \dots$$

for some constants  $0 < \rho, s_0 < 1$

- **Exact recovery** when  $\mathbf{w} = \mathbf{0}$  but with a constant fraction of outliers  $s \asymp 1$
- **Stable recovery** with additional bounded noise
- Resist outliers **obliviously**: no prior knowledge of outliers (except sparsity)

# Numerical experiment with both dense noise and sparse outliers



Median-TWF with outliers achieves almost identical accuracy as TWF without outliers

# **Outline**

---

- Part I: Overview
- Part II: Phase retrieval: a case study
  - Spectral initialization
  - Local refinement: algorithm and analysis
- Part III: Low-rank matrix estimation
- Part IV: Closing remarks

# Motivation

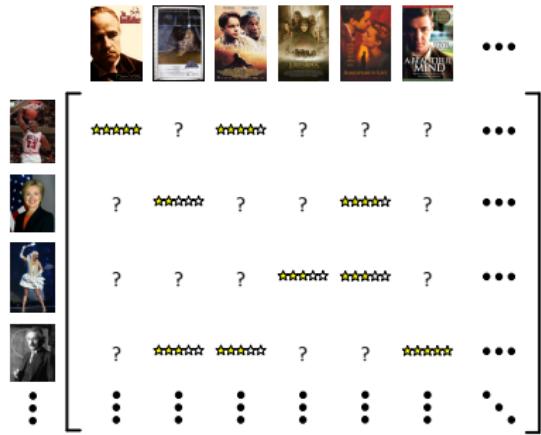
Low-rank matrix estimation problems arise in many applications

A popular example is **recommendation systems**: how to predict unseen user ratings for movies?

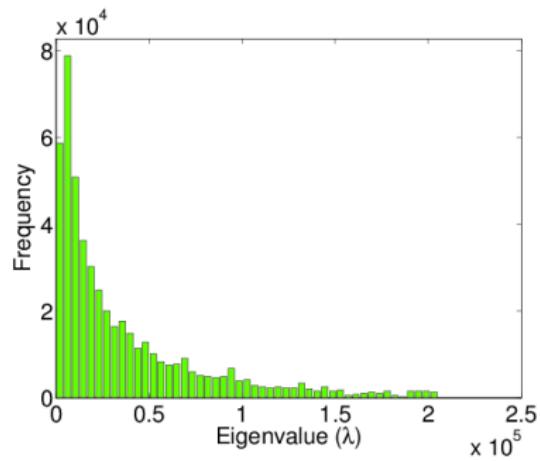


*figure credit: E. Candès*

# Low-rank modeling



*figure credit: E. Candès*



A few factors explain most of the data

# Low-rank modeling

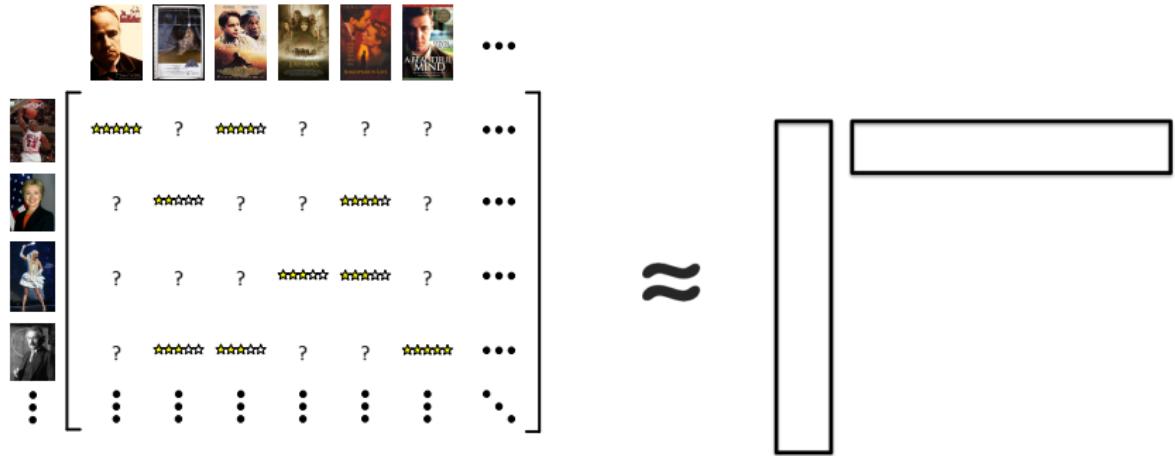


figure credit: E. Candès

A few factors explain most of the data → **low-rank** approximation

How to exploit (approx.) low-rank structure in prediction?

## Other problems with low-rank matrices

---

- sensor network localization
- structure from motion
- system identification and time series analysis
- spatial-temporal data modeling, e.g. video, network traffic, ..
- face recognition
- quantum state tomography
- community detection
- ...

# Rank-constrained optimization

---

**Rank-constrained optimization:**

$$\text{minimize}_{M \in \mathbb{R}^{n \times n}} \quad F(M) \quad \text{s.t.} \quad \text{rank}(M) \leq r,$$

where  $F(M)$  is convex in  $M$ , and  $r \ll n$

- useful model for many low-rank estimation problems;
- computationally intractable.

# Convex relaxation

**Convex relaxation:**

$$\text{minimize}_{M \in \mathbb{R}^{n \times n}} \quad F(M) \quad \text{s.t.} \quad \|M\|_* \leq \zeta$$

where  $\|\cdot\|_*$  is nuclear norm — convex relaxation of rank

- **Pros:** mature theory; versatile to incorporate other constraints
- **Cons:** run-time in  $O(n^3)$ ; even  $M$  itself takes  $O(n^2)$  storage

**Question:** can we develop algorithms that work with computational and memory complexities nearly linear in  $n$ ?

# Burer-Monteiro factorization

---

**Matrix factorization:**

$$\text{minimize}_{\mathbf{U}, \mathbf{V}} f(\mathbf{U}, \mathbf{V}) := F(\mathbf{U}\mathbf{V}^\top)$$

where  $\mathbf{M} = \mathbf{U}\mathbf{V}^\top$ , where  $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{n \times r}$ .

- pioneered by Burer, Monteiro '03
- highly non-convex
- global ambiguity: for any orthonormal  $\mathbf{R} \in \mathbb{R}^{r \times r}$  and  $\alpha \neq 0$ ,

$$\mathbf{U}\mathbf{V}^\top = (\alpha \mathbf{U}\mathbf{R})(\alpha^{-1}\mathbf{V}\mathbf{R})^\top$$

i.e. if  $(\mathbf{U}, \mathbf{V})$  is a global minimizer, so does  $(\alpha \mathbf{U}\mathbf{R}, \alpha^{-1}\mathbf{V}\mathbf{R})$

# Revisiting PCA

---

Given PSD  $M \in \mathbb{R}^{n \times n}$  (not necessarily low-rank), solve *low-rank approximation problem* (best rank- $r$  approximation):

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

Solution is truncated eigen-decomposition ([Eckart-Young theorem](#))

- let  $M = \sum_{i=1}^n \sigma_i u_i u_i^\top$  be EVD of  $M$  ( $\sigma_1 \geq \dots \geq \sigma_n$ ), then

$$\widehat{M} = \sum_{i=1}^r \sigma_i u_i u_i^\top$$

— *nonconvex, but tractable*

## Optimization viewpoint

---

Factorize  $Z = \mathbf{X}\mathbf{X}^\top$  with  $\mathbf{X} \in \mathbb{R}^{n \times r}$ . We're interested in the landscape of

$$f(\mathbf{X}) := \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

# Optimization viewpoint

---

Factorize  $Z = \mathbf{X}\mathbf{X}^\top$  with  $\mathbf{X} \in \mathbb{R}^{n \times r}$ . We're interested in the landscape of

$$f(\mathbf{X}) := \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

To simplify exposition: set  $r = 1$ .

$$f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2$$

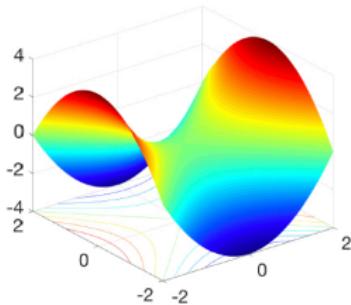
## Definition 10 (critical points)

A first-order critical point (stationary point) of  $f$  satisfies

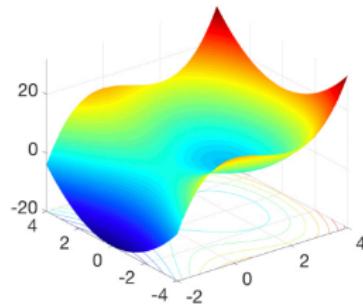
$$\nabla f(\mathbf{x}) = \mathbf{0}$$

# Several types of critical points

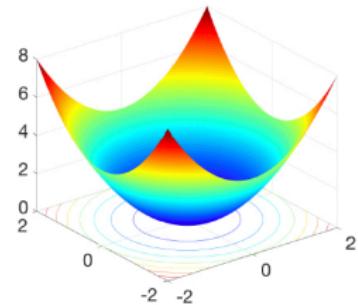
---



(a) strict saddle



(b) local minimum



(c) global minimum

*figure credit: Li et al. '16*

## Critical points of $f(\mathbf{x})$

---

$\mathbf{x}$  is critical point, i.e.  $\nabla f(\mathbf{x}) = (\mathbf{x}\mathbf{x}^\top - \mathbf{M})\mathbf{x} = \mathbf{0}$

$\Updownarrow$

$$\mathbf{M}\mathbf{x} = \|\mathbf{x}\|_2^2 \mathbf{x}$$

$\Updownarrow$

$\mathbf{x}$  aligns with eigenvectors of  $\mathbf{M}$  or  $\mathbf{x} = \mathbf{0}$

Since  $\mathbf{M}\mathbf{u}_i = \sigma_i \mathbf{u}_i$ , set of critical points is given by

$$\{\mathbf{0}\} \cup \{\sqrt{\sigma_i} \mathbf{u}_i, i = 1, \dots, n\}$$

# Categorization of critical points

---

Critical points can be further categorized based on **Hessians**:

$$\nabla^2 f(\mathbf{x}) := 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}$$

- For any non-zero critical points  $\mathbf{x}_k := \sqrt{\sigma_k} \mathbf{u}_k$ :

$$\begin{aligned}\nabla^2 f(\mathbf{x}_k) &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \mathbf{I} - \mathbf{M} \\ &= 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top + \sigma_k \left( \sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top \right) - \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{u}_i^\top \\ &= \sum_{i:i \neq k} (\sigma_k - \sigma_i) \mathbf{u}_i \mathbf{u}_i^\top + 2\sigma_k \mathbf{u}_k \mathbf{u}_k^\top\end{aligned}$$

# Categorization of critical points

---

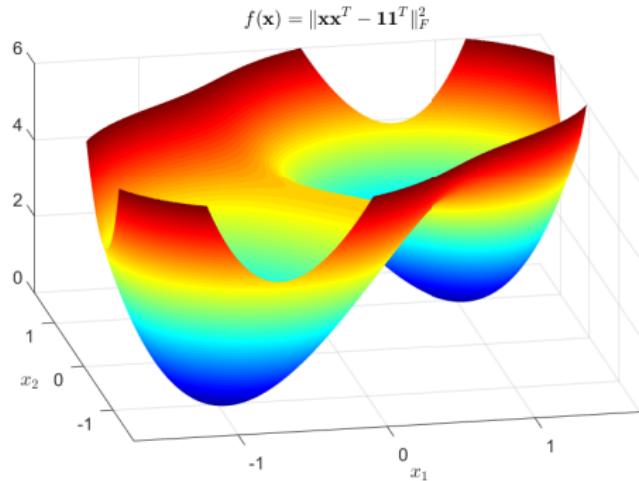
Critical points can be further categorized based on **Hessians**:

$$\nabla^2 f(\mathbf{x}) := 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}$$

- If  $\sigma_1 > \sigma_2 \geq \dots \geq \sigma_n \geq 0$ , then
  - $k = 1$ :  $\nabla^2 f(\mathbf{x}_1) \succ \mathbf{0}$  → local minima
  - $1 < k \leq n$ :  $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) < 0, \lambda_{\max}(\nabla^2 f(\mathbf{x}_k)) > 0$   
→ strict saddle
  - $\mathbf{x} = \mathbf{0}$ :  $\nabla^2 f(\mathbf{0}) \preceq \mathbf{0}$  → local maxima (or strict saddle)

## Good news: benign landscape

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



global minima  $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$  & strict saddle  $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , and  $\pm \begin{bmatrix} 1 \\ -1 \end{bmatrix}$   
— No “spurious” local minima!

## Key messages from landscape analysis

---

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

If  $\sigma_r > \sigma_{r+1}$ :

- **all local minima are global:**  $\mathbf{X}$  contains top- $r$  eigenvectors (up to orthonormal transformation)
- **strict saddle points:** all stationary points are saddle points except global optimum

# Low-rank recovery with few measurements

---

Consider linear measurements:

$$\mathbf{y} = \mathcal{A}(\mathbf{M}), \quad \mathbf{y} \in \mathbb{R}^m, \quad m \ll n^2$$

where  $\mathbf{M} \in \mathbb{R}^{n \times n}$  is rank- $r$  ( $r \ll n$ ) and PSD (for simplicity).

- Consider least-squares loss function:

$$f(\mathbf{X}) := \frac{1}{4} \|\mathcal{A}(\mathbf{X}\mathbf{X}^\top - \mathbf{M})\|_{\text{F}}^2$$

- If  $\mathcal{A}$  is isotropic (i.e.  $\mathbb{E}[\mathcal{A}^*\mathcal{A}] = \mathcal{I}$ ), then

$$\mathbb{E}[f(\mathbf{X})] = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

- Does  $f(\mathbf{X})$  inherit benign landscape?

# Landscape preserving under RIP

---

## Definition 11

Rank- $r$  restricted isometry constants  $\delta_r$  is smallest quantity obeying

$$(1 - \delta_r) \|\mathbf{M}\|_{\text{F}}^2 \leq \|\mathcal{A}(\mathbf{M})\|_{\text{F}}^2 \leq (1 + \delta_r) \|\mathbf{M}\|_{\text{F}}^2, \quad \forall \mathbf{M} : \text{rank}(\mathbf{M}) \leq r$$

# Landscape preserving under RIP

---

## Definition 11

Rank- $r$  restricted isometry constants  $\delta_r$  is smallest quantity obeying

$$(1 - \delta_r) \|\mathbf{M}\|_{\text{F}}^2 \leq \|\mathcal{A}(\mathbf{M})\|_{\text{F}}^2 \leq (1 + \delta_r) \|\mathbf{M}\|_{\text{F}}^2, \quad \forall \mathbf{M} : \text{rank}(\mathbf{M}) \leq r$$

**Key message:** benign landscape is preserved when  $\mathcal{A}$  satisfies RIP  
e.g., when  $\mathcal{A}$  follows the Gaussian design

# Landscape preserving under RIP

## Definition 11

Rank- $r$  restricted isometry constants  $\delta_r$  is smallest quantity obeying

$$(1 - \delta_r) \|\mathbf{M}\|_{\text{F}}^2 \leq \|\mathcal{A}(\mathbf{M})\|_{\text{F}}^2 \leq (1 + \delta_r) \|\mathbf{M}\|_{\text{F}}^2, \quad \forall \mathbf{M} : \text{rank}(\mathbf{M}) \leq r$$

## Theorem 12 (Bhojanapalli et al. '16, Ge et al. '17)

If  $\mathcal{A}$  satisfies RIP with  $\delta_{2r} < \frac{1}{10}$ , then

- all local min are global: any local minimum  $\mathbf{X}$  of  $f(\cdot)$  satisfies  $\mathbf{X}\mathbf{X}^\top = \mathbf{M}$
- strict saddle points: any non-local min critical point  $\mathbf{X}$  of  $f(\cdot)$  satisfies  $\lambda_{\min}[\nabla^2 f(\mathbf{X})] \leq -\frac{2}{5}\sigma_r$

# Landscape without RIP

---

**Matrix completion:**

Complete  $\mathbf{M}$  from partial entries  $M_{i,j}, (i,j) \in \Omega$

where  $(i,j)$  is included in  $\Omega$  independently with prob.  $p$

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

In matrix completion, RIP does not hold

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

# Incoherence for matrix completion

## Definition 13 (Incoherence for matrix completion)

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

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

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

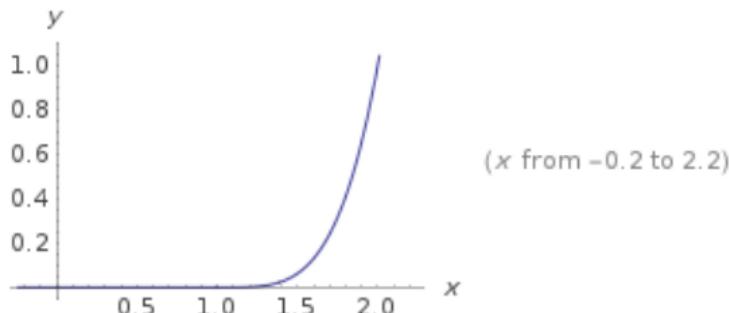
# Regularization

---

One possible regularizer:

$$Q(\mathbf{X}) = \sum_{i=1}^n (\underbrace{\|\mathbf{e}_i^\top \mathbf{X}\|_2}_{\text{row norm}} - \alpha)_+^4 := \sum_{i=1}^n Q_i(\|\mathbf{e}_i^\top \mathbf{X}\|_2)$$

where  $\alpha$  is regularization parameter, and  $z_+ = \max\{z, 0\}$



# MC has no spurious local minima under proper regularization

Consider *regularized* loss function

$$f_{\text{reg}}(\mathbf{X}) = \frac{1}{p} \|\mathcal{P}_{\Omega}(\mathbf{X}\mathbf{X}^{\top} - \mathbf{M})\|_{\text{F}}^2 + \underbrace{\lambda Q(\mathbf{X})}_{\text{promote incoherence}}$$

where  $\lambda$ : regularization parameter

## Theorem 14 (Ge et al, 2016)

If sample size  $n^2 p \gtrsim \mu^4 n r^6 \log n$  and if  $\alpha$  and  $\lambda$  are chosen properly, then with high prob.,

- all local min are global: any local minimum  $\mathbf{X}$  of  $f_{\text{reg}}(\cdot)$  satisfies  $\mathbf{X}\mathbf{X}^{\top} = \mathbf{M}$
- saddle points that are not local minima are strict saddles

# Initialization-free theory

---

## Implications:

- Under benign landscape, local search algorithms that can find local minima are often sufficient, *regardless of initialization*
- Key algorithm issue: how to escape saddle points

# Saddle-point escaping algorithms

---

- *Vanilla GD with random initialization*: converges to global minimizers almost surely, but no rates are known (Lee et al. '16)
- *Second-order algorithms (Hessian-based)*: trust-region methods, ... (Sun et al. '16)
- *First-order algorithms*: (perturbed) gradient descent, stochastic gradient descent, ... (Jin et al. '17)

**Open problem:** does MC converge fast with random initialization?

# 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  $\mathbb{P}((i, j) \in \Omega) = p$  and  $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ <sup>1</sup>

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

---

<sup>1</sup>can be relaxed to sub-Gaussian noise and 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  $\mathbb{P}((i, j) \in \Omega) = p$  and  $E_{i,j} \sim \mathcal{N}(0, \sigma^2)$ <sup>1</sup>

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

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

---

<sup>1</sup>can be relaxed to sub-Gaussian noise and 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 optimal transform from the  $t$ th iterate  $\mathbf{X}^t$  to  $\mathbf{X}^*$  as

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

## Theorem 15 (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

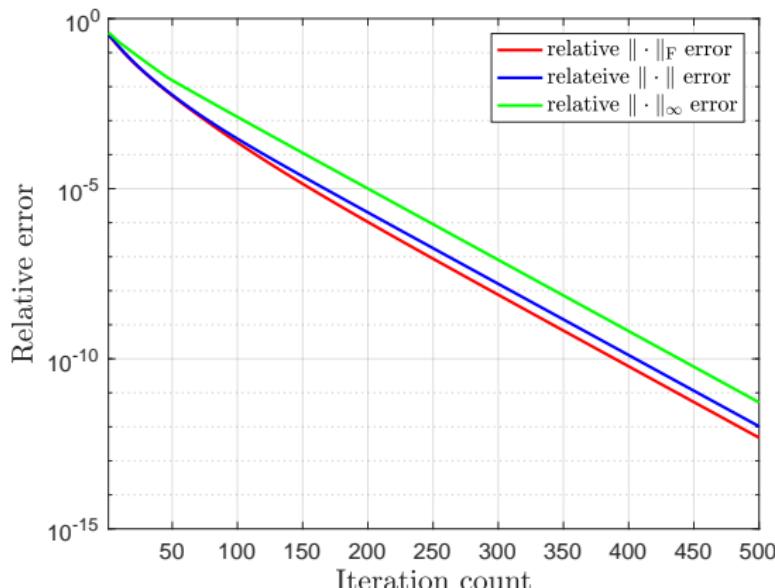
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{\text{F}} \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|_{\text{F}},$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\| \lesssim \rho^t \mu r \frac{1}{\sqrt{np}} \|\mathbf{X}^*\|, \quad (\text{spectral})$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{2,\infty}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

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

- vanilla gradient descent 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 MC, where  $n = 1000$ ,  $r = 10$ ,  $p = 0.1$ , and  $\eta_t = 0.2$

# Noisy matrix completion

## Theorem 16 (Noisy MC, Ma, Wang, Chi, Chen '17)

Under sample complexity of Theorem 15, if noise satisfies

$$\sigma \sqrt{\frac{n}{p}} \ll \frac{\sigma_{\min}}{\sqrt{\kappa^3 \mu r \log^3 n}}, \text{ then GD iterates satisfy}$$

$$\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{\text{F}} \lesssim \left( \rho^t \mu r \frac{1}{\sqrt{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|\mathbf{X}^*\|_{\text{F}},$$

$$\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{2,\infty} \lesssim \left( \rho^t \mu r \sqrt{\frac{\log n}{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n \log n}{p}} \right) \|\mathbf{X}^*\|_{2,\infty},$$

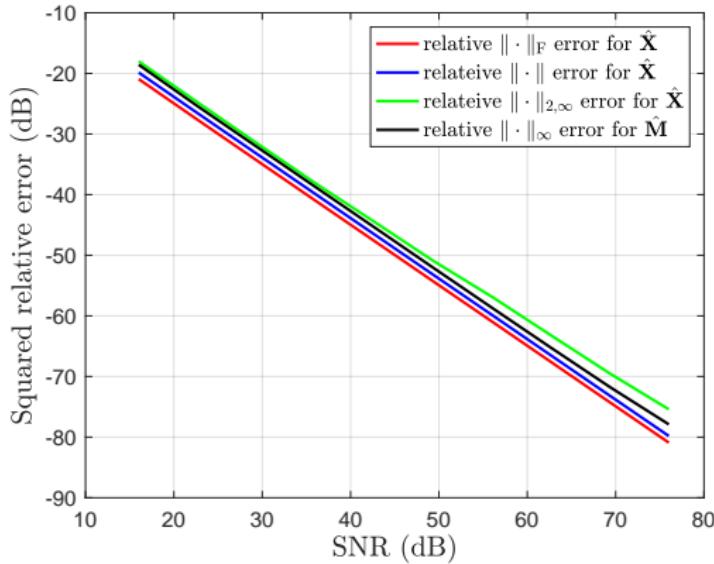
$$\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\| \lesssim \left( \rho^t \mu r \frac{1}{\sqrt{np}} + \frac{\sigma}{\sigma_{\min}} \sqrt{\frac{n}{p}} \right) \|\mathbf{X}^*\|$$

- minimax entrywise error control in  $\|\mathbf{X}^t \mathbf{X}^{t\top} - \mathbf{X}^* \mathbf{X}^{*\top}\|_\infty$

# Numerical evidence for noisy data

---

$$\text{Set SNR} := \frac{\|M\|_F^2}{n^2\sigma^2}$$



Squared relative error of the estimate  $\widehat{X}$  (measured by  $\|\cdot\|_F$ ,  $\|\cdot\|$ ,  $\|\cdot\|_{2,\infty}$ ) and  $\widehat{M} = \widehat{X}\widehat{X}^\top$  (measured by  $\|\cdot\|_\infty$ ) vs. SNR, where  $n = 500$ ,  $r = 10$ ,  $p = 0.1$ , 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

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

# Quadratic sampling

$$\begin{array}{c} \textbf{\textit{A}} \\ \left\{ \begin{array}{c} m \\ \text{---} \\ n \end{array} \right. \end{array} \quad \begin{array}{c} \textbf{\textit{X}} \\ \underbrace{\hspace{1cm}}_r \end{array} \quad = \quad \begin{array}{c} \textbf{\textit{AX}} \\ \begin{array}{|ccc|} \hline 1 & 1 & 1 \\ -3 & 0 & 1 \\ 2 & 2 & 0 \\ -1 & -1 & -1 \\ 4 & 1 & -1 \\ 2 & 2 & 2 \\ -2 & 0 & 1 \\ -1 & 0 & -1 \\ 3 & 3 & 3 \\ -1 & 4 & 1 \\ \hline \end{array} \end{array} \quad \begin{array}{c} y_i = \|\textbf{\textit{a}}_i^\top \textbf{\textit{X}}\|_2^2 \\ \longrightarrow \\ \begin{array}{|c|} \hline 3 \\ 10 \\ 8 \\ 3 \\ 18 \\ 12 \\ 5 \\ 2 \\ 27 \\ 18 \\ \hline \end{array} \end{array}$$

Recover  $\textbf{\textit{X}}^* \in \mathbb{R}^{n \times r}$  from  $m$  random quadratic measurements

$$y_i = \|\textbf{\textit{a}}_i^\top \textbf{\textit{X}}^*\|_2^2, \quad i = 1, \dots, m$$

*Applications: quantum state tomography, covariance sketching, ...*

# Gradient descent with spectral initialization

---

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \|\mathbf{a}_k^\top \mathbf{X}\|_2^2 - y_k \right)^2$$

# Gradient descent with spectral initialization

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{k=1}^m \left( \|\mathbf{a}_k^\top \mathbf{X}\|_2^2 - y_k \right)^2$$

## Theorem 17 (Quadratic sampling)

Under i.i.d. Gaussian designs  $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$ , GD (with spectral initialization) achieves

- $\max_l \|\mathbf{a}_l^\top (\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*)\|_2 \lesssim \sqrt{\log n} \frac{\sigma_r^2(\mathbf{X}^*)}{\|\mathbf{X}^*\|_{\text{F}}} \text{ (incoherence)}$
- $\|\mathbf{X}^t \mathbf{Q}^t - \mathbf{X}^*\|_{\text{F}} \lesssim \left(1 - \frac{\sigma_r^2(\mathbf{X}^*)\eta}{2}\right)^t \|\mathbf{X}^*\|_{\text{F}} \text{ (linear convergence)}$

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

# Demixing sparse and low-rank matrices

---

Suppose we are given a matrix

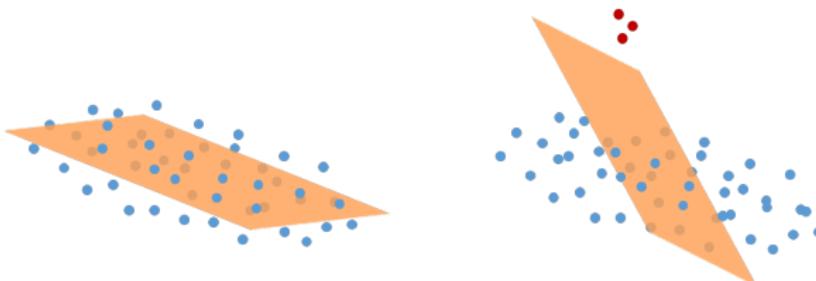
$$M = \underbrace{L}_{\text{low-rank}} + \underbrace{S}_{\text{sparse}} \in \mathbb{R}^{n \times n}$$

**Question:** can we hope to recover both  $L$  and  $S$  from  $M$ ?

# Applications

---

- Robust PCA



- Video surveillance: separation of background and foreground



## Nonconvex approach

- $\text{rank}(\mathbf{L}) \leq r$ ; if we write the SVD of  $\mathbf{L} = \mathbf{U}\Sigma\mathbf{V}^\top$ , set

$$\mathbf{X}^* = \mathbf{U}_L \Sigma^{1/2}; \quad \mathbf{Y}^* = \mathbf{V} \Sigma^{1/2}$$

- non-zero entries of  $\mathbf{S}$  are “spread out” (no more than  $s$  fraction of non-zeros per row/column), but otherwise arbitrary

$$\mathcal{S}_s = \{\mathbf{S} \in \mathbb{R}^{n \times n} : \|\mathbf{S}_{i,:}\|_0 \leq s \cdot n; \|\mathbf{S}_{:,j}\|_0 \leq s \cdot n\}$$

$$\underset{\mathbf{X}, \mathbf{Y}, \mathbf{S} \in \mathcal{S}_s}{\text{minimize}} F(\mathbf{X}, \mathbf{Y}, \mathbf{S}) := \underbrace{\|\mathbf{M} - \mathbf{X}\mathbf{Y}^\top - \mathbf{S}\|_{\text{F}}^2}_{\text{least-squares loss}} + \underbrace{\frac{1}{4} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_{\text{F}}^2}_{\text{fix scaling ambiguity}}$$

where  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$ .

# Gradient descent and hard thresholding

---

$$\text{minimize}_{\mathbf{X}, \mathbf{Y}, \mathbf{S} \in \mathcal{S}_s} \quad F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$$

- **Spectral initialization:** Set  $\mathbf{S}^0 = \mathcal{H}_{\gamma_s}(\mathbf{M})$ . Let  $\mathbf{U}^0 \boldsymbol{\Sigma}^0 \mathbf{V}^{0\top}$  be rank- $r$  SVD of  $\mathbf{M}^0 := \mathcal{P}_\Omega(\mathbf{M} - \mathbf{S}^0)$ ; set  $\mathbf{X}^0 = \mathbf{U}^0 (\boldsymbol{\Sigma}^0)^{1/2}$  and  $\mathbf{Y}^0 = \mathbf{V}^0 (\boldsymbol{\Sigma}^0)^{1/2}$

# Gradient descent and hard thresholding

---

$$\text{minimize}_{\mathbf{X}, \mathbf{Y}, \mathbf{S} \in \mathcal{S}_s} \quad F(\mathbf{X}, \mathbf{Y}, \mathbf{S})$$

- **Spectral initialization:** Set  $\mathbf{S}^0 = \mathcal{H}_{\gamma_s}(\mathbf{M})$ . Let  $\mathbf{U}^0 \boldsymbol{\Sigma}^0 \mathbf{V}^{0\top}$  be rank- $r$  SVD of  $\mathbf{M}^0 := \mathcal{P}_\Omega(\mathbf{M} - \mathbf{S}^0)$ ; set  $\mathbf{X}^0 = \mathbf{U}^0 (\boldsymbol{\Sigma}^0)^{1/2}$  and  $\mathbf{Y}^0 = \mathbf{V}^0 (\boldsymbol{\Sigma}^0)^{1/2}$
- **for**  $t = 0, 1, 2, \dots$ 
  - **Hard thresholding:**  $\mathbf{S}^{t+1} = \mathcal{H}_{\gamma_s}(\mathbf{M} - \mathbf{X}^t \mathbf{Y}^{t\top})$
  - **Gradient updates:**

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

## Efficient nonconvex recovery

### Theorem 18 (Nonconvex RPCA, Yi et al. '16)

Set  $\gamma = 2$  and  $\eta = 1/(36\sigma_{\max})$ . Suppose that

$$s \lesssim \min \left\{ \frac{1}{\mu\sqrt{\kappa r^3}}, \frac{1}{\mu\kappa^2 r} \right\}$$

Then GD+HT satisfies

$$\|\mathbf{X}^t \mathbf{Y}^{t\top} - \mathbf{L}\|_{\text{F}}^2 \lesssim \left(1 - \frac{1}{288\kappa}\right)^t \mu^2 \kappa r^3 s^2 \sigma_{\max}$$

- $O(\kappa \log \frac{1}{\varepsilon})$  iterations to reach  $\varepsilon$  accuracy
- for adversarial outliers, optimal fraction is  $s = O(1/\mu r)$ ;  
Theorem 18 is suboptimal by a factor of  $\sqrt{r}$
- extendable to partial observation models

# Outline

---

- Part I: Overview
- Part II: Phase retrieval: a case study
  - Spectral initialization
  - Local refinement: algorithm and analysis
- Part III: Low-rank matrix estimation
- Part IV: Closing remarks

## A growing list of “benign” nonconvex problems

---

- blind deconvolution / self-calibration
- dictionary learning
- tensor decomposition
- robust PCA
- mixed linear regression
- Gaussian mixture models
- etc...

## Topics we did not cover

---

- **other algorithms:** alternating minimization, stochastic gradient descent, mirror descent, singular value projection, etc...
- **additional structures:** e.g. sparsity, piece-wise smoothness
- **saddle-point escaping algorithms**

# Advertisement

---

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

— Y. Chi, Y. Lu and Y. Chen, arXiv: 1809.09573

# Reference

---

- [1] Dr. Ju Sun's webpage: "<http://sunju.org/research/nonconvex/>".
- [2] "*Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation,*" Y. Chen, and Y. Chi, *arXiv preprint arXiv:1802.08397*, IEEE Signal Processing Magazine, to appear.
- [3] "*Phase retrieval via Wirtinger flow: Theory and algorithms,*" E. Candes, X. Li, M. Soltanolkotabi, *IEEE Transactions on Information Theory*, 2015.
- [4] "*Solving random quadratic systems of equations is nearly as easy as solving linear systems,*" Y. Chen, E. Candes, *Communications on Pure and Applied Mathematics*, 2017.
- [5] "*Provable non-convex phase retrieval with outliers: Median truncated Wirtinger flow,*" H. Zhang, Y. Chi, and Y. Liang, *ICML 2016*.
- [6] "*Implicit Regularization in Nonconvex Statistical Estimation: Gradient Descent Converges Linearly for Phase Retrieval, Matrix Completion and Blind Deconvolution,*" C. Ma, K. Wang, Y. Chi and Y. Chen, *arXiv preprint arXiv:1711.10467*, 2017.

# Reference

---

- [7] "Gradient Descent with Random Initialization: Fast Global Convergence for Nonconvex Phase Retrieval," Y. Chen, Y. Chi, J. Fan, C. Ma, *arXiv preprint arXiv:1803.07726*, 2018.
- [8] "Solving systems of random quadratic equations via truncated amplitude flow," G. Wang, G. Giannakis, and Y. Eldar, *IEEE Transactions on Information Theory*, 2017.
- [9] "Matrix completion from a few entries," R. Keshavan, A. Montanari, and S. Oh, *IEEE Transactions on Information Theory*, 2010.
- [10] "Guaranteed matrix completion via non-convex factorization," R. Sun, T. Luo, *IEEE Transactions on Information Theory*, 2016.
- [11] "Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees," Y. Chen and M. Wainwright, *arXiv preprint arXiv:1509.03025*, 2015.
- [12] "Fast Algorithms for Robust PCA via Gradient Descent," X. Yi, D. Park, Y. Chen, and C. Caramanis, *NIPS*, 2016.

# Reference

---

- [13] "Matrix completion has no spurious local minimum," R. Ge, J. Lee, and T. Ma, *NIPS*, 2016.
- [14] "No Spurious Local Minima in Nonconvex Low Rank Problems: A Unified Geometric Analysis," R. Ge, C. Jin, and Y. Zheng, *ICML*, 2017.
- [15] "Symmetry, Saddle Points, and Global Optimization Landscape of Nonconvex Matrix Factorization," X. Li et al., *arXiv preprint arxiv:1612.09296*, 2016.
- [16] "Phase Transitions of Spectral Initialization for High-Dimensional Nonconvex Estimation," Y. M. Lu and G. Li, *Information and Inference*, to appear, *arXiv:1702.06435*, 2018.
- [17] "Kaczmarz Method for Solving Quadratic Equations," Y. Chi and Y. M. Lu, *IEEE Signal Processing Letters*, vol. 23, no. 9, pp. 1183-1187, 2016.
- [18] "Scaling Limit: Exact and Tractable Analysis of Online Learning Algorithms with Applications to Regularized Regression and PCA," C. Wang, J. Mattingly and Y. M. Lu, *arXiv preprint arXiv:1712.04332*, 2017.
- [19] "A Geometric Analysis of Phase Retrieval," S. Ju, Q. Qu, and J. Wright, to appear, *Foundations of Computational Mathematics*, 2016.

# Reference

---

- [20] "Gradient descent converges to minimizers," J. Lee, M. Simchowitz, M. Jordan, B. Recht, *Conference on Learning Theory*, 2016.
- [21] "Fundamental limits of weak recovery with applications to phase retrieval," M. Mondelli, and A. Montanari, arXiv:1708.05932, 2017.
- [22] "Phase retrieval using alternating minimization," P. Netrapalli, P. Jain, and S. Sanghavi, *NIPS*, 2013.
- [23] "Optimization-based AMP for Phase Retrieval: The Impact of Initialization and  $\ell_2$ -regularization," J. Ma, J. Xu, and A. Maleki, arXiv:1801.01170, 2018.
- [24] "How to escape saddle points efficiently," C. Jin, R. Ge, P. Netrapalli, S. Kakade, M. Jordan, arXiv:1703.00887, 2017.
- [25] "Complete dictionary recovery over the sphere," J. Sun, Q. Qu, J. Wright, *IEEE Transactions on Information Theory*, 2017.
- [26] "A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization," S. Burer, and R. Monteiro, *Mathematical Programming*, 2003.

# Reference

---

- [27] “*Memory-efficient Kernel PCA via Partial Matrix Sampling and Nonconvex Optimization: a Model-free Analysis of Local Minima,*” J. Chen, X. Li, arXiv:1711.01742, 2017.
- [28] “*Rapid, robust, and reliable blind deconvolution via nonconvex optimization,*” X. Li, S. Ling, T. Strohmer, K. Wei, arXiv:1606.04933, 2016.
- [29] “*Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming,*” E. Candes, T. Strohmer, V. Voroninski, *Communications on Pure and Applied Mathematics*, 2012.
- [30] “*Exact matrix completion via convex optimization,*” E. Candes, B. Recht, *Foundations of Computational mathematics*, 2009.
- [31] “*Low-rank solutions of linear matrix equations via procrustes flow,*” S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, B. Recht, arXiv:1507.03566, 2015.
- [32] “*Global optimality of local search for low rank matrix recovery,*” S. Bhojanapalli, B. Neyshabur, and N. Srebro, NIPS, 2016.

# Reference

---

- [33] "Phase retrieval via matrix completion," E. Candes, Y. Eldar, T. Strohmer, and V. Voroninski, *SIAM Journal on Imaging Sciences*, 2013.
- [34] "Optimal rates of convergence for noisy sparse phase retrieval via thresholded Wirtinger flow," T. Cai, X. Li, Z. Ma, *The Annals of Statistics*, 2016.
- [35] "The landscape of empirical risk for non-convex losses," S. Mei, Y. Bai, and A. Montanari, arXiv:1607.06534, 2016.
- [36] "Non-convex robust PCA," P. Netrapalli, U. Niranjan, S. Sanghavi, A. Anandkumar, and P. Jain, NIPS, 2014.
- [37] "Non-square matrix sensing without spurious local minima via the Burer-Monteiro approach," D. Park, A. Kyriolidis, C. Caramanis, and S. Sanghavi, arXiv:1609.03240, 2016.
- [38] "Solving almost all systems of random quadratic equations" G. Wang, G. Giannakis, Y. Saad, and J. Chen, arXiv:1705.10407, 2017.
- [39] "A Nonconvex Approach for Phase Retrieval: Reshaped Wirtinger Flow and Incremental Algorithms," H. Zhang, Y. Zhou, Y. Liang and Y. Chi, *Journal of Machine Learning Research*, 2017.

# Reference

---

- [40] “*Nonconvex Matrix Factorization from Rank-One Measurements,*” Y. Li, C. Ma, Y. Chen and Y. Chi, arXiv:1802.06286, 2018.
- [41] “*Structured signal recovery from quadratic measurements: Breaking sample complexity barriers via nonconvex optimization,*”, M. Soltanolkotabi, arXiv:1702.06175, 2017.

**Thanks!**