# The Power of Preconditioning in Overparameterized Low-Rank Matrix Sensing

Cong Ma

Department of Statistics, UChicago
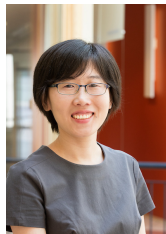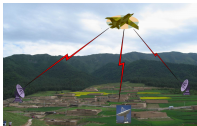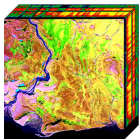
Xingyu Xu
CMU

Yandi Shen
UChicago → CMU

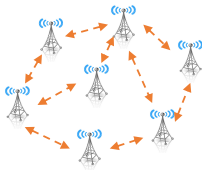Yuejie Chi
CMU

# Low-rank matrices in data science



radar imaging



hyperspectral imaging



recommendation systems



localization



community detection



bioinformatics

# Low-rank matrix recovery



$M \in \mathbb{R}^{n_1 \times n_2}$
$\text{rank}(M) = r$

$\mathcal{A}(\cdot)$
linear map

$y \in \mathbb{R}^m$

$$y = \mathcal{A}(M)$$

**Goal:** recover $M$ in the sample-starved regime

$$\underbrace{(n_1 + n_2)r}_{\text{degrees of freedom}} \quad \lesssim \quad \underbrace{m}_{\text{no. of measurements}} \quad \ll \quad \underbrace{n_1 n_2}_{\text{ambient dimension}}$$

# Convex relaxation via nuclear norm minimization

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \quad \text{rank}(\boldsymbol{Z}) \qquad \text{s.t.} \qquad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

# Convex relaxation via nuclear norm minimization

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\boldsymbol{Z}) \qquad \text{s.t.} \qquad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

⇩ cvx surrogate

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \|\boldsymbol{Z}\|_*$$
$$\text{s.t.} \quad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

where $\| \cdot \|_*$ is the nuclear norm

# Convex relaxation via nuclear norm minimization

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \quad \text{rank}(\boldsymbol{Z}) \qquad \text{s.t.} \qquad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

⇓ cvx surrogate

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \quad \|\boldsymbol{Z}\|_*$$
$$\text{s.t.} \quad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

where $\|\cdot\|_*$ is the nuclear norm

## Significant developments in the last decade:

Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09, Candès, Tao '10, Cai et al. '10, Gross '10, Negahban,

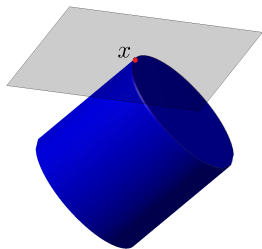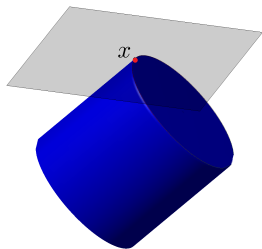Wainwright '11, Sanghavi et al. '13, Chen, Chi '14, ...

# Convex relaxation via nuclear norm minimization

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\boldsymbol{Z}) \qquad \text{s.t.} \qquad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

⇩ cvx surrogate

$$\min_{\boldsymbol{Z} \in \mathbb{R}^{n_1 \times n_2}} \|\boldsymbol{Z}\|_*$$
$$\text{s.t.} \quad \boldsymbol{y} \approx \mathcal{A}(\boldsymbol{Z})$$

where $\|\cdot\|_*$ is the nuclear norm



$x$

## Significant developments in the last decade:

Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09, Candès, Tao '10, Cai et al. '10, Gross '10, Negahban,

Wainwright '11, Sanghavi et al. '13, Chen, Chi '14, ...
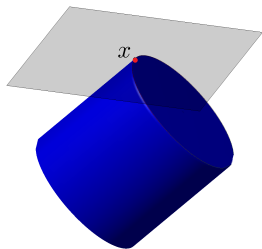
**Poor scalability:** operate in the *ambient* matrix space

# Low-rank matrix factorization

$$\min_{\text{rank}(\boldsymbol{Z})=r} \quad \frac{1}{2} \left\| \boldsymbol{y} - \mathcal{A}(\boldsymbol{Z}) \right\|_2^2$$

# Low-rank matrix factorization

$$\min_{\mathrm{rank}(\boldsymbol{Z})=r} \quad \frac{1}{2}\left\|\boldsymbol{y} - \mathcal{A}(\boldsymbol{Z})\right\|_2^2$$



$$\min_{\boldsymbol{X}\in\mathbb{R}^{n_1\times r}, \boldsymbol{Y}\in\mathbb{R}^{n_2\times r}} \quad f(\boldsymbol{X}, \boldsymbol{Y}) = \frac{1}{2}\left\|\boldsymbol{y} - \mathcal{A}(\boldsymbol{X}\boldsymbol{Y}^\top)\right\|_2^2$$

# Low-rank matrix factorization

$$\min_{\mathrm{rank}(\boldsymbol{Z})=r} \quad \frac{1}{2}\left\|\boldsymbol{y}-\mathcal{A}(\boldsymbol{Z})\right\|_2^2$$



$$\min_{\boldsymbol{X}\in\mathbb{R}^{n_1\times r},\boldsymbol{Y}\in\mathbb{R}^{n_2\times r}} \quad f(\boldsymbol{X},\boldsymbol{Y})=\frac{1}{2}\left\|\boldsymbol{y}-\mathcal{A}(\boldsymbol{X}\boldsymbol{Y}^\top)\right\|_2^2$$

# Nonconvex problems are hard (in theory)

# Statistics meets optimization



Statistical model

worst case                    average case

# Statistics meets optimization



worst case    →  **Statistical model**    average case

Simple algorithms can be efficient for nonconvex problems

# Matrix sensing: GD with balancing regularization

$$\min_{\boldsymbol{X}, \boldsymbol{Y}} \quad f_{\mathrm{reg}}(\boldsymbol{X}, \boldsymbol{Y}) = \frac{1}{2} \left\| \boldsymbol{y} - \mathcal{A}(\boldsymbol{X}\boldsymbol{Y}^{\top}) \right\|_2^2 + \frac{1}{8} \left\| \boldsymbol{X}^{\top}\boldsymbol{X} - \boldsymbol{Y}^{\top}\boldsymbol{Y} \right\|_{\mathrm{F}}^2$$

"Basin of attraction"

- **Spectral initialization:** find an initial point in the "basin of attraction"

$$(\boldsymbol{X}_0, \boldsymbol{Y}_0) \leftarrow \mathsf{SVD}_r(\mathcal{A}^*(\boldsymbol{y}))$$

- **Gradient iterations:** for $t = 0, 1, \ldots,$

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \, \nabla_{\boldsymbol{X}} f_{\mathrm{reg}}(\boldsymbol{X}_t, \boldsymbol{Y}_t)$$
$$\boldsymbol{Y}_{t+1} = \boldsymbol{Y}_t - \eta \, \nabla_{\boldsymbol{Y}} f_{\mathrm{reg}}(\boldsymbol{X}_t, \boldsymbol{Y}_t)$$

# Prior theory for vanilla GD

Condition number $\kappa = \frac{\sigma_{\max}(\boldsymbol{M})}{\sigma_{\min}(\boldsymbol{M})}$

**Theorem 1 (Tu et al., ICML 2016)**

*For low-rank matrix sensing with i.i.d. Gaussian design, vanilla GD (with spectral initialization) achieves*

$$\|\boldsymbol{X}_t \boldsymbol{Y}_t^\top - \boldsymbol{M}\|_{\mathrm{F}} \leq \varepsilon \cdot \sigma_{\min}(\boldsymbol{M})$$

- **Computational:** *within* $O(\kappa \log \frac{1}{\varepsilon})$ *iterations;*
- **Statistical:** *as long as the sample size satisfies*

$$m \gtrsim (n_1 + n_2) r^2 \kappa^2$$

Similar results hold for many other low-rank problems

# GD slows down for ill-conditioned matrices

Condition number $\kappa = \frac{\sigma_{\max}(\boldsymbol{M})}{\sigma_{\min}(\boldsymbol{M})}$



Vanilla GD converges in $O\left(\kappa \log \frac{1}{\varepsilon}\right)$ iterations

# Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots

power-law spectrum

# Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots

rank-5 approximation

# Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots



rank-10 approximation

# Condition number can be large



$96\%$

$\kappa \approx 60$

chlorine concentration levels
120 junctions, 180 time slots

rank-10 approximation

*Can we accelerate the convergence rate of GD to $O(\log \frac{1}{\varepsilon})$?*

Data source: www.epa.gov/water-research/epanet

# A recipe: scaled gradient descent (ScaledGD)

—joint work with Tian Tong, and Yuejie Chi

$f(\boldsymbol{X}, \boldsymbol{Y}) = \|\boldsymbol{y} - \mathcal{A}(\boldsymbol{X}\boldsymbol{Y}^{\top})\|_2^2$



- **Spectral initialization:** find an initial point in the "basin of attraction"

- **Scaled gradient iterations:** for $t = 0, 1, \ldots,$

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \nabla_{\boldsymbol{X}} f(\boldsymbol{X}_t, \boldsymbol{Y}_t) \underbrace{(\boldsymbol{Y}_t^{\top} \boldsymbol{Y}_t)^{-1}}_{\text{preconditioner}}$$

$$\boldsymbol{Y}_{t+1} = \boldsymbol{Y}_t - \eta \nabla_{\boldsymbol{Y}} f(\boldsymbol{X}_t, \boldsymbol{Y}_t) \underbrace{(\boldsymbol{X}_t^{\top} \boldsymbol{X}_t)^{-1}}_{\text{preconditioner}}$$

# A recipe: scaled gradient descent (ScaledGD)

—joint work with Tian Tong, and Yuejie Chi

$f(\boldsymbol{X}, \boldsymbol{Y}) = \|\boldsymbol{y} - \mathcal{A}(\boldsymbol{X}\boldsymbol{Y}^\top)\|_2^2$

- **Spectral initialization:** find an initial point in the "basin of attraction"

- **Scaled gradient iterations:** for $t = 0, 1, \ldots,$

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \nabla_{\boldsymbol{X}} f(\boldsymbol{X}_t, \boldsymbol{Y}_t) \underbrace{(\boldsymbol{Y}_t^\top \boldsymbol{Y}_t)^{-1}}_{\text{preconditioner}}$$

$$\boldsymbol{Y}_{t+1} = \boldsymbol{Y}_t - \eta \nabla_{\boldsymbol{Y}} f(\boldsymbol{X}_t, \boldsymbol{Y}_t) \underbrace{(\boldsymbol{X}_t^\top \boldsymbol{X}_t)^{-1}}_{\text{preconditioner}}$$

ScaledGD is a *preconditioned* gradient method
without balancing regularization

# ScaledGD for low-rank matrix completion



**Huge computational saving:** ScaledGD converges in a $\kappa$-independent manner with minimal overhead

## A closer look at ScaledGD

**Connection to quasi-Newton method :**

Define $\boldsymbol{F}_t = [\boldsymbol{X}_t^\top, \boldsymbol{Y}_t^\top]^\top \in \mathbb{R}^{(n_1+n_2)\times r}$. One can write update rule as

$$\text{vec}(\boldsymbol{F}_{t+1})$$
$$= \text{vec}(\boldsymbol{F}_t) - \eta \underbrace{\begin{bmatrix} (\boldsymbol{Y}_t^\top \boldsymbol{Y}_t) \otimes \boldsymbol{I}_{n_1} & \boldsymbol{0} \\ \boldsymbol{0} & (\boldsymbol{X}_t^\top \boldsymbol{X}_t) \otimes \boldsymbol{I}_{n_2} \end{bmatrix}^{-1}}_{=:\boldsymbol{H}_t^{-1}} \text{vec}(\nabla_{\boldsymbol{F}} \mathcal{L}(\boldsymbol{F}_t))$$

# A closer look at ScaledGD

**Invariance to invertible transforms:**



$(\boldsymbol{X}_t, \boldsymbol{Y}_t)$

$(\boldsymbol{X}_t \boldsymbol{Q}, \boldsymbol{Y}_t \boldsymbol{Q}^{-\top})$

$\boldsymbol{M}_t = \boldsymbol{X}_t \boldsymbol{Y}_t^{\top}$

$\boldsymbol{M}_{t+1} = \boldsymbol{X}_{t+1} \boldsymbol{Y}_{t+1}^{\top}$

$(\boldsymbol{X}_{t+1}, \boldsymbol{Y}_{t+1})$

$(\boldsymbol{X}_{t+1} \boldsymbol{Q}, \boldsymbol{Y}_{t+1} \boldsymbol{Q}^{-\top})$

— not true for GD

# Theoretical guarantees of ScaledGD

**Theorem 2 (Tong, Ma and Chi, JMLR 2021)**

*For low-rank matrix sensing with i.i.d. Gaussian design, ScaledGD with spectral initialization achieves*

$$\|\boldsymbol{X}_t \boldsymbol{Y}_t^\top - \boldsymbol{M}\|_{\mathrm{F}} \lesssim \varepsilon \cdot \sigma_{\min}(\boldsymbol{M})$$

- **Computational:** *within $O\left(\log \frac{1}{\varepsilon}\right)$ iterations*
- **Statistical:** *the sample complexity satisfies*

$$m \gtrsim (n_1 + n_2) r^2 \kappa^2$$

# Theoretical guarantees of ScaledGD

**Theorem 2 (Tong, Ma and Chi, JMLR 2021)**

*For low-rank matrix sensing with i.i.d. Gaussian design, ScaledGD with spectral initialization achieves*

$$\|\boldsymbol{X}_t \boldsymbol{Y}_t^\top - \boldsymbol{M}\|_{\mathrm{F}} \lesssim \varepsilon \cdot \sigma_{\min}(\boldsymbol{M})$$

- **Computational:** within $O\left(\log \frac{1}{\varepsilon}\right)$ iterations
- **Statistical:** the sample complexity satisfies

$$m \gtrsim (n_1 + n_2) r^2 \kappa^2$$

**Strict improvement over Tu et al.:** ScaledGD provably accelerates vanilla GD with the same sample complexity

# Key ingredient in analysis

**Scaled distance metric:**

$$\text{dist}^2 \left( \begin{bmatrix} \boldsymbol{X}_t \\ \boldsymbol{Y}_t \end{bmatrix}, \begin{bmatrix} \boldsymbol{X} \\ \boldsymbol{Y} \end{bmatrix} \right) = \inf_{\boldsymbol{Q} \in \text{GL}(r)} \left\| (\boldsymbol{X}_t \boldsymbol{Q} - \boldsymbol{X}) \boldsymbol{\Sigma}^{1/2} \right\|_{\text{F}}^2$$
$$+ \left\| (\boldsymbol{Y}_t \boldsymbol{Q}^{-\top} - \boldsymbol{Y}) \boldsymbol{\Sigma}^{1/2} \right\|_{\text{F}}^2$$

where $\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}$ is SVD of $\boldsymbol{M}$, $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}^{1/2}$, $\boldsymbol{U} = \boldsymbol{V}\boldsymbol{\Sigma}^{1/2}$

# Key ingredient in analysis

**Scaled distance metric:**

$$\mathsf{dist}^2\left(\begin{bmatrix} \boldsymbol{X}_t \\ \boldsymbol{Y}_t \end{bmatrix}, \begin{bmatrix} \boldsymbol{X} \\ \boldsymbol{Y} \end{bmatrix}\right) = \inf_{\boldsymbol{Q}\in\mathrm{GL}(r)} \left\|(\boldsymbol{X}_t\boldsymbol{Q} - \boldsymbol{X})\boldsymbol{\Sigma}^{1/2}\right\|_{\mathrm{F}}^2$$
$$+ \left\|(\boldsymbol{Y}_t\boldsymbol{Q}^{-\top} - \boldsymbol{Y})\boldsymbol{\Sigma}^{1/2}\right\|_{\mathrm{F}}^2$$

where $\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}$ is SVD of $\boldsymbol{M}$, $\boldsymbol{X} = \boldsymbol{U}\boldsymbol{\Sigma}^{1/2}$, $\boldsymbol{U} = \boldsymbol{V}\boldsymbol{\Sigma}^{1/2}$

- Account for ambiguity arising from invertible transforms
- Fidelity to reconstruction loss: locally, we have

$$\mathsf{dist}^2\left(\begin{bmatrix} \boldsymbol{X}_t \\ \boldsymbol{Y}_t \end{bmatrix}, \begin{bmatrix} \boldsymbol{X} \\ \boldsymbol{Y} \end{bmatrix}\right) \asymp \|\boldsymbol{X}_t\boldsymbol{Y}_t^{\top} - \boldsymbol{M}\|_{\mathrm{F}}^2$$

# ScaledGD works more broadly



| Algorithms | Robust PCA | | Matrix completion | |
|---|---|---|---|---|
| | corruption fraction | iteration complexity | sample complexity | iteration complexity |
| GD | $\frac{1}{\mu r^{3/2}\kappa^{3/2}\vee\mu r\kappa^2}$ | $\kappa\log\frac{1}{\varepsilon}$ | $(\mu\vee\log n)\mu n r^2\kappa^2$ | $\kappa\log\frac{1}{\varepsilon}$ |
| ScaledGD | $\frac{1}{\mu r^{3/2}\kappa}$ | $\log\frac{1}{\varepsilon}$ | $(\mu\kappa^2\vee\log n)\mu n r^2\kappa^2$ | $\log\frac{1}{\varepsilon}$ |

Huge computational saving at comparable sample complexities

# What if we do not know the exact rank?

So far we have assumed the exact rank is given.... what if we do not know the exact rank?

# What if we do not know the exact rank?

So far we have assumed the exact rank is given.... what if we do not know the exact rank?

**Misspecification by overparameterization:**

$$\boldsymbol{M} = \boldsymbol{X}\boldsymbol{X}^\top, \qquad \boldsymbol{X} \in \mathbb{R}^{n \times \tilde{r}}, \qquad \tilde{r} > r$$

# What if we do not know the exact rank?

So far we have assumed the exact rank is given.... what if we do not know the exact rank?

**Misspecification by overparameterization:**

$$\boldsymbol{M} = \boldsymbol{X}\boldsymbol{X}^\top, \qquad \boldsymbol{X} \in \mathbb{R}^{n \times \tilde{r}}, \qquad \tilde{r} > r$$

**ScaledGD:**

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \, \nabla_{\boldsymbol{X}} f(\boldsymbol{X}_t) \underbrace{(\boldsymbol{X}_t^\top \boldsymbol{X}_t)^{-1}}_{\text{preconditioner}}$$

*analysis break down and might be unstable...*

# What if we do not know the exact rank?

So far we have assumed the exact rank is given.... what if we do not know the exact rank?

**Misspecification by overparameterization:**

$$\boldsymbol{M} = \boldsymbol{X}\boldsymbol{X}^\top, \qquad \boldsymbol{X} \in \mathbb{R}^{n \times \tilde{r}}, \qquad \tilde{r} > r$$

**ScaledGD($\lambda$):**

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta\,\nabla_{\boldsymbol{X}} f(\boldsymbol{X}_t)\underbrace{(\boldsymbol{X}_t^\top \boldsymbol{X}_t + \lambda \boldsymbol{I})^{-1}}_{\texttt{preconditioner}}$$

*add regularization to stablize the preconditioner*

# Does preconditioning hurt generalization?

- Infinitely many global minima, not all generalize
- Can we still guarantee generalization?



optimization

generalization

## WHEN DOES PRECONDITIONING HELP OR HURT GENERALIZATION?

*Shun-ichi Amari[1], Jimmy Ba[2,3], Roger Grosse[2,3], Xuechen Li[4], Atsushi Nitanda[5,6], Taiji Suzuki[5,6], Denny Wu[2,3], Ji Xu[7]

[1]RIKEN CBS, [2]University of Toronto, [3]Vector Institute, [4]Google Research, Brain Team, [5]University of Tokyo, [6]RIKEN AIP, [7]Columbia University
amari@brain.riken.jp, {jba,rgrosse,lxuechen,dennywu}@cs.toronto.edu, {nitanda,taiji}@mist.i.u-tokyo.ac.jp, jixu@cs.columbia.edu

# Theoretical guarantees

## Theorem 3 (Xu, Shen, Chi, Ma, ICML 2023)

*For low-rank matrix sensing with i.i.d. Gaussian design, overparameterized ScaledGD($\lambda$) with $\lambda \asymp \sigma_{\min}(\boldsymbol{M})$, $\eta \asymp 1$, and a sufficiently small random initialization achieves*

$$\|\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}\|_{\mathrm{F}} \lesssim \varepsilon \cdot \sigma_{\min}(\boldsymbol{M})$$

- **Computational:** *within* $O\big(\log \kappa \log(\kappa n) + \log \frac{1}{\varepsilon}\big)$ *iterations;*
- **Statistical:** *the sample complexity satisfies*

$$m \gtrsim n r^2 \mathsf{poly}(\kappa)$$

- Our analysis also enables exact convergence under random initialization with correct rank specification

# Comparison with overparameterized GD

# Comparison with overparameterized GD

# Comparison with overparameterized GD



*ScaledGD picks up the signal component much faster than GD even from small random initialization*

# Comparisons with prior art

Comparison with Zhang, Fattahi, and Zhang '21

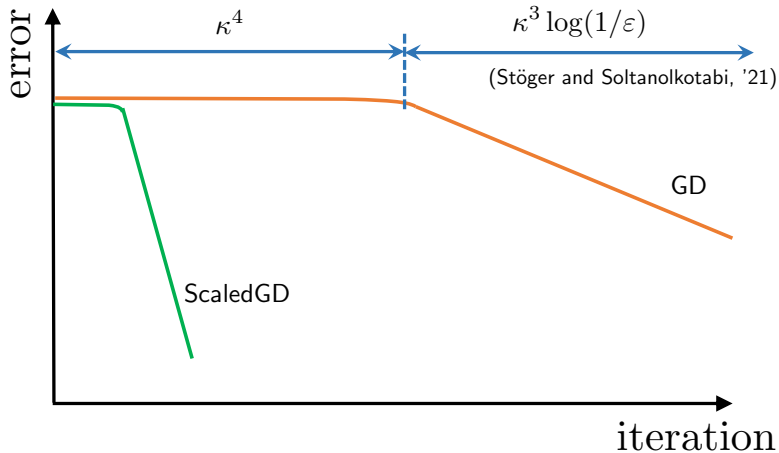$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \, \nabla_{\boldsymbol{X}} f(\boldsymbol{X}_t) \underbrace{(\boldsymbol{X}_t^\top \boldsymbol{X}_t + \lambda_t \boldsymbol{I})^{-1}}_{\texttt{preconditioner}}$$

where $\lambda_t = \|\mathcal{A}(\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M})\|$

- Local analysis: require spectral initialization
- Large sample complexity: sample complexity is $n\tilde{r}^2 \operatorname{poly}(\kappa)$, depending on the overparameterized rank $\tilde{r}$ instead of the true rank $r$

# Robustness to noise

Consider the noisy setting

$$y_i = \langle A_i, \boldsymbol{M} \rangle + \xi_i, \quad \text{where} \quad \xi_i \sim \mathcal{N}(0, \sigma^2)$$

---

**Theorem 4 (Xu, Shen, Chi, Ma, '23)**

*For low-rank matrix sensing with i.i.d. Gaussian design, overparameterized ScaledGD($\lambda$) with the same configuration as before achieves*
$$\|\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}\|_{\mathrm{F}} \lesssim \kappa^2 \sigma \sqrt{nr}$$

# ScaledGD($\lambda$) is nearly optimal

ScaledGD($\lambda$) achieves

$$\|\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}\|_F \lesssim \kappa^2 \sigma \sqrt{nr}$$

- ScaledGD($\lambda$) is minimax optimal (up to $\kappa^2$) for recovering rank-$r$ matrices, cf. Candès and Plan '09
- Both the rate and sample size requirement improve over prior art (e.g., Zhuo et al., '21, Zhang et al., '23) as ours depend on true rank $r$

A little analysis

# Phase I: approximating power method

Recall update rule of ScaledGD$(\lambda)$

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \mathcal{A}^* \mathcal{A}(\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}) \boldsymbol{X}_t (\boldsymbol{X}_t^\top \boldsymbol{X}_t + \lambda \boldsymbol{I})^{-1}$$

Since initialization is small, i.e., $\boldsymbol{X}_t \approx \boldsymbol{0}$, we have

$$\boldsymbol{X}_{t+1} \approx \boldsymbol{X}_t + \eta \mathcal{A}^* \mathcal{A}(\boldsymbol{M}) \boldsymbol{X}_t \lambda^{-1}$$
$$= \underbrace{\left( \boldsymbol{I} + \frac{\eta}{\lambda} \mathcal{A}^* \mathcal{A}(\boldsymbol{M}) \right) \boldsymbol{X}_t}_{\text{power method iterates}}$$

# Phase I: approximating power method

Indeed, we show that

$$\boldsymbol{X}_t \approx \left( I + \frac{\eta}{\lambda} \mathcal{A}^* \mathcal{A}(\boldsymbol{M}) \right)^t \boldsymbol{X}_0, \qquad \text{when } t \lesssim \frac{1}{\eta}$$

Consequently, ScaledGD($\lambda$) has three nice properties after phase I

- subspace misalignment is small
- signal strength is mildly large
- overparameterization error remains small

# Phase II: exponential amplification of the signal

In phase II, equipped with the three properties, signal is exponentially amplified in the sense that

$$\sigma_{\min}(\mathbf{\Sigma}^{-1/2}\mathbf{U}^\top \mathbf{X}_t) \quad \text{grows at rate } 1 + \eta$$

until a constant level

$$\mathbf{U}^\top \mathbf{X}_t \mathbf{X}_t^\top \mathbf{U} \succeq 0.1\mathbf{\Sigma}$$

> Scaled signal strength $\sigma_{\min}(\mathbf{\Sigma}^{-1/2}\mathbf{U}^\top \mathbf{X}_t)$ is the key

# Phase II: comparison with GD

Note that signal is amplified in a scale-independent way

$$\sigma_{\min}(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{U}^{\top}\boldsymbol{X}_t) \quad \text{grows with rate } 1 + \eta$$
$$\implies \sigma_i^2(\boldsymbol{U}^{\top}\boldsymbol{X}_t)/\sigma_i(\boldsymbol{M}) \quad \text{grows uniformly with rate } 1 + \eta$$

In contrast, for GD the growth of different singular values are different:
$$\sigma_i^2(\boldsymbol{U}^{\top}\boldsymbol{X}_t^{\mathsf{GD}}) \text{ grows with rate } 1 + \eta\sigma_i(\boldsymbol{M}),$$

## Phase II: comparison with GD

Note that signal is amplified in a scale-independent way

$$\sigma_{\min}(\boldsymbol{\Sigma}^{-1/2}\boldsymbol{U}^{\top}\boldsymbol{X}_t) \quad \text{grows with rate } 1+\eta$$
$$\implies \sigma_i^2(\boldsymbol{U}^{\top}\boldsymbol{X}_t)/\sigma_i(\boldsymbol{M}) \quad \text{grows uniformly with rate } 1+\eta$$

In contrast, for GD the growth of different singular values are different:

$$\sigma_i^2(\boldsymbol{U}^{\top}\boldsymbol{X}_t^{\mathsf{GD}}) \text{ grows with rate } 1+\eta\sigma_i(\boldsymbol{M}),$$

**Issue:** GD requires $\eta\sigma_{\max}(\boldsymbol{M}) \lesssim 1$ to stay in control, but then the growth rate for $\sigma_r^2(\boldsymbol{U}^{\top}\boldsymbol{X}_t^{\mathsf{GD}})$ would only be $1+O(\kappa^{-1})$

# Phase III: local convergence

Recall update rule of ScaledGD($\lambda$)

$$\boldsymbol{X}_{t+1} = \boldsymbol{X}_t - \eta \mathcal{A}^* \mathcal{A}(\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}) \boldsymbol{X}_t (\boldsymbol{X}_t^\top \boldsymbol{X}_t + \lambda \boldsymbol{I})^{-1}$$

When signal is at constant level, $\boldsymbol{X}_t^\top \boldsymbol{X}_t$ dominates $\lambda \boldsymbol{I}$, which yields

$$\boldsymbol{X}_{t+1} \approx \boldsymbol{X}_t - \eta \mathcal{A}^* \mathcal{A}(\boldsymbol{X}_t \boldsymbol{X}_t^\top - \boldsymbol{M}) \boldsymbol{X}_t (\boldsymbol{X}_t^\top \boldsymbol{X}_t)^{-1}$$

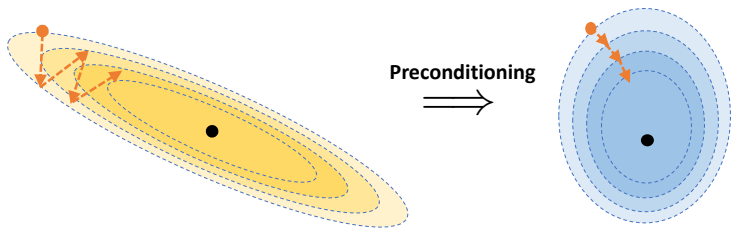ScaledGD($\lambda$) is similar to ScaledGD locally

*Concluding remarks*

# Preconditioning helps!



Preconditioning can dramatically increase the computational efficiency of vanilla gradient methods without hurting statistical efficiency

# Preconditioning helps!



Preconditioning

Preconditioning can dramatically increase the computational efficiency of vanilla gradient methods without hurting statistical efficiency

**Future directions:**

- streaming/stochastic variants of ScaledGD
- generalizing the idea of ScaledGD to other learning problems

**Papers:**

"The power of preconditioning in overparameterized low-rank matrix sensing,"
X. Xu, Y. Shen, Y. Chi, and C. Ma, ICML 2023

"Accelerating ill-conditioned low-rank matrix estimation via scaled gradient
descent," T. Tong, C. Ma, and Y. Chi, JMLR 2021