

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



Cong Ma

Department of Statistics, UChicago

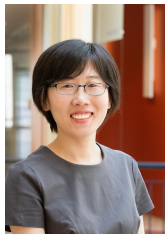
*BDML, Oct. 2023*



Xingyu Xu  
CMU



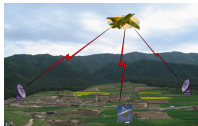
Yandi Shen  
UChicago → CMU



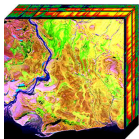
Yuejie Chi  
CMU

# Low-rank matrices in data science

---



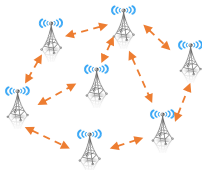
radar imaging



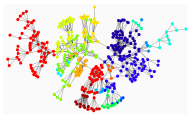
hyperspectral imaging



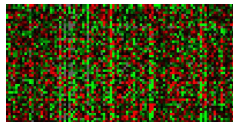
recommendation systems



localization



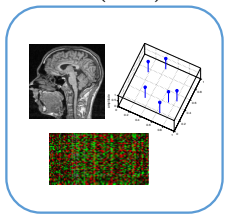
community detection



bioinformatics

# Low-rank matrix recovery

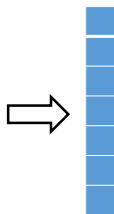
$$\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$$
$$\text{rank}(\mathbf{M}) = r$$



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



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



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

**Goal:** recover  $\mathbf{M}$  in the sample-starved regime

$$\underbrace{(n_1 + n_2)r}_{\text{degrees of freedom}} \lesssim \underbrace{m}_{\text{sensing budget}} \ll \underbrace{n_1 n_2}_{\text{ambient dimension}}$$

# Low-rank matrix factorization

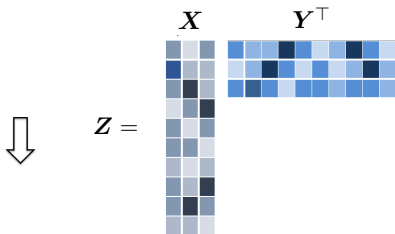
---

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

# Low-rank matrix factorization

---

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

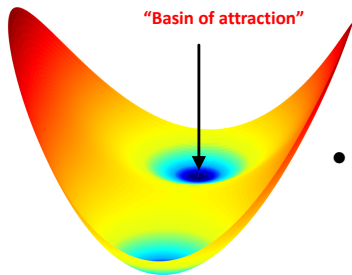


$$\min_{\mathbf{X} \in \mathbb{R}^{n_1 \times r}, \mathbf{Y} \in \mathbb{R}^{n_2 \times r}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top)\|_2^2$$



# Prior art: GD with balancing regularization

$$\min_{\mathbf{X}, \mathbf{Y}} f_{\text{reg}}(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top)\|_2^2 + \frac{1}{8} \|\mathbf{X}^\top \mathbf{X} - \mathbf{Y}^\top \mathbf{Y}\|_F^2$$



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

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

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

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

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



# Prior theory for vanilla GD

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

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

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

$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \leq \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

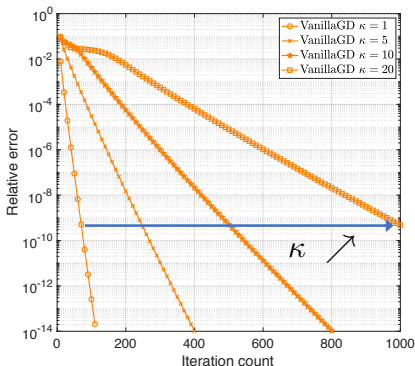
- **Computational:** within  $O(\kappa \log \frac{1}{\varepsilon})$  iterations;
- **Statistical:** as long as the sample size satisfies

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

Similar results hold for many low-rank problems

# Convergence of vanilla gradient descent

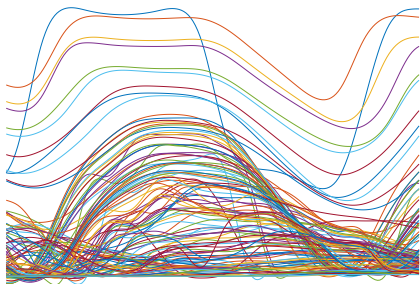
$$\text{Condition number } \kappa = \frac{\sigma_{\max}(\mathbf{M})}{\sigma_{\min}(\mathbf{M})}$$



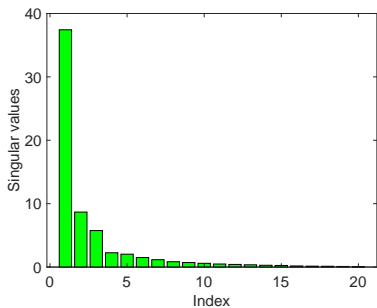
Vanilla GD converges in  $O(\kappa \log \frac{1}{\epsilon})$  iterations

# Condition number can be large

---



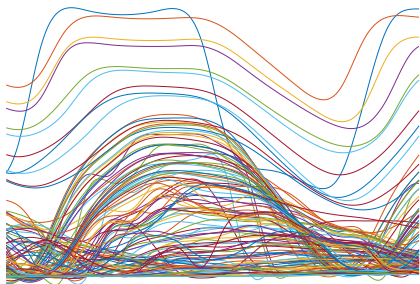
chlorine concentration levels  
120 junctions, 180 time slots



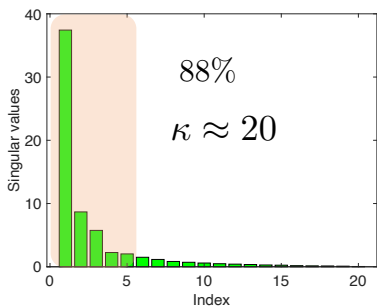
power-law spectrum

# Condition number can be large

---



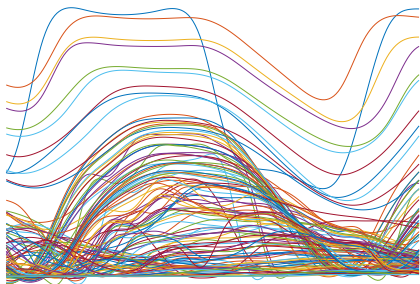
chlorine concentration levels  
120 junctions, 180 time slots



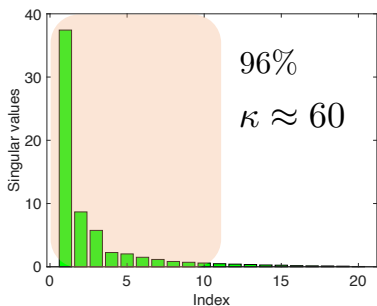
rank-5 approximation

# Condition number can be large

---

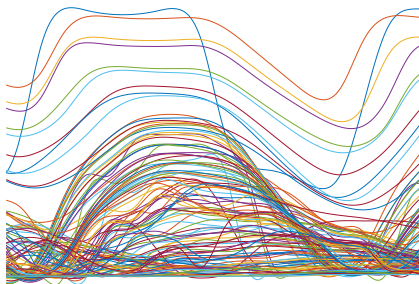


chlorine concentration levels  
120 junctions, 180 time slots

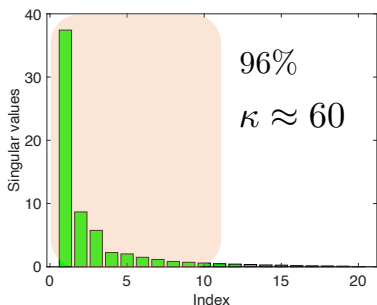


rank-10 approximation

# Condition number can be large



chlorine concentration levels  
120 junctions, 180 time slots



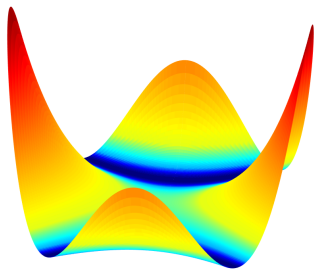
rank-10 approximation

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

# A recipe: scaled gradient descent (ScaledGD)

---

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



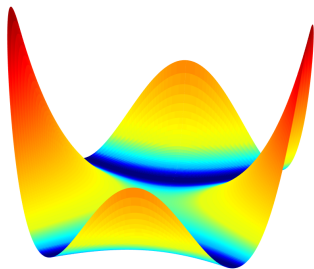
- **Spectral initialization:** find an initial point in the “basin of attraction”
- **Scaled gradient iterations:** for  $t = 0, 1, \dots$ ,

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

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

# A recipe: scaled gradient descent (ScaledGD)

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



- **Spectral initialization:** find an initial point in the “basin of attraction”
- **Scaled gradient iterations:** for  $t = 0, 1, \dots$ ,

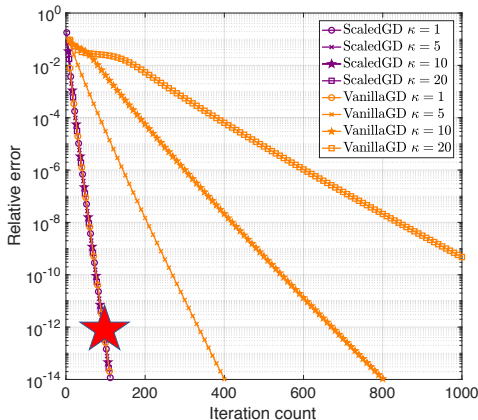
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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 :

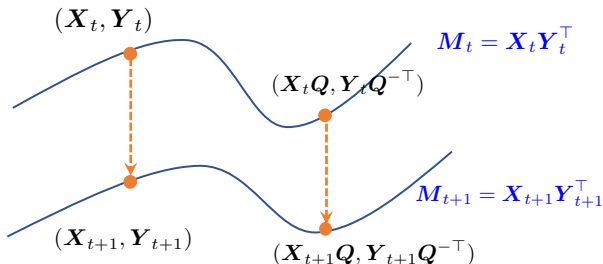
Define  $\mathbf{F}_t = [\mathbf{X}_t^\top, \mathbf{Y}_t^\top]^\top \in \mathbb{R}^{(n_1+n_2) \times r}$ . One can write update rule as

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

# A closer look at ScaledGD

---

Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



— 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*

$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **Computational:** within  $O(\log \frac{1}{\varepsilon})$  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

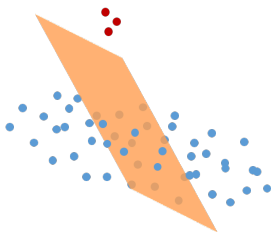
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **Computational:** within  $O(\log \frac{1}{\varepsilon})$  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

# ScaledGD works more broadly



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

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

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:**

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



# 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:**

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

**ScaledGD:**

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

*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:**

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

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

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

*add regularization to stabilize the preconditioner*

# 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  $\lambda \asymp \sigma_{\min}(\mathbf{M})$ ,  $\eta \asymp 1$ , and a sufficiently small random initialization achieves

$$\|\mathbf{X}_t \mathbf{X}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

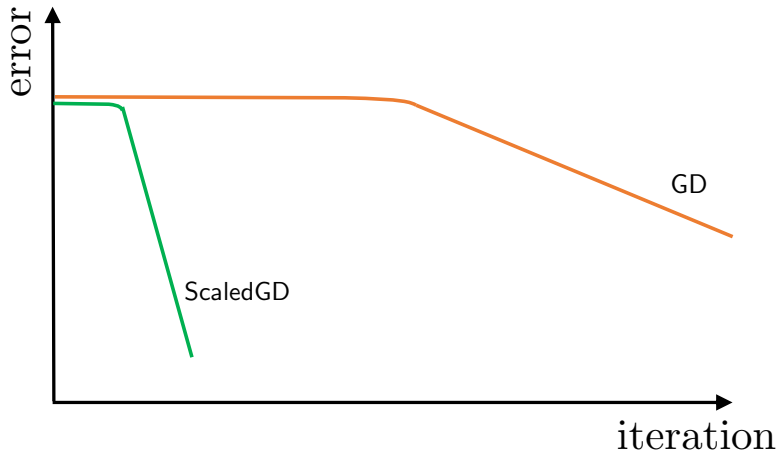
- **Computational:** within  $O(\log \kappa \log(\kappa n) + \log \frac{1}{\varepsilon})$  iterations;
- **Statistical:** the sample complexity satisfies

$$m \gtrsim nr^2 \text{poly}(\kappa)$$

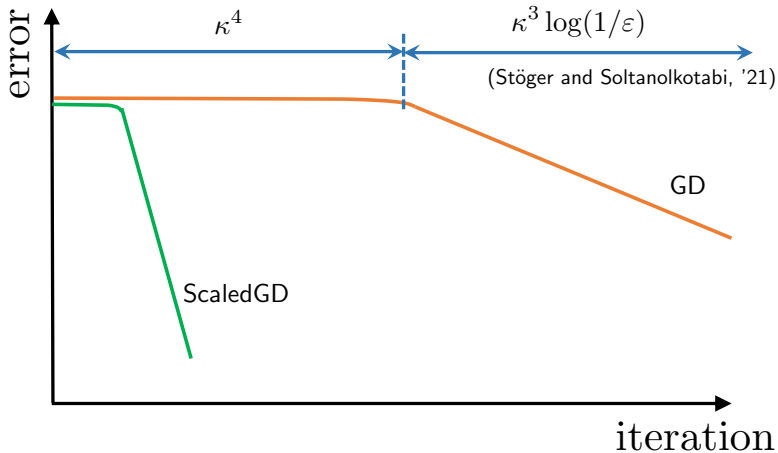
- Our analysis also enables exact convergence under random initialization with correct rank specification

## Comparison with overparameterized GD

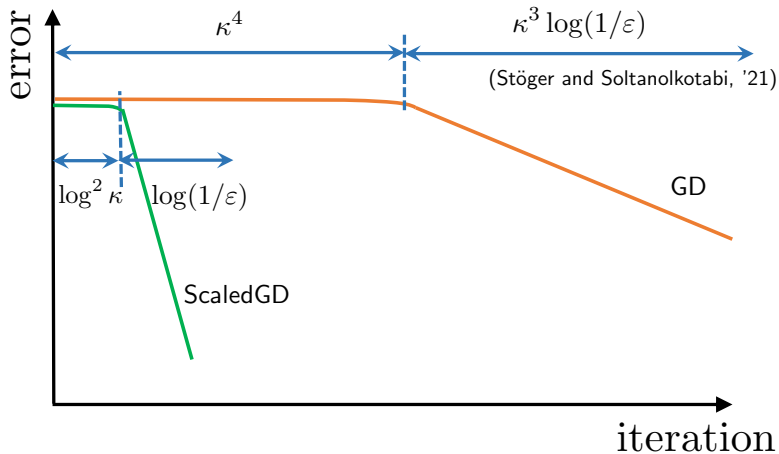
---



# 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

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

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

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

## Extension to noisy case

---

Consider the noisy setting

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

### 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*

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



# ScaledGD( $\lambda$ ) is nearly optimal

---

ScaledGD( $\lambda$ ) achieves

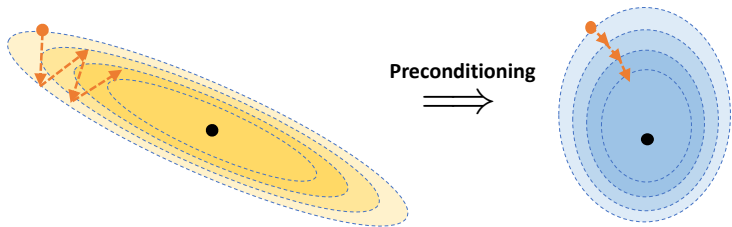
$$\|\mathbf{X}_t \mathbf{X}_t^\top - \mathbf{M}\|_F \lesssim \kappa^2 \sigma \sqrt{nr}$$

- 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  $r$

*Concluding remarks*

# Preconditioning helps!

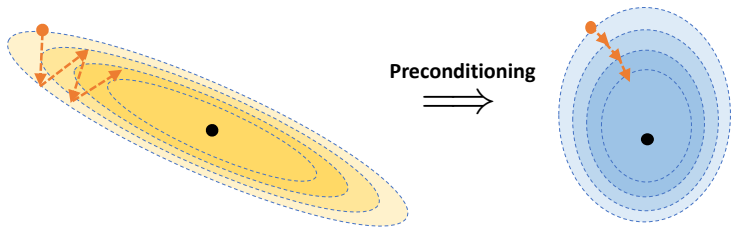
---



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

# Preconditioning helps!

---



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