

A tale of preconditioning and overparameterization in ill-conditioned low-rank estimation

Yuejie Chi

Carnegie Mellon University

CAMDA Conference
May 2023

Sensing, computing, and imaging advances

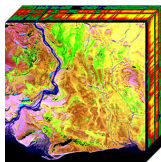
New imaging/sensing modalities allow us to probe the nature in unprecedented manners.



healthcare



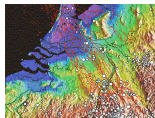
Radio astronomy



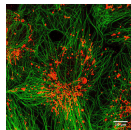
hyperspectral



Internet traffic



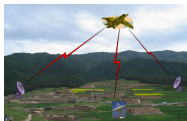
seismic imaging



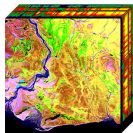
microscopy

The large amount of data brings exciting opportunities that call for new tools that are **scalable in computation and memory**.

Low-rank matrices in data science



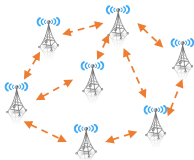
radar imaging



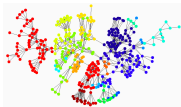
hyperspectral imaging



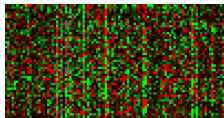
recommendation systems



localization



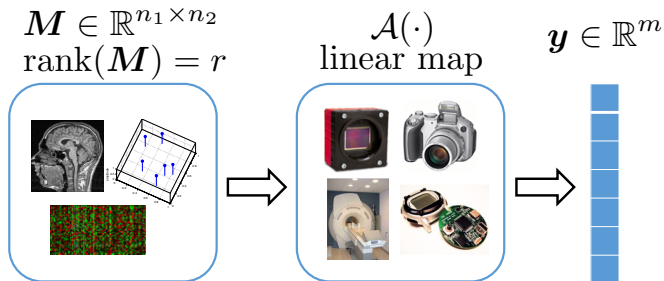
community detection



bioinformatics

Low-rank representations encode latent structures

A canonical problem: low-rank matrix sensing



$$y = \mathcal{A}(M) + \text{noise}$$

Recover M in the sample-starved regime:

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

Convex relaxation via nuclear norm minimization

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

Convex relaxation via nuclear norm minimization

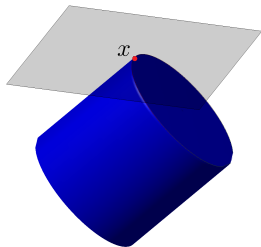
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ cvx surrogate

$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \|\mathbf{Z}\|_*$$

s.t. $\mathbf{y} \approx \mathcal{A}(\mathbf{Z})$

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



Convex relaxation via nuclear norm minimization

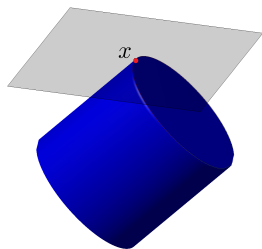
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ cvx surrogate

$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \|\mathbf{Z}\|_*$$

s.t. $\mathbf{y} \approx \mathcal{A}(\mathbf{Z})$

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

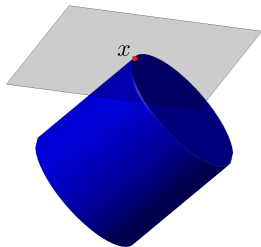
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ cvx surrogate

$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \|\mathbf{Z}\|_*$$

s.t. $\mathbf{y} \approx \mathcal{A}(\mathbf{Z})$

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, ...

Poor scalability: operate in the *ambient* matrix space

Low-rank matrix factorization

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

Low-rank matrix factorization

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



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

Low-rank matrix factorization

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



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



$$\mathbf{Z} = \begin{matrix} & \mathbf{X} & & \mathbf{Y}^\top \\ \begin{matrix} \text{[Vertical grid of 12x12 cells]} \end{matrix} & & & \begin{matrix} \text{[Horizontal grid of 12x8 cells]} \end{matrix} \end{matrix}$$

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

Low-rank matrix factorization

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



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

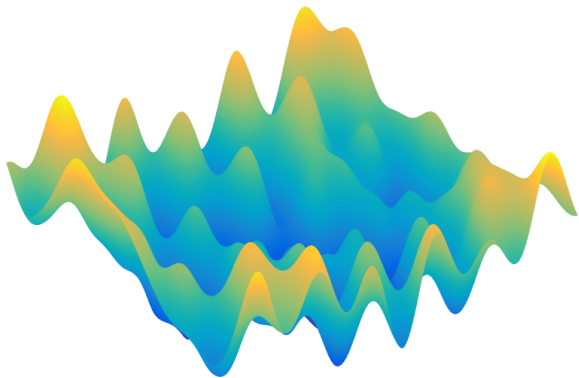
**more scalable,
but nonconvex!**



$$\mathbf{Z} = \begin{matrix} & \mathbf{X} & \mathbf{Y}^\top \\ \begin{matrix} \text{[Vertical grid of blue and grey squares]} \end{matrix} & & \begin{matrix} \text{[Horizontal grid of blue and dark blue squares]} \end{matrix} \end{matrix}$$

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

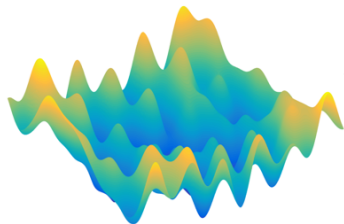
Nonconvex problems are hard (in theory)!



Nonconvex problems are hard (in theory)!

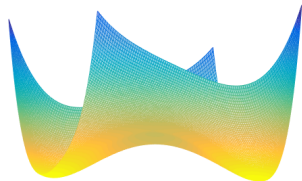


Statistics meets optimization



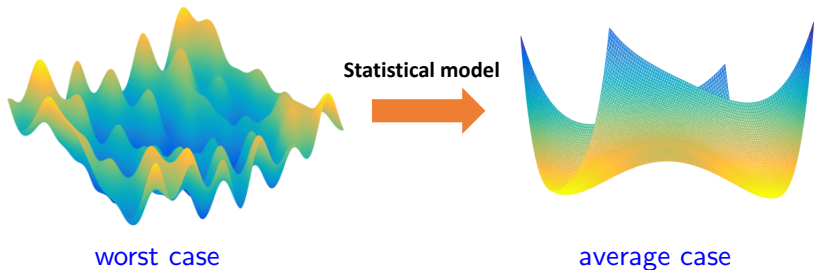
worst case

Statistical model



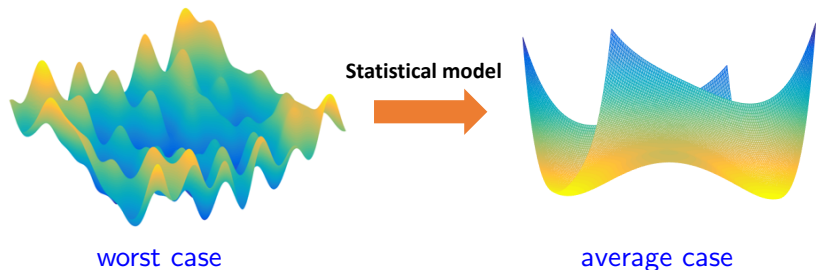
average case

Statistics meets optimization



Simple algorithms can be efficient for nonconvex problems!

Statistics meets optimization



Simple algorithms can be efficient for nonconvex problems!

Vanilla gradient descent (GD):

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

for $t = 0, 1, \dots$

Low-rank matrix sensing: GD with balancing regularization

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_2^2 = \frac{1}{2} \left\| \mathcal{A}(\mathbf{M} - \mathbf{X}\mathbf{Y}^\top) \right\|_2^2$$

Low-rank matrix sensing: GD with balancing regularization

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

Low-rank matrix sensing: GD with balancing regularization

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

- **Spectral initialization:** find an initial point in the “basin of attraction”.

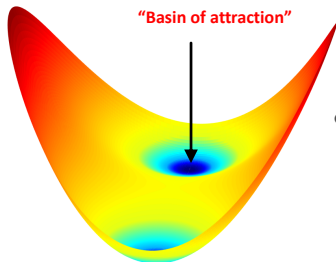
$$(\mathbf{X}_0, \mathbf{Y}_0) \leftarrow \text{SVD}_r(\mathcal{A}^*(\mathbf{y}))$$

- **Gradient iterations:**

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

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

for $t = 0, 1, \dots$



Prior art: GD for asymmetric low-rank matrix sensing

Theorem (Tu et al., ICML 2016)

Suppose $\mathbf{M} = \mathbf{X}_* \mathbf{Y}_*^\top$ is rank- r and has a *condition number* $\kappa = \sigma_{\max}(\mathbf{M})/\sigma_{\min}(\mathbf{M})$. 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 complexity satisfies

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

Prior art: GD for asymmetric low-rank matrix sensing

Theorem (Tu et al., ICML 2016)

Suppose $M = X_* Y_*^\top$ is rank- r and has a *condition number* $\kappa = \sigma_{\max}(M)/\sigma_{\min}(M)$. For low-rank matrix sensing with i.i.d. Gaussian design, vanilla GD (with spectral initialization) achieves

$$\|X_t Y_t^\top - M\|_F \leq \varepsilon \cdot \sigma_{\min}(M)$$

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

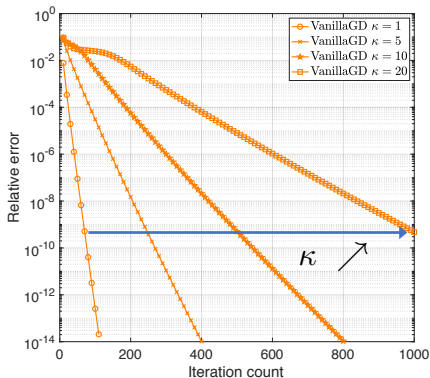
$$m \gtrsim (n_1 + n_2)r^2\kappa^2.$$

Similar results hold for many low-rank problems: matrix completion, robust PCA, etc...

(Netrapalli et al. '13, Candès, Li, Soltanolkotabi '14, Sun and Luo '15, Chen and Wainwright '15, Zheng and Lafferty '15, Ma et al. '17,)

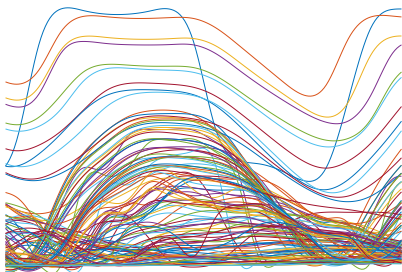
Convergence slows down for ill-conditioned matrices

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}) \right\|_{\text{F}}^2$$

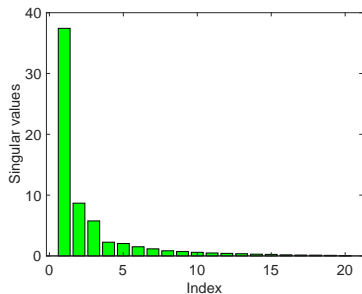


Vanilla GD converges in $O(\kappa \log \frac{1}{\epsilon})$ iterations.

Condition number can be large



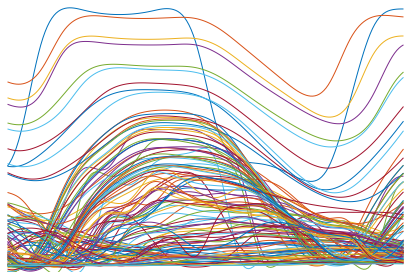
chlorine concentration levels
120 junctions, 180 time slots



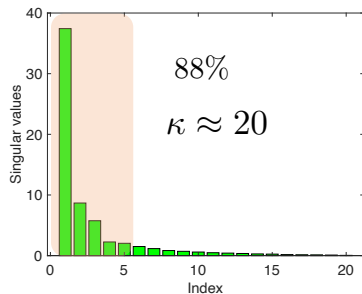
power-law spectrum

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

Condition number can be large



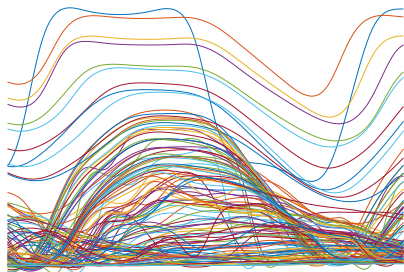
chlorine concentration levels
120 junctions, 180 time slots



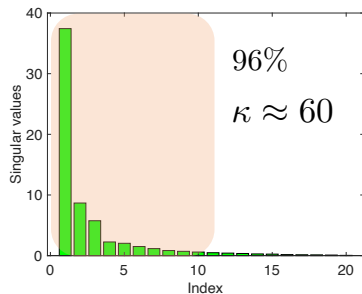
rank-5 approximation

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

Condition number can be large



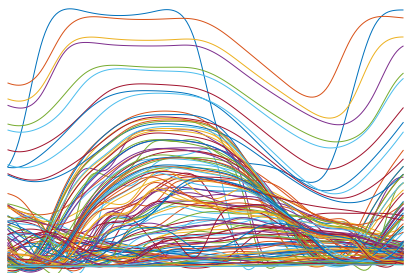
chlorine concentration levels
120 junctions, 180 time slots



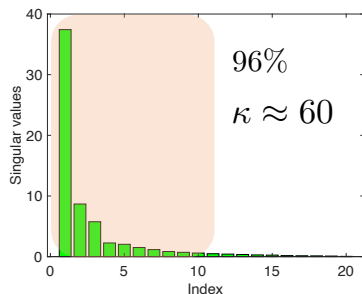
rank-10 approximation

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

Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots

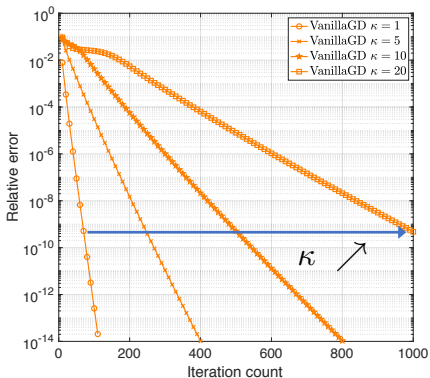


rank-10 approximation

Must mind the condition number!

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

Getting rid of the condition number?



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

This talk: the power of preconditioning

Acceleration for ill-conditioned matrix estimation:

Can we design provably fast gradient algorithms that are insensitive to the condition number of low-rank matrices?

This talk: the power of preconditioning

Acceleration for ill-conditioned matrix estimation:

Can we design provably fast gradient algorithms that are insensitive to the condition number of low-rank matrices?

Robustness to adversarial outliers:

Can we design provably robust variants that are simultaneously oblivious to the presence of outliers?

This talk: the power of preconditioning

Acceleration for ill-conditioned matrix estimation:

Can we design provably fast gradient algorithms that are insensitive to the condition number of low-rank matrices?

Robustness to adversarial outliers:

Can we design provably robust variants that are simultaneously oblivious to the presence of outliers?

Generalization to tensors:

Can we generalize to higher-dimensional objects?

This talk: the power of preconditioning

Acceleration for ill-conditioned matrix estimation:

Can we design provably fast gradient algorithms that are insensitive to the condition number of low-rank matrices?

Robustness to adversarial outliers:

Can we design provably robust variants that are simultaneously oblivious to the presence of outliers?

Generalization to tensors:

Can we generalize to higher-dimensional objects?

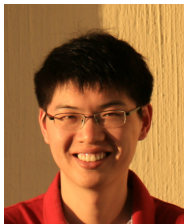
Going beyond spectral initialization and exact parameterization:

Can we still succeed with a misspecified rank?

*Accelerating gradient descent for ill-conditioned
low-rank matrix estimation*



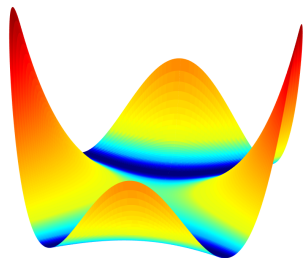
Tian Tong
CMU→Amazon



Cong Ma
UChicago

Our recipe: scaled gradient descent (ScaledGD)

$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_2^2$$



- **Spectral initialization:** find an initial point in the “basin of attraction”.
- **Scaled gradient iterations:**

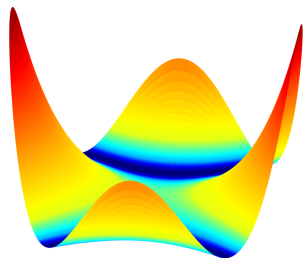
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

Our recipe: scaled gradient descent (ScaledGD)

$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_2^2$$



- **Spectral initialization:** find an initial point in the “basin of attraction”.
- **Scaled gradient iterations:**

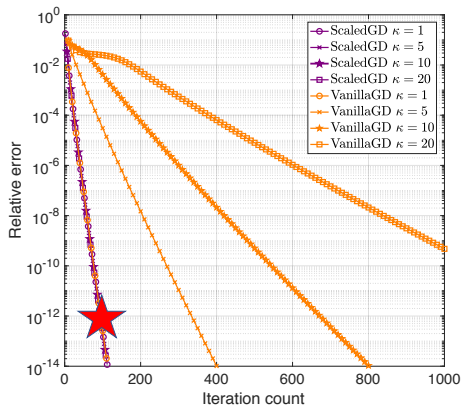
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

ScaledGD is a *preconditioned* gradient method
without balancing regularization!

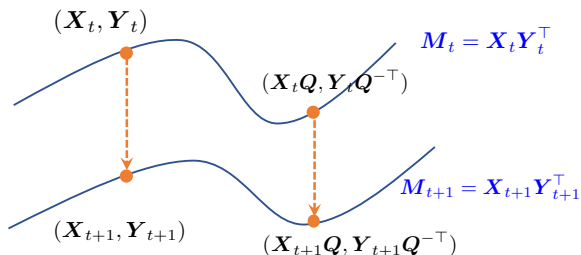
ScaledGD for low-rank matrix completion



Huge computational saving: ScaledGD converges in an κ -independent manner with a minimal overhead!

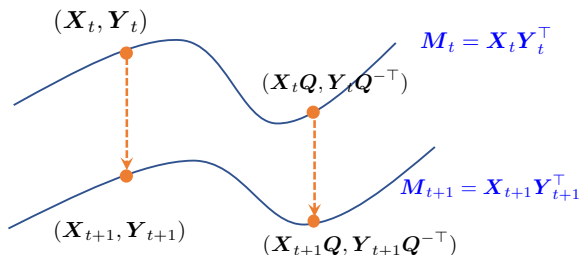
A closer look at ScaledGD

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



A closer look at ScaledGD

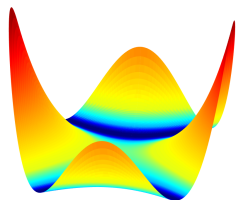
Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



New distance metric as Lyapunov function:

$$\text{dist}^2 \left(\begin{bmatrix} X \\ Y \end{bmatrix}, \begin{bmatrix} X_* \\ Y_* \end{bmatrix} \right) = \inf_{Q \in \text{GL}(r)} \left\| (XQ - X_*) \Sigma_*^{1/2} \right\|_F^2 + \left\| (YQ^{-T} - Y_*) \Sigma_*^{1/2} \right\|_F^2$$

+ a careful trajectory-based analysis



Theoretical guarantees of ScaledGD

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

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

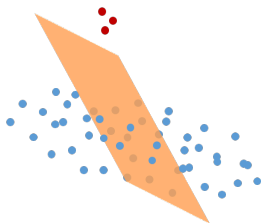
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **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 at the same sample complexity!

ScaledGD works more broadly



| | | | | |
|---|---|---|---|---|
| ✓ | ? | ? | ? | ✓ |
| ? | ? | ✓ | ✓ | ? |
| ✓ | ? | ? | ✓ | ? |
| ? | ? | ✓ | ? | ? |
| ✓ | ? | ? | ? | ? |
| ? | ✓ | ? | ? | ✓ |

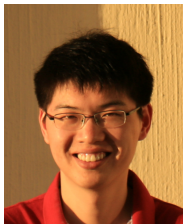
| | Robust PCA | | Matrix completion | |
|------------|--|----------------------------------|---|----------------------------------|
| Algorithms | corruption fraction | iteration complexity | sample complexity | iteration complexity |
| GD | $\frac{1}{\mu r^{3/2} \kappa^{3/2} \sqrt{\mu r \kappa^2}}$ | $\kappa \log \frac{1}{\epsilon}$ | $(\mu \vee \log n) \mu n r^2 \kappa^2$ | $\kappa \log \frac{1}{\epsilon}$ |
| 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 computation savings at comparable sample complexities!

Robustness to outliers and corruptions?

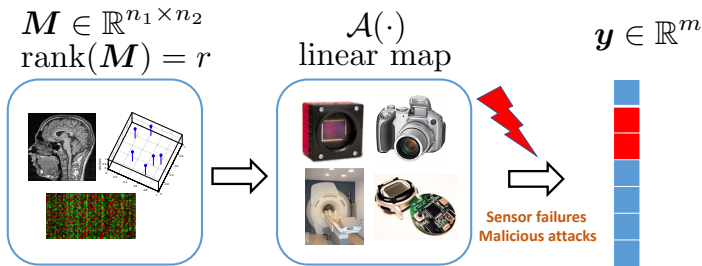


Tian Tong
CMU→Amazon



Cong Ma
UChicago

Outlier-corrupted low-rank matrix sensing



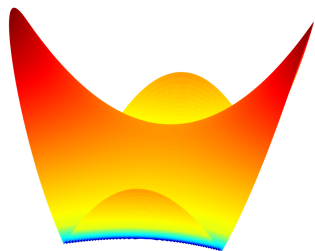
$$y = \underbrace{\mathcal{A}(M)}_{\text{outliers}} + \underbrace{s}_{\text{outliers}}, \quad \mathcal{A}(M) = \{\langle A_i, M \rangle\}_{i=1}^m$$

Arbitrary but sparse outliers: $\|s\|_0 \leq \alpha \cdot m$, where $0 \leq \alpha < 1$ is fraction of outliers.

Dealing with outliers: subgradient methods

Least absolute deviation (LAD):

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$



- **Median-truncated spectral initialization:** (Li et.al.'19).
- **Subgradient iterations:** (Charisopoulos et.al.'19; Li et al'18)

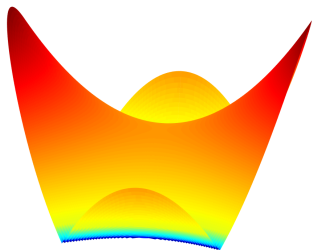
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

Dealing with outliers: subgradient methods

Least absolute deviation (LAD):

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$



- **Median-truncated spectral initialization:** (Li et.al.'19).
- **Subgradient iterations:** (Charisopoulos et.al.'19; Li et al'18)

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

Suffer from similar slow down due to ill-conditioning.

Dealing with outliers: scaled subgradient methods

Least absolute deviation (LAD):

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$

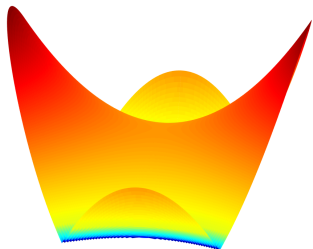
- **Median-truncated spectral initialization:** (Li et.al.'19).

- **Scaled subgradient iterations:**

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

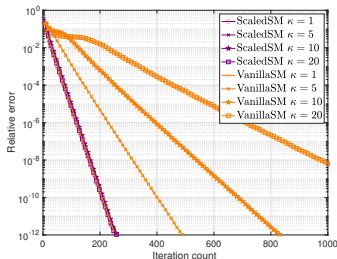
$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{X}_t^\top \mathbf{X}_t)^{-1}}_{\text{preconditioner}}$$

where η_t is set as Polyak's or geometric decaying stepsize.



Performance guarantees

| | matrix sensing | quadratic sensing |
|--|--|---|
| Subgradient Method (Charisopoulos et al, '19) | $\frac{\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ | $\frac{r\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ |
| ScaledSM (Tong, Ma, Chi, TSP '21) | $\frac{1}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ | $\frac{r}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ |



Robustness to both ill-conditioning and adversarial corruptions!

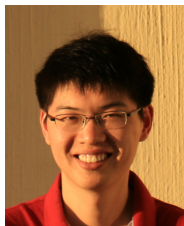
Generalization to tensors



Tian Tong
CMU→Amazon

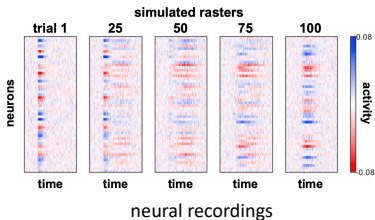


Harry Dong
CMU

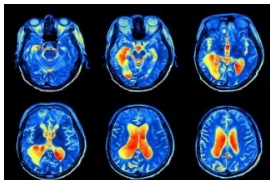


Cong Ma
UChicago

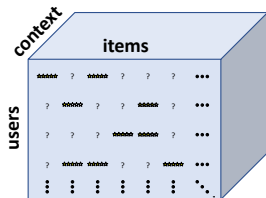
Capturing multi-way interactions by tensors



video surveillance



neuroimaging



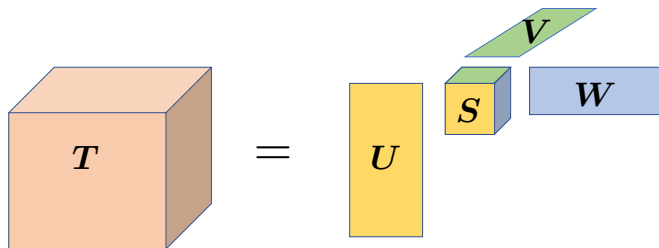
recommendation system

High-order tensors capture multi-way interactions across modalities.

Low-rank tensor under Tucker decomposition

Low-rank Tucker decomposition of a tensor:

$$T(i_1, i_2, i_3) = \sum_{j_1, j_2, j_3} S(j_1, j_2, j_3) U(i_1, j_1) V(i_2, j_2) W(i_3, j_3)$$



$$T = (U, V, W) \cdot S,$$

where $U \in \mathbb{R}^{n_1 \times r_1}$, $V \in \mathbb{R}^{n_2 \times r_2}$, $W \in \mathbb{R}^{n_3 \times r_3}$ and $S \in \mathbb{R}^{r_1 \times r_2 \times r_3}$.

Evidence that tensor problems are more challenging

Low-rank tensor recovery

Recover low-rank \mathbf{T} from $\mathbf{y} = \mathcal{A}(\mathbf{T})$.

Evidence that tensor problems are more challenging

Low-rank tensor recovery

Recover low-rank \mathbf{T} from $\mathbf{y} = \mathcal{A}(\mathbf{T})$.

- **Computation hardness:** the nuclear norm of a tensor is NP-hard to compute (Hillar and Lim, '13);

Evidence that tensor problems are more challenging

Low-rank tensor recovery

Recover low-rank \mathbf{T} from $\mathbf{y} = \mathcal{A}(\mathbf{T})$.

- **Computation hardness:** the nuclear norm of a tensor is NP-hard to compute (Hillar and Lim, '13);
- **Computational barrier:** polynomial-time algorithm exists when the sample size is above $\Omega(n^{3/2})$ (Barak and Moitra, '16);

Evidence that tensor problems are more challenging

Low-rank tensor recovery

Recover low-rank \mathbf{T} from $\mathbf{y} = \mathcal{A}(\mathbf{T})$.

- **Computation hardness:** the nuclear norm of a tensor is NP-hard to compute (Hillar and Lim, '13);
- **Computational barrier:** polynomial-time algorithm exists when the sample size is above $\Omega(n^{3/2})$ (Barak and Moitra, '16);
- **Little existing results for the Tucker case:** no provably efficient first-order algorithm for low-rank tensor completion (Han, Zhang, Willett, '20).

How to construct scaled gradients for tensors?

$$\min_{\mathbf{F}=(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \|\mathcal{A}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{y}\|_2^2$$

How to construct scaled gradients for tensors?

$$\min_{\mathbf{F}=(\mathbf{U},\mathbf{V},\mathbf{W},\mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \|\mathcal{A}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{y}\|_2^2$$

Step 1: unfolding the tensor along mode-1:

$$\mathcal{M}_1(\mathbf{T}) = \mathbf{U} \underbrace{\mathcal{M}_1(\mathbf{S})(\mathbf{V} \otimes \mathbf{W})^\top}_{\check{\mathbf{U}}^\top}$$

How to construct scaled gradients for tensors?

$$\min_{\mathbf{F}=(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \|\mathcal{A}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{y}\|_2^2$$

Step 1: unfolding the tensor along mode-1:

$$\mathcal{M}_1(\mathbf{T}) = \mathbf{U} \underbrace{\mathcal{M}_1(\mathbf{S})(\mathbf{V} \otimes \mathbf{W})^\top}_{\check{\mathbf{U}}^\top}$$

Step 2: Treat this as a matrix problem for updating factor \mathbf{U} :

$$\mathbf{U}_{t+1} = \mathbf{U}_t - \eta \nabla_{\mathbf{U}} f(\mathbf{F}_t) (\check{\mathbf{U}}_t^\top \check{\mathbf{U}}_t)^{-1}$$

How to construct scaled gradients for tensors?

$$\min_{\mathbf{F}=(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \|\mathcal{A}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{y}\|_2^2$$

Step 1: unfolding the tensor along mode-1:

$$\mathcal{M}_1(\mathbf{T}) = \mathbf{U} \underbrace{\mathcal{M}_1(\mathbf{S})(\mathbf{V} \otimes \mathbf{W})^\top}_{\check{\mathbf{U}}^\top}$$

Step 2: Treat this as a matrix problem for updating factor \mathbf{U} :

$$\mathbf{U}_{t+1} = \mathbf{U}_t - \eta \nabla_{\mathbf{U}} f(\mathbf{F}_t) (\check{\mathbf{U}}_t^\top \check{\mathbf{U}}_t)^{-1}$$

Step 3: update the core tensor \mathbf{S} :

$$\mathbf{S}_{t+1} = \mathbf{S}_t - \eta \left((\mathbf{U}_t^\top \mathbf{U}_t)^{-1}, (\mathbf{V}_t^\top \mathbf{V}_t)^{-1}, (\mathbf{W}_t^\top \mathbf{W}_t)^{-1} \right) \cdot \nabla_{\mathbf{S}} f(\mathbf{F}_t)$$

ScaledGD for ill-conditioned low-rank tensor estimation

$$\min_{\mathbf{F}=(\mathbf{U},\mathbf{V},\mathbf{W},\mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \|\mathcal{A}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{y}\|_2^2$$

Scaled gradient iterations:

$$\mathbf{U}_{t+1} = \mathbf{U}_t - \eta \nabla_{\mathbf{U}} f(\mathbf{F}_t) (\check{\mathbf{U}}_t^\top \check{\mathbf{U}}_t)^{-1},$$

$$\mathbf{V}_{t+1} = \mathbf{V}_t - \eta \nabla_{\mathbf{V}} f(\mathbf{F}_t) (\check{\mathbf{V}}_t^\top \check{\mathbf{V}}_t)^{-1},$$

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \nabla_{\mathbf{W}} f(\mathbf{F}_t) (\check{\mathbf{W}}_t^\top \check{\mathbf{W}}_t)^{-1},$$

$$\mathbf{S}_{t+1} = \mathbf{S}_t - \eta ((\mathbf{U}_t^\top \mathbf{U}_t)^{-1}, (\mathbf{V}_t^\top \mathbf{V}_t)^{-1}, (\mathbf{W}_t^\top \mathbf{W}_t)^{-1}) \cdot \nabla_{\mathbf{S}} f(\mathbf{F}_t),$$

where $\check{\mathbf{U}}_t := (\mathbf{V}_t \otimes \mathbf{W}_t) \mathcal{M}_1(\mathbf{S}_t)^\top$, $\check{\mathbf{V}}_t := (\mathbf{U}_t \otimes \mathbf{W}_t) \mathcal{M}_2(\mathbf{S}_t)^\top$, and $\check{\mathbf{W}}_t := (\mathbf{U}_t \otimes \mathbf{V}_t) \mathcal{M}_3(\mathbf{S}_t)^\top$. Here, $\mathcal{M}_k(\mathbf{S})$ is the matricization of \mathbf{S} along the k -th mode.

Key property: invariance to parameterization.

ScaledGD for low-rank tensor completion

Theorem (Tong et. al., JMLR 2022)

For low-rank tensor completion under Bernoulli sampling, assume $n = n_1 = n_2 = n_3$, ScaledGD with spectral initialization and projection achieves

$$\|(\mathbf{U}_t, \mathbf{V}_t, \mathbf{W}_t) \cdot \mathbf{S}_t - \mathbf{T}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{T})$$

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

$$n^3 p \gtrsim \mu^{3/2} r^{5/2} n^{3/2} \kappa^3 \log n.$$

ScaledGD for low-rank tensor completion

Theorem (Tong et. al., JMLR 2022)

For low-rank tensor completion under Bernoulli sampling, assume $n = n_1 = n_2 = n_3$, ScaledGD with spectral initialization and projection achieves

$$\|(\mathbf{U}_t, \mathbf{V}_t, \mathbf{W}_t) \cdot \mathbf{S}_t - \mathbf{T}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{T})$$

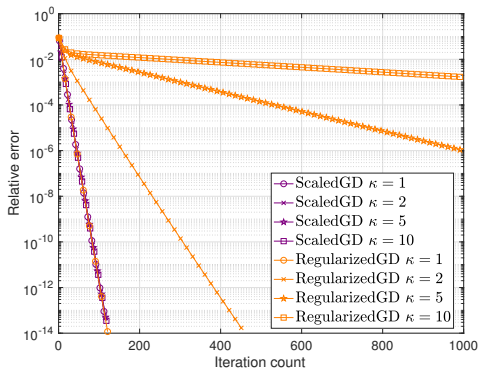
- **Computational:** within $O(\log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** as long as the sample complexity satisfies

$$n^3 p \gtrsim \mu^{3/2} r^{5/2} n^{3/2} \kappa^3 \log n.$$

First provable linear convergence at a near-optimal sample complexity for low-Tucker-rank tensor completion!

Numerical evidence

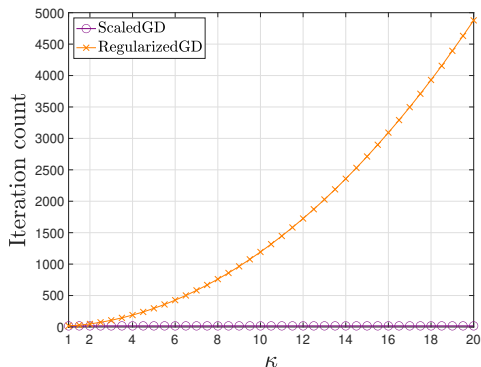
$$\min_{\mathbf{F}=(\mathbf{U},\mathbf{V},\mathbf{W},\mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{T} \right\|_{\mathbf{F}}^2$$



The benefit of ScaledGD is even more evident for tensors!

Numerical evidence

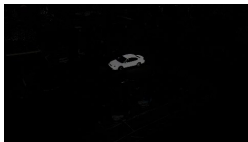
$$\min_{\mathbf{F}=(\mathbf{U},\mathbf{V},\mathbf{W},\mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{T} \right\|_{\mathbf{F}}^2$$



The benefit of ScaledGD is even more evident for tensors!

Tensor robust principal component analysis

Data = Sparse + Low-rank



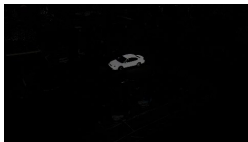
Theorem (Dong, Tong, Ma, Chi, 2022)

For a low-rank plus sparse tensor, ScaledGD with spectral initialization and *iteration-varying* thresholding converges at a constant rate, as long as *the corruption level* per fiber satisfies

$$\alpha \lesssim \frac{1}{\mu^2 r^3 \kappa}.$$

Tensor robust principal component analysis

Data = Sparse + Low-rank



Theorem (Dong, Tong, Ma, Chi, 2022)

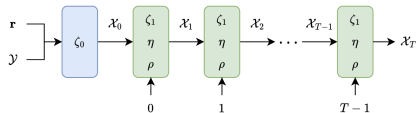
For a low-rank plus sparse tensor, ScaledGD with spectral initialization and *iteration-varying* thresholding converges at a constant rate, as long as *the corruption level per fiber* satisfies

$$\alpha \lesssim \frac{1}{\mu^2 r^3 \kappa}.$$

Can use selective mode updates to accelerate computation!

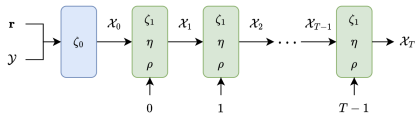
Hyperparameter tuning via self-supervised learning

unfolding + self-supervised learning



Hyperparameter tuning via self-supervised learning

unfolding + self-supervised learning

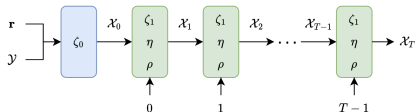


some materials data



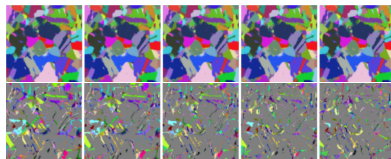
Hyperparameter tuning via self-supervised learning

unfolding + self-supervised learning



low-rank + sparse decomposition

some materials data



“Deep Unfolded Tensor Robust PCA with Self-supervised Learning”, Dong, Shah, Donegan, and Chi, ICASSP 2023.

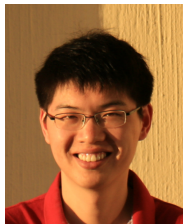
Overparameterizing (Misspecified) ScaledGD?



Xingyu Xu
CMU



Yandi Shen
UChicago



Cong Ma
UChicago

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:

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

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:

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

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:

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

ScaledGD(λ):

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

add regularization to stabilize the preconditioner

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 Amari¹, Jimmy Ba^{2,3}, Roger Grosse^{2,3}, Xuechen Li⁴, Atsushi Nitanda^{5,6},
Taiji Suzuki^{5,6}, Denny Wu^{2,3}, Ji Xu⁷

¹RIKEN CBS, ²University of Toronto, ³Vector Institute, ⁴Google Research, Brain Team,

⁵University of Tokyo, ⁶RIKEN AIP, ⁷Columbia 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 (Xu, Shen, Ma, Chi, ICML 2023)

For low-rank matrix sensing with i.i.d. Gaussian design, overparameterized ScaledGD(λ) with $\lambda \asymp \sigma_{\min}(\mathbf{M})$, $\eta \asymp 1$, and $\mathbf{X}_0 \sim \alpha \mathcal{N}(0, 1/n)$ with sufficiently small α achieves

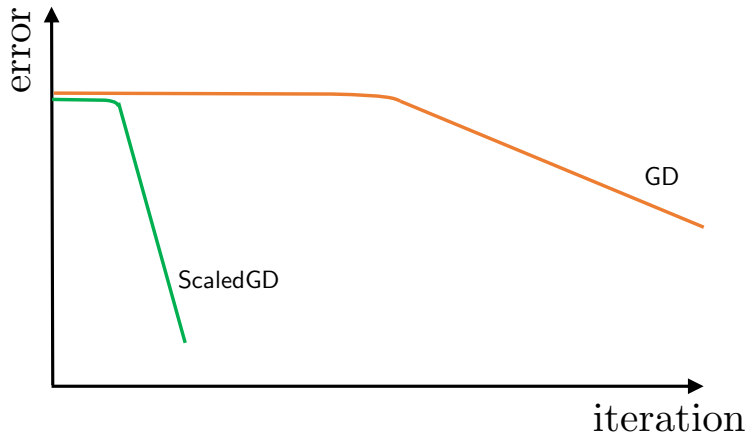
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

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

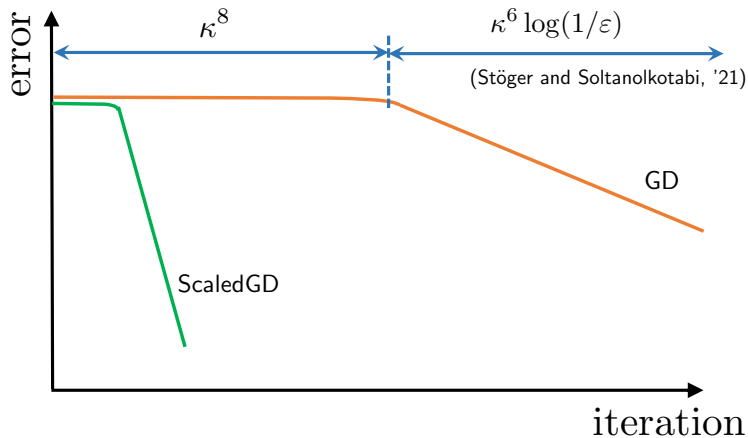
$$m \gtrsim nr^2 \text{poly}(\kappa).$$

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

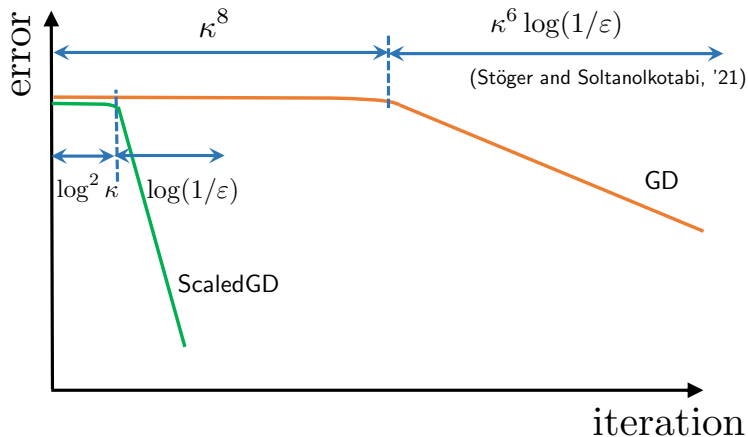
Comparison with overparameterized GD



Comparison with overparameterized GD



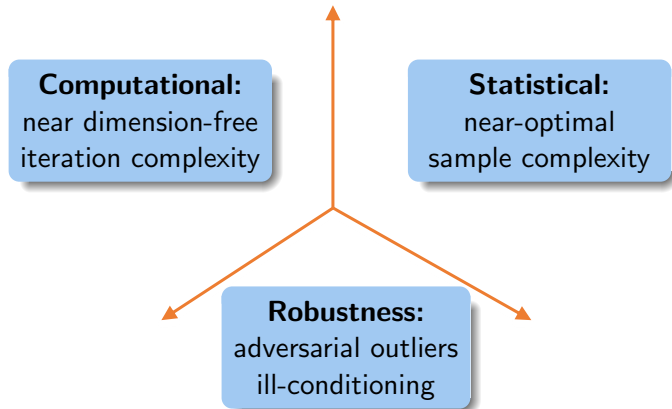
Comparison with overparameterized GD



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

Concluding remarks

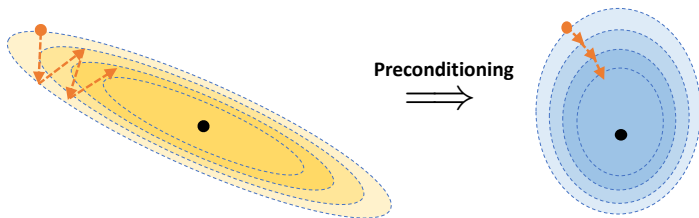
Bridging the theory-practice gap



Nonconvex low-rank matrix and tensor estimation:

- identification and exploitation of benign geometric properties;
- analyzing iterate trajectories beyond black-box optimization;
- simple variants of GD lead to robust and accelerated convergence.

Preconditioning helps!



Preconditioning dramatically increases the efficiency of vanilla gradient methods even for challenging nonconvex problems!

Ongoing directions:

- asymmetric ScaledGD with overparameterization.
- Generalizing the idea of ScaledGD to other learning and estimation problems.

Selected References

Overview:

- Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview, *IEEE Trans. on Signal Processing*, 2019.

ScaledGD for low-rank matrix estimation:

- The Power of Preconditioning in Overparameterized Low-Rank Matrix Sensing, *arXiv preprint arXiv:2302.01186*, 2023. Short version at ICML 2023.
- Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent, *Journal of Machine Learning Research*, 2021.
- Low-rank matrix recovery with scaled subgradient methods: Fast and robust convergence without the condition number, *IEEE Trans. on Signal Processing*, 2021.

ScaledGD for low-rank tensor estimation:

- Scaling and scalability: Provable nonconvex low-rank tensor estimation from incomplete measurements, *Journal of Machine Learning Research*, 2022.
- Fast and provable tensor robust principal component analysis via scaled gradient descent, *Information and Inference*, accepted.

Thanks!



<https://users.ece.cmu.edu/~yuejiec/>