

Scaling and Scalability: Accelerating Ill-conditioned Low-rank Estimation

Yuejie Chi

Carnegie Mellon University

IEEE SAM TC Webinar, December 2021



Tian Tong
Carnegie Mellon



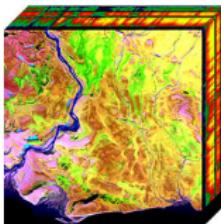
Cong Ma
University of Chicago

Sensing and imaging advances

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



healthcare



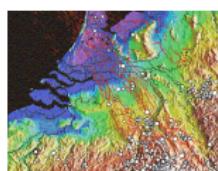
hyperspectral



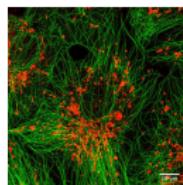
Radio astronomy



Internet traffic



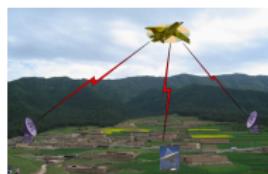
seismic imaging



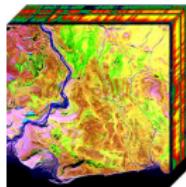
microscopy

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



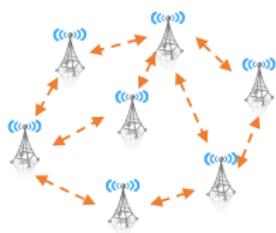
radar imaging



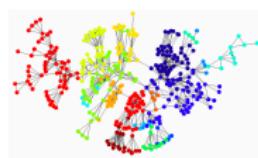
hyperspectral imaging



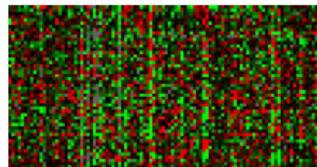
recommendation systems



localization



community detection



bioinformatics

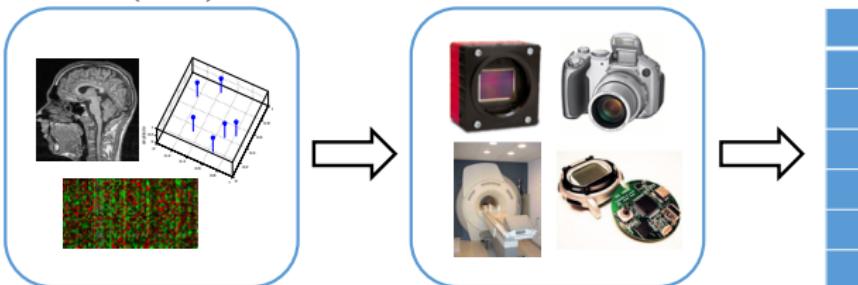
Low-rank representations encode latent structures

A canonical problem: low-rank matrix sensing

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

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

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



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

Recover \mathbf{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

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

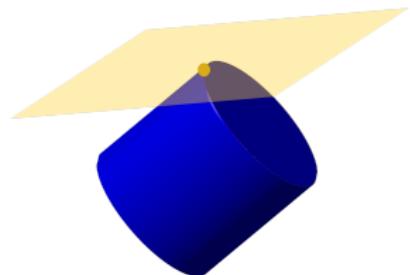
Convex relaxation via nuclear norm minimization

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

↓ cvx surrogate

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

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



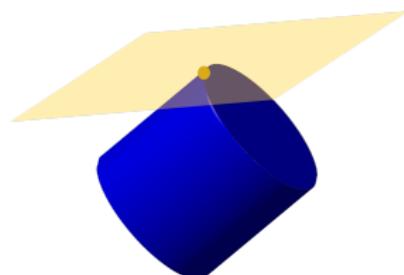
Convex relaxation via nuclear norm minimization

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

↓ cvx surrogate

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

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, ...

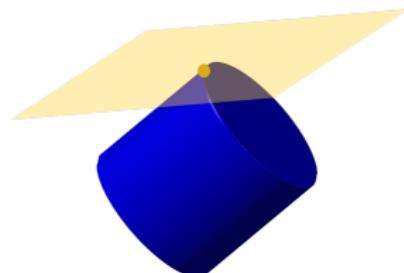
Convex relaxation via nuclear norm minimization

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

↓ cvx surrogate

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

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

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

Low-rank matrix factorization

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



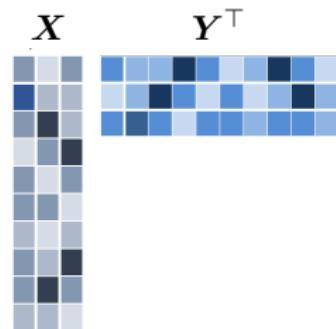
$$\min_{\text{rank}(\mathbf{Z})=r} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{Z})\|_2^2$$

Low-rank matrix factorization

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



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



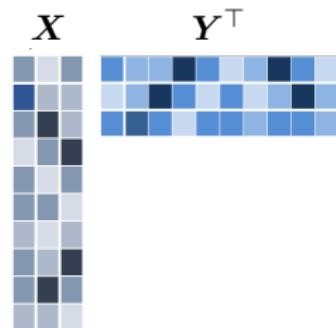
$$\min_{X \in \mathbb{R}^{n_1 \times r}, Y \in \mathbb{R}^{n_2 \times r}} f(X, Y) = \frac{1}{2} \|y - \mathcal{A}(XY^T)\|_2^2$$

Low-rank matrix factorization

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



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



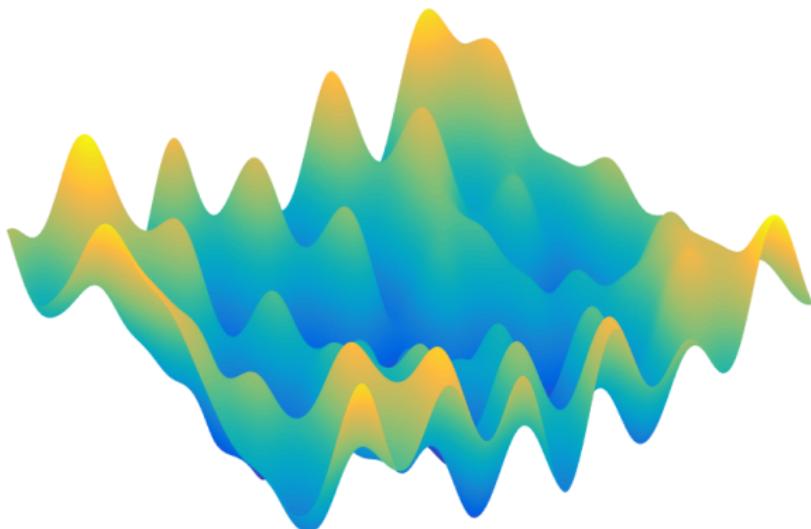
**more scalable,
but nonconvex!**



$$Z =$$

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

Nonconvex problems are hard (in theory)!



“...in fact, the great watershed in optimization isn’t between linearity and nonlinearity, but convexity and nonconvexity.

R. T. Rockafellar, in SIAM Review, 1993

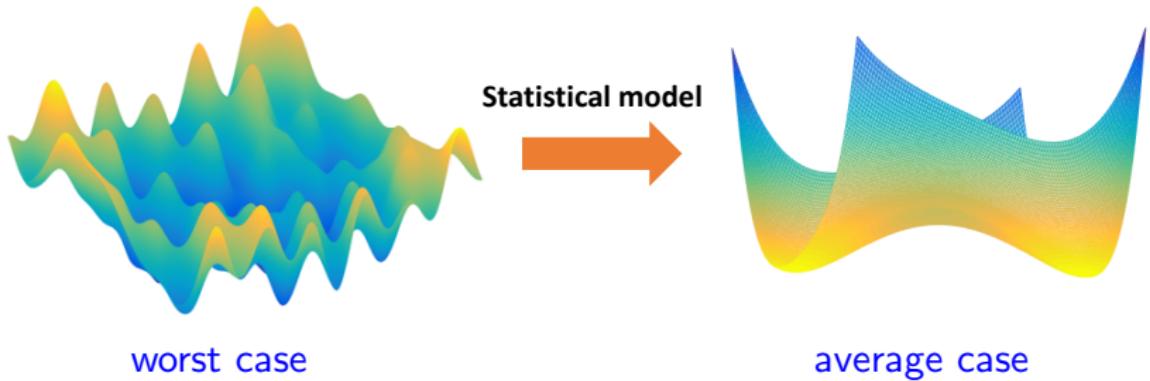
Nonconvex problems are hard (in theory)!



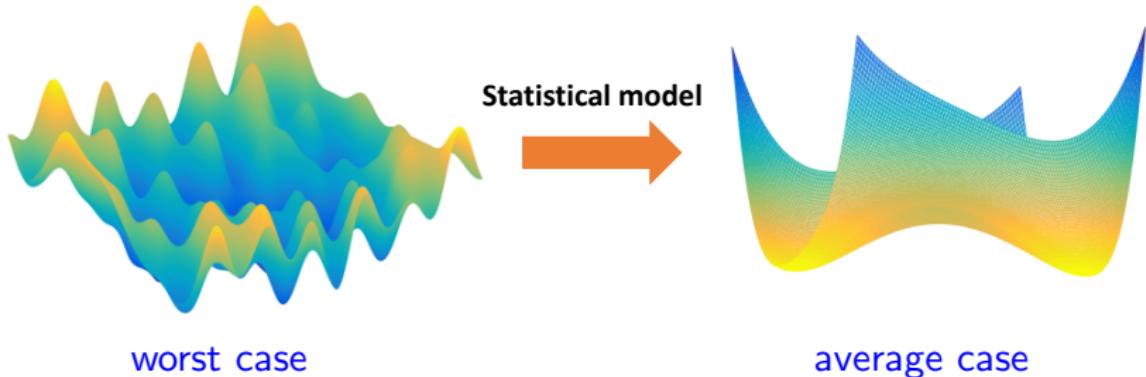
“...in fact, the great watershed in optimization isn’t between linearity and nonlinearity, but convexity and nonconvexity.

R. T. Rockafellar, in SIAM Review, 1993

Statistics meets optimization

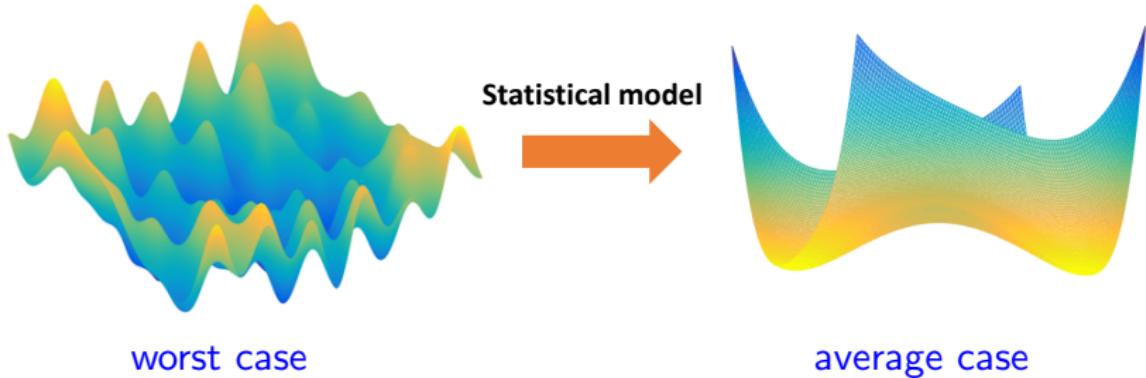


Statistics meets optimization



Simple algorithms can be efficient for nonconvex problems!

Statistics meets optimization



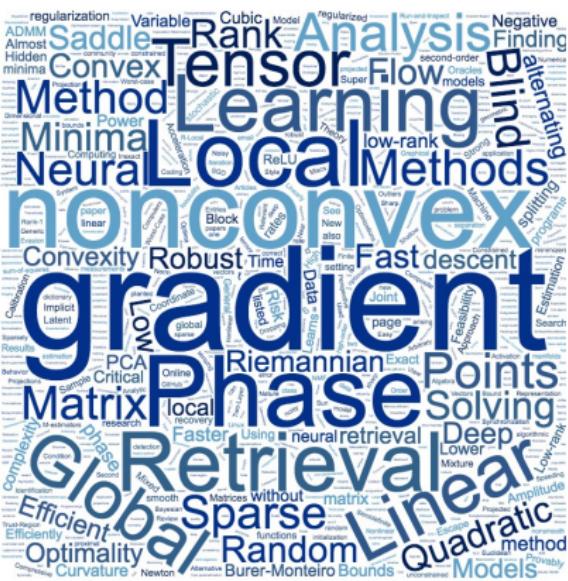
Simple algorithms can be efficient for nonconvex problems!

Vanilla gradient descent (GD):

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

for $t = 0, 1, \dots$

Recent developments: provable nonconvex optimization



"Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview," Chi, Lu, Chen, TSP 2019

Phase retrieval: Netrapalli et al. '13, Candès, Li, Soltanolkotabi '14, Chen, Candès '15, Cai, Li, Ma '15, Zhang et al. '16, Wang et al. '16, Sun, Qu, Wright '16, Ma et al. '17, Chen et al. '18, Soltani, Hegde '18, Ruan and Duchi. '18, ...

Matrix sensing/completion: Keshavan et al. '09, Jain et al. '09, Hardt '13, Jain et al. '13, Sun, Luo '15, Chen, Wainwright '15, Tu et al. '15, Zheng, Lafferty '15, Bhojanapalli et al. 16, Ge, Lee, Ma '16, Jin et al. '16, Ma et al. '17, Chen and Li '17, Cai et al. '18, Li, Zhu, Tang, Wakin '18, Charisopoulos et al. '19, ...

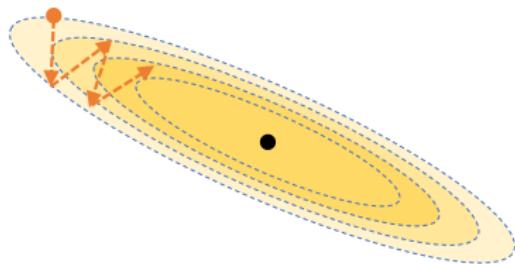
Blind deconvolution/demixing: Li et al.'16, Lee et al.'16, Cambareri, Jacques'16, Ling, Strohmer'16, Huang, Hand'16, Ma et al.'17, Zhang et al.'18, Li, Bresler'18, Dong, Shi'18, Shi, Chi'19, Qu et al.'19...

Dictionary learning: Arora et al. '14, Sun et al. '15, Chatterji, Bartlett '17, Bai et al. '18, Gilboa et al. '18, Rambhatla et al. '19, Qu et al. '19, ...

Robust principal component analysis: Netrapalli et al. '14, Yi et al. '16, Gu et al. '16, Ge et al. '17, Cherapanamjeri et al. '17, Vaswani et al. '18, Maunu et al. '19, ...

Deep learning: Zhong et al. '17, Bai, Mei, Montanari '17, Du et al. '17, Ge, Lee, Ma '17, Du et al. '18, Soltanolkotabi and Oymak, '18...

Acceleration via preconditioning

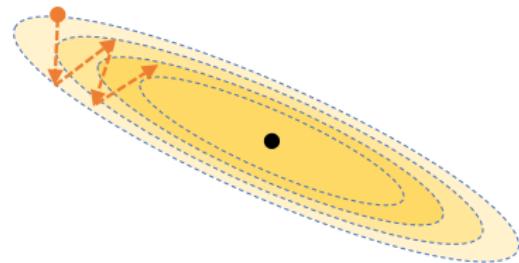


Vanilla GD:

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

⌚ Slows down with ill-conditioning.

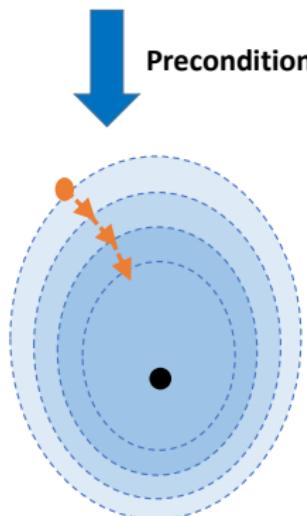
Acceleration via preconditioning



Vanilla GD:

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

⌚ Slows down with ill-conditioning.

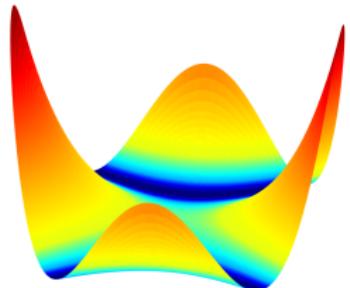


Preconditioned GD:

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

⌚ Preconditioning helps!

Robustness via nonsmooth optimization

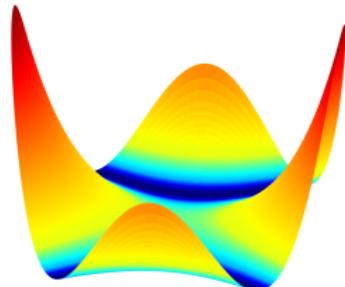


Least squares:

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

⌚ Sensitive to outliers.

Robustness via nonsmooth optimization



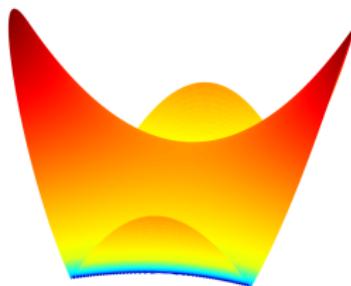
Least squares:

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

☺ Sensitive to outliers.



Nonsmooth



Least absolute deviation:

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

☺ Nonsmoothness helps!

This talk

Optimization geometry:

When and why does simple gradient descent work well for low-rank matrix estimation?

This talk

Optimization geometry:

When and why does simple gradient descent work well for low-rank matrix estimation?

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

Optimization geometry:

When and why does simple gradient descent work well for low-rank matrix estimation?

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

Optimization geometry:

When and why does simple gradient descent work well for low-rank matrix estimation?

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:

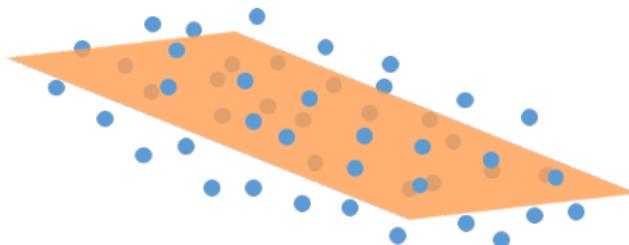
Can we design provably robust variants that are simultaneously oblivious to the presence of outliers?

Generalization to tensors:

Can we generalize to higher-dimensional objects?

*Warm-up: understanding the geometry
of low-rank matrix factorization*

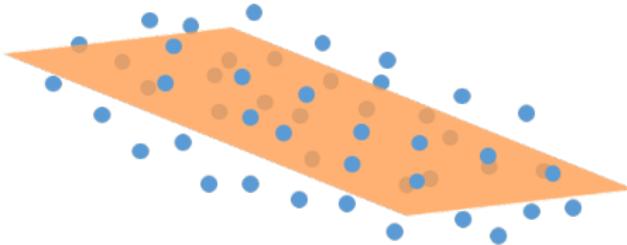
Revisiting PCA: in search of low-rank representation



Given $M \succeq 0 \in \mathbb{R}^{n \times n}$ (e.g. sample covariance matrix), find its best rank- r approximation:

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

Revisiting PCA: in search of low-rank representation



This problem admits a closed-form solution:

- let $\mathbf{M} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ be eigen-decomposition of \mathbf{M} ($\lambda_1 \geq \dots \lambda_r > \lambda_{r+1} \dots \geq \lambda_n$), then

$$\widehat{\mathbf{M}} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

— *nonconvex, but tractable*

An optimization viewpoint

Low-rank factorization: if we factorize $Z = XX^\top$ with $X \in \mathbb{R}^{n \times r}$, then it leads to a nonconvex problem:

$$\text{minimize}_{X \in \mathbb{R}^{n \times r}} \quad f(X) = \|XX^\top - M\|_F^2$$

An optimization viewpoint

Low-rank factorization: if we factorize $\mathbf{Z} = \mathbf{X}\mathbf{X}^\top$ with $\mathbf{X} \in \mathbb{R}^{n \times r}$, then it leads to a nonconvex problem:

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

Theorem (Baldi and Hornik, 1989)

Suppose \mathbf{M} has a strict eigen-gap between λ_r and λ_{r+1} , the critical points of $f(\mathbf{X})$ can be categorized into

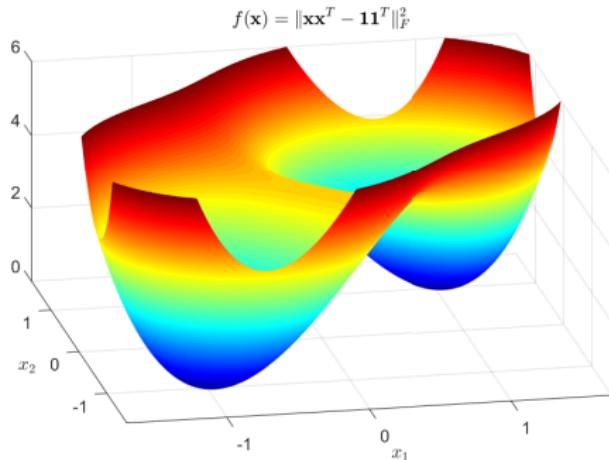
- global minima;
- strict saddle points, from which there exist directions to strictly decrease $f(\mathbf{X})$.

In other words, *all local minima are global minima!*

Baldi and Hornik. "Neural networks and principal component analysis: Learning from examples without local minima." Neural networks 2.1 (1989): 53-58.

Benign landscape of PCA

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



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

Local geometry: the hidden convexity

$$f(\mathbf{X}) := \left\| \mathbf{X} \mathbf{X}^\top - \mathbf{X}_\star \mathbf{X}_\star^\top \right\|_{\text{F}}^2, \quad \mathbf{X} \in \mathbb{R}^{n \times r}$$

A modified distance metric:

$$\text{dist}^2(\mathbf{X}, \mathbf{X}_\star) = \|\mathbf{X} \mathbf{H}_{\mathbf{X}} - \mathbf{X}_\star\|_{\text{F}}^2,$$

where $\mathbf{H}_{\mathbf{X}} := \operatorname{argmin}_{\mathbf{H} \in \mathcal{O}^{r \times r}} \|\mathbf{X} \mathbf{H} - \mathbf{X}_\star\|_{\text{F}}^2$.

Local geometry: the hidden convexity

$$f(\mathbf{X}) := \left\| \mathbf{X} \mathbf{X}^\top - \mathbf{X}_\star \mathbf{X}_\star^\top \right\|_{\text{F}}^2, \quad \mathbf{X} \in \mathbb{R}^{n \times r}$$

A modified distance metric:

$$\text{dist}^2(\mathbf{X}, \mathbf{X}_\star) = \|\mathbf{X} \mathbf{H}_{\mathbf{X}} - \mathbf{X}_\star\|_{\text{F}}^2,$$

where $\mathbf{H}_{\mathbf{X}} := \operatorname{argmin}_{\mathbf{H} \in \mathcal{O}^{r \times r}} \|\mathbf{X} \mathbf{H} - \mathbf{X}_\star\|_{\text{F}}^2$.

Restricted strong convexity:

$$\operatorname{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{X}) \operatorname{vec}(\mathbf{V}) \gtrsim \|\mathbf{V}\|_{\text{F}}^2, \quad \mathbf{V} := \mathbf{X} \mathbf{H}_{\mathbf{X}} - \mathbf{X}_\star.$$

Local geometry: the hidden convexity

$$f(\mathbf{X}) := \left\| \mathbf{X} \mathbf{X}^\top - \mathbf{X}_\star \mathbf{X}_\star^\top \right\|_{\text{F}}^2, \quad \mathbf{X} \in \mathbb{R}^{n \times r}$$

A modified distance metric:

$$\text{dist}^2(\mathbf{X}, \mathbf{X}_\star) = \|\mathbf{X} \mathbf{H}_{\mathbf{X}} - \mathbf{X}_\star\|_{\text{F}}^2,$$

where $\mathbf{H}_{\mathbf{X}} := \operatorname{argmin}_{\mathbf{H} \in \mathcal{O}^{r \times r}} \|\mathbf{X} \mathbf{H} - \mathbf{X}_\star\|_{\text{F}}^2$.

Restricted strong convexity:

$$\operatorname{vec}(\mathbf{V})^\top \nabla^2 f(\mathbf{X}) \operatorname{vec}(\mathbf{V}) \gtrsim \|\mathbf{V}\|_{\text{F}}^2, \quad \mathbf{V} := \mathbf{X} \mathbf{H}_{\mathbf{X}} - \mathbf{X}_\star.$$

Linear convergence: The GD iterates obey

$$\text{dist}^2(\mathbf{X}_t, \mathbf{X}_\star) \leq \left(1 - \frac{c}{\kappa}\right)^t \text{dist}^2(\mathbf{X}_0, \mathbf{X}_\star), \quad t \geq 0,$$

as long as $\text{dist}^2(\mathbf{X}_0, \mathbf{X}_\star) \lesssim \lambda_1$. Here, $\kappa := \lambda_1/\lambda_r$.

Gradient descent for low-rank matrix sensing

Low-rank matrix sensing: GD with balancing regularization

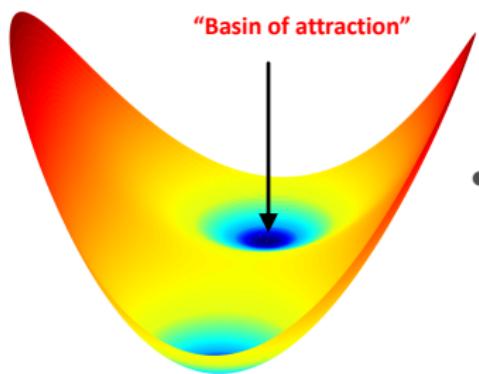
$$\min_{\mathbf{X}, \mathbf{Y}} \quad f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X} \mathbf{Y}^\top) \right\|_2^2 = \frac{1}{2} \left\| \mathcal{A}(\mathbf{M} - \mathbf{X} \mathbf{Y}^\top) \right\|_2^2$$

Low-rank matrix sensing: GD with balancing regularization

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

Low-rank matrix sensing: GD with balancing regularization

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



- **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 $M = X_\star Y_\star^\top$ 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

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

Prior art: GD for asymmetric low-rank matrix sensing

Theorem (Tu et al., ICML 2016)

Suppose $M = X_\star Y_\star^\top$ 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

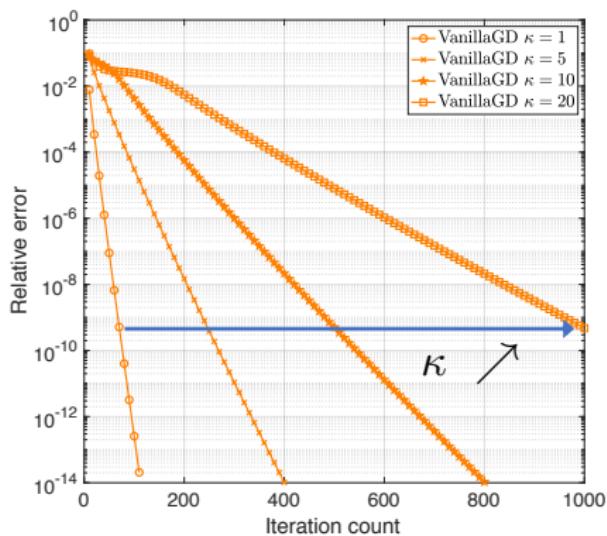
$$m \gtrsim (n_1 + n_2)r^2\kappa^2.$$

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,)

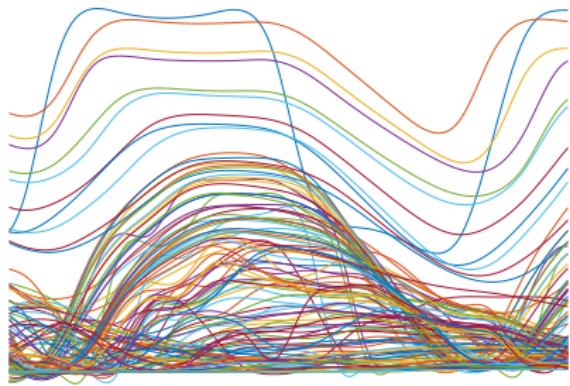
Convergence slows down for ill-conditioned matrices

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

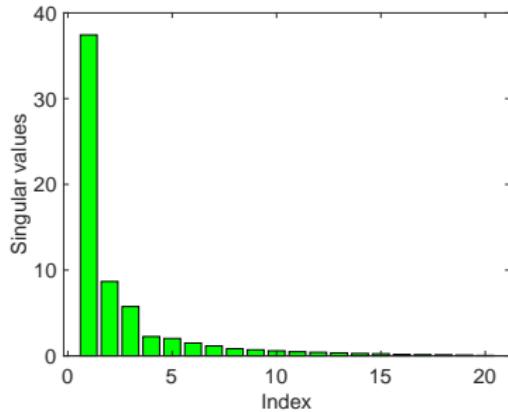


Vanilla GD converges in $O\left(\kappa \log \frac{1}{\varepsilon}\right)$ iterations.

Condition number can be large

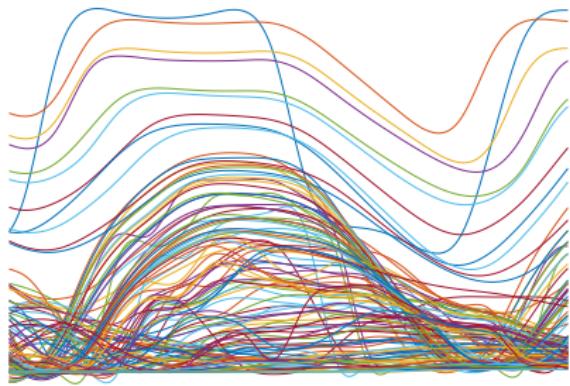


chlorine concentration levels
120 junctions, 180 time slots

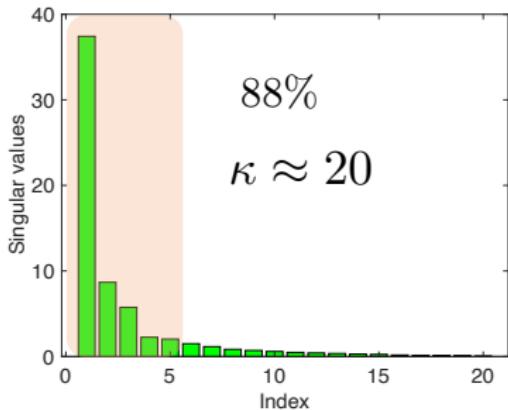


power-law spectrum

Condition number can be large



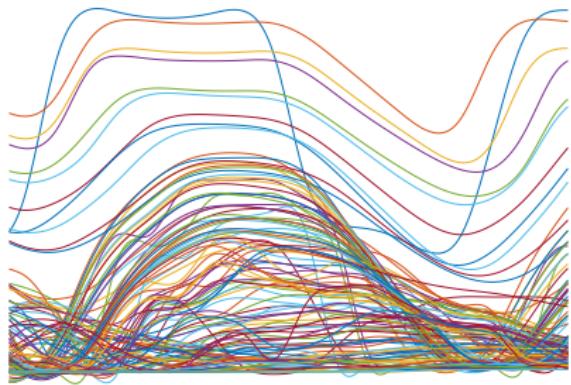
chlorine concentration levels
120 junctions, 180 time slots



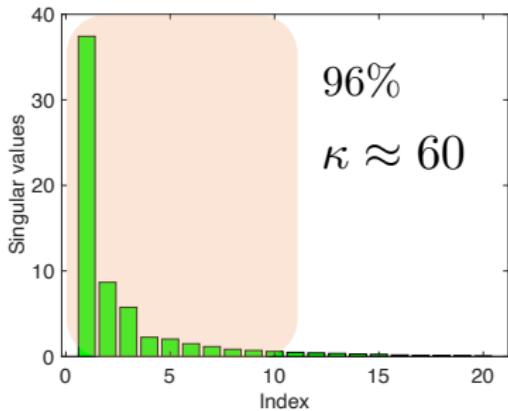
rank-5 approximation

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

Condition number can be large

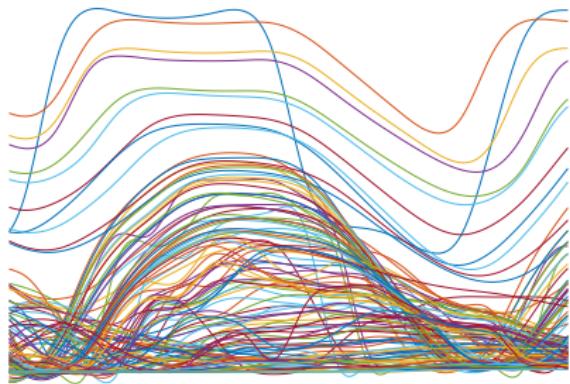


chlorine concentration levels
120 junctions, 180 time slots

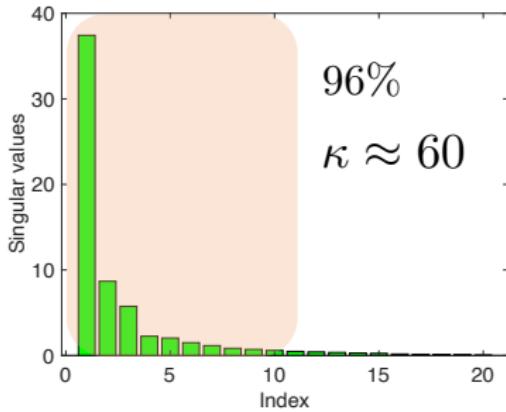


rank-10 approximation

Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots



rank-10 approximation

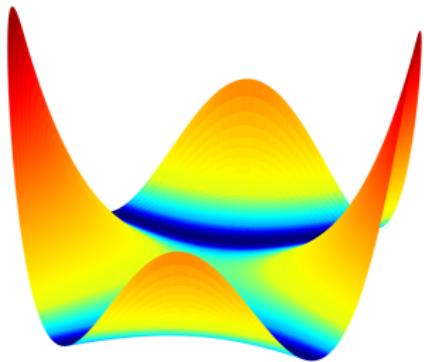
Can we accelerate the convergence rate of GD to $O(\log \frac{1}{\epsilon})$?

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

Accelerating ill-conditioned matrix estimation

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:**

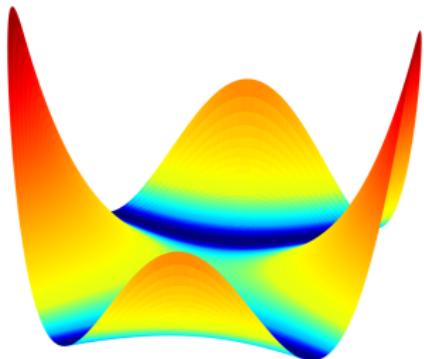
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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:**

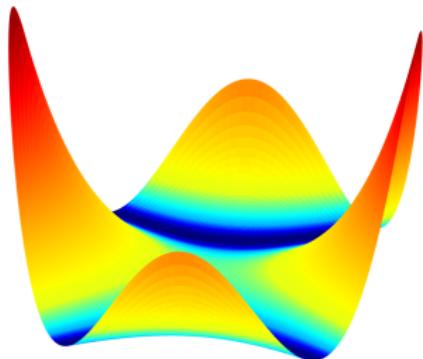
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

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:**

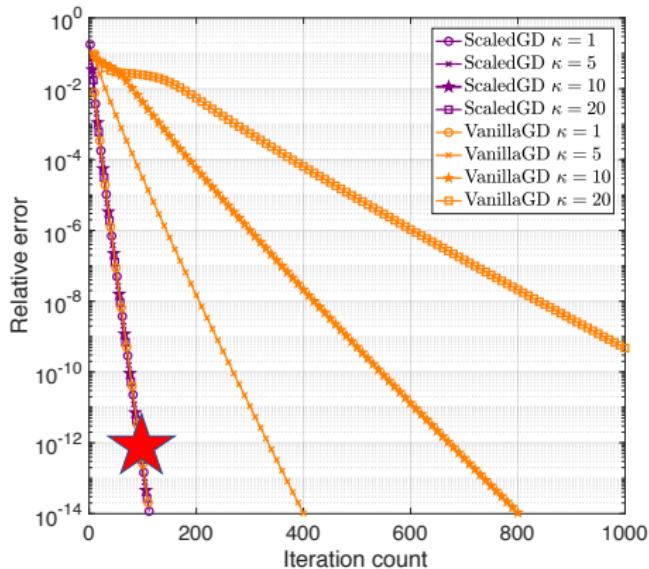
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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!

ScaledGD as a quasi-Newton method

$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{XY}^\top - \mathbf{X}_\star \mathbf{Y}_\star^\top \right\|_{\text{F}}^2$$

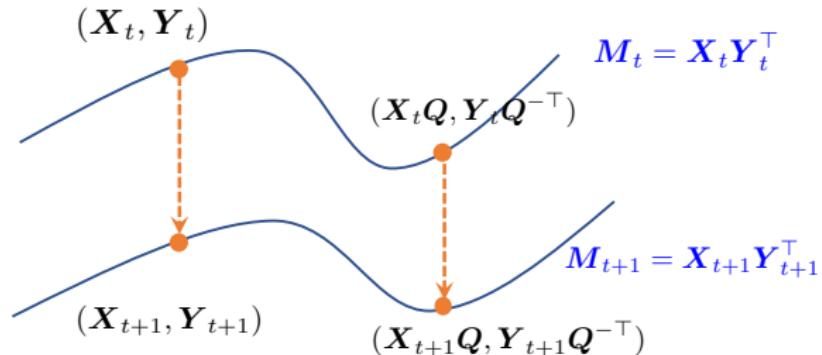
For low-rank matrix factorization, ScaledGD is equivalent to:

$$\begin{bmatrix} \text{vec}(\mathbf{X}) \\ \text{vec}(\mathbf{Y}) \end{bmatrix} \Longleftarrow \begin{bmatrix} \text{vec}(\mathbf{X}) \\ \text{vec}(\mathbf{Y}) \end{bmatrix} - \eta \begin{bmatrix} \nabla_{\mathbf{X}, \mathbf{X}}^2 f & \mathbf{0} \\ \mathbf{0} & \nabla_{\mathbf{Y}, \mathbf{Y}}^2 f \end{bmatrix}^{-1} \begin{bmatrix} \text{vec}(\nabla_{\mathbf{X}} f) \\ \text{vec}(\nabla_{\mathbf{Y}} f) \end{bmatrix}.$$

The preconditioners are chosen as the inverse of the block diagonal approximation of the Hessian to low-rank matrix factorization.

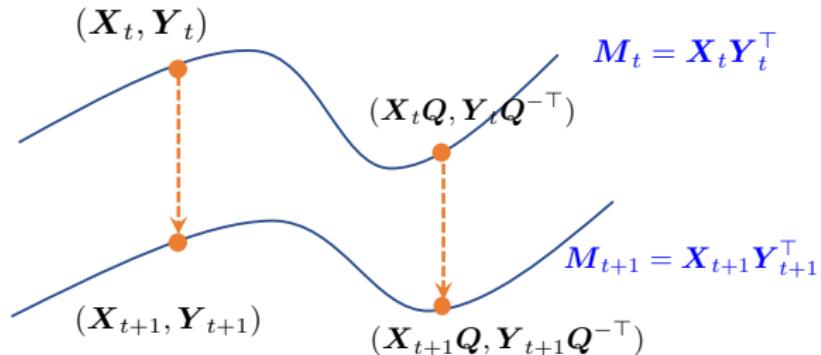
A closer look at ScaledGD

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



A closer look at ScaledGD

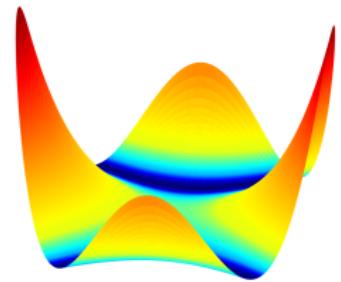
Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



New distance metric as Lyapunov function:

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

+ a careful trajectory-based analysis



Theoretical guarantees of ScaledGD

Theorem (Tong, Ma and Chi, 2020)

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

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

- **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 (Tong, Ma and Chi, 2020)

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

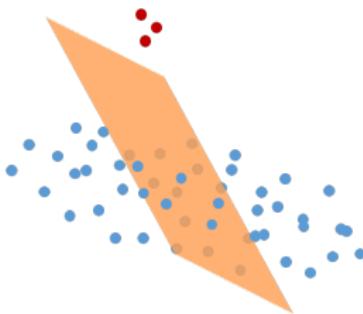
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_{\text{F}} \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **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 at the same sample complexity!

ScaledGD works more broadly



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

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

Code available at <https://github.com/Titan-Tong/ScaledGD>

What about the run time?

The run time of ScaledGD is rather competitive, with additional suitability for parallel implementation.

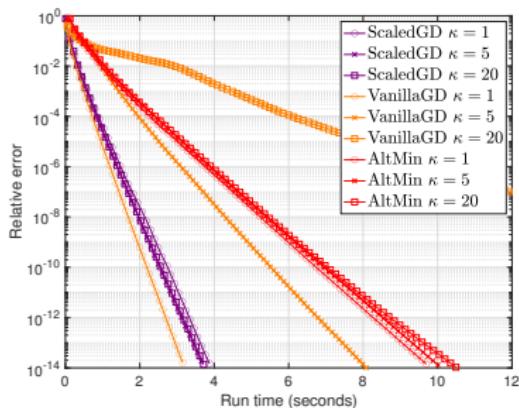
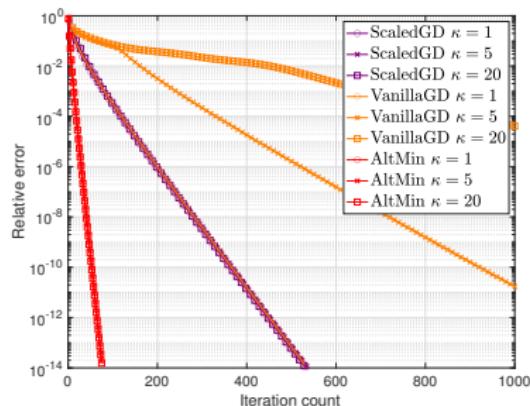
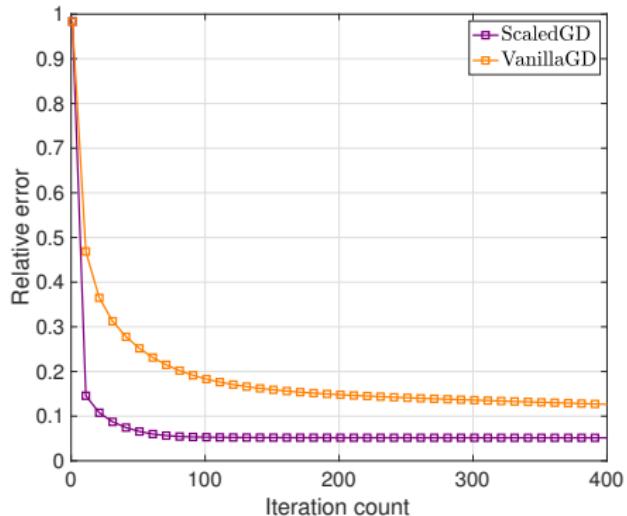


Figure: Run time for matrix completion with $n = 1000$, $p = 0.2$, $r = 50$.

Numerical stability

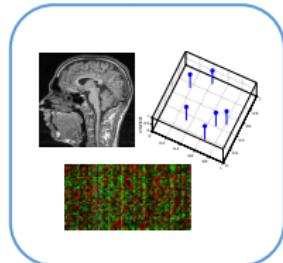
ScaledGD converges faster than GD in a small number of iterations (they eventually reach the same accuracy).



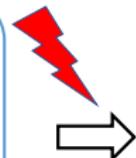
Robustness to outliers and corruptions?

Outlier-corrupted low-rank matrix sensing

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



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



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



Sensor failures
Malicious attacks

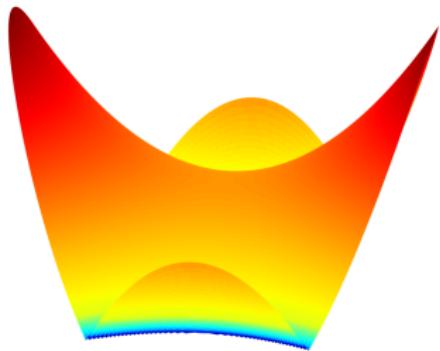
$$\mathbf{y} = \mathcal{A}(\mathbf{M}) + \underbrace{\mathbf{s}}_{\text{outliers}}, \quad \mathcal{A}(\mathbf{M}) = \{\langle \mathbf{A}_i, \mathbf{M} \rangle\}_{i=1}^m$$

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

Dealing with outliers: subgradient methods

Least absolute deviation (LAD):

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



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

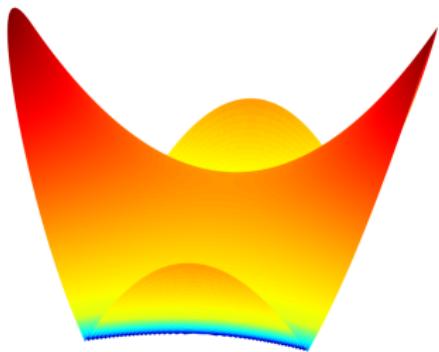
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

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

Dealing with outliers: subgradient methods

Least absolute deviation (LAD):

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



- **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

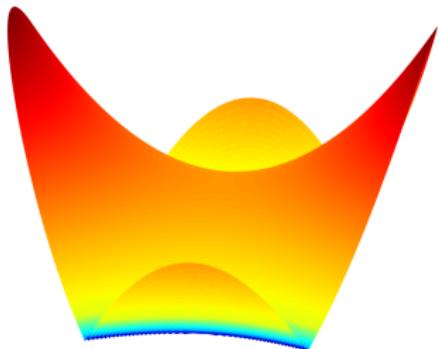
Least absolute deviation (LAD):

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

- 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.

Stepsize schedule

Polyak's stepsize:

$$\eta_t = \frac{f(\mathbf{X}_t \mathbf{Y}_t^\top) - f(\mathbf{M})}{\|\partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1/2}\|_F^2 + \|\partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{X}_t^\top \mathbf{X}_t)^{-1/2}\|_F^2}.$$

- Use the distance concerted with preconditioners.
- Require the knowledge of the optimal value $f(\mathbf{M})$.

Stepsize schedule

Polyak's stepsize:

$$\eta_t = \frac{f(\mathbf{X}_t \mathbf{Y}_t^\top) - f(\mathbf{M})}{\|\partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1/2}\|_F^2 + \|\partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{X}_t^\top \mathbf{X}_t)^{-1/2}\|_F^2}.$$

- Use the distance concerted with preconditioners.
- Require the knowledge of the optimal value $f(\mathbf{M})$.

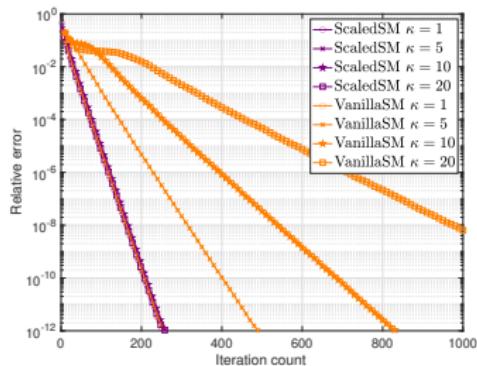
Geometrically decaying stepsize:

$$\eta_t = \frac{\lambda q^t}{\sqrt{\|\partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1/2}\|_F^2 + \|\partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)(\mathbf{X}_t^\top \mathbf{X}_t)^{-1/2}\|_F^2}}$$

- Parameters λ, q need to be tuned.
- Perform similarly as Polyak's stepsize under well-tuned λ, q .

Performance guarantees

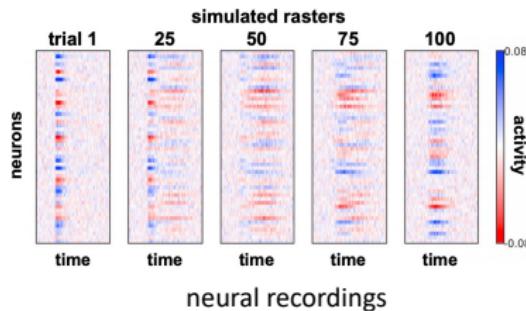
	matrix sensing	quadratic sensing
Subgradient Method (Charisopoulos et al, '19)	$\frac{\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$	$\frac{r\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$
ScaledSM (Tong, Ma, Chi, '20)	$\frac{1}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$	$\frac{r}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$



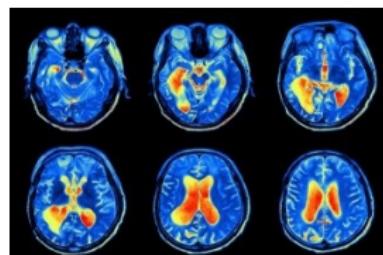
Robustness to both ill-conditioning and adversarial corruptions!

Generalization to tensors

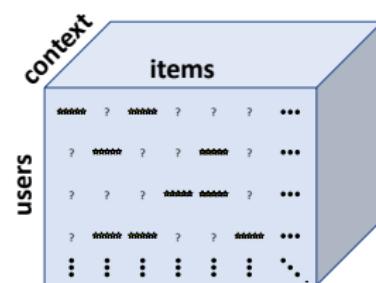
Capturing multi-way interactions by tensors



video surveillance



neuroimaging



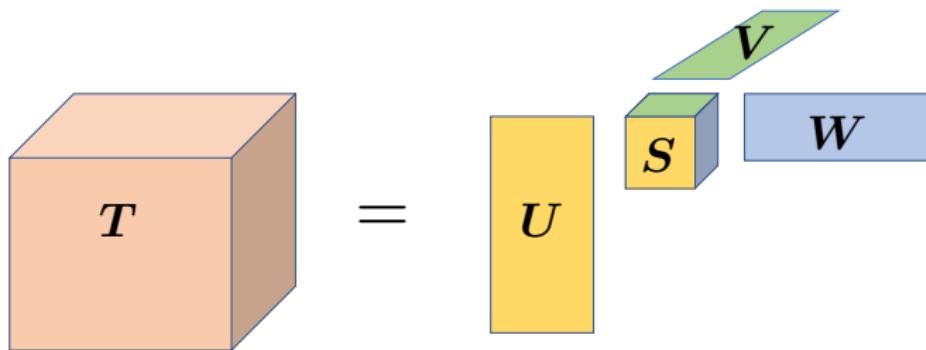
recommendation system

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

Low-rank tensor under Tucker decomposition

Low-rank Tucker decomposition of a tensor:

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



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

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

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:

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

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:

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

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

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

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

where $\breve{\mathbf{U}}_t := (\mathbf{V}_t \otimes \mathbf{W}_t) \mathcal{M}_1(\mathbf{S}_t)^\top$, $\breve{\mathbf{V}}_t := (\mathbf{U}_t \otimes \mathbf{W}_t) \mathcal{M}_2(\mathbf{S}_t)^\top$, and $\breve{\mathbf{W}}_t := (\mathbf{U}_t \otimes \mathbf{V}_t) \mathcal{M}_3(\mathbf{S}_t)^\top$. 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., 2021)

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

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

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

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

ScaledGD for low-rank tensor completion

Theorem (Tong et. al., 2021)

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

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

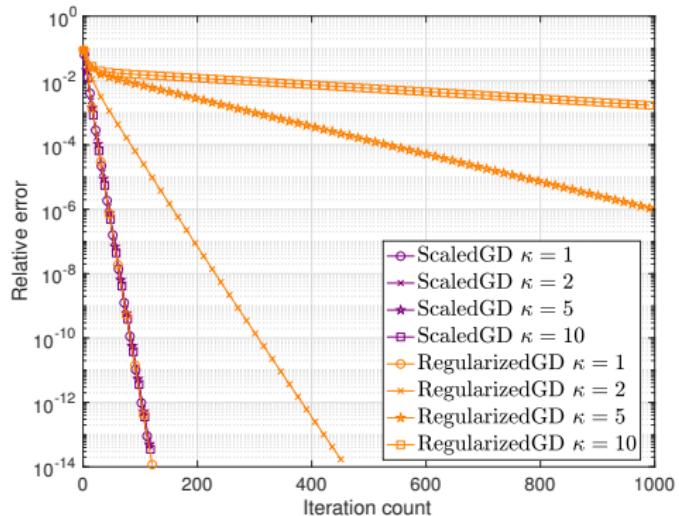
- **Computational:** within $O(\log \frac{1}{\varepsilon})$ iterations;
- **Statistical:** as long as the sample complexity satisfies

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

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

Numerical evidence

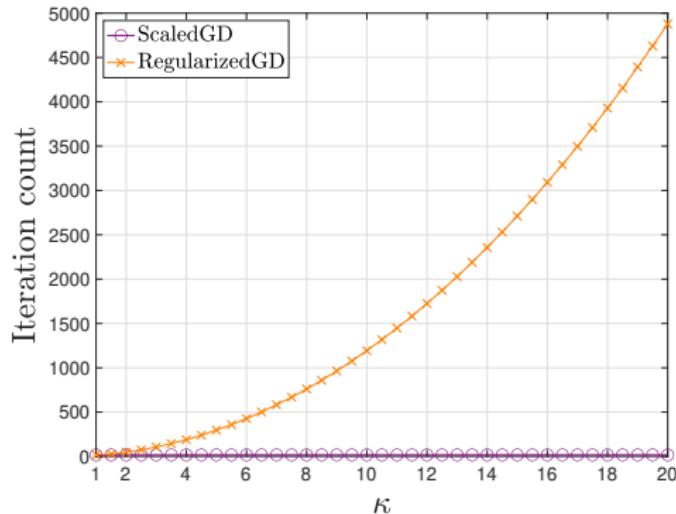
$$\min_{\mathbf{F}=(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{T} \right\|_{\text{F}}^2$$



The benefit of ScaledGD is even more evident for tensors!

Numerical evidence

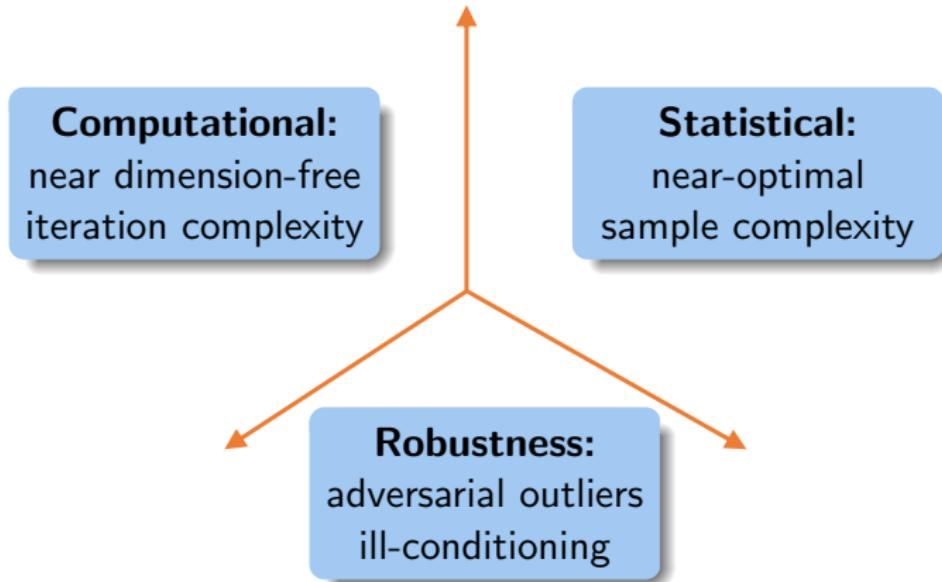
$$\min_{\mathbf{F}=(\mathbf{U}, \mathbf{V}, \mathbf{W}, \mathbf{S})} f(\mathbf{F}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}((\mathbf{U}, \mathbf{V}, \mathbf{W}) \cdot \mathbf{S}) - \mathbf{T} \right\|_{\text{F}}^2$$



The benefit of ScaledGD is even more evident for tensors!

Concluding remarks

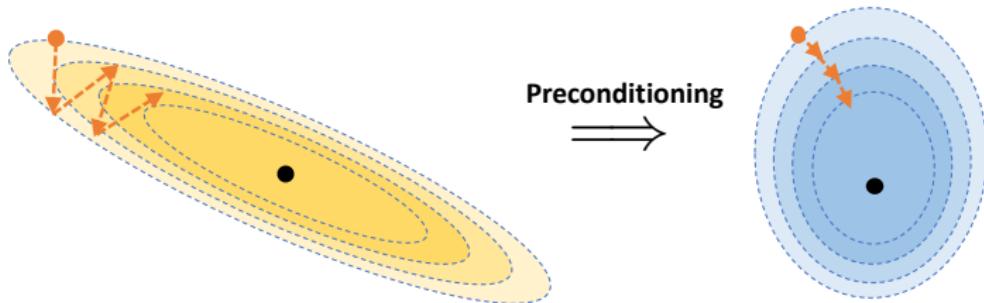
Bridging the theory-practice gap



Nonconvex low-rank matrix and tensor estimation:

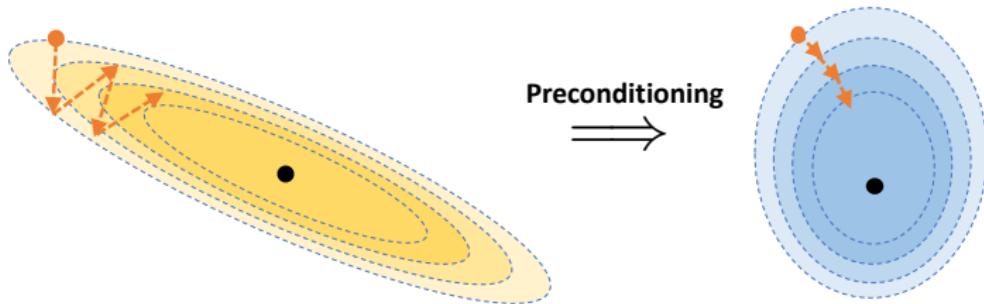
- identification and exploitation of benign geometric properties;
- analyzing iterate trajectories beyond black-box optimization;
- simple variants of GD lead to robust and accelerated convergence.

Preconditioning helps!



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

Preconditioning helps!



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

Promising directions:

- streaming/stochastic variants of ScaledGD;
- applications of ScaledGD to other problems.

Selected References

Overview:

1. Nonconvex Optimization Meets Low-Rank Matrix Factorization: An Overview, Y. Chi, Y. M. Lu and Y. Chen, *IEEE Trans. on Signal Processing*, 2019.
2. Spectral Methods for Data Science: A Statistical Perspective", Y. Chen, Y. Chi, J. Fan and C. Ma, *Foundations and Trends in Machine Learning*, 2021.
3. Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation, Y. Chen and Y. Chi, *IEEE Signal Processing Magazine*, 2018.

Geometry of factored gradient descent:

1. Implicit Regularization for Nonconvex Statistical Estimation, C. Ma, K. Wang, Y. Chi and Y. Chen, *Foundations of Computational Mathematics*, 2020.
2. Beyond Procrustes: Balancing-free Gradient Descent for Asymmetric Low-Rank Matrix Sensing, C. Ma, Y. Li and Y. Chi, *IEEE Trans. on Signal Processing*, 2021.
3. Nonconvex Matrix Factorization from Rank-One Measurements, Y. Li, C. Ma, Y. Chen, and Y. Chi, *IEEE Trans. on Information Theory*, 2020.

Selected References

Robustness to ill-conditioning:

1. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent, T. Tong, C. Ma, and Y. Chi, *Journal of Machine Learning Research*, 2021.
2. Scaling and scalability: Provable nonconvex low-rank tensor estimation from incomplete measurements, T. Tong, C. Ma, A. Prater-Bennette, E. Tripp, and Y. Chi, *arXiv preprint arXiv:2104.14526*, 2021.

Robustness to adversarial outliers:

1. Low-rank matrix recovery with scaled subgradient methods: Fast and robust convergence without the condition number, T. Tong, C. Ma, and Y. Chi, *IEEE Trans. on Signal Processing*, 2021.
2. Non-convex low-rank matrix recovery with arbitrary outliers via median-truncated gradient descent, Y. Li, Y. Chi, H. Zhang and Y. Liang, *Information and Inference: A Journal of the IMA*, 2020.

Thanks!



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

Advertisement: 2022 IEEE SPS Distinguished Lecturer

- nonconvex optimization for statistical estimation
- convex optimization for superresolution
- reinforcement learning
- communication-efficient distributed learning

Please consider inviting me to your local chapter!