

Scalable and Robust Nonconvex Approaches for Low-rank Structure Estimation

Yuejie Chi

Carnegie Mellon University

International Workshop on Intelligent Signal Processing
Zhejiang University, September 2021

Sensing and imaging advances

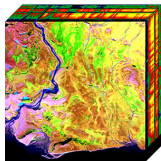
New imaging/sensing modalities allow us to probe the nature in unprecedented manners.



healthcare



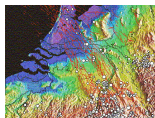
Radio astronomy



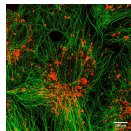
hyperspectral



Internet traffic



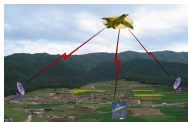
seismic imaging



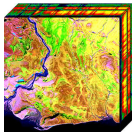
microscopy

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



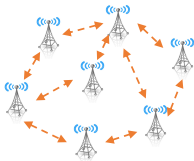
radar imaging



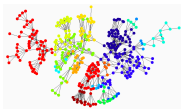
hyperspectral imaging



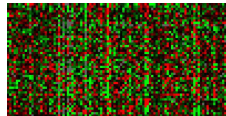
recommendation systems



localization



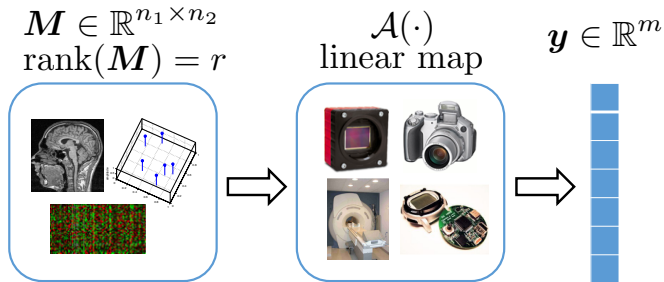
community detection



bioinformatics

Low-rank matrices are redundant representations of latent information

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

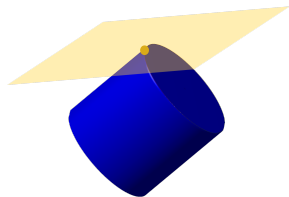
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

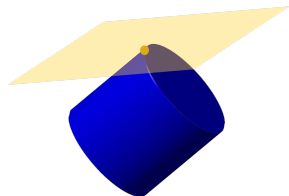
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

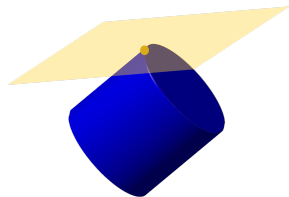
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

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

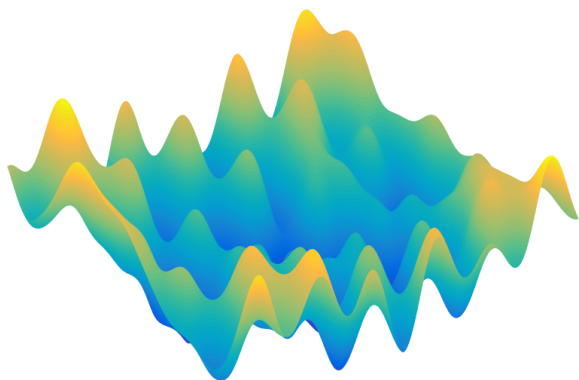
Low-rank matrix factorization

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



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

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

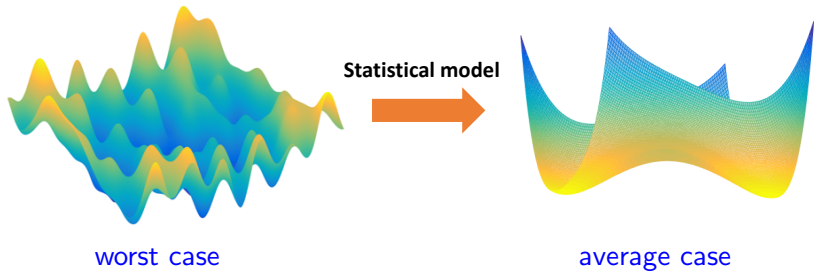
Nonconvex problems are hard (in theory)!



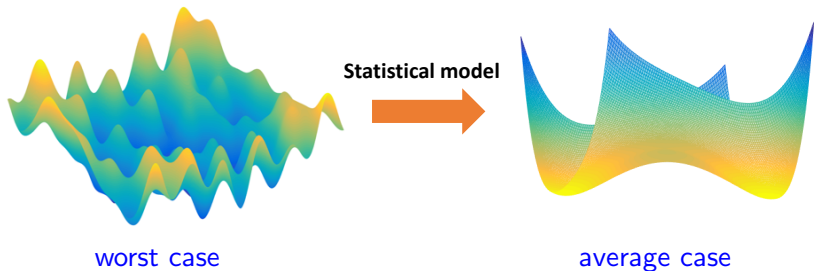
“...in fact, the great watershed in optimization isn’t between linearity and nonlinearity, but convexity and nonconvexity.”

R. T. Rockafellar, in SIAM Review, 1993

Statistics meets optimization

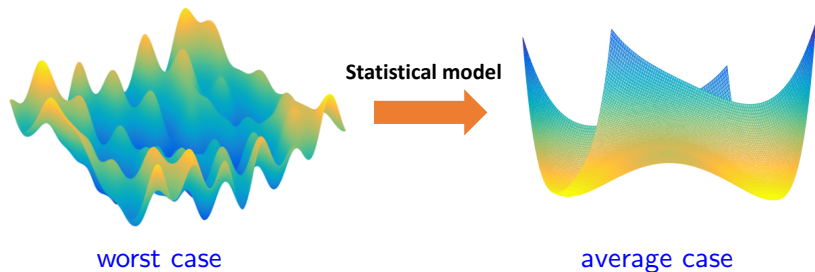


Statistics meets optimization



Simple algorithms can be efficient for nonconvex learning!

Statistics meets optimization



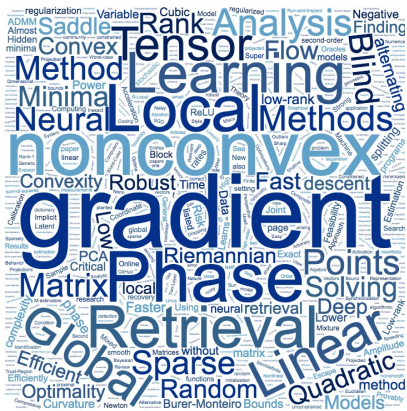
Simple algorithms can be efficient for nonconvex learning!

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, ...

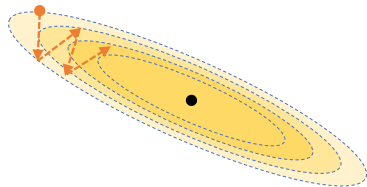
Blind deconvolution/demixing: Li et al. '16, Lee et al. '16, Cambareri, Jacques '16, Ling, Strohmer '16, Huang, Hand '16, Ma et al. '17, Zhang et al. '18, Li, Bresler '18, Dong, Shi '18, Shi, Chi '19, Qu et al. '19...

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

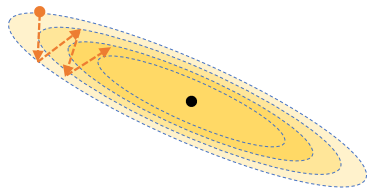


Vanilla GD:

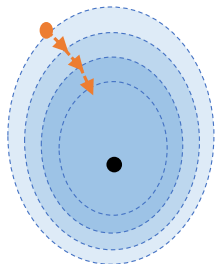
$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

☹ **Slows down with ill-conditioning.**

Acceleration via preconditioning



Preconditioning



Vanilla GD:

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

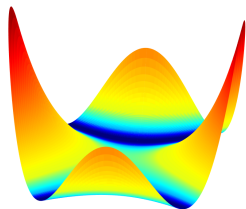
☹ **Slows down with ill-conditioning.**

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

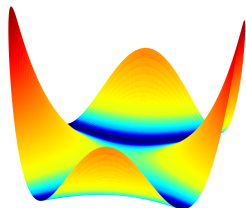


Least squares:

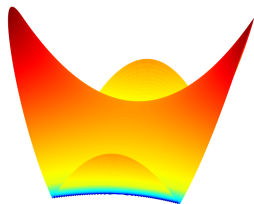
$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_2^2$$

☹ **Sensitive to outliers.**

Robustness via nonsmooth optimization



Nonsmooth



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.**

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 gradient algorithms that are 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 gradient algorithms that are oblivious to the presence of outliers?

Generalization to tensors:

Can we generalize to higher-dimensional objects?

A bit preliminaries of optimization

Unconstrained optimization

Consider an unconstrained optimization problem

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x})$$

Definition (first-order critical points)

A first-order critical point of f satisfies

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

Unconstrained optimization

Consider an unconstrained optimization problem

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x})$$

Definition (second-order critical points)

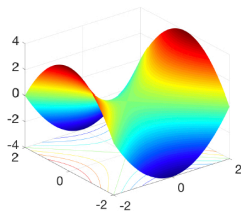
A second-order critical point \boldsymbol{x} satisfies

$$\nabla f(\boldsymbol{x}) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\boldsymbol{x}) \succeq \mathbf{0}$$

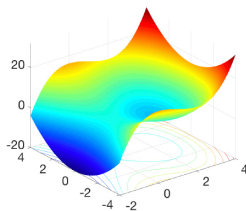
Several types of critical points

For any first-order critical point \mathbf{x} :

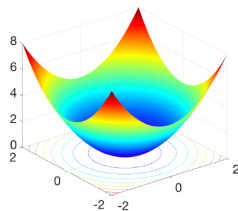
- $\nabla^2 f(\mathbf{x}) \prec \mathbf{0}$ \rightarrow local maximum
- $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ \rightarrow local minimum
- $\lambda_{\min}(\nabla^2 f(\mathbf{x})) < 0$ \rightarrow *strict* saddle point



(a) strict saddle



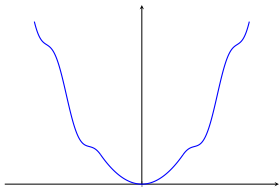
(b) local minimum



(c) global minimum

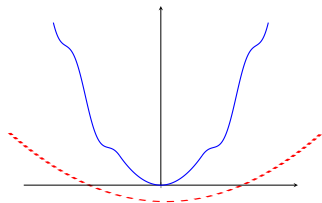
figure credit: Li et al. '16

Gradient descent theory



Two standard conditions that enable geometric convergence of GD

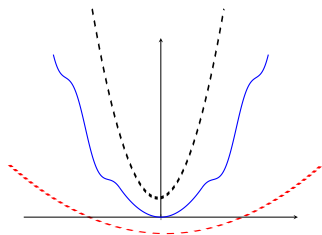
Gradient descent theory



Two standard conditions that enable geometric convergence of GD

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

Gradient descent theory



Two standard conditions that enable geometric convergence of GD

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

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

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$) with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$) with $\eta = 1/\beta$ obeys

$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

- Condition number β/α determines rate of convergence

Gradient descent theory revisited

f is said to be α -strongly convex and β -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

ℓ_2 error contraction: GD ($\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$) with $\eta = 1/\beta$ obeys

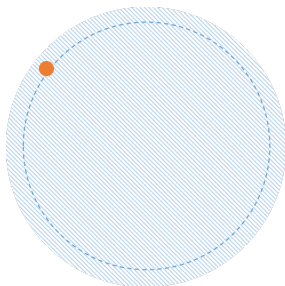
$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

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

Gradient descent theory revisited

$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

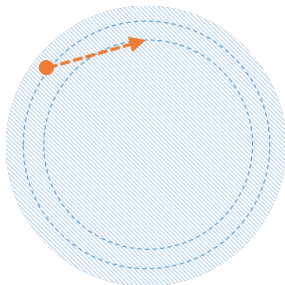
- region of local strong convexity + smoothness



Gradient descent theory revisited

$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

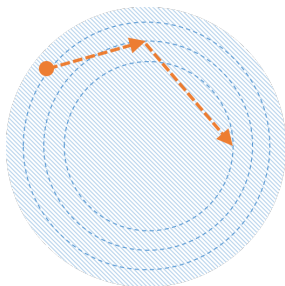
- region of local strong convexity + smoothness



Gradient descent theory revisited

$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

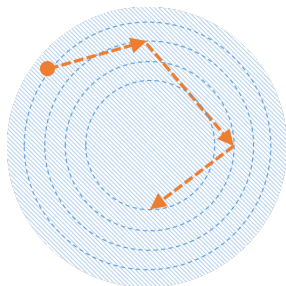
- region of local strong convexity + smoothness



Gradient descent theory revisited

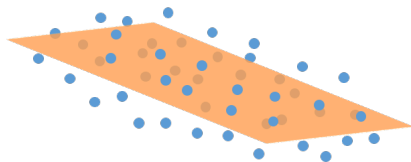
$$\|\mathbf{x}_{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}_t - \mathbf{x}_{\text{opt}}\|_2$$

- region of local strong convexity + smoothness



Warm-up: understanding the geometry of PCA

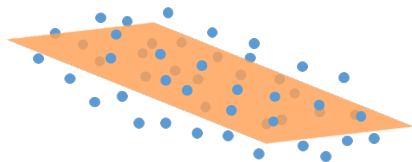
Revisiting PCA



Given $\mathbf{M} \succeq \mathbf{0} \in \mathbb{R}^{n \times n}$ (not necessarily low-rank), find its best rank- r approximation:

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

Revisiting PCA



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} ($\lambda_1 \geq \dots \geq \lambda_n$), 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 $Z = \mathbf{X}\mathbf{X}^\top$ with $\mathbf{X} \in \mathbb{R}^{n \times r}$, then it leads to a nonconvex problem:

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

An optimization viewpoint

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

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

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

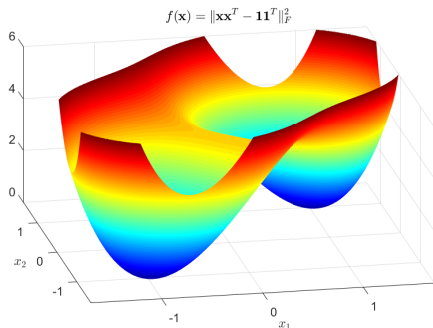
- global minima;
- strict saddle points, from which there exist directions to strictly decrease $f(\mathbf{X})$.

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 strong convexity and local linear convergence

- The global minimizers: $\mathbf{x}_{\text{opt}} = \pm\sqrt{\lambda_1}\mathbf{u}_1$
- For all \mathbf{x} obeying $\underbrace{\|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}}_{\text{basin of attraction}}$, one has

$$0.25(\lambda_1 - \lambda_2)\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq 4.5\lambda_1\mathbf{I}_n$$

Local strong convexity and local linear convergence

- The global minimizers: $\mathbf{x}_{\text{opt}} = \pm\sqrt{\lambda_1}\mathbf{u}_1$
- For all \mathbf{x} obeying $\underbrace{\|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}}_{\text{basin of attraction}}$, one has

$$0.25(\lambda_1 - \lambda_2)\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq 4.5\lambda_1\mathbf{I}_n$$

ℓ_2 **error contraction:** The GD iterates obey

$$\|\mathbf{x}_t - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \left(1 - \frac{\lambda_1 - \lambda_2}{18\lambda_1}\right)^t \|\mathbf{x}_0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2, \quad t \geq 0,$$

as long as $\|\mathbf{x}_0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$

Extension to the low-rank case

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

Cannot be uniquely determined \mathbf{X} up to orthogonal transform.

Extension to the low-rank case

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

Cannot be uniquely determined \mathbf{X} up to orthogonal transform.

- A modified distance metric:

$$\text{dist}^2(\mathbf{X}, \mathbf{X}_\star) = \min_{\mathbf{H} \in \mathcal{O}^{r \times r}} \|\mathbf{X}\mathbf{H} - \mathbf{X}_\star\|_{\mathbb{F}}^2.$$

- Optimal alignment matrix (the Procrustes problem):

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

Restricted strong convexity

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

f satisfies α -**restricted strong convexity** and β -smoothness:

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

where $\beta \asymp \lambda_1$ and $\alpha \asymp \lambda_r$.

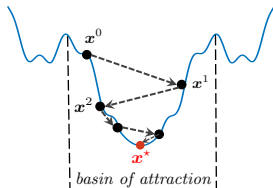
ℓ_2 error contraction: The GD iterates obey

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

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

Two vignettes

Two-stage approach:



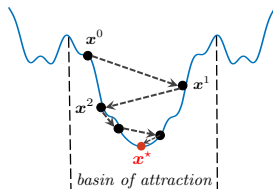
smart initialization

+

local refinement

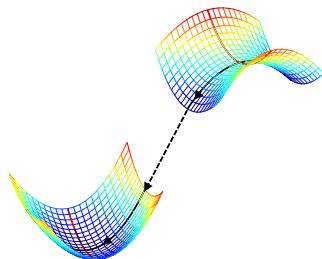
Two vignettes

Two-stage approach:



smart initialization
+
local refinement

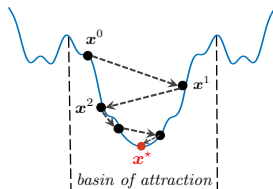
Global landscape:



benign landscape
+
saddle-point escaping

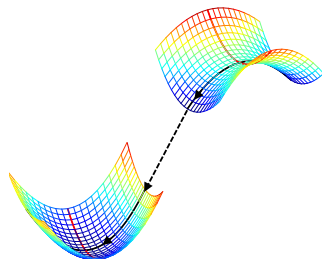
Two vignettes

Two-stage approach:



smart initialization
+
local refinement

Global landscape:



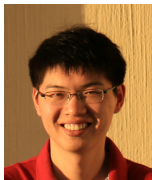
benign landscape
+
saddle-point escaping

This tutorial will mostly focus on the two-stage approach.

*Geometry and implicit regularization
in nonconvex low-rank matrix estimation*



Yuxin Chen
Princeton



Cong Ma
Chicago



Kaizheng Wang
Columbia

Low-rank matrix completion: dealing with missing data



Given partial samples of a *low-rank* matrix $M = X_* X_*^T \in \mathbb{R}^{n \times n}$ in an index set Ω , fill in missing entries.

A natural least-squares formulation

given: $\mathcal{P}_\Omega(\mathbf{M})$

↓

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \left\| \mathcal{P}_\Omega(\mathbf{X}\mathbf{X}^\top - \mathbf{M}) \right\|_F^2$$

A natural least-squares formulation

given: $\mathcal{P}_\Omega(\mathbf{M})$

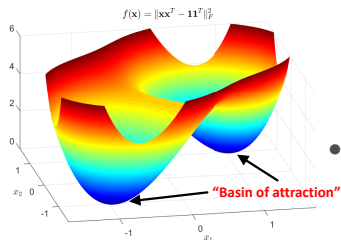
↓

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \left\| \mathcal{P}_\Omega(\mathbf{X}\mathbf{X}^\top - \mathbf{M}) \right\|_F^2$$

- **Bernoulli sampling:** Assume every entry is observed i.i.d. with $0 < p \leq 1$:

$$\mathbb{E}[f(\mathbf{X})] = p \left\| \mathbf{X}\mathbf{X}^\top - \mathbf{M} \right\|_F^2.$$

Two-stage approach



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

$$\mathbf{X}_0 = \text{SVD}_r(\mathcal{P}_\Omega(\mathbf{M}))$$

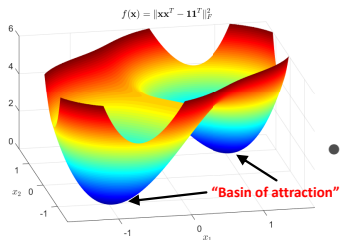
- **Gradient iterations:**

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

for $t = 0, 1, \dots$

“Spectral methods for data science: A statistical perspective”, Y. Chen, Y. Chi, J. Fan, C. Ma, FnT ML, 2021.

Two-stage approach



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

$$\mathbf{X}_0 = \text{SVD}_r(\mathcal{P}_\Omega(\mathbf{M}))$$

- **Gradient iterations:**

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

for $t = 0, 1, \dots$

Question: Does vanilla GD still work with partial observations?

“Spectral methods for data science: A statistical perspective”, Y. Chen, Y. Chi, J. Fan, C. Ma, FnT ML, 2021.

Incoherence

Which is easier to complete?

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}$$

vs.

$$\underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}$$

Incoherence

Which is easier to complete?

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}}_{\text{hard}} \quad \text{vs.} \quad \underbrace{\begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & & \\ 1 & 1 & 1 & \cdots & 1 \end{bmatrix}}_{\text{easy}}$$

Definition (Incoherence for matrix completion)

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

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

Incoherence

Which is easier to complete?

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

Definition (Incoherence for matrix completion)

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

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

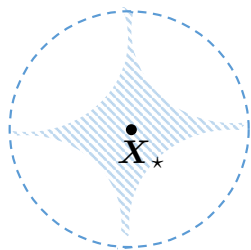
Which region has benign geometry?

Finite-sample level ($p \asymp \frac{\text{polylog}n}{n}$) :

$f(\mathbf{X})$ restricted strongly convex and smooth


along descent direction \mathbf{V} **only when \mathbf{X} is incoherent:**

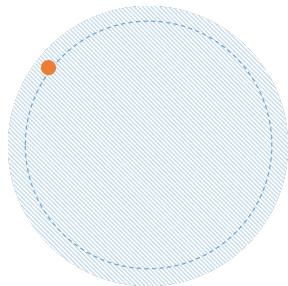
$$\|\mathbf{X}\mathbf{H}\mathbf{X} - \mathbf{X}_\star\|_{2,\infty} \ll \|\mathbf{X}_\star\|_{2,\infty}$$



region of local strong convexity + smoothness

Vanilla gradient descent is at risk

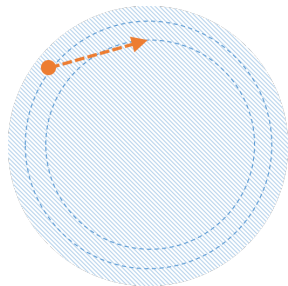
 region of local strong convexity + smoothness



GD on the pop. loss

Vanilla gradient descent is at risk

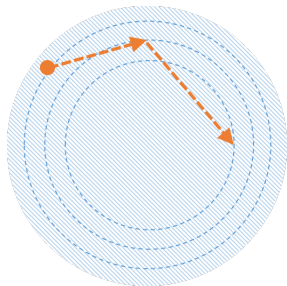
● region of local strong convexity + smoothness



GD on the pop. loss

Vanilla gradient descent is at risk

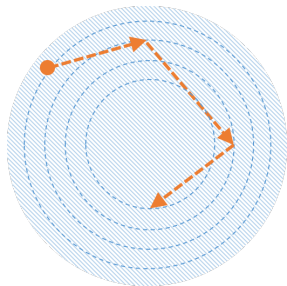
● region of local strong convexity + smoothness



GD on the pop. loss

Vanilla gradient descent is at risk

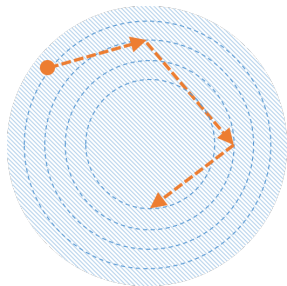
● region of local strong convexity + smoothness



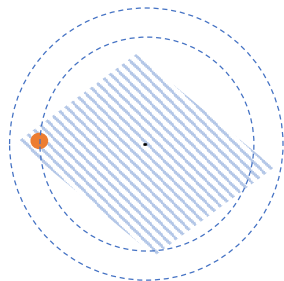
GD on the pop. loss

Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



GD on the pop. loss

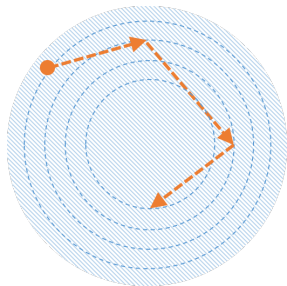


GD on the emp. loss

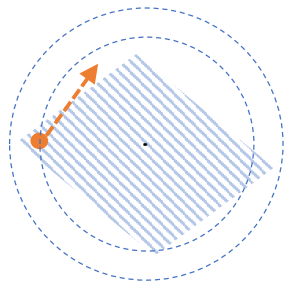
- Generic optimization theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



GD on the pop. loss

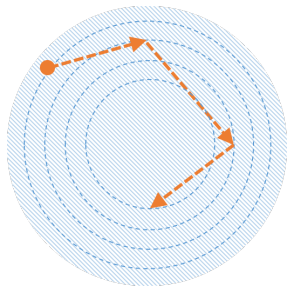


GD on the emp. loss

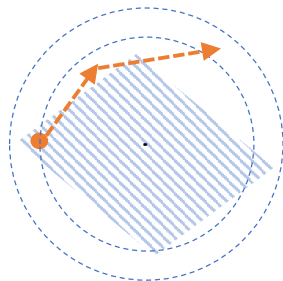
- Generic optimization theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



GD on the pop. loss

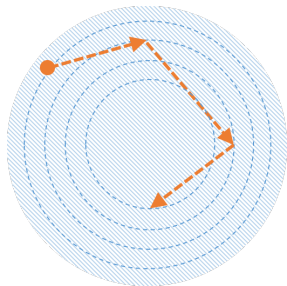


GD on the emp. loss

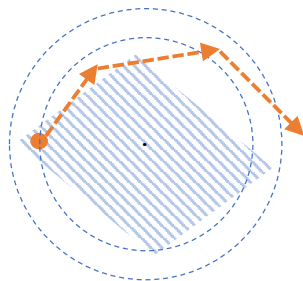
- Generic optimization theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



GD on the pop. loss

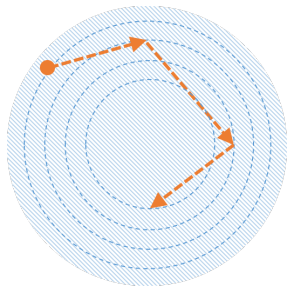


GD on the emp. loss

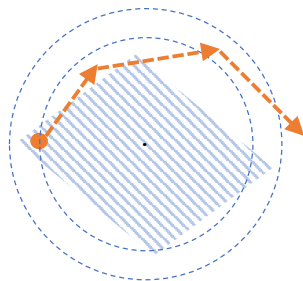
- Generic optimization theory only ensures that iterates remain in ℓ_2 ball but not incoherence region

Vanilla gradient descent is at risk

● region of local strong convexity + smoothness



GD on the pop. loss

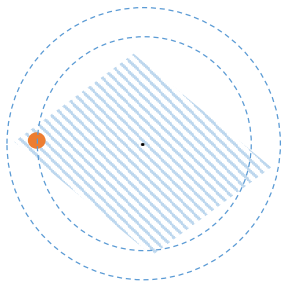


GD on the emp. loss

- Generic optimization theory only ensures that iterates remain in ℓ_2 ball but not incoherence region
- Existing algorithms enforce regularization, or apply sample splitting to promote incoherence

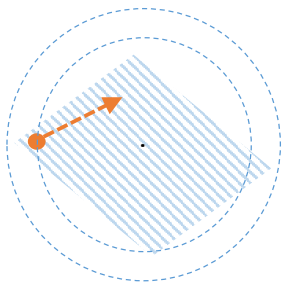
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



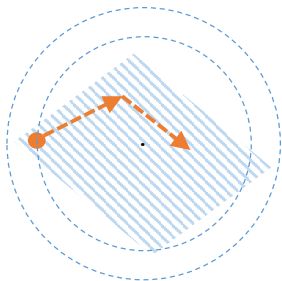
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



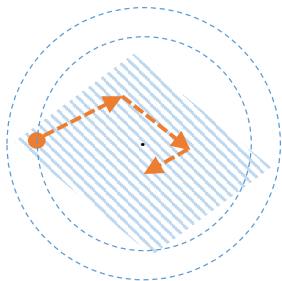
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



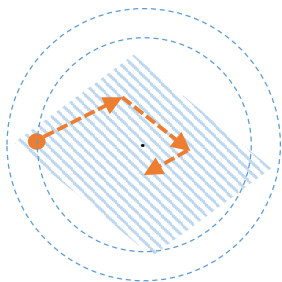
Our findings: GD is implicitly regularized

- region of local strong convexity + smoothness



Our findings: GD is implicitly regularized

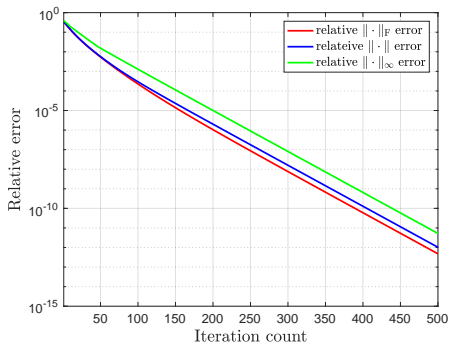
- region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**
even without regularization

Matrix completion via vanilla GD

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



Vanilla GD converges fast without regularization!

Theoretical guarantees - noise-free case

Theorem (Ma, Wang, Chi, Chen, FoCM 2020)

Suppose $M = X_* X_*^\top$ is rank- r , μ -incoherent and has a condition number $\kappa = \sigma_{\max}(M)/\sigma_{\min}(M)$. Vanilla GD (with spectral initialization) achieves

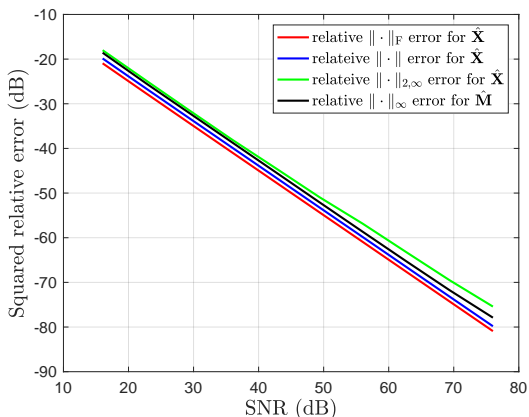
$$\|X_t X_t^\top - M\|_F \leq \varepsilon \cdot \sigma_{\min}(M)$$

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

$$n^2 p \gtrsim nr^3 \text{poly}(\mu, \kappa, \log n).$$

First convergence guarantee of vanilla GD for matrix completion

Noisy matrix completion via vanilla GD



Near-optimal entrywise error control:

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

The phenomenon is quite general

| | Prior theory | | Our theory | |
|----------------------------|-------------------|--|---------------------------|---|
| | sample complexity | iteration complexity | sample complexity | iteration complexity |
| Phase retrieval | $n \log n$ | $n \log \left(\frac{1}{\epsilon}\right)$ | $n \log n$ | $\log n \log \left(\frac{1}{\epsilon}\right)$ |
| Quadratic sensing | $nr^6 \log^2 n$ | $n^4 r^2 \log \left(\frac{1}{\epsilon}\right)$ | $nr^4 \log n$ | $r^2 \log \left(\frac{1}{\epsilon}\right)$ |
| Matrix completion | n/a | n/a | $nr^3 \text{poly} \log n$ | $\log \left(\frac{1}{\epsilon}\right)$ |
| Blind deconvolution | n/a | n/a | $K \text{poly} \log m$ | $\log \left(\frac{1}{\epsilon}\right)$ |



Huge computational savings!

An aside: minimax stability of nuclear norm minimization

convex



nonconvex

$$\min_{\mathbf{Z} \in \mathbb{R}^{n \times n}} \sum_{(i,j) \in \Omega} (Z_{i,j} - M_{i,j})^2 + \lambda \|\mathbf{Z}\|_*$$

$$\min_{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}} \sum_{(i,j) \in \Omega} \left[(\mathbf{X}\mathbf{Y}^\top)_{i,j} - M_{i,j} \right]^2 + \frac{\lambda}{2} \|\mathbf{X}\|_F^2 + \frac{\lambda}{2} \|\mathbf{Y}\|_F^2$$

Theorem (Chen, Chi, Fan, Ma, Yan '19)

With high prob., any minimizer $\widehat{\mathbf{M}}_{\text{cvx}}$ of convex program is nearly rank- r and is minimax near-optimal:

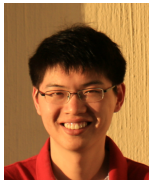
$$\|\widehat{\mathbf{M}}_{\text{cvx}} - \mathbf{M}\|_F \lesssim \sigma \sqrt{\frac{n}{p}}, \quad \|\widehat{\mathbf{M}}_{\text{cvx}} - \mathbf{M}\|_\infty \lesssim \sigma \sqrt{\frac{n \log n}{p}} \cdot \frac{1}{n}$$

Noisy Matrix Completion: Understanding Statistical Guarantees for Convex Relaxation via Nonconvex Optimization, SIAM Journal on Optimization.

Accelerating ill-conditioned matrix estimation



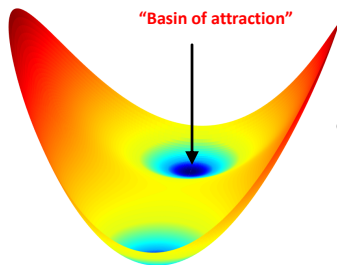
Tian Tong
CMU



Cong Ma
Chicago

The asymmetric case: 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 $t = 0, 1, \dots$

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

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

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

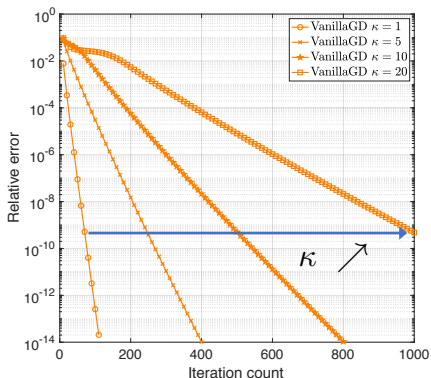
$$m \gtrsim (n_1 + n_2)r^2\kappa^2.$$

Similar results hold for many low-rank problems.

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

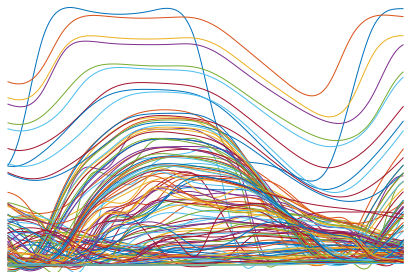
Convergence slows down for ill-conditioned matrices

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

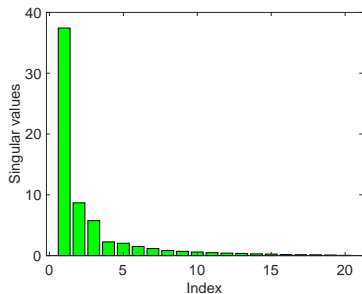


Vanilla GD converges in $O(\kappa \log \frac{1}{\epsilon})$ iterations.

Condition number can be large



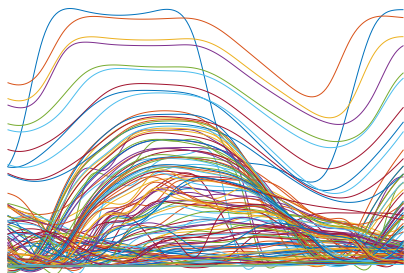
chlorine concentration levels
120 junctions, 180 time slots



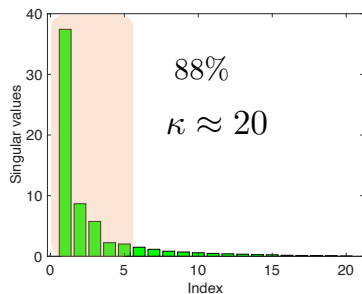
power-law spectrum

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

Condition number can be large



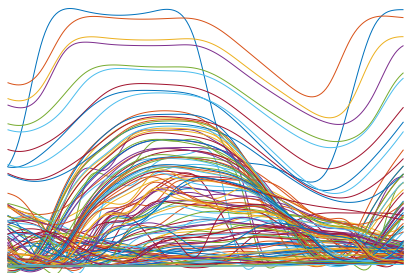
chlorine concentration levels
120 junctions, 180 time slots



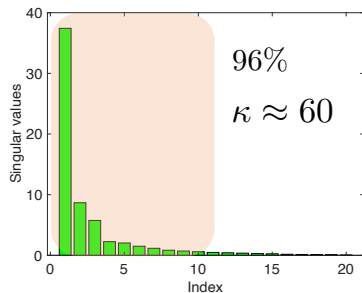
rank-5 approximation

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

Condition number can be large



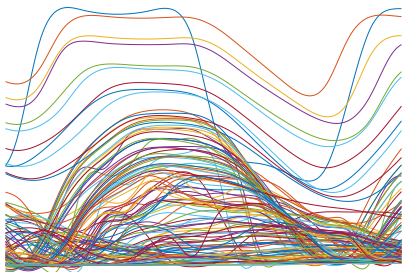
chlorine concentration levels
120 junctions, 180 time slots



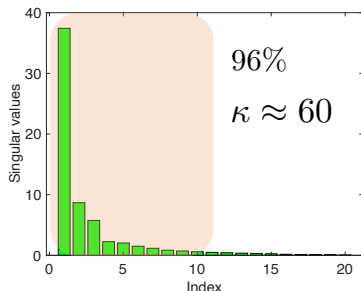
rank-10 approximation

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

Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots



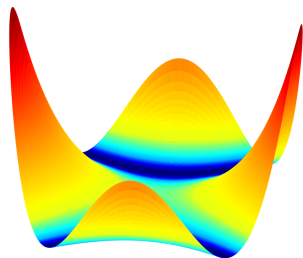
rank-10 approximation

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

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

A new algorithm: 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:**

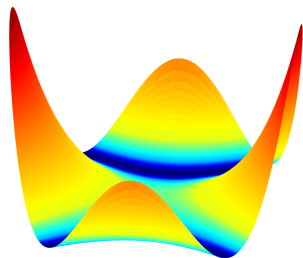
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

A new algorithm: 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:**

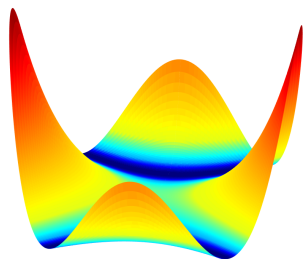
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

A new algorithm: 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:**

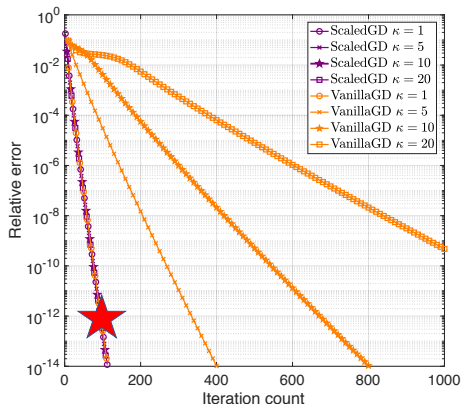
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

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

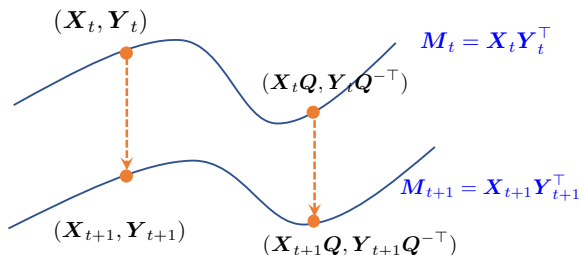
ScaledGD for low-rank matrix completion



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

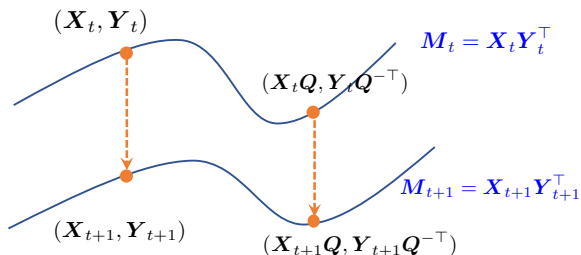
A closer look at ScaledGD

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



A closer look at ScaledGD

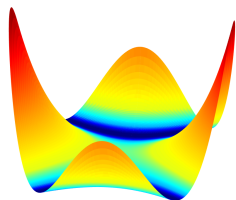
Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



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, 2020)

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

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

- **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, 2020)

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

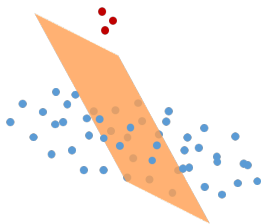
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **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!

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.

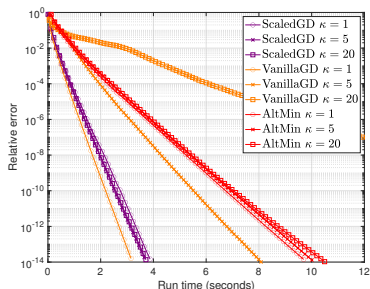
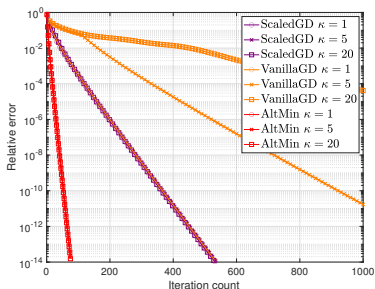
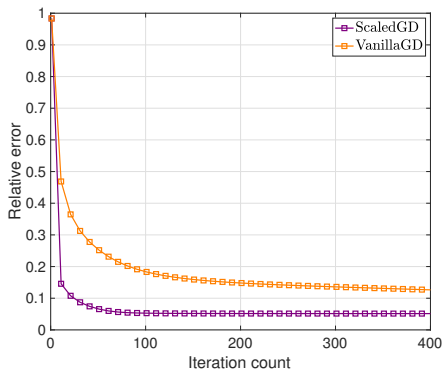


Figure: Run time for matrix completion with $n = 1000$, $p = 0.2$, $r = 50$.

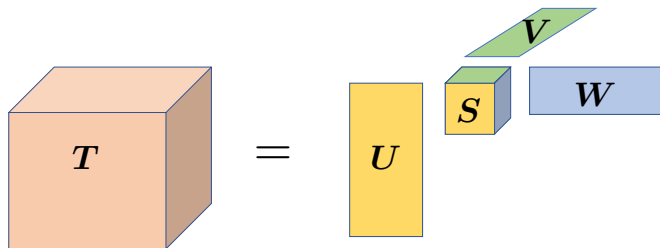
Numerical stability

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



Generalization to tensors

Low-rank tensor under Tucker decomposition



Low-rank Tucker decomposition of a tensor:

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

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 $S \in \mathbb{R}^{r_1 \times r_2 \times r_3}$.

Applications in fMRI imaging, recommendation systems, etc...

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., 2021)

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

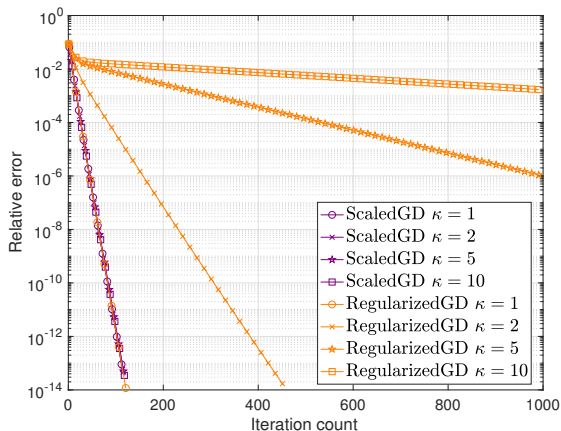
$$\|(\mathbf{U}_t, \mathbf{V}_t, \mathbf{W}_t) \cdot \mathbf{S}_t - \mathbf{T}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{T})$$

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

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

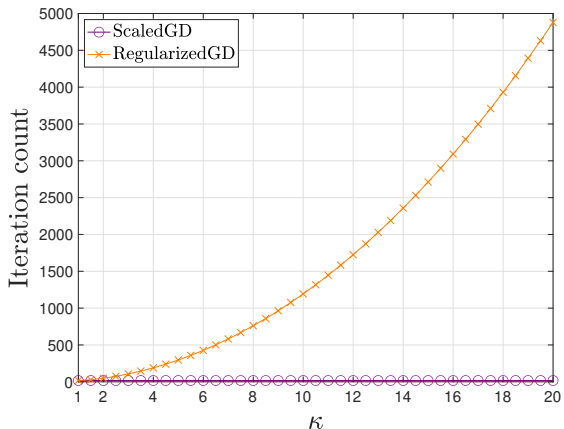
First provable linear convergence at a near-optimal sample complexity for low-Tucker-rank tensor completion!

Numerical evidence



The benefit of ScaledGD is even more evident for tensors!

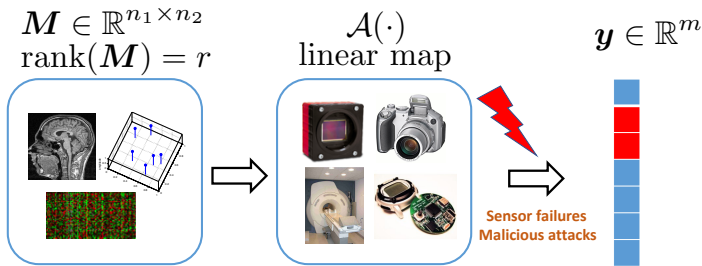
Numerical evidence



The benefit of ScaledGD is even more evident for tensors!

Robustness to outliers and corruptions?

Outlier-corrupted low-rank matrix sensing



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

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

Existing approaches fail

- **Spectral initialization would fail:**
 $\mathbf{X}_0 \leftarrow$ top- r SVD of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{A}_i$$



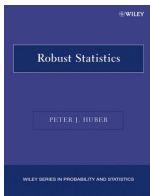
- **Gradient iterations would fail:**

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\eta}{m} \sum_{i=1}^m \nabla l_i(y_i; \mathbf{X}_t)$$

for $t = 0, 1, \dots$

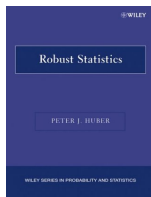
Even a single outlier can fail the algorithm!

Median-truncated gradient descent



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

Median-truncated gradient descent

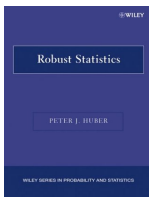


Key idea: “median-truncation” — discard samples *adaptively* based on how large sample gradients / values deviate from median

- **Robustify spectral initialization:** $X_0 \leftarrow$ top- r SVD of

$$Y = \frac{1}{m} \sum_{i: |y_i| \lesssim \text{median}(|y_i|)} y_i A_i$$

Median-truncated gradient descent



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

- **Robustify spectral initialization:** $\mathbf{X}_0 \leftarrow$ top- r SVD of

$$\mathbf{Y} = \frac{1}{m} \sum_{i: |y_i| \lesssim \text{median}(|y_i|)} y_i \mathbf{A}_i$$

- **Robustify gradient descent:**

$$\mathbf{X}_{t+1} = \mathbf{X}_t - \frac{\eta}{m} \sum_{i: |r_t^i| \lesssim \text{median}(|r_t^i|)} \nabla \ell_i(y_i; \mathbf{X}_t), \quad t = 0, 1, \dots$$

where $r_t^i := |y_i - \langle \mathbf{A}_i, \mathbf{X}_t \mathbf{X}_t^\top \rangle|$ is the size of the gradient.

Theoretical guarantees

Theorem (Li, Chi, Zhang, and Liang, IMIAI 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, median-truncated GD (with robust spectral initialization) achieves

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

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

$$m \gtrsim nr^2 \text{poly}(\kappa, \log n);$$

- **Robustness:** *and the fraction of outliers*

$$\alpha \lesssim 1/\sqrt{r}.$$

Theoretical guarantees

Theorem (Li, Chi, Zhang, and Liang, IMIAI 2020)

For low-rank matrix sensing with i.i.d. Gaussian design, median-truncated GD (with robust spectral initialization) achieves

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

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

$$m \gtrsim nr^2 \text{poly}(\kappa, \log n);$$

- **Robustness:** and the fraction of outliers

$$\alpha \lesssim 1/\sqrt{r}.$$

Median-truncated GD adds robustness to GD *obliviously*.

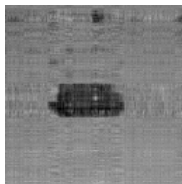
Numerical example

Low-rank matrix sensing:

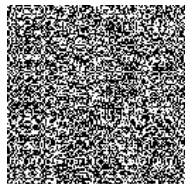
$$y_i = \langle \mathbf{A}_i, \mathbf{M} \rangle + s_i, \quad i = 1, \dots, m$$



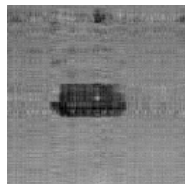
Ground truth



GD
no outliers



GD
1% outliers



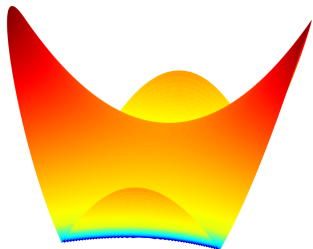
median-TGD
1% outliers

Median-truncated GD achieves similar performance as if performing GD on the clean data.

Dealing with outliers: subgradient methods

Least absolute deviation (LAD): (Charisopoulos et.al.'19; Li et al'18)

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



Subgradient iterations:

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

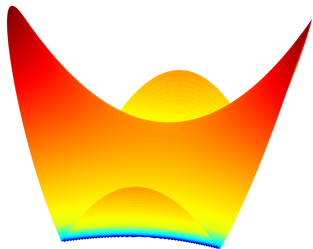
$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t)$$

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

Dealing with outliers: scaled subgradient methods

Least absolute deviation (LAD):

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



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:

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

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

Stepsize schedule

Polyak's stepsize:

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

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

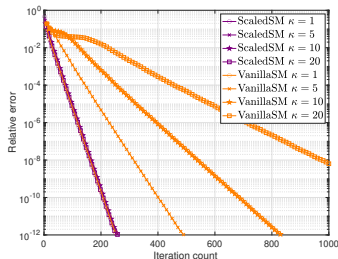
Geometrically decaying stepsize:

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

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

Performance guarantees

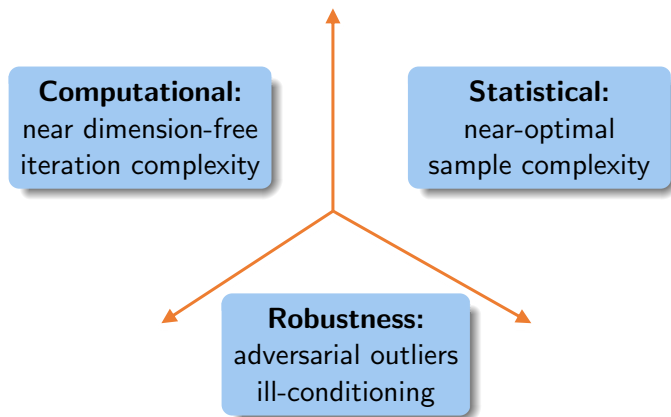
| | matrix sensing | quadratic sensing |
|--|--|---|
| Subgradient Method (Charisopoulos et al, '19) | $\frac{\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ | $\frac{r\kappa}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$ |
| ScaledSM (Tong, Ma, Chi, '20) | $\frac{1}{(1-2\alpha)^2} \log \frac{1}{\epsilon}$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



Robustness to both ill-conditioning and adversarial corruptions!

Concluding remarks

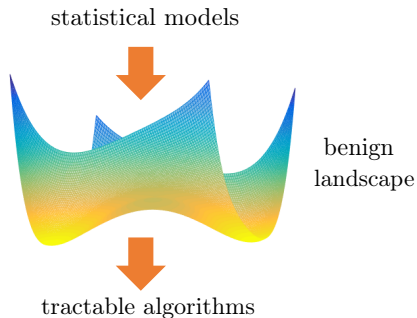
Bridging the theory-practice gap



Nonconvex low-rank matrix 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.

Statistical thinking + Optimization efficiency



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

A growing list of “benign” nonconvex problems

- phase retrieval
- matrix sensing
- matrix completion
- blind deconvolution / self-calibration
- dictionary learning
- tensor decomposition / completion
- robust PCA
- mixed linear regression
- learning one-layer neural networks
- ...

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/>