

Nonconvex Low-Rank Matrix Estimation: Geometry, Robustness, and Acceleration

Yuejie Chi

Carnegie Mellon University

SIAM Conference on Imaging Science
July 9, 2020

Acknowledgements

Our research is supported by National Science Foundation, Office of Naval Research and Army Research Office.



Sensing and imaging advances

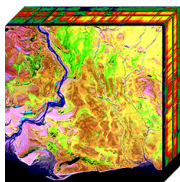
New imaging/sensing modalities allow us to probe the nature in unprecedented manners.



healthcare



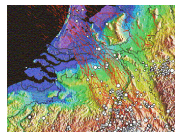
Radio astronomy



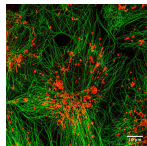
hyperspectral



Internet traffic



seismic imaging



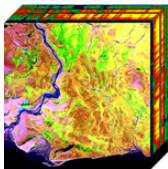
microscopy

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

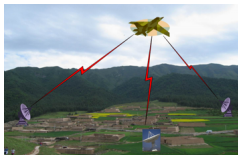
Low-rank matrices in imaging science

Why low-rank images?

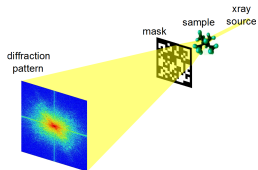
- *redundant representations of latent information;*
- *a small number of sources of interest;*
- *“lifting” of indirect correlation measurements.*



hyperspectral imaging



radar imaging

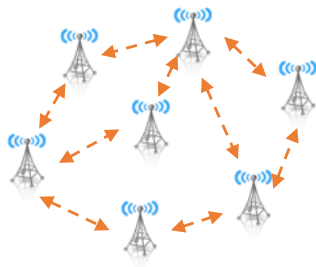


optical imaging

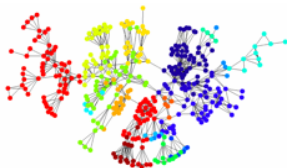
Beyond imaging science



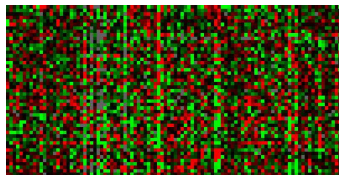
recommendation systems



localization

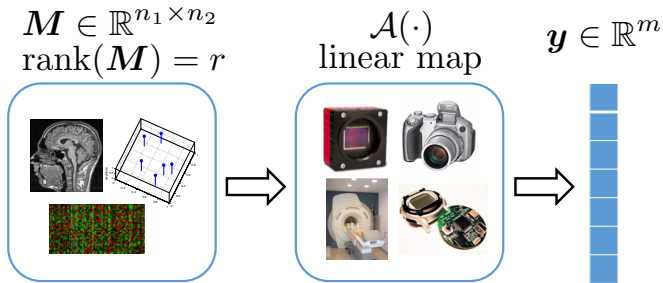


community detection



bioinformatics

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

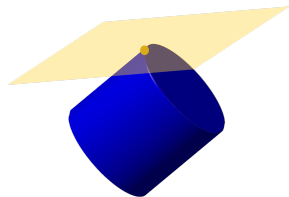
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

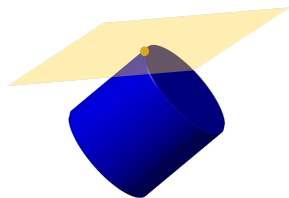
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

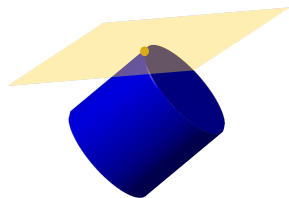
$$\min_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \text{rank}(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{y} \approx \mathcal{A}(\mathbf{Z})$$

↓ 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

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

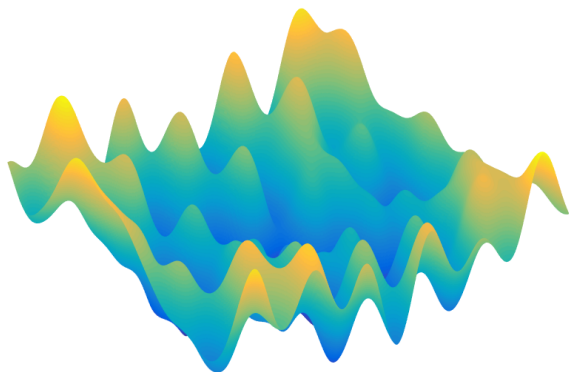
Low-rank matrix factorization

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



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

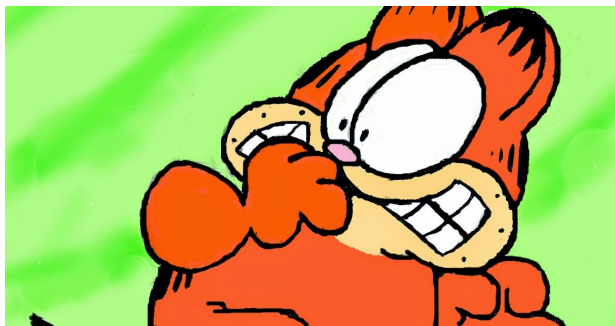
Nonconvex problems are hard (in theory)!



“...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

This talk: geometry, robustness, acceleration

Optimization geometry:

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

Robustness to adversarial outliers:

Can we design provably robust gradient algorithms that are oblivious to the presence of outliers?

Acceleration for ill-conditioned matrix estimation:

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

Geometry and implicit regularization in nonconvex low-rank matrix estimation



Yuxin Chen
Princeton



Cong Ma
Princeton

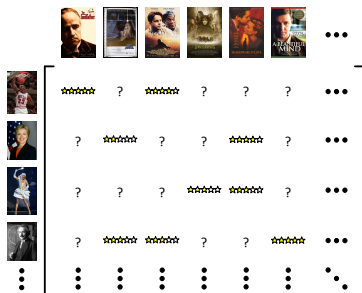


Kaizheng Wang
Princeton



Yuanxin Li
CMU

Low-rank matrix completion: dealing with missing data



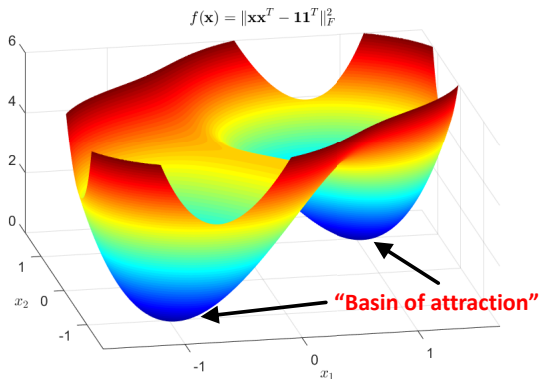
Given partial samples of a *low-rank* matrix $M = X_* X_*^\top \in \mathbb{R}^{n \times n}$ in an index set Ω , fill in missing entries.

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{2} \left\| \mathcal{P}_\Omega(\mathbf{X}\mathbf{X}^\top - M) \right\|_F^2$$

What might the loss function look like?

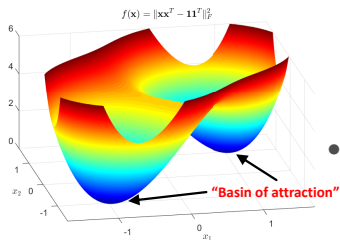
Full observation = PCA: $f(\mathbf{X}) = \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_F^2$.

$f(\mathbf{X})$ restricted strongly convex and smooth
along descent direction \mathbf{V} when \mathbf{X} is close to \mathbf{X}_* .



Parameter recovery via gradient descent (GD)

a two-step recovery strategy:



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

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

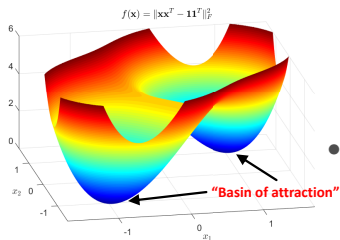
- **Gradient iterations:**

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

for $t = 0, 1, \dots$

Parameter recovery via gradient descent (GD)

a two-step recovery strategy:



- **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?

Which region has benign geometry?

Finite-sample level ($p \asymp \frac{\text{polylog}n}{n}$) : *assume every entry is observed i.i.d. with probability $0 < p \leq 1$.*

Question: which matrix is easier to complete?

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

Which region has benign geometry?

Finite-sample level ($p \asymp \frac{\text{polylog}n}{n}$) : assume every entry is observed i.i.d. with probability $0 < p \leq 1$.

Question: which matrix is easier to complete?

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

coherent

vs.

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

incoherent

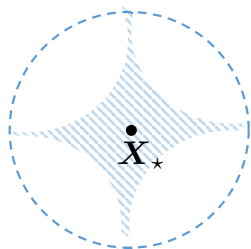
Low-rank matrix completion is only well-defined for “incoherent” matrices whose energies are spread evenly across the entries.

Which region has benign geometry?

Finite-sample level ($p \asymp \frac{\text{polylog}n}{n}$) : assume every entry is observed i.i.d. with probability $0 < p \leq 1$.

$f(\mathbf{X})$ restricted strongly convex and smooth along descent direction \mathbf{V} **only when \mathbf{X} is incoherent:**

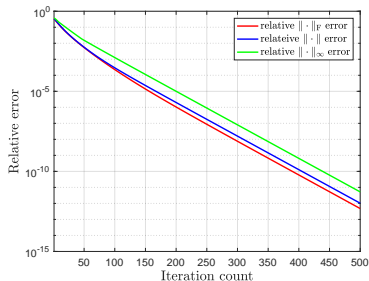
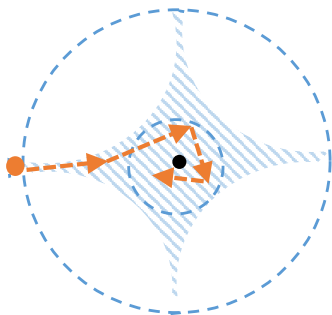
$$\|\mathbf{X} - \mathbf{X}_\star\|_{2,\infty} \ll \|\mathbf{X}_\star\|_{2,\infty}$$



region of local strong convexity + smoothness

Our findings: gradient descent is implicitly regularized

● region of local strong convexity + smoothness



Gradient descent implicitly forces iterates to remain **incoherent** even without regularization

Theoretical guarantees - noise-free case

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

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

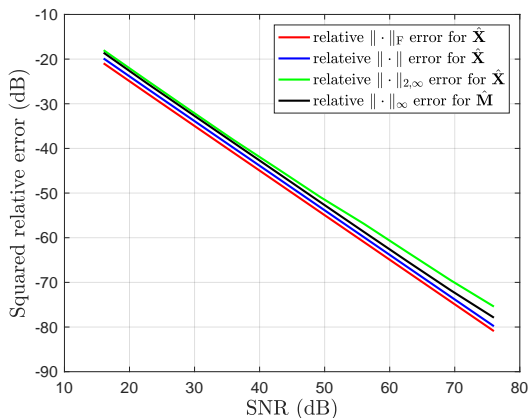
$$\|\mathbf{X}_t \mathbf{X}_t^\top - \mathbf{M}\|_{\text{F}} \leq \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

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

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

Key idea: the iterates are implicitly regularized

Noisy matrix completion via vanilla GD



Near-optimal entry-wise 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!

Towards robustness to adversarial outliers



Yuanxin Li
CMU

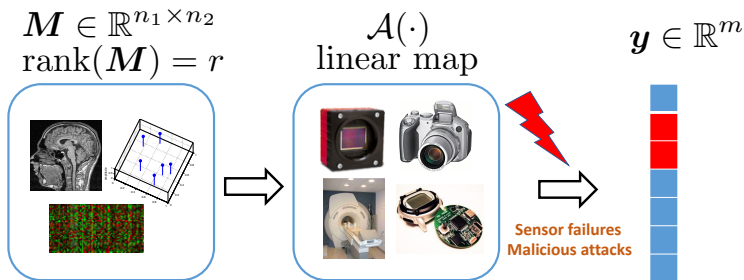


Yingbin Liang
OSU



Huishuai Zhang
MSRA

Outlier-corrupted low-rank matrix sensing



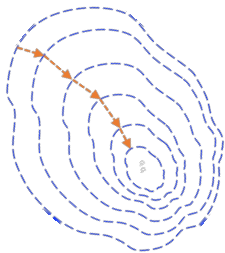
$$y = \mathcal{A}(M) + \underbrace{s}_{\text{outliers}}, \quad \mathcal{A}(M) = \{\langle A_i, M \rangle\}_{i=1}^m$$

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

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



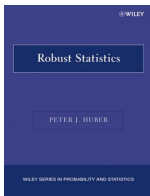
- **Gradient iterations would fail:**

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

for $t = 0, 1, \dots$

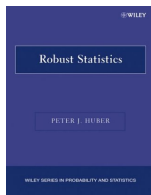
Even a single outlier can fail the algorithm!

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

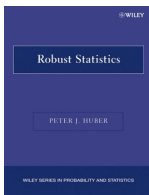


Key idea: “median-truncation” — discard samples *adaptively* based on how large sample gradients / values deviate from median

- **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 $r_t^i := |y_i - \langle \mathbf{A}_i, \mathbf{X}_t \rangle|$ 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*.

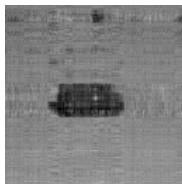
Numerical example

Low-rank matrix sensing:

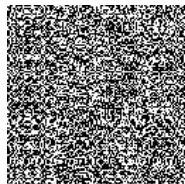
$$y_i = \langle \mathbf{A}_i, \mathbf{M} \rangle + s_i, \quad i = 1, \dots, m$$



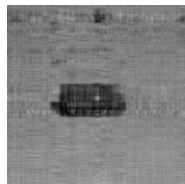
Ground truth



GD
no outliers



GD
1% outliers



median-TGD
1% outliers

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

Accelerating ill-conditioned matrix estimation



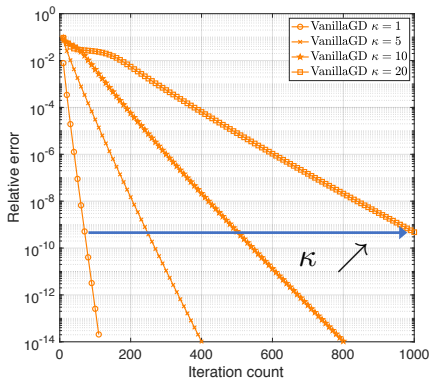
Tian Tong
CMU



Cong Ma
Princeton

Convergence slows down for ill-conditioned matrices

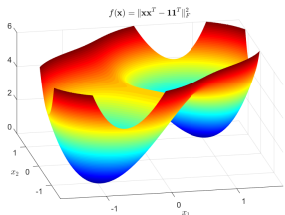
$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathcal{P}_{\Omega}(\mathbf{X}\mathbf{Y}^{\top} - \mathbf{M}) \right\|_{\text{F}}^2$$



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

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

A new algorithm: scaled gradient descent (ScaledGD)

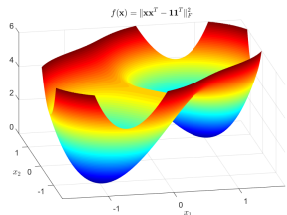


- Spectral initialization.
- Scaled gradient iterations:

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

for $t = 0, 1, \dots$

A new algorithm: scaled gradient descent (ScaledGD)



- Spectral initialization.
- Scaled gradient iterations:

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

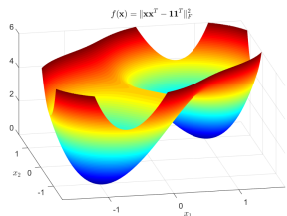
for $t = 0, 1, \dots$

For the asymmetric case:

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

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

A new algorithm: scaled gradient descent (ScaledGD)



- Spectral initialization.
- Scaled gradient iterations:

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

for $t = 0, 1, \dots$

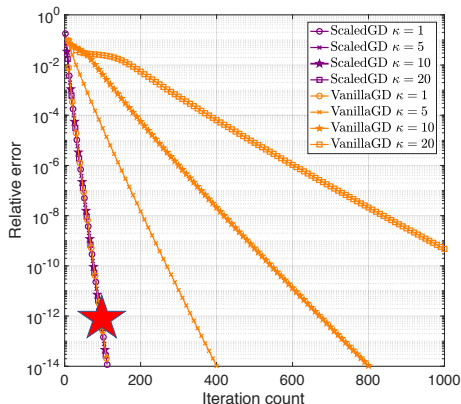
For the asymmetric case:

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

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

ScaledGD is a *preconditioned* gradient method.

ScaledGD for low-rank matrix completion



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

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

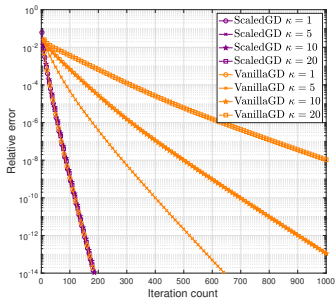
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

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

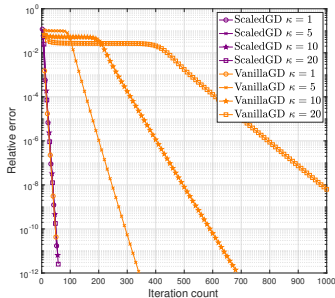
$$m \gtrsim nr^2 \kappa^2.$$

Acceleration for ill-conditioning: ScaledGD provably accelerates vanilla GD for low-rank matrix sensing.

ScaledGD works more broadly



Robust PCA



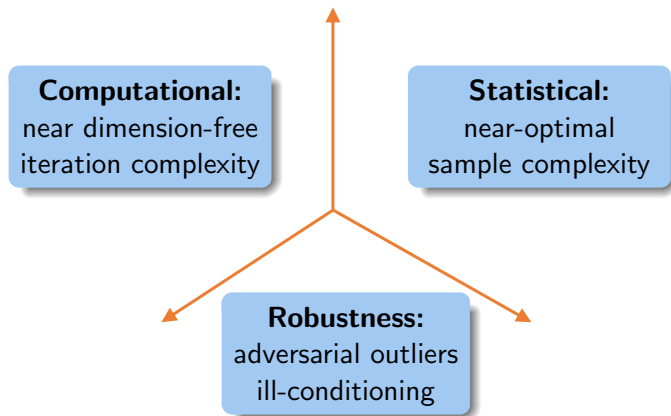
Hankel matrix completion

ScaledGD is more efficient when the low-rank matrix is ill-conditioned.

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

Final remarks

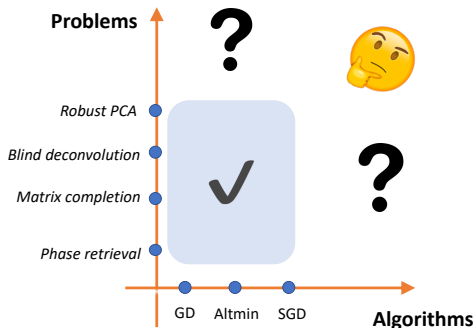
Bridging the theory-practice gap



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.

Future directions



Limitations of current framework:

- largely case-by-case: lengthy proofs, somewhat similar recipes;
- somewhat strong assumptions, e.g. Gaussian measurements, uniformly sampling...

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. Harnessing Structures in Big Data via Guaranteed Low-Rank Matrix Estimation, Y. Chen and Y. Chi, *IEEE Signal Processing Magazine*, 2018.

Geometry and implicit regularization:

1. Implicit Regularization for Nonconvex Statistical Estimation, C. Ma, K. Wang, Y. Chi and Y. Chen, *Foundations of Computational Mathematics*, 2020.
2. Nonconvex Matrix Factorization from Rank-One Measurements, Y. Li, C. Ma, Y. Chen, and Y. Chi, AISTATS 2019.

Accelerating ill-conditioned low-rank matrix estimation:

1. Accelerating Ill-Conditioned Low-Rank Matrix Estimation via Scaled Gradient Descent, T. Tong, C. Ma, and Y. Chi, preprint, 2020.

Robustness to adversarial outliers:

1. 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.
2. Median-Truncated Nonconvex Approach for Phase Retrieval with Outliers, H. Zhang, Y. Chi and Y. Liang, *IEEE Trans. on Information Theory*, 2019.

Thanks!



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