

How to Stop Worrying about Ill-Conditioning in Low-Rank Matrix Estimation

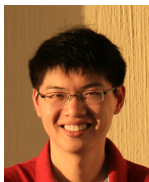
Yuejie Chi

Carnegie Mellon University

November 2020

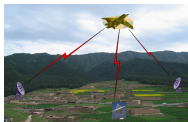


Tian Tong
CMU

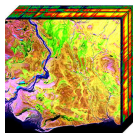


Cong Ma
Berkeley

Low-rank matrices in data science



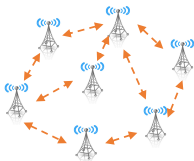
radar imaging



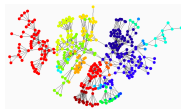
hyperspectral imaging



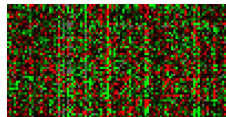
recommendation systems



localization



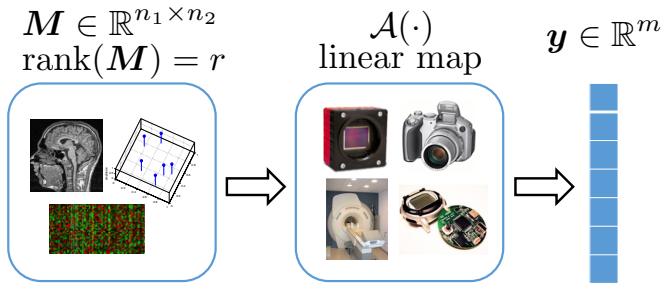
community detection



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:

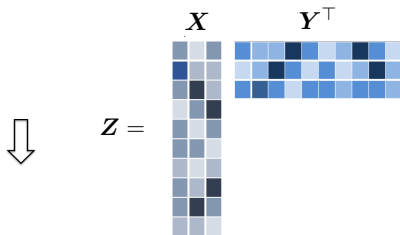
$$\underbrace{(n_1 + n_2)r}_{\text{degree of freedom}} \lesssim \underbrace{m}_{\text{sensing budget}} \ll \underbrace{n_1 n_2}_{\text{ambient dimension}}$$

Low-rank matrix factorization

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

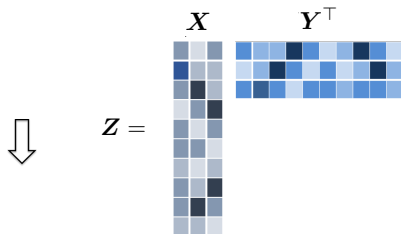
Low-rank matrix factorization

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



Low-rank matrix factorization

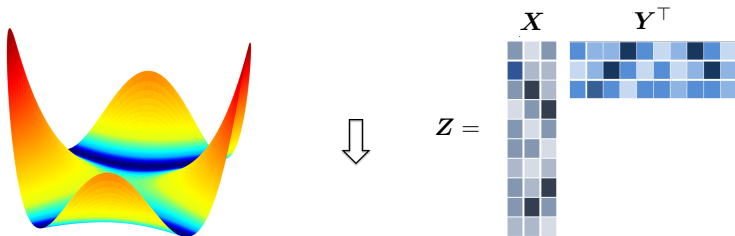
$$\min_{\text{rank}(\mathbf{Z})=r} \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{Z})\|_2^2$$



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

Low-rank matrix factorization

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

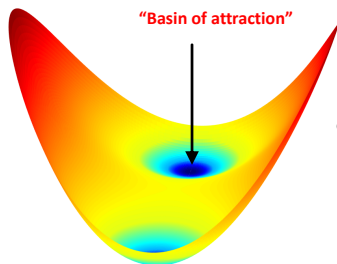


$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^T)\|_2^2$$

Saves memory and computation but introduces nonconvexity!

Prior art: GD with balancing regularization

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



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

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

- **Gradient iterations:**

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

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

for $t = 0, 1, \dots$

Prior theory for vanilla GD

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

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

Prior theory for vanilla GD

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

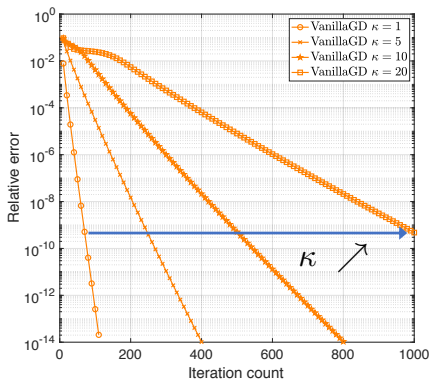
$$m \gtrsim (n_1 + n_2)r^2\kappa^2.$$

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,)

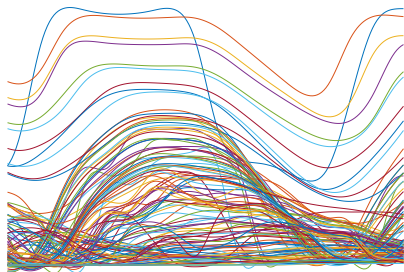
Convergence slows down for ill-conditioned matrices

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

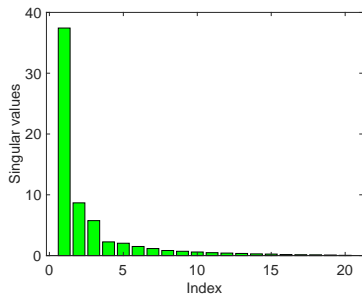


Vanilla GD converges in $O(\kappa \log \frac{1}{\epsilon})$ iterations.

Condition number can be large



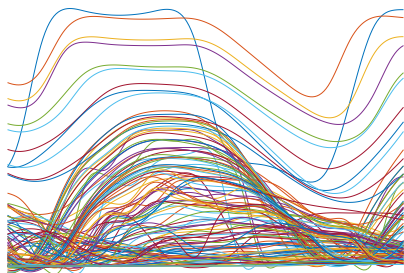
chlorine concentration levels
120 junctions, 180 time slots



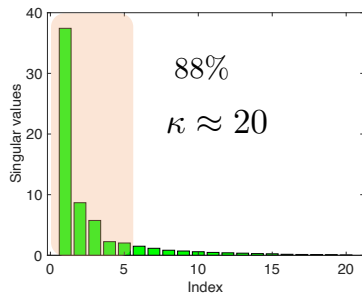
power-law spectrum

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

Condition number can be large



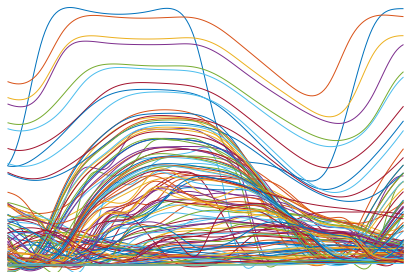
chlorine concentration levels
120 junctions, 180 time slots



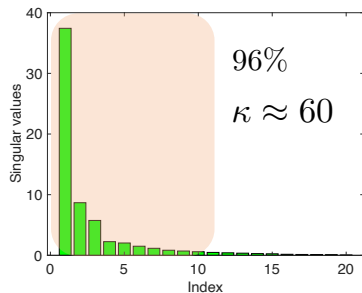
rank-5 approximation

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

Condition number can be large



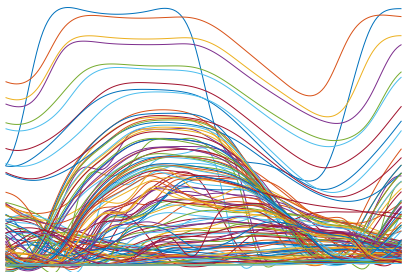
chlorine concentration levels
120 junctions, 180 time slots



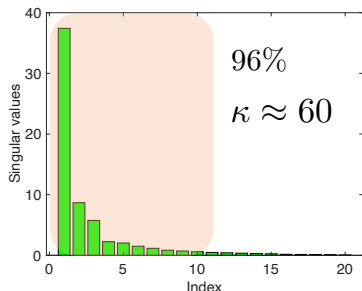
rank-10 approximation

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

Condition number can be large



chlorine concentration levels
120 junctions, 180 time slots



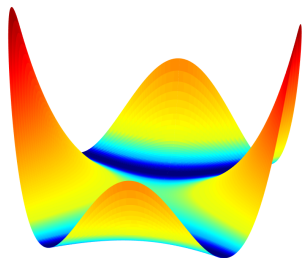
rank-10 approximation

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:**

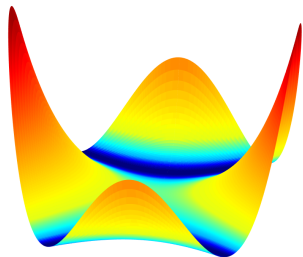
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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:**

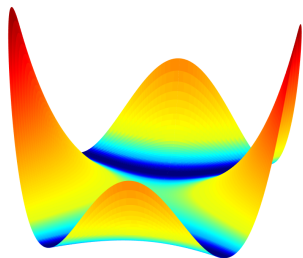
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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:**

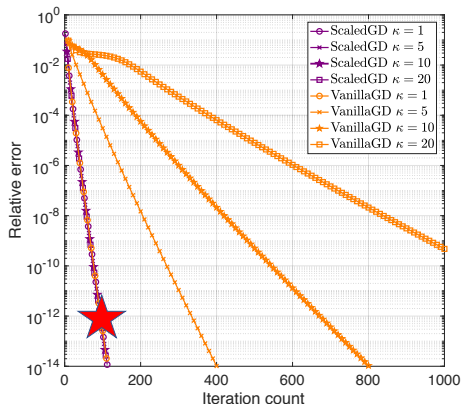
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

for $t = 0, 1, \dots$

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

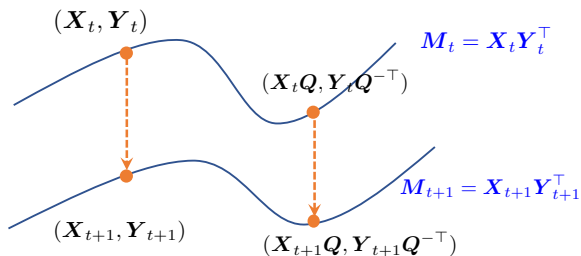
ScaledGD for low-rank matrix completion



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

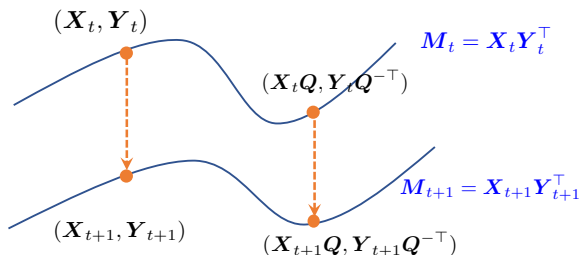
A closer look at ScaledGD

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



A closer look at ScaledGD

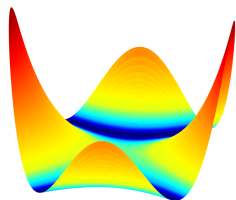
Invariance to invertible transforms: (Tanner and Wei, '16; Mishra '16)



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

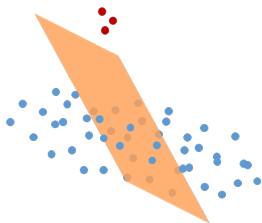
$$\|\mathbf{X}_t \mathbf{Y}_t^\top - \mathbf{M}\|_F \lesssim \varepsilon \cdot \sigma_{\min}(\mathbf{M})$$

- **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



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

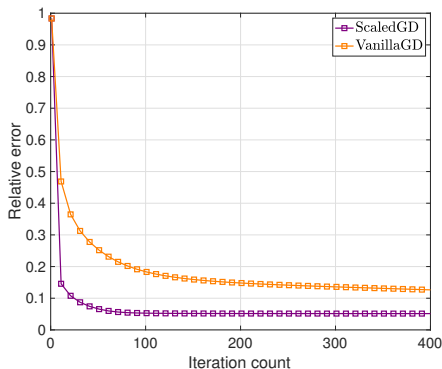
	Robust PCA		Matrix completion	
Algorithms	corruption fraction	iteration complexity	sample complexity	iteration complexity
GD	$\frac{1}{\mu r^{3/2} \kappa^{3/2} \sqrt{\mu r \kappa^2}}$	$\kappa \log \frac{1}{\epsilon}$	$(\mu \vee \log n) \mu n r^2 \kappa^2$	$\kappa \log \frac{1}{\epsilon}$
ScaledGD	$\frac{1}{\mu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

Huge computation savings at comparable sample complexities!

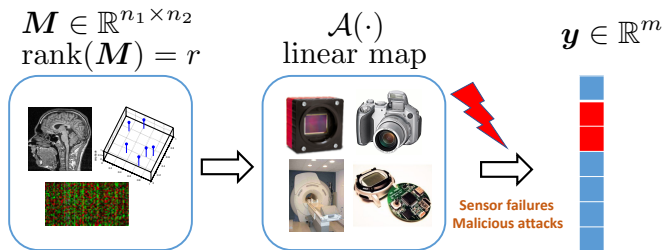
Code available at <https://github.com/Titan-Tong/ScaledGD>

Numerical stability

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

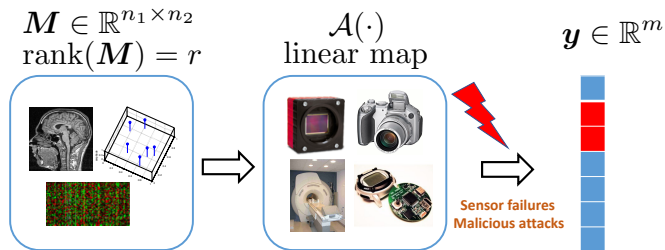


Outlier-corrupted low-rank matrix sensing



$$y = \mathcal{A}(M) + \underbrace{\text{sparse outliers}}_{\text{a small fraction (e.g. } p_s \approx 5\%)}$$

Outlier-corrupted low-rank matrix sensing

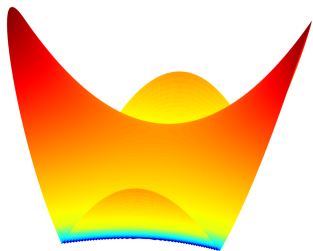


$$y = \mathcal{A}(M) + \underbrace{\text{sparse outliers}}_{\text{a small fraction (e.g. } p_s \approx 5\%)}$$

Least absolute deviation (LAD)

$$\min_{\mathbf{X}, \mathbf{Y}} f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \left\| \mathbf{y} - \mathcal{A}(\mathbf{X}\mathbf{Y}^\top) \right\|_1$$

Scaled subgradient methods



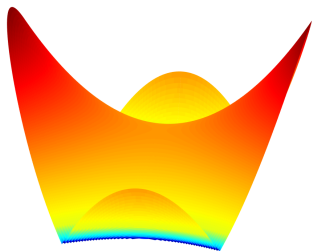
Scaled subgradient iterations:

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

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

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

Scaled subgradient methods



Scaled subgradient iterations:

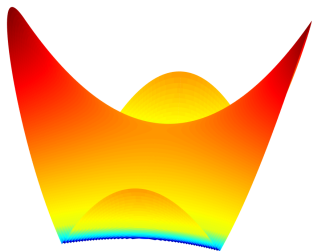
$$\mathbf{X}_{t+1} = \mathbf{X}_t - \eta_t \partial_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{Y}_t^\top \mathbf{Y}_t)^{-1}}_{\text{preconditioner}}$$

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

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

	matrix sensing	quadratic sensing
Subgradient Method (Charisopoulos et al, '19)	$\frac{\kappa}{(1-2p_s)^2} \log \frac{1}{\varepsilon}$	$\frac{r\kappa}{(1-2p_s)^2} \log \frac{1}{\varepsilon}$
ScaledSM	$\frac{1}{(1-2p_s)^2} \log \frac{1}{\varepsilon}$	$\frac{r}{(1-2p_s)^2} \log \frac{1}{\varepsilon}$

Scaled subgradient methods



Scaled subgradient iterations:

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

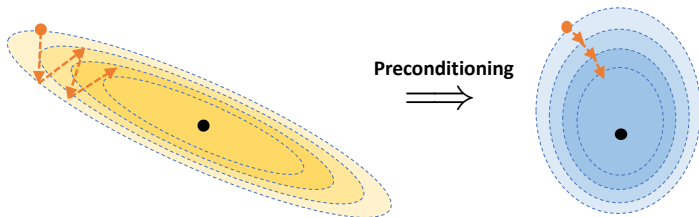
$$\mathbf{Y}_{t+1} = \mathbf{Y}_t - \eta_t \partial_{\mathbf{Y}} f(\mathbf{X}_t, \mathbf{Y}_t) \underbrace{(\mathbf{X}_t^\top \mathbf{X}_t)^{-1}}_{\text{preconditioner}}$$

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

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

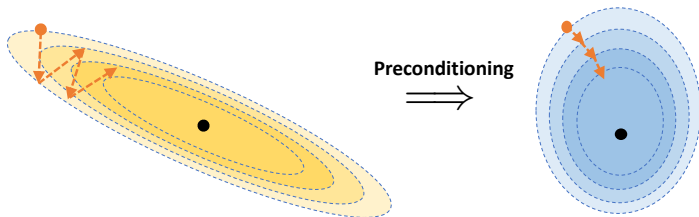
Robustness to both ill-conditioning and adversarial corruptions!

Concluding remarks



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

Concluding remarks



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

Promising directions: unveiling the power of preconditioning in

- Statistical learning
- Reinforcement learning
- Many more ...

Thanks!

- Accelerating Ill-Conditioned Low-Rank Matrix Estimation via **Scaled Gradient Descent**, arXiv 2005.08898.
- Low-Rank Matrix Recovery with **Scaled Subgradient Methods**: Fast and Robust Convergence Without the Condition Number, arXiv 2010.13364.

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