

Noisy Bilinear Low-Rank Matrix Sketching

Wenbin Wang, Xindi Ping, Cheng Cheng, Ziping Zhao

ShanghaiTech University, Shanghai, China
2025 IEEE Information Theory Workshop

October 2nd, 2025

Compressed Sensing

Classic Framework:

$$y_i = \mathbf{a}_i^\top \mathbf{x}, \quad i = 1, \dots, m,$$

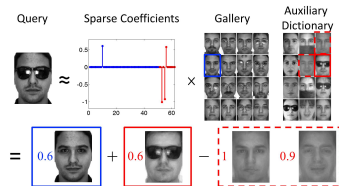
where $\{\mathbf{a}_i\}_{i=1}^m$ are the sketching/sensing vectors, \mathbf{x} is the data/signal, $\{y_i\}_{i=1}^m$ are the measurements.



(a) Magnetic Resonance Imaging



(b) Image Denoising



(c) Robust Face Recognition

Challenges in Modern Data Acquisition

Data generation at unprecedented rate: data samples are

- high-dimensional (dimension \gg data number);
- not observable due to privacy or security constraints;
- distributed at multiple locations.

Limited processing power at sensor platforms:

- **time-sensitive:** impossible to obtain a complete snapshot of the system;
- **storage-limited:** cannot store the whole data set;
- **power-hungry:** minimize the number of observations.

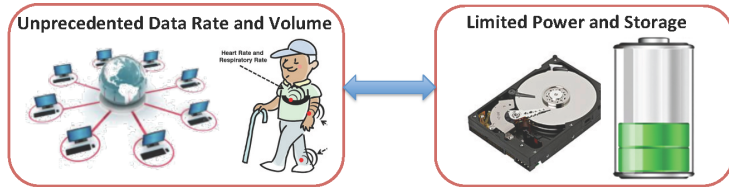


Figure 1: Mismatch streaming¹

¹<https://yuejiechi.github.io/GeometricConstraints.html>

Bilinear Matrix Sketching

$$\underbrace{Y}_{\text{observation}} = \underbrace{A}_{\text{measurement matrix}} \underbrace{X}_{\text{unknown}} \underbrace{B^\top}_{\text{measurement matrix}} + \underbrace{E}_{\text{noise}}.$$

$$Y, E \in \mathbb{R}^{m \times m}, \quad X \in \mathbb{R}^{d \times d}, \quad A, B \in \mathbb{R}^{m \times d}, \quad m \ll d.$$

Bilinear Matrix Sketching

$$\underbrace{Y}_{\text{observation}} = \underbrace{A}_{\text{measurement matrix}} \underbrace{X}_{\text{unknown}} \underbrace{B^\top}_{\text{measurement matrix}} + \underbrace{E}_{\text{noise}}.$$

$$Y, E \in \mathbb{R}^{m \times m}, \quad X \in \mathbb{R}^{d \times d}, \quad A, B \in \mathbb{R}^{m \times d}, \quad m \ll d.$$

Why Bilinear? Why Matrix? Why low-rank?

Bilinear Matrix Sketching

$$\underbrace{Y}_{\text{observation}} = \underbrace{A}_{\text{measurement matrix}} \underbrace{X}_{\text{unknown}} \underbrace{B^T}_{\text{measurement matrix}} + \underbrace{E}_{\text{noise}}.$$

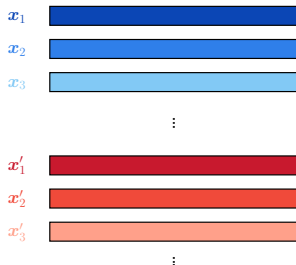
$$Y, E \in \mathbb{R}^{m \times m}, \quad X \in \mathbb{R}^{d \times d}, \quad A, B \in \mathbb{R}^{m \times d}, \quad m \ll d.$$

Why Bilinear? Why Matrix? Why low-rank?

Application-driven

Covariance Sketching

Consider two data \mathbf{x}, \mathbf{x}' possibly distributedly observed at m sensors:



Bilinear Sketching:

- two sketching matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times d}$ with specific distribution;
- two observations $\mathbf{z} = \mathbf{A}\mathbf{x}$ and $\mathbf{z}' = \mathbf{B}\mathbf{x}'$ with the cross-covariance matrix of the sketches:

$$\mathbb{E}(\mathbf{z}\mathbf{z}'^\top) = \mathbf{A} \underbrace{\mathbf{X}^\star}_{\mathbb{E}(\mathbf{x}\mathbf{x}'^\top)} \mathbf{B}$$

Graph Sketching

Consider a directed graph \mathcal{G} with d nodes with adjacency matrix \mathbf{X} .

- First, we consider $\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{A}^\top$
- Define $\mathbf{A} \in \mathbb{R}^{m \times d}$ as composed of i.i.d. Bernoulli entries such as

$$\mathbf{A}_{u,i} = \begin{cases} 1, & \text{if } i \in u, \\ 0, & \text{otherwise.} \end{cases}$$

- Then,

$$Y_{u,v} = \sum_{i=0}^{d-1} \sum_{j=0}^{d-1} \mathbf{A}_{u,i} \mathbf{X}_{i,j} \mathbf{A}_{v,j}$$

Illustration of Graph Sketching

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

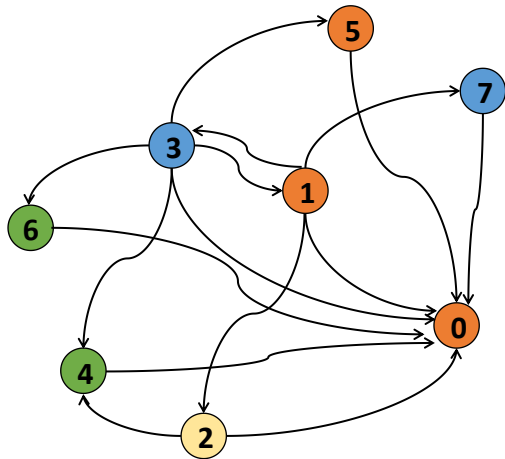


Figure 2: Original Graph: \mathbf{X}

Illustration of Graph Sketching

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$
$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

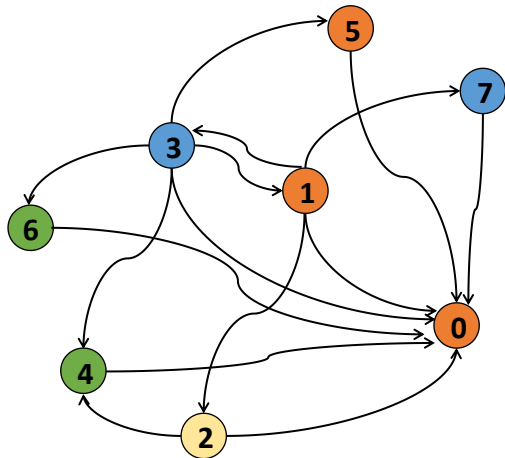


Figure 3: Original Graph: \mathbf{X}

Illustration of Graph Sketching

$$\mathbf{Y} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 4 & 2 \\ 1 & 2 & 2 & 0 \\ 0 & 0 & 2 & 0 \end{bmatrix}$$

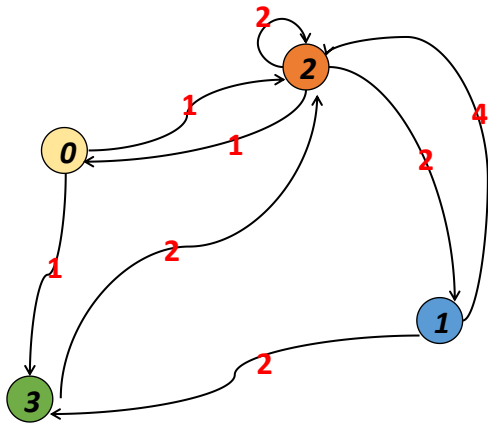


Figure 4: Compressed Graph: \mathbf{Y}

Graph Sketching

Consider a directed graph \mathcal{G} with d nodes with adjacency matrix \mathbf{X} .

- Now, we consider $\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{B}^\top$
- Define $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times d}$ as composed of i.i.d. Bernoulli entries such as

$$\mathbf{A}_{u,i} = \begin{cases} 1, & \text{if } i \in u, \\ 0, & \text{otherwise.} \end{cases} \quad \mathbf{B}_{v,j} = \begin{cases} 1, & \text{if } j \in v, \\ 0, & \text{otherwise.} \end{cases}$$

- Then,

$$Y_{u,v} = \sum_{i=0}^{d-1} \sum_{j=0}^{d-1} \mathbf{A}_{u,i} \mathbf{X}_{i,j} \mathbf{B}_{v,j} = \sum_{i \in u} \sum_{j \in v} \mathbf{X}_{i,j}$$

The sketching matrices \mathbf{A} and \mathbf{B} respectively partition the original graph \mathcal{G} in two different dimensions.

Illustration of Graph Sketching

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

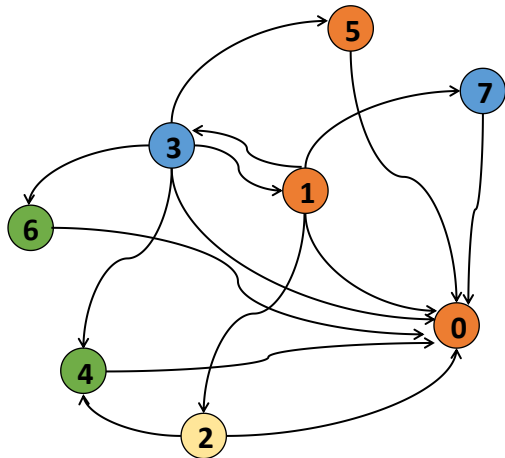


Figure 5: Original Graph: \mathbf{Y}

Illustration of Graph Sketching

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$
$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

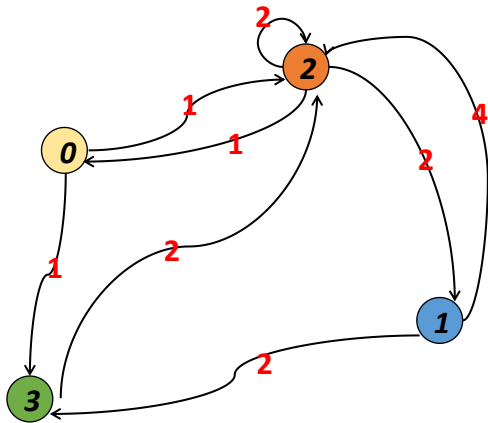


Figure 6: Compressed Graph: $\mathbf{Y} = \mathbf{A}\mathbf{X}\mathbf{A}^\top$

Illustration of Graph Sketching

$$X = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$

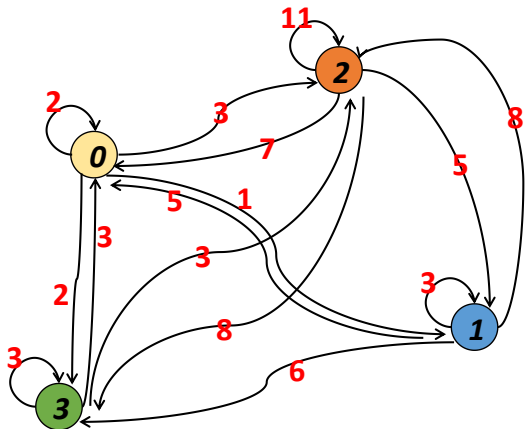


Figure 7: Compressed Graph: $Y = BXB^T$

Illustration of Graph Sketching

$$\begin{aligned}
 AXA^\top &= \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 0 & 4 & 2 \\ 1 & 2 & 2 & 0 \\ 0 & 0 & 2 & 0 \end{bmatrix} \\
 BXB^\top &= \begin{bmatrix} 2 & 1 & 3 & 2 \\ 5 & 3 & 8 & 6 \\ 7 & 5 & 11 & 8 \\ 3 & 0 & 3 & 3 \end{bmatrix} \\
 AXB^\top &= \begin{bmatrix} 1 & 1 & 2 & 1 \\ 3 & 2 & 5 & 4 \\ 3 & 2 & 4 & 3 \\ 2 & 0 & 2 & 2 \end{bmatrix}
 \end{aligned}$$

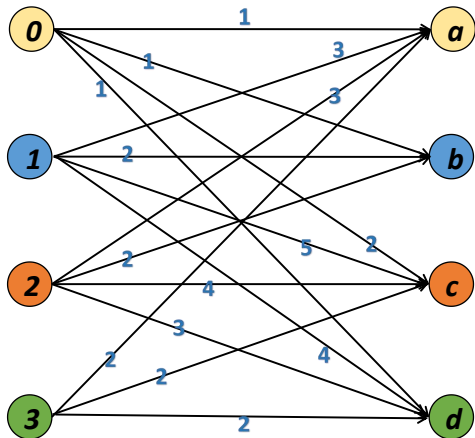






Figure 8: Compressed Graph: $Y = AXB^\top$

Related Work

The *distributed sparsity* assumption: the matrix \mathbf{X}^* is called d -distributed sparse if each row/column of \mathbf{X} cannot have more than d non-zeros.

- Convex optimization program (Dasarathy 2012)

$$\widehat{\mathbf{X}} = \arg \min_{\mathbf{X}} \left\{ \left\| \mathbf{A} \mathbf{X} \mathbf{B}^{\top} - \mathbf{Y} \right\|_{\text{F}}^2 + \lambda \left\| \mathbf{X} \right\|_1 \right\}$$

-  is easy to obtain by iterative algorithm due to its convexity,
-  introduces a non-negligible bias.
- **Our work:**
 -  no bias
 -  low-rank

Problem Formulation

We propose to estimate the low-rank matrices from bilinear measurements using the non-convex penalty

$$\underset{\mathbf{X}}{\text{minimize}} \quad \frac{1}{2m^2} \left\| \mathbf{Y} - \mathbf{A} \mathbf{X} \mathbf{B}^\top \right\|_{\text{F}}^2 + P_\lambda(\mathbf{X})$$

where $P_\lambda(\mathbf{X}) = \sum_{i=1}^d p_\lambda(\sigma_i(\mathbf{X}))$ is a decomposable nonconvex penalty imposed on the singular values of \mathbf{X} such as $P_\lambda(\mathbf{X}) = \lambda \|\mathbf{X}\|_* + \sum_{i=1}^d q_\lambda(\sigma_i(\mathbf{X}))$.

Assumption 1

- There exists $\nu > 0$ such that the derivative satisfies $p'_\lambda(t) = 0$ for all $t \geq \nu$;
- Both $p_\lambda(t)$ and $q_\lambda(t)$ are symmetric about zero, i.e., $p_\lambda(t) = p_\lambda(-t)$, $q_\lambda(t) = q_\lambda(-t)$;
- The derivative $q'_\lambda(t)$ is monotonic and Lipschitz continuous in the interval $[0, \infty)$.
Explicitly, for $t_2 \geq t_1 \geq 0$, there exist constants $\zeta^- \geq \zeta^+ > 0$ such that
$$-\zeta^- \leq \frac{q'_\lambda(t_2) - q'_\lambda(t_1)}{(t_2 - t_1)} \leq -\zeta^+.$$
- Both $q_\lambda(t)$ and its derivative vanish at zero, i.e., $q_\lambda(0) = q'_\lambda(0) = 0$;
- There exists a constant $\lambda > 0$ bounding the magnitude of the derivative, i.e., $|q'_\lambda(t)| \leq \lambda$.

Optimization Algorithm

Algorithm 1: Proximal Gradient Algorithm

Input $\lambda_0 > 0, \epsilon > 0, L_{\min} > 0, \eta \in (0, 1), \delta \in (0, 1)$

Initialize $\mathbf{X}^0 = \mathbf{0}, L_0 = L_{\min}$

for $t = 0, 1, \dots, T - 1$ **do**

$\lambda_{t+1} = \eta \lambda_t; \epsilon_{t+1} = \lambda_t/4;$

$k = 0; \mathbf{X}^k = \mathbf{X}^t;$

while $\omega_{\lambda_{t+1}}(\mathbf{X}^k) > \epsilon_{t+1}$ **do**

$k = k + 1;$

$\mathbf{X}^k = \arg \min_{\mathbf{X}} \tilde{F}_{L, \lambda}(\mathbf{X}; \mathbf{X}^{k-1});$

if $F(\mathbf{X}^k) > \tilde{F}(\mathbf{X}^k; \mathbf{X}^{k-1})$ **then**

$L_{k-1} = 2L_{k-1};$

end if

$L_k = \max\{L_{\min}, L_{k-1}/2\}$

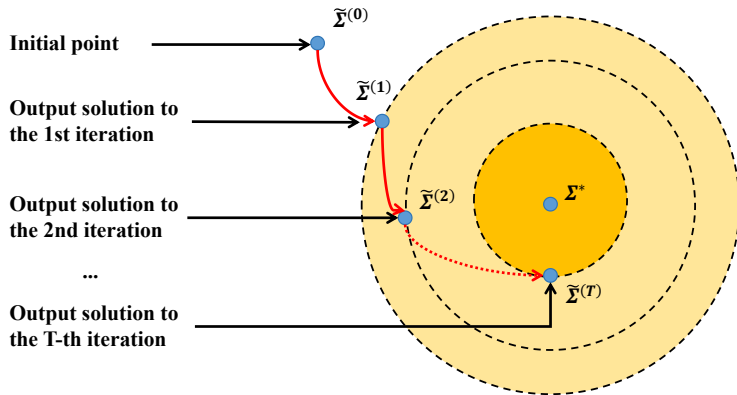
end while

$\mathbf{X}^{t+1} = \mathbf{X}^k; L_{t+1} = L_k$

end for

Output $\{\mathbf{X}^t\}_{t=1}^T$

Algorithm Illustration



Preliminaries

Consider the singular value decomposition $\mathbf{X}^* = \mathbf{U}^* \mathbf{\Sigma}^* \mathbf{V}^{*\top}$, where $\mathbf{U}^*, \mathbf{V}^* \in \mathbb{R}^{d \times r}$, and $\mathbf{\Sigma}^* = \text{diag}(\sigma_1^*, \dots, \sigma_r^*)$. We introduce the subspace \mathcal{F} and \mathcal{F}^\perp , which are defined in terms of the row and column spaces of the matrices:

$$\begin{aligned}\mathcal{F}(\mathbf{U}^*, \mathbf{V}^*) &:= \{\mathbf{\Delta} \mid \text{row}(\mathbf{\Delta}) \subseteq \mathbf{V}^*, \text{col}(\mathbf{\Delta}) \subseteq \mathbf{U}^*\}, \\ \mathcal{F}^\perp(\mathbf{U}^*, \mathbf{V}^*) &:= \{\mathbf{\Delta} \mid \text{row}(\mathbf{\Delta}) \perp \mathbf{V}^*, \text{col}(\mathbf{\Delta}) \perp \mathbf{U}^*\}.\end{aligned}$$

Restricted Region

Define a local region \mathcal{R} as

$$\mathcal{R} = \{\mathbf{\Delta} \mid \|\Pi_{\mathcal{F}^\perp}(\mathbf{\Delta})\|_* \leq 5 \|\Pi_{\mathcal{F}}(\mathbf{\Delta})\|_*\},$$

where $\Pi_{\mathcal{F}(\cdot)}$ is the projection operator that projects matrices into the subspace \mathcal{F} .

Essential Assumptions

Assumption 2 (RSC & RSM)

- The empirical loss function $f(\cdot)$ is ρ^- -strongly convex and ρ^+ -smooth over \mathcal{R} with $\infty > \rho^+ \geq \rho^- > 0$. Specifically, for all $\mathbf{X} - \mathbf{X}' \in \mathcal{R}$, we have:

$$\begin{aligned}\langle \mathbf{X} - \mathbf{X}', \nabla f(\mathbf{X}) - \nabla f(\mathbf{X}') \rangle &\geq \rho^- \|\mathbf{X} - \mathbf{X}'\|_{\text{F}}^2, \\ \langle \mathbf{X} - \mathbf{X}', \nabla f(\mathbf{X}) - \nabla f(\mathbf{X}') \rangle &\geq \frac{\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|_{\text{F}}^2}{\rho^+}.\end{aligned}$$

Assumption 3 (Minimal Signal Strength)

- The singular value of the ground truth \mathbf{X}^* satisfies:

$$\min_{i \in \mathcal{S}_1 \cup \mathcal{S}_2} |\sigma_i^*| \geq \nu + 2\sqrt{s_1 + s_2} \|\mathbf{A}^\top \mathbf{E} \mathbf{B}\|_{\text{F}} / (m^2 \rho).$$

Theoretical Results I

Define $\mathcal{S}_1 = \{i \mid \sigma_i^* \geq \nu\}$, $\mathcal{S}_2 = \{i \mid \nu > \sigma_i^* > 0\}$ with their corresponding cardinalities given by $s_1 = |\mathcal{S}_1|$ and $s_2 = |\mathcal{S}_2|$.

Theorem 1

Suppose Assumptions 1 and 2 hold, if $\rho^- > \zeta^-$, $\lambda \gtrsim \|\mathbf{A}^\top \mathbf{E} \mathbf{B}\|_F / m^2$, we have:

$$\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F \lesssim \tau \sqrt{s_1} + \sqrt{s_2}$$

where $\tau = \|\Pi_{\mathcal{F}_{S_1}}(\nabla f(\mathbf{X}^*))\|_F$ and \mathcal{F}_{S_1} is a subspace of \mathcal{F} associated with S_1 .

- The *oracle rate* refers to the statistical convergence rate of the *oracle estimator*, defined as $\widehat{\mathbf{X}}^O = \arg \min_{\mathbf{X} \in \mathcal{F}(\mathbf{U}^*, \mathbf{V}^*)} f(\mathbf{X})$, which knows the true rank spaces in advance. By the mean value theorem, it is easy to obtain that $\|\widehat{\mathbf{X}}^O - \mathbf{X}^*\|_F \lesssim \|\Pi_{\mathcal{F}}(\nabla f(\mathbf{X}^*))\|_F$.

Theoretical Results II

Theorem 2 (oracle property)

Suppose Assumptions 1, 2 and 3 hold.

If $\rho > \zeta^-$, and

$$\lambda \geq \frac{(\rho^- + \sqrt{s_1 + s_2}\rho^+) \|\mathbf{A}^\top \mathbf{E} \mathbf{B}\|_F}{2m^2 \rho^-},$$

we have

$$\text{rank}(\widehat{\mathbf{X}}) = \text{rank}(\widehat{\mathbf{X}}^O) = \text{rank}(\mathbf{X}^*)$$

and

$$\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F \lesssim \sqrt{s_1} \tau,$$

where $\tau = \|\Pi_{\mathcal{F}}(\nabla f(\mathbf{X}^))\|_F$.*

Theoretical Results III

Corollary 3

Consider the noise entries are i.i.d. sub-Gaussian random variables with variance κ , and the vectorized sketching matrices $\text{vec}(\mathbf{A})$ and $\text{vec}(\mathbf{B})$ follow sub-Gaussian distributions $\text{vec}(\mathbf{A}) \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Theta}_1)$, $\text{vec}(\mathbf{B}) \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Theta}_2)$. We term these as $\boldsymbol{\Theta}_1$ -ensemble and $\boldsymbol{\Theta}_2$ -ensemble, and $\varpi_1(\boldsymbol{\Theta}_1) = \sqrt{\sup_{\|u\|_2=1, \|v\|_2=1} \text{Var}(u^\top \mathbf{A}v)}$, $\varpi_2(\boldsymbol{\Theta}_2) = \sqrt{\sup_{\|u\|_2=1, \|v\|_2=1} \text{Var}(u^\top \mathbf{B}v)}$. Assuming Assumptions 1 and 2 hold, and \mathbf{A} and \mathbf{B} are sampled from $\boldsymbol{\Theta}_1$ -ensemble and $\boldsymbol{\Theta}_2$ -ensemble, respectively, if $\rho \gtrsim \sqrt{\lambda_{\min}(\boldsymbol{\Theta}_1)\lambda_{\min}(\boldsymbol{\Theta}_2)} > \zeta^-$ and $\lambda \gtrsim \kappa\sqrt{\varpi_1\varpi_2 d}/m$, then with probability at least $1 - \exp(-d)$. Additionally, according to Theorem 1, we have:

$$\|\widehat{\mathbf{X}} - \mathbf{X}^\star\|_{\text{F}} \lesssim \mathcal{O}\left(\sqrt{\frac{\varpi_1\varpi_2}{\lambda_{\min}(\boldsymbol{\Theta}_1)\lambda_{\min}(\boldsymbol{\Theta}_2)}} \frac{\kappa(\sqrt{s_2 d} + s_1)}{m}\right).$$

With Assumption 3, the convergence rate improves to

$$\mathcal{O}\left(\sqrt{\frac{\varpi_1\varpi_2}{\lambda_{\min}(\boldsymbol{\Theta}_1)\lambda_{\min}(\boldsymbol{\Theta}_2)}} \frac{\kappa s_1}{m}\right).$$

Experiment I

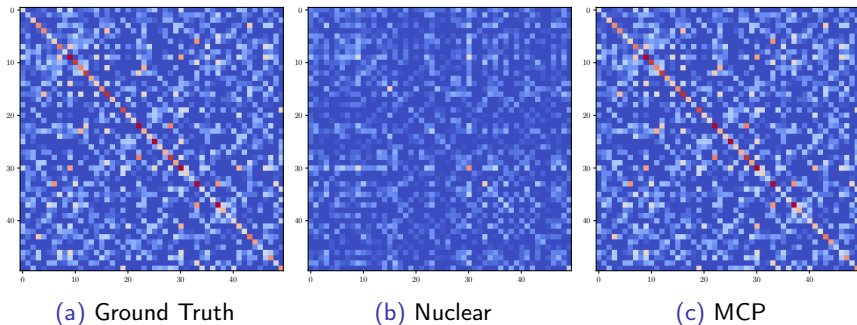


Figure 9: Heatmaps show the recovery of a 50×50 low-rank matrix (rank = 10, Gaussian generated) from noisy bilinear measurements with $\mathcal{N}(0, 0.01)$ noise. MCP achieves near-perfect recovery, while nuclear-norm minimization exhibits excessive smoothing and weakens low-rank features.

Experiment II

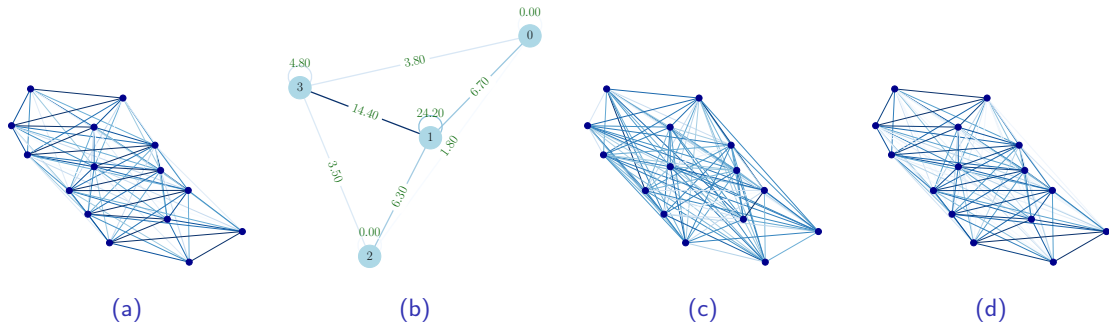


Figure 10: An illustrative example of graph sketching is shown as follows: (a) The original graph \mathcal{G} with 15 nodes; (b) The sketch of the graph \mathcal{G} , where the nodes represent the partitions and the edges represent the total number of edges of \mathcal{G} that cross these partitions; (c) The graph recovered using least squares error minimization; (d) The graph recovered using the SCAD penalty.

Experiment III

Dataset	Nuclear	Weighted Nuclear	SCAD	MCP
Fashion-MNIST	0.7683 ± 0.1076	0.7286 ± 0.0630	0.0124 ± 0.0033	0.0108 ± 0.0012
Places365	0.4472 ± 0.0827	0.4647 ± 0.0484	0.0079 ± 0.0013	0.0066 ± 0.0021
ImageNet-O	0.4574 ± 0.1502	0.5069 ± 0.1149	0.0137 ± 0.0079	0.0138 ± 0.0056

Table 1: Low-rank recovery experiments on three real-world image datasets: Fashion-MNIST ($d = 28, r = 10$), Places365 ($d = 256, r = 100$), and ImageNet-O ($d = 512, r = 200$). Observations are formed using bilinear sketching with $m = 5, 50, 80$ and additive $\mathcal{N}(0, 0.01)$ noise. We compare nuclear norm, weighted nuclear norm, and nonconvex methods (SCAD and MCP).

Conclusions

- We have proposed a novel approach for low-rank matrix estimation from bilinear measurements using the non-convex penalty.
- We have presented both the theoretical and empirical results.

Thank you!

Website: <https://wenbinwang12.github.io>

Email: wangwb2023@shanghaitech.edu.cn