Factor Group-Sparse Regularization for Efficient Low-Rank Matrix Recovery

Jicong Fan, Lijun Ding, Yudong Chen, and Madeleine Udell Cornell University

Low-rank matrix completion (LRMC)

Recover a low-rank $M \in \mathbb{R}^{m \times n}$ from partial observation $\mathcal{P}_{\Omega}(M)$.

<i>M</i>	A	lacksquare	$P_{\Omega}(M)$	M
0.19 1.33 1.18 -0.77 2.06	1.53 0.14	0.06 0.74 0.82 -0.35 1.31	0.19 ? 1.18 -0.77 ?	0.19 1.33 1.18 -0.77 2.06
0.43 0.85 -0.49 -1.05 0.09	-0.19 0.67	0.66 1.47 -0.51 -1.66 0.49	? 0.85 -0.49 ? 0.09	0.43 0.85 -0.49 -1.05 0.09
0.17 0.32 -0.21 -0.40 0.01	-0.09 0.26		0.17 0.32 ? -0.40 0.01	0.17 0.32 -0.21 -0.40 0.01
0.08 0.57 0.51 -0.33 0.88	0.65 0.06		? 0.57 0.51 -0.33 ?	0.08 0.57 0.51 -0.33 0.88
-0.37 -0.99 0.03 0.98 -0.63	-0.29 -0.53		-0.37 -0.99 ? 0.98 -0.63	-0.37 -0.99 0.03 0.98 -0.63

Applications

ocollaborative filtering (recommendation system)

(a) Low-rank matrix

- classification (especially on incomplete data)
- image inpainting
- other data pre-processing tasks

Rank minimization and Schatten-p "norm" surrogate

Rank Minimization:

minimize rank(
$$X$$
), subject to $\mathcal{P}_{\Omega}(X) = \mathcal{P}_{\Omega}(M)$

② Need certain surrogate R(X) of rank(X):

minimize
$$R(X)$$
, subject to $\mathcal{P}_{\Omega}(X) = \mathcal{P}_{\Omega}(M)$.

3 Schatten-p "norm", 0 :

$$R(\boldsymbol{X}) = \|\boldsymbol{X}\|_{\mathrm{Sp}} := \Big(\sum_{i=1}^{\min(m,n)} \sigma_i^p(\boldsymbol{X})\Big)^{1/p}.$$

- **10 Approximation advantage**: smaller value p leads to better approximation of rank(X)
- Computational challenge: full SVD computation

Factorization Approaches for $||X||_{Sp}$

Factorize X = AB, where $A \in \mathbb{R}^{m \times d}$, $B \in \mathbb{R}^{d \times n}$

- ① Factored nuclear (F-nuclear) norm^[1]: $||X||_* = \min_{AB=X} \frac{1}{2} (||A||_F^2 + ||B||_F^2)$
- 2 Bi-nuclear norm^[2]: $||X||_{S_{1/2}} = \min_{AB-X} ||A||_* ||B||_* = \min_{AB-X} \left(\frac{||A||_* + ||B||_*}{2}\right)^2$
- 3 F²+nuclear norm^[3]: $||X||_{S_{\overline{3}}^2} = \min_{AB=X} \left[\frac{2}{3} (||A(||_* + \frac{1}{2}||B)||_F^2) \right]^{3/2}$

Problematic for larger d

- Tull rankness of factors *A*, *B* during optimization procedures
- 2 (Possible) SVD computation of factors A, B

Notations

We factor $X \in \mathbb{R}^{m \times n}$ as $A = [a_1, a_2, \dots, a_d] \in \mathbb{R}^{m \times d}$ and $B = [b_1, b_2, \dots, b_d]^T \in \mathbb{R}^{d \times n}$, where $d \ge r := \operatorname{rank}(X)$, and a_i and b_i are column vectors. WLOG, we assume $m \le n$.

Our Methods: (i) group-sparse regularizers (FGSR)

We suggest using the following surrogates R(X):

$$FGSR_{1/2}(X) := \frac{1}{2} \min_{AB=X} ||A||_{2,1} + ||B^T||_{2,1}, \quad FGSR_{2/3}(X) := \frac{2}{3\alpha^{1/3}} \min_{AB=X} ||A||_{2,1} + \frac{\alpha}{2} ||B||_F^2,$$
 where $||A||_{2,1} := \sum_{j=1}^d ||a_j||_{2,j}$.

Theorem 1 (Relation to $||X||_{Sv}$)

$$FGSR_{1/2}(X) = \sum_{j=1}^{r} \sigma_j^{1/2}(X) = \|X\|_{S_{1/2}'}^{1/2} \qquad FGSR_{2/3}(X) = \sum_{j=1}^{r} \sigma_j^{2/3}(X) = \|X\|_{S_{2/3}}^{2/3}$$

In general, we have

Theorem 2 Fix $\alpha > 0$, and choose $q \in \{1, \frac{1}{2}, \frac{1}{4}, \cdots\}$. Set $p = \frac{2q}{(2+q)}$. For any matrix $X \in \mathbb{R}^{m \times n}$ with rank $(X) = r \leq d \leq \min(m, n)$, we have

$$FGSR_{p}(X) := \frac{1}{(1/2 + 1/q)\alpha^{q/(q+2)}} \min_{X = \sum_{j=1}^{d} a_{j} \boldsymbol{b}_{j}^{T}} \sum_{j=1}^{d} \frac{1}{q} \|\boldsymbol{a}_{j}\|^{q} + \frac{\alpha}{2} \|\boldsymbol{b}_{j}\|^{2} = \|X\|_{Sp}^{p}.$$
 (1)

Our Methods: (ii) optimization templates and algorithms

LRMC

(b) Matrix completion^[1]

Optimization template:

minimize
$$FGSR_p(X)$$
, subject to $P_{\Omega}(X) = P_{\Omega}(M)$. (2)

Take $FGSR_{2/3}$ as an example, and rewrite (2) as

minimize
$$||A||_{2,1} + \frac{\alpha}{2}||B||_F^2$$
, subject to $X = AB$, $P_{\Omega}(X) = P_{\Omega}(M)$.

• Algorithm: Alternating Direction Method of Multipliers (ADMM) [4] with linearization. **Noisy LRMC**

Observe $P_{\Omega}(M_e)$ with $M_e = M + E$ with noise E.

Optimization template:

minimize
$$\frac{1}{2} \|\mathcal{P}_{\Omega}(\mathbf{M}_e - \mathbf{A}\mathbf{B})\|_F^2 + \gamma \left(\frac{1}{q} \|\mathbf{A}\|_{2,q}^q + \frac{\alpha}{2} \|\mathbf{B}^T\|_F^2\right),$$
 (3)

where $q \in \{1, \frac{1}{2}, \frac{1}{4}, \cdots\}$ and $||A||_{2,q} := \left(\sum_{j=1}^{d} ||a_j||^q\right)^{1/q}$.

• Algorithm: Proximal alternating linearized minimization (PALM) [5] for better efficiency. Our methods also extends to Robust PCA (see section 6 of the paper for detail)

Generalization error bound for LRMC

Each entry of a rank-r ($r \le d$) matrix M is observed with probability ρ plus some noise $E_{ij} \sim \mathcal{N}(0, \epsilon^2)$. Observation is $P_{\Omega}(M + E)$. Consider solving

$$\underset{\|X\|_{S_n}^p \leqslant R_p, \operatorname{rank}(X) \leqslant d}{\operatorname{minimize}} \|\mathcal{P}_{\Omega}(M + E - X)\|_F^2. \tag{4}$$

This is equivalent to (3) for certain choice of γ and R_p .

Theorem 3

Suppose $||M||_{S_n}^p \le R_p$, \hat{M} is the optimal solution of (4), and $|\Omega| \ge \frac{32}{3}n \log^2 n$. Denote

 $\zeta := \max\{\|M\|_{\infty}, \|\hat{M}\|_{\infty}\}$. Then there exist numerical constants c_1 and c_2 such that the following inequality holds with probability at least $1 - 5n^{-2}$

$$\|\mathbf{M} - \hat{\mathbf{M}}\|_{F}^{2} \leq \max \left\{ c_{1} \zeta^{2} \frac{n \log n}{|\Omega|}, (5.5 + \sqrt{10}) R_{p} \left((4\sqrt{3}\epsilon_{0} + c_{2}\zeta)^{2} \frac{n \log n}{|\Omega|} \right)^{1-p/2} \right\}.$$
 (5)

In sum, Theorem 3 shows it is possible to reduce the matrix completion error by using a smaller p in (4) or a smaller q in (3) (as the second term in the brace of (5) is the dominant term for sufficiently large $|\Omega|$, which decreases as p decreases).

Advantage of our methods

- **Tighter rank approximation.** Compared to the nuclear norm, the spectral quantities in Theorem 1 are tighter approximations to the rank of *X*.
- •Robust to rank initialization. The iterative algorithms we proposed to minimize $FGSR_{2/3}$ and $FGSR_{1/2}$ quickly force some of the columns of A and B^T to zero, where they remain. Hence the number of nonzero columns is reduced dynamically, and converges to r quickly in experiments: these methods are rank-revealing.
- Low computational cost and SVD-free! The per iteration complexity of algorithms LRMC can be as low as $O(d'\operatorname{card}(\Omega))$.

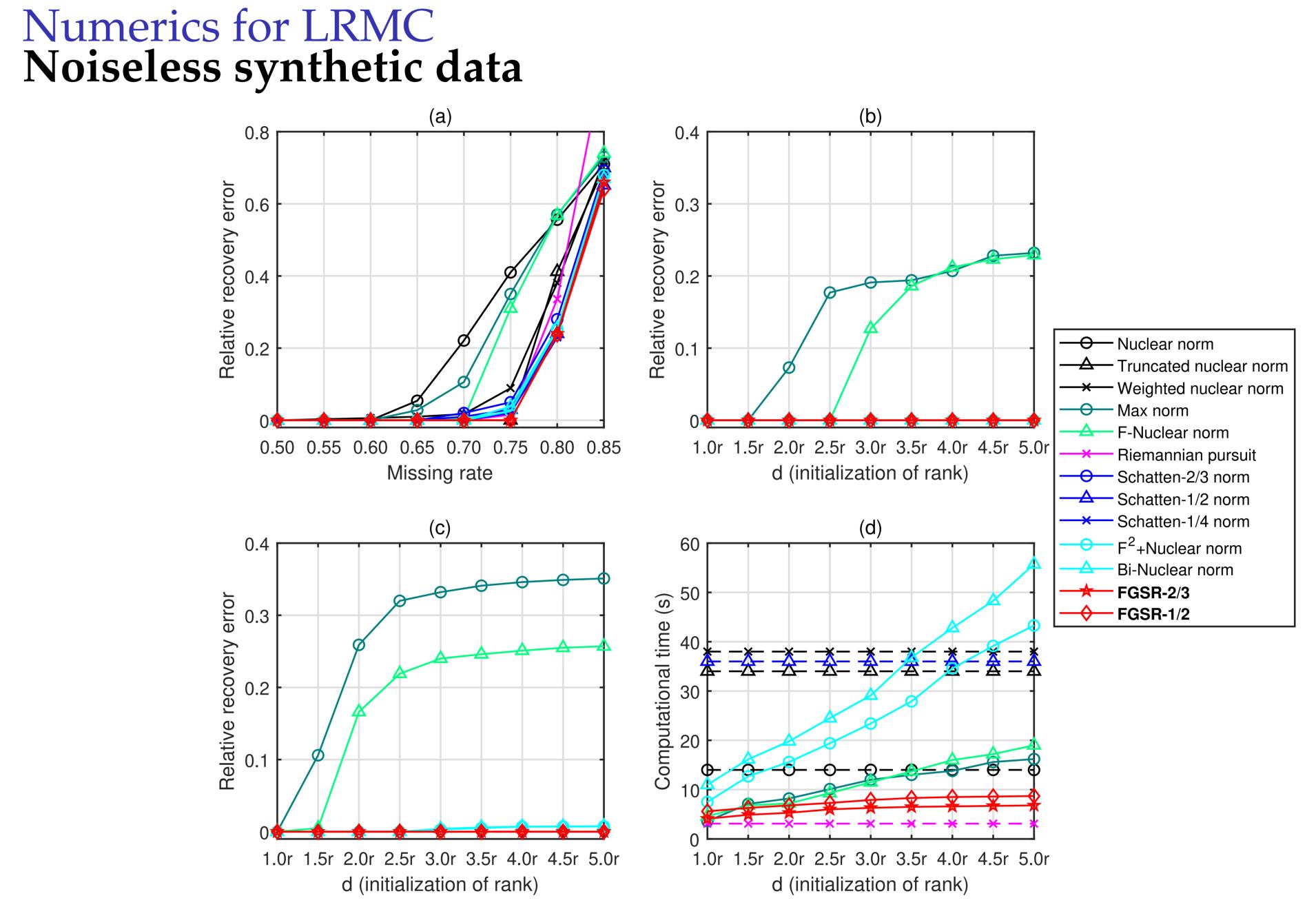


Figure: Matrix completion on noiseless synthetic data (r = 50): (a) the effect of missing rate on recovery error; (b)(c) the effect of rank initialization on recovery error (missing rate = 0.6 or 0.7); (d) the effect of rank initialization on computational cost (missing rate = 0.7)

Noisy synthetic data

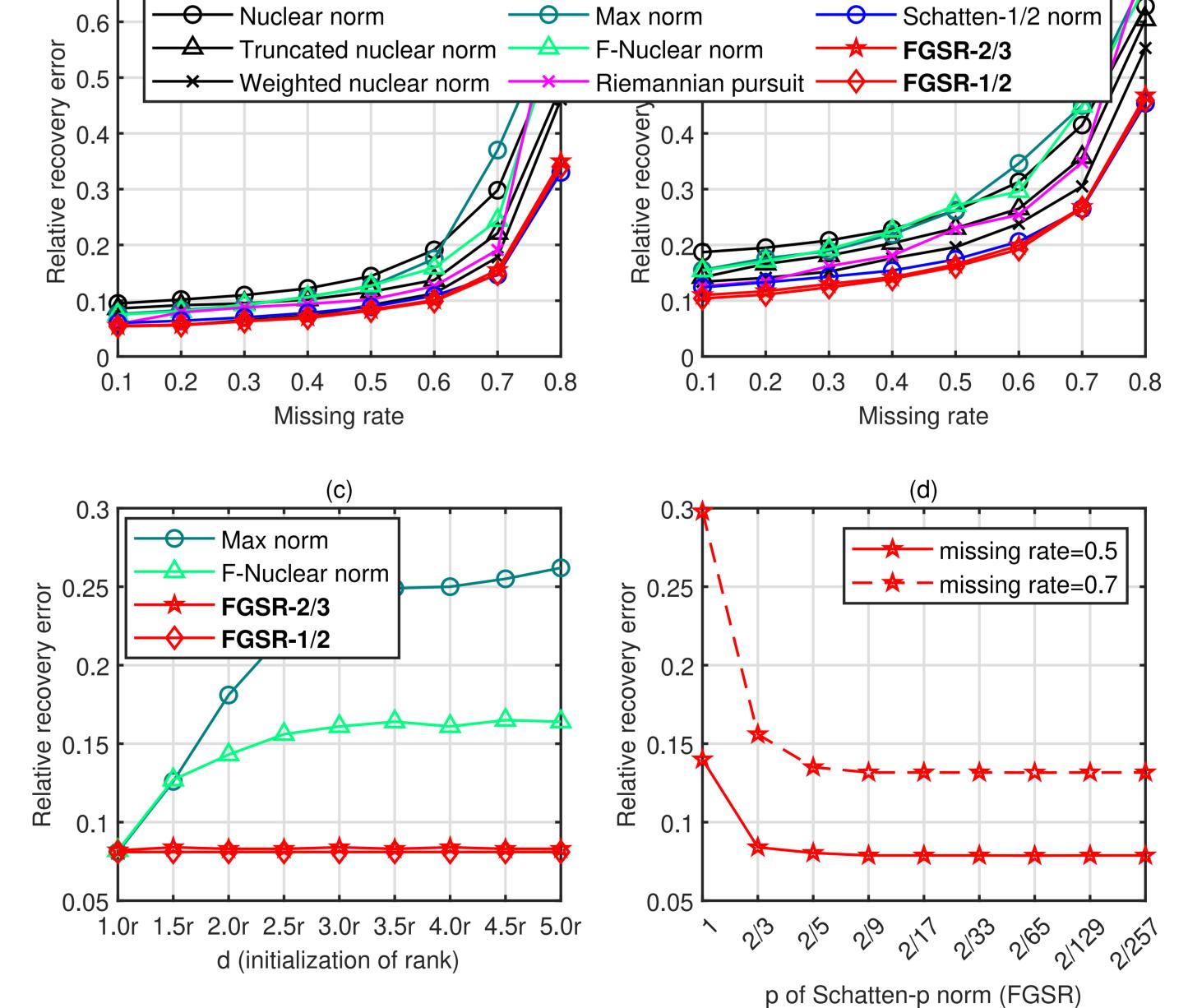


Figure: Matrix completion on noisy synthetic data: (a)(b) recovery error when SNR = 10 or 5; (c) the effect of rank initialization on recovery error (SNR = 10, missing rate = 0.5); (d) the effect of p in Schatten-p norm (using FGSR when p < 1).

Similar results for the MovieLens-1M dataset (shown in the paper).

Reference

- [1] Srebro, Nathan, Jason Rennie, and Tommi S. Jaakkola. NIPS, 2005.
- [2] Shang, Fanhua, Yuanyuan Liu, and James Cheng. AISTATS, 2016.
- [3] Shang, Fanhua, James Cheng, Yuanyuan Liu, Zhi-Quan Luo, and Zhouchen Lin. IEEE TPAMI, 2017
- [4] Qinghua Liu, Xinyue Shen, and Yuantao Gu. arXiv preprint arXiv, 2017
- [5] Bolte, Jérôme, Shoham Sabach, and Marc Teboulle. Mathematical Programming 146, no. 1-2 (2014): 459-494.