Yuejie Chi

# Low-Rank Matrix Completion

Imagine one observes a small subset of entries in a large matrix and aims to recover the entire matrix. Without a priori knowledge of the matrix, this problem is highly ill-posed. Fortunately, data matrices often exhibit low-dimensional structures that can be used effectively to regularize the solution space. The celebrated effectiveness of principal component analysis (PCA) in science and engineering suggests that most variability of real-world data can be accounted for by projecting the data onto a few directions known as the principal components. Correspondingly, the data matrix can be modeled as a low-rank matrix, at least approximately. Is it possible to complete a partially observed matrix if its rank, i.e., its maximum number of linearly independent row or column vectors, is small?

Low-rank matrix completion arises in a variety of applications in recommendation systems, computer vision, and signal processing. As a motivating example, consider users' ratings of products arranged in a rating matrix. Each rating may only be affected by a small number of factors—such as price, quality, and utility—and how they are reflected on the products' specifications and users' expectations. Naturally, this suggests that the rating matrix is low rank, since the numbers of users and products are much higher than the number of factors. Often, the rating matrix is sparsely observed, and it is of great interest to predict the missing ratings to make targeted recommendations.

## Relevance

The theory and algorithms of low-rank matrix completion have been significantly expanded in the last decade with converging efforts from signal processing, applied mathematics, statistics, optimization, and machine learning. This lecture note provides an introductory exposition of some key results in this rapidly developing field.

## Prerequisites

We expect the readers to be familiar with basic concepts in linear algebra, optimization, and probability.

## Problem statement

Let $M \in \mathbb{R}^{n_1 \times n_2}$ be a rank-$r$ matrix, whose thin singular value decomposition (SVD) is given as

$$M = U\Sigma V^{\top}, \qquad (1)$$

where $U \in \mathbb{R}^{n_1 \times r}, V \in \mathbb{R}^{n_2 \times r}$ are composed of orthonormal columns, and $\Sigma$ is an $r$-dimensional diagonal matrix with the singular values arranged in a nonincreasing order, i.e., $\sigma_1 \geq \cdots \geq \sigma_r > 0$. The "degrees of freedom" of $M$ is $(n_1 + n_2 - r)r$, which is the total number of parameters we need to uniquely specify $M$.

Assume we are given partial observations of $M$ over an index set $\Omega \subset \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, n_2\}$. To concisely put it, define the observation operator $\mathcal{P}_\Omega: \mathbb{R}^{n_1 \times n_2} \to \mathbb{R}^{n_1 \times n_2}$ as

$$[\mathcal{P}_\Omega(M)]_{ij} = \begin{cases} M_{ij}, & (i,j) \in \Omega \\ 0, & \text{otherwise} \end{cases}.$$

Our goal is to recover $M$ from $\mathcal{P}_\Omega(M)$, when the number of observation $m = |\Omega| \ll n_1 n_2$ is much smaller than the number of entries in $M$, under the assumption that $M$ is low rank, i.e., $r \ll \min\{n_1, n_2\}$. For notational simplicity in the sequel, let $n = \max\{n_1, n_2\}$.

## Solution

### Which low-rank matrices can we complete?

To begin, we ask the following question: What kind of low-rank matrices can we complete? As motivation, consider the following $4 \times 4$ rank-1 matrices $M_1$ and $M_2$, given as

$$M_1 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad M_2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}.$$

The matrix $M_1$ is more difficult to complete, since most of its entries are zero, and we need to collect more measurements to make sure enough mass comes from its nonzero entries. In contrast, the mass of $M_2$ is more uniformly distributed across all entries, making it easier to propagate information from one entry to another.

> To put it differently, a low-rank matrix is easier to complete if its energy spreads evenly across different coordinates.

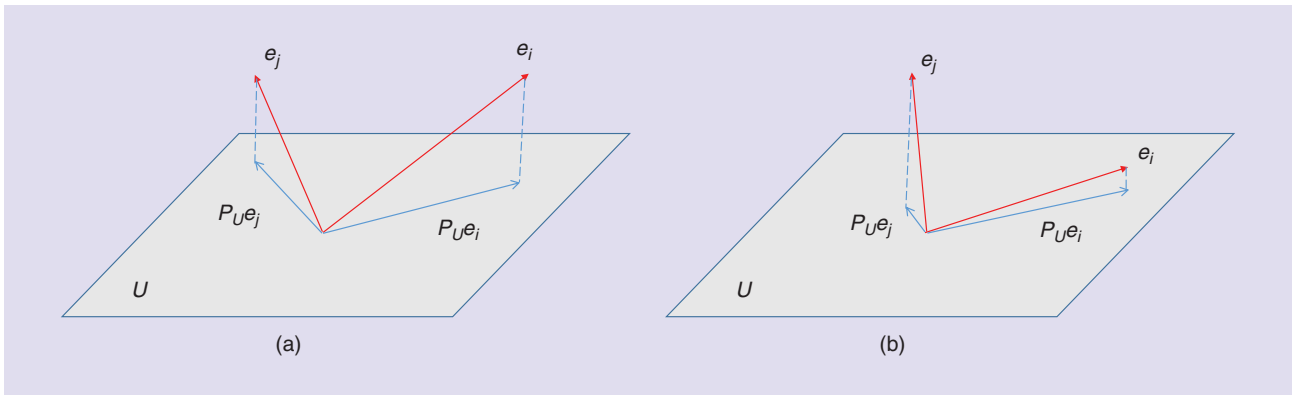To put it differently, a low-rank matrix is easier to complete if its energy spreads evenly across different coordinates. This property is captured by the notion of *coherence* [1], which measures the alignment between the column/row spaces of the low-rank matrix with standard basis vectors. For a matrix $U \in \mathbb{R}^{n_1 \times r}$ with orthonormal columns, let $P_U$ be the orthogonal projection onto the column space of $U$. The coherence parameter of $U$ is defined as

**FIGURE 1.** An illustration of the coherence parameter $\mu(U)$. $\mu(U)$ is small when all the standard basis vectors $e_i$ have approximately the same projections onto the subspace $U$, as shown in (a); $\mu(U)$ is large if $U$ is too aligned with certain standard basis vector, as shown in (b).

$$\mu(U) = \frac{n_1}{r} \max_{1 \le i \le n_1} \| P_U e_i \|_2^2$$
$$= \frac{n_1}{r} \max_{1 \le i \le n_1} \| U^\top e_i \|_2^2, \qquad (2)$$

where $e_i$ is the $i$th standard basis vector. Figure 1 provides a geometric illustration of the coherence parameter $\mu(U)$.

For a low-matrix $M$ whose SVD is given in (1), the coherence of $M$ is defined as

$$\mu = \max \{ \mu(U), \mu(V) \}. \qquad (3)$$

Notably, the coherence $\mu$ is determined by the singular vectors of $M$ and independent of its singular values. Since $1 \le \mu(U) \le n_1/r$ and $1 \le \mu(V) \le n_2/r$, we have $1 \le \mu \le n/r$. In the previous example, the coherence of $M_1$ matches the upper bound $n/r$, while the coherence of $M_2$ matches the lower bound one. The smaller $\mu$ is, the easier it is to complete the matrix.

### Which observation patterns can we handle?

Low-rank matrix completion can still be hopeless even when most of the entries are revealed. Consider, for example, the following observation pattern for a $4 \times 4$ matrix:

$$\begin{bmatrix} \star & \star & \star & ? \\ \star & \star & \star & ? \\ \star & \star & \star & ? \\ \star & \star & \star & ? \end{bmatrix},$$

where $\star$ indicates an observed entry, and ? indicates a missing entry. The last column of the matrix cannot be recovered since it can lie anywhere in the column space of the low-rank matrix. Therefore, we require at least $r$ observations per column/row. To bypass such pessimistic observation patterns, it is useful to think of random observation patterns. A popular choice is the Bernoulli model, where each entry is observed independently and identically with probability $p := m/(n_1 n_2)$. By a coupon-collecting argument [2], under the Bernoulli model, it is impossible to recover a low-rank matrix with less than some constant times $\mu n r \log n$ measurements using any algorithm, which is referred to as the *information-theoretic lower bound*. Compared with the degrees of freedom, which is on the order of $nr$, we pay a price in sample complexity by a factor of $\mu \log n$, highlighting again the role of coherence in low-rank matrix completion.

### Matrix completion via convex optimization

We present the first algorithm based on convex optimization. To promote the low-rank structure of the solution, a natural heuristic is to find the matrix with the minimum rank that is consistent with the observations, leading to

$$\min_{\Phi \in \mathbb{R}^{n_1 \times n_2}} \text{rank} (\Phi)$$
$$\text{s.t.} \quad \mathcal{P}_\Omega(\Phi) = \mathcal{P}_\Omega(M). \qquad (4)$$

However, since rank minimization is NP-hard, the above formulation is intractable. Motivated by the success of $\ell_1$ norm minimization for sparse recovery in compressed sensing [3], we consider convex relaxation for the rank heuristic. Observing that the rank of $\Phi$ equals to the number of its nonzero singular values, we replace rank ($\Phi$) by the sum of its singular values, denoted as the nuclear norm:

$$\| \Phi \|_* \triangleq \sum_{i=1}^{\min \{n_1, n_2\}} \sigma_i(\Phi),$$

where $\sigma_i(\Phi)$ is the $i$th singular value of $\Phi$. The nuclear norm is the tightest convex relaxation of the rank constraint, i.e., the nuclear norm ball $\{ \Phi : \| \Phi \|_* \le 1 \}$ is the convex hull of the collection of unit-norm rank-1 matrices: $\{ uv^\top : \| u \|_2 = \| v \|_2 = 1 \}$. Notably, the nuclear norm is also unitarily invariant, and can be represented as the solution to a semidefinite program,

$$\| \Phi \|_* = \min_{W_1, W_2} \frac{1}{2}(\text{Tr}(W_1) + \text{Tr}(W_2))$$
$$\text{s.t.} \quad \begin{bmatrix} W_1 & \Phi \\ \Phi^\top & W_2 \end{bmatrix} \succeq 0.$$

Hence, instead of solving (4) directly, we solve nuclear norm minimization, which searches for a matrix with the minimum nuclear norm that satisfies all the measurements:

$$\min_{\Phi \in \mathbb{R}^{n_1 \times n_2}} \| \Phi \|_* \quad \text{s.t.} \quad \mathcal{P}_\Omega(\Phi) = \mathcal{P}_\Omega(M). \qquad (5)$$

This gives a convex program that can be solved efficiently in polynomial time. Moreover, it doesn't require knowledge of the rank a priori.

The performance of nuclear norm minimization has been investigated in a recent line of elegant works [2]–[5], which suggests it can exactly recover a low-rank matrix as soon as the number
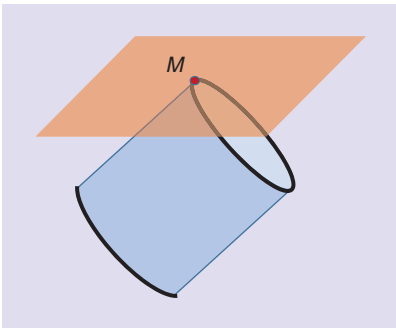
**FIGURE 2.** A geometric illustration of nuclear norm minimization: the cylinder represents level sets of the nuclear norm, and the hyperplane represents the measurement constraint. The two sets intersect at the thickened edges, which correspond to low-rank solutions.

of measurements is slightly larger than the information-theoretic lower bound by a logarithmic factor. Suppose that each entry of $M$ is observed independently with probability $p \in (0,1)$. If $p$ satisfies

$$p \geq C \frac{\mu r \log^2 n}{n},$$

for some large enough constant $C > 0$, then with high probability, the nuclear norm minimization algorithm (5) exactly recovers $M$ as the unique optimal solution of (5). Figure 2 illustrates the geometry of nuclear norm minimization when the number of measurements is sufficiently large. When both $\mu$ and $r$ are much smaller than $n$, this means we can recover a low-rank matrix even when the proportion of observations is vanishingly small.

### Matrix completion via nonconvex optimization

The computational and memory complexities of nuclear norm minimization can be quite expensive for large-scale problems, even with first-order methods, due to optimizing over and storing the matrix variable $\Phi$. Therefore, it is necessary to consider alternative approaches whose complexities scale more favorably in $n$. This leads to the second algorithm based on gradient descent using a proper initialization. If the rank of the matrix $M$ is known, it is natural to incorporate this knowledge and consider a rank-constrained least-squares problem

$$\min_{\Phi \in \mathbb{R}^{n_1 \times n_2}} \| \mathcal{P}_\Omega(\Phi - M) \|_{\mathrm{F}}^2,$$
$$\text{s.t rank}(\Phi) \leq r, \tag{6}$$

where $\| \cdot \|_{\mathrm{F}}$ is the Frobenius norm of a matrix. Invoking the low-rank factorization $\Phi = XY^\top$, where $X \in \mathbb{R}^{n_1 \times r}$ and $Y \in \mathbb{R}^{n_2 \times r}$, we can rewrite (6) as an unconstrained, yet nonconvex optimization problem:

$$\min_{X,Y} f(X,Y) := \| \mathcal{P}_\Omega(XY^\top - M) \|_{\mathrm{F}}^2. \tag{7}$$

On one end, the memory complexities of $X$ and $Y$ are linear in $n$ instead of quadratic in $n$ when dealing with $\Phi$. On the other end, we can only determine $X$ and $Y$ up to invertible transforms in (7), since for any invertible matrix $Q \in \mathbb{R}^{r \times r}$, we have $XY^\top = (XQ)(YQ^{-\top})^\top$. To fix the scaling ambiguity, it is useful to consider a modified loss function

$$F(X,Y) = \frac{1}{4p} f(X,Y)$$
$$+ \frac{1}{16} \| X^\top X - Y^\top Y \|_{\mathrm{F}}^2,$$

where the second term is introduced to motivate solutions where $X$ and $Y$ have balanced norms. The observation probability $p$, if not known, can be faithfully estimated by the sample proportion $|\Omega|/(n_1 n_2)$.

How do we optimize the nonconvex loss $F(X,Y)$? A plausible strategy proceeds in two steps.

1) The first step aims to find an initialization that is close to the ground truth, which can be provided via the so-called spectral method [6]. Consider the partially observed matrix $(1/p)\mathcal{P}_\Omega(M)$, which is an unbiased estimate of $M$ with expectation $\mathbb{E}[(1/p)\mathcal{P}_\Omega(M)] = M$. Therefore, its best rank-$r$ approximation produces a reasonably good initial guess. Let $U_0 \Sigma_0 V_0^\top$ be the best rank-$r$ approximation of $(1/p)\mathcal{P}_\Omega(M)$, where $U_0 \in \mathbb{R}^{n_1 \times r}, V_0 \in \mathbb{R}^{n_2 \times r}$ contain orthonormal columns and $\Sigma_0$ is an $r \times r$ diagonal matrix. The spectral initialization sets $X_0 = U_0 \Sigma_0^{1/2}$ and $Y_0 = V_0 \Sigma_0^{1/2}$.

2) The second step aims to refine the initial estimate locally via simple iterative methods, such as gradient descent [7], [8], following the update rule

$$\begin{bmatrix} X_{t+1} \\ Y_{t+1} \end{bmatrix} = \begin{bmatrix} X_t \\ Y_t \end{bmatrix} - \eta_t \begin{bmatrix} \nabla_X F(X_t, Y_t) \\ \nabla_Y F(X_t, Y_t) \end{bmatrix}, \tag{8}$$

where $\eta_t$ is the step size, and $\nabla_X F(X, Y), \nabla_Y F(X,Y)$ are the partial derivatives with respect to $X$ and $Y$ that can be derived easily.
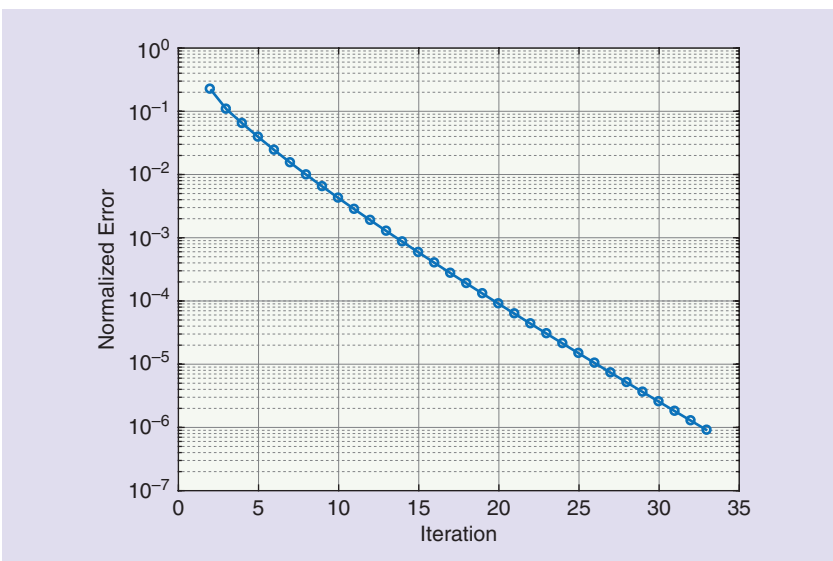


**FIGURE 3.** The normalized error of low-rank matrix completion with respect to the iteration count via gradient descent with the spectral initialization for a $10^4 \times 10^4$ matrix of rank-10 using about 5% observations.

Recall the SVD of $M$ in (1), and denote $X^\natural = U\Sigma^{1/2}$ and $Y^\natural = V\Sigma^{1/2}$; this allows us to write the factorization as $M = X^\natural Y^{\natural\top}$ and call $Z^\natural = [X^{\natural\top}, Y^{\natural\top}]^\top \in \mathbb{R}^{(n_1+n_2)\times r}$ the ground truth. Since $Z^\natural$ is only identifiable up to orthonormal transforms, let the optimal transform between the $t$th iterate $Z_t = [X_t^\top, Y_t^\top]^\top \in \mathbb{R}^{(n_1+n_2)\times r}$ and $Z^\natural$ as

$$H_t := \operatorname*{argmin}_{R\in\mathbb{R}^{r\times r}, RR^\top=I} \left\| Z_t R - Z^\natural \right\|_F.$$

Assume the condition number $\kappa := \sigma_1/\sigma_r$ of $M$ is a bounded constant, then as long as

$$p \geq C_1 \frac{\mu^3 r^3 \log^3 n}{n}$$

for some sufficiently large constant $C_1 > 0$, with high probability, the iterates satisfy [8]

$$\left\| Z_t H_t - Z^\natural \right\|_F \leq C_2 \rho^t \mu r \frac{1}{\sqrt{np}} \left\| Z^\natural \right\|_F,$$
$$\forall t \geq 0,$$

where $C_2 > 0, 0 < \rho < 1$ are some constants, provided that the step size $0 < \eta_t \equiv \eta \leq 2/(25\kappa\sigma_1)$. Hence, gradient descent converges at a geometric rate, as soon as the number of measurements is on the order of $\mu^3 r^3 n \log^3 n$, which scales linearly in $n$ up to logarithmic factors. To reach $\epsilon$-accuracy, i.e., $\left\| Z_t H_t - Z^\natural \right\|_F / \left\| Z^\natural \right\|_F \leq \epsilon$, gradient descent needs an order of $\log(1/\epsilon)$ iterations. The number of iterations is independent of the problem size and therefore the computational cost is much cheaper in conjunction with low cost per iteration.

## Summary

Table 1 summarizes the figures-of-merit of the discussed algorithms using state-of-the-art theory.

## Numerical example

Let $M$ be a rank-10 matrix of size $10^4 \times 10^4$ with about 5% of observed entries, i.e., $p = 0.05$, where $X^\natural$ and $Y^\natural$ are generated with i.i.d. standard Gaussian entries. We implement gradient descent with spectral initialization to recover $M$. Figure 3 plots the normalized error $\left\| X_t Y_t^\top - M \right\|_F / \left\| M \right\|_F$ with respect to the iteration counts, which verifies the geometric convergence predicted by the theory. Indeed, the normalized error is below $10^{-5}$ within 30 iterations!

## What we have learned

Under mild statistical models, low-rank matrix completion admits efficient algorithms with provable near-optimal performance guarantees, using both convex and nonconvex optimization techniques. The theory and algorithms discussed herein can be extended to recover matrices that are approximately low rank using noisy measurements. Low-rank matrix completion can be viewed as a special case of low-rank matrix estimation using an underdetermined set of linear equations. Other linear measurement patterns are also actively studied, motivated by applications such as sensor network localization, phase retrieval, quantum state tomography, and so on. Furthermore, low-rank matrix completion can be made robust even when many of the observations are corrupted by outliers of arbitrary magnitudes, known as the *sparse* and *low-rank decomposition problem* [9].

Low-rank structures are ubiquitous in modern data science problems and becoming increasingly popular as a modeling tool. Understanding the algorithmic and theoretical properties of estimation of low-rank structures is still an active area of research that will have a growing impact in future years. For a recent survey on low-rank matrix estimation, please see [10].

## Acknowledgments

## Author

*Yuejie Chi* (yuejiechi@cmu.edu) received her B.E. (Hon.) degree in electrical engineering from Tsinghua University, Beijing, China, in 2007 and her Ph.D. degree in electrical engineering from Princeton University, New Jersey, in 2012. She is currently an associate professor with the Department of Electrical and Computer Engineering at Carnegie Mellon University, Pittsburgh, Pennsylvania. Her research interests include statistical signal processing, machine learning, and large-scale optimization and their applications in data science, inverse problems, imaging, and sensing systems. She is a Senior Member of the IEEE.

## References

[1] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Found. Comput. Math.*, vol. 9, no. 6, pp. 717–772, Apr. 2009.

[2] E. J. Candès and T. Tao, "The power of convex relaxation: Near-optimal matrix completion," *IEEE Trans. Inf. Theory*, vol. 56, no. 5, pp. 2053–2080, 2010.

[3] D. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.

[4] D. Gross, "Recovering low-rank matrices from few coefficients in any basis," *IEEE Trans. Inf. Theory*, vol. 57, no. 3, pp. 1548–1566, Mar. 2011.

[5] Y. Chen, "Incoherence-optimal matrix completion," *IEEE Trans. Inf. Theory*, vol. 61, no. 5, pp. 2909–2923, 2015.

[6] R. H. Keshavan, A. Montanari, and S. Oh, "Matrix completion from a few entries," *IEEE Trans. Inf. Theory*, vol. 56, no. 6, pp. 2980–2998, 2010.

[7] R. Sun and Z.-Q. Luo, "Guaranteed matrix completion via non-convex factorization," *IEEE Trans. Inf. Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.

[8] C. Ma, K. Wang, Y. Chi, and Y. Chen, "Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion and blind deconvolution," arXiv Preprint, arXiv:1711.10467, 2017.

[9] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *J. ACM*, vol. 58, no. 3, pp. 11, 2011.

[10] Y. Chen and Y. Chi, "Harnessing structures in big data via guaranteed low-rank matrix estimation," *IEEE Signal Process. Mag.*, vol. 35, no. 4, pp. 14–31, July 2018.

**SP**

**Table 1. Figure-of-merits for low-rank matrix completion in terms of order-wise sample complexity and computational complexity.**

| | Sample Complexity | Computational Complexity |
|---|---|---|
| Information-theoretic lower bound | $\mu n r \log n$ | NP-hard |
| Nuclear norm minimization | $\mu n r \log^2 n$ | Polynomial time |
| Gradient descent with spectral initialization | $\mu^3 n r^3 \log^3 n$ | Linear time |