

GLOBAL OPTIMALITY IN LOW-RANK MATRIX OPTIMIZATION

Zhihui Zhu, Qiuwei Li, Gongguo Tang, and Michael B. Wakin

Department of Electrical Engineering, Colorado School of Mines, Golden, CO USA

ABSTRACT

This paper considers the minimization of a general objective function $f(\mathbf{X})$ over the set of non-square $n \times m$ matrices where the optimal solution \mathbf{X}^* is low-rank. To reduce the computational burden, we factorize the variable \mathbf{X} into a product of two smaller matrices and optimize over these two matrices instead of \mathbf{X} . We analyze the global geometry for a general and yet well-conditioned objective function $f(\mathbf{X})$ whose restricted strong convexity and restricted strong smoothness constants are comparable. In particular, we show that the reformulated objective function has no spurious local minima and obeys the strict saddle property. These geometric properties imply that a number of iterative optimization algorithms (such as gradient descent) can provably solve the factored problem with global convergence.

Index Terms— Low-rank matrix optimization, matrix sensing, nonconvex optimization, optimization geometry, strict saddle

1. INTRODUCTION

Consider the minimization of a general objective function $f(\mathbf{X})$ over all $n \times m$ matrices:

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times m}}{\text{minimize}} f(\mathbf{X}), \quad (1)$$

which we suppose admits a low-rank solution $\mathbf{X}^* \in \mathbb{R}^{n \times m}$ with $\text{rank}(\mathbf{X}^*) = r^*$. Low-rank matrix optimizations of the form (1) appear in a wide variety of applications, including quantum tomography [1], projection matrix design for compressive sensing [2], collaborative filtering [3], low-rank matrix recovery from compressive measurements [4], and matrix completion [5]. In order to find a low-rank solution, the nuclear norm is widely used in matrix inverse problems [4] arising in machine learning [6], signal processing [7], and control [8]. Although nuclear norm minimization enjoys strong statistical guarantees [5], its computational complexity is very high (as most algorithms require performing an expensive singular value decomposition (SVD) in each iteration), prohibiting it from scaling to practical problems.

This work was supported by NSF grant CCF-1409261, NSF grant CCF-1464205, and NSF CAREER grant CCF-1149225. Email: {zzhu, qiuli, gtang, mwakin}@mines.edu.

To relieve the computational bottleneck, recent studies propose to factorize the matrix variable into the Burer-Monteiro type decomposition [9, 10] with $\mathbf{X} = \mathbf{U}\mathbf{V}^T$, and optimize over the $n \times r$ and $m \times r$ ($r \geq r^*$) matrices \mathbf{U} and \mathbf{V} . With this parameterization of \mathbf{X} , we can recast (1) into the following program:

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}, \mathbf{V} \in \mathbb{R}^{m \times r}}{\text{minimize}} h(\mathbf{U}, \mathbf{V}) := f(\mathbf{U}\mathbf{V}^T). \quad (2)$$

The bilinear nature of the parameterization renders the objective function of (2) nonconvex even when $f(\mathbf{X})$ is a convex function. Hence, the objective function in (2) can potentially have spurious local minima (i.e., local minimizers that are not global minimizers) or “bad” saddle points that prevent a number of iterative algorithms from converging to the global solution. By analyzing the landscape of nonconvex functions, several recent works have shown that with an exact factorization ($r = r^*$), the factored objective function $h(\mathbf{U}, \mathbf{V})$ in matrix inverse problems has no spurious local minima [11–17].

We generalize this line of work by focusing on a general objective function $f(\mathbf{X})$ in the optimization (1), not necessarily a quadratic loss function coming from a matrix inverse problem. We provide a geometric analysis for the factored program (2) and show that all the critical points of the objective function are well-behaved. Our characterization of the geometry of the objective function ensures a number of iterative optimization algorithms converge to a global minimum.

This paper continues in Section 2 with formal definitions for strict saddles and the strict saddle property. We present the main results in Section 3 and their implications in matrix recovery in Section 4.

2. PRELIMINARIES

2.1. Notation

To begin, we first briefly introduce some notation used throughout the paper. The symbols \mathbf{I} and $\mathbf{0}$ respectively represent the identity matrix and zero matrix with appropriate sizes. The set of $r \times r$ orthonormal matrices is denoted by $\mathcal{O}_r := \{\mathbf{R} \in \mathbb{R}^{r \times r} : \mathbf{R}^T \mathbf{R} = \mathbf{I}\}$. If a function $h(\mathbf{U}, \mathbf{V})$ has two arguments, $\mathbf{U} \in \mathbb{R}^{n \times r}$ and $\mathbf{V} \in \mathbb{R}^{m \times r}$, we occasionally use the notation $h(\mathbf{W})$ when we put these two arguments into a new one as $\mathbf{W} = \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}$. For a scalar func-

tion $f(\mathbf{Z})$ with a matrix variable $\mathbf{Z} \in \mathbb{R}^{n \times m}$, its gradient is an $n \times m$ matrix whose (i, j) -th entry is $[\nabla f(\mathbf{Z})]_{ij} = \frac{\partial f(\mathbf{Z})}{\partial \mathbf{Z}_{ij}}$ for all $i \in \{1, 2, \dots, n\}, j \in \{1, 2, \dots, m\}$. The Hessian of $f(\mathbf{Z})$ can be represented with a bilinear form defined via $[\nabla^2 f(\mathbf{Z})](\mathbf{A}, \mathbf{B}) = \sum_{i,j,k,l} \frac{\partial^2 f(\mathbf{Z})}{\partial \mathbf{Z}_{ij} \partial \mathbf{Z}_{kl}} \mathbf{A}_{ij} \mathbf{B}_{kl}$ for any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times m}$.

2.2. Strict Saddle Property

Suppose $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable objective function. We begin with the notion of strict saddles and the strict saddle property.

Definition 1 (Strict saddles). *A critical point \mathbf{x} is a strict saddle if the Hessian matrix evaluated at this point has a strictly negative eigenvalue, i.e., $\lambda_{\min}(\nabla^2 h(\mathbf{x})) < 0$.*

Definition 2 (Strict saddle property [18]). *A twice differentiable function satisfies the strict saddle property if each critical point either corresponds to a local minimum or is a strict saddle.*

Intuitively, the strict saddle property requires a function to have a directional negative curvature at all the critical points but local minima. This property allows a number of iterative algorithms such as noisy gradient descent [18] and the trust region method [19] to further decrease the function value at all the strict saddles and thus converge to a local minimum.

Theorem 1. [18, 20, 21] (informal) *For a twice continuously differentiable objective function satisfying the strict saddle property, a number of iterative optimization algorithms (such as gradient descent and the trust region method) can find a local minimum.*

3. PROBLEM FORMULATION AND MAIN RESULTS

3.1. Problem Formulation

This paper considers the problem (1) of minimizing a general function $f(\mathbf{X})$ admitting a low-rank solution \mathbf{X}^* with $\text{rank}(\mathbf{X}^*) = r^* \leq r$. We factorize the variable $\mathbf{X} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U} \in \mathbb{R}^{n \times r}, \mathbf{V} \in \mathbb{R}^{m \times r}$ and transform (1) into its factored counterpart (2). Throughout the paper, \mathbf{X}, \mathbf{W} and $\widehat{\mathbf{W}}$ are matrices depending on \mathbf{U} and \mathbf{V} :

$$\mathbf{W} = \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}, \quad \widehat{\mathbf{W}} = \begin{bmatrix} \mathbf{U} \\ -\mathbf{V} \end{bmatrix}, \quad \mathbf{X} = \mathbf{U}\mathbf{V}^T.$$

Although the new variable \mathbf{W} has much smaller size than \mathbf{X} when $r \ll \min\{n, m\}$, the objective function in the factored problem (2) may have a much more complicated landscape due to the bilinear form about \mathbf{U} and \mathbf{V} . The reformulated objective function $h(\mathbf{U}, \mathbf{V})$ could introduce spurious local minima or degenerate saddle points even when $f(\mathbf{X})$ is convex. Our goal is to guarantee that this does not happen.

Let $\mathbf{X}^* = \mathbf{Q}_{\mathbf{U}^*} \mathbf{\Sigma}^* \mathbf{Q}_{\mathbf{V}^*}^T$ denote an SVD of \mathbf{X}^* , where $\mathbf{Q}_{\mathbf{U}^*} \in \mathbb{R}^{n \times r}$ and $\mathbf{Q}_{\mathbf{V}^*} \in \mathbb{R}^{m \times r}$ are orthonormal matrices of appropriate sizes, and $\mathbf{\Sigma}^* \in \mathbb{R}^{r \times r}$ is a diagonal matrix with non-negative diagonals (but with some zero diagonals if $r > r^* = \text{rank}(\mathbf{X}^*)$). We denote

$$\mathbf{U}^* = \mathbf{Q}_{\mathbf{U}^*} \mathbf{\Sigma}^{*1/2}, \quad \mathbf{V}^* = \mathbf{Q}_{\mathbf{V}^*} \mathbf{\Sigma}^{*1/2},$$

where $\mathbf{X}^* = \mathbf{U}^* \mathbf{V}^{*T}$ forms a balanced factorization of \mathbf{X}^* since \mathbf{U}^* and \mathbf{V}^* have the same singular values. Throughout the paper, we utilize the following two ways to stack \mathbf{U}^* and \mathbf{V}^* together:

$$\mathbf{W}^* = \begin{bmatrix} \mathbf{U}^* \\ \mathbf{V}^* \end{bmatrix}, \quad \widehat{\mathbf{W}}^* = \begin{bmatrix} \mathbf{U}^* \\ -\mathbf{V}^* \end{bmatrix}.$$

Before moving on, we note that for any solution (\mathbf{U}, \mathbf{V}) to (2), $(\mathbf{U}\mathbf{\Psi}, \mathbf{V}\mathbf{\Phi})$ is also a solution to (2) for any $\mathbf{\Psi}, \mathbf{\Phi} \in \mathbb{R}^{r \times r}$ such that $\mathbf{U}\mathbf{\Psi}\mathbf{\Phi}^T\mathbf{V}^T = \mathbf{U}\mathbf{V}^T$. In order to address this ambiguity (i.e., to reduce the search space of \mathbf{W} for (2)), we utilize the trick in [13, 22, 23] by introducing a regularizer

$$g(\mathbf{U}, \mathbf{V}) = \frac{\mu}{4} \|\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V}\|_F^2 \quad (3)$$

and solving the following problem

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}, \mathbf{V} \in \mathbb{R}^{m \times r}}{\text{minimize}} \quad \rho(\mathbf{U}, \mathbf{V}) := f(\mathbf{U}\mathbf{V}^T) + g(\mathbf{U}, \mathbf{V}), \quad (4)$$

where $\mu > 0$ controls the weight for the term $\|\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V}\|_F^2$, which will be discussed soon.

We remark that \mathbf{W}^* is still a global minimizer to the factored problem (4) since $f(\mathbf{X})$ and $g(\mathbf{W})$ achieve their global minimum at \mathbf{X}^* and \mathbf{W}^* , respectively. The regularizer $g(\mathbf{W})$ is applied to force the difference between the two Gram matrices of \mathbf{U} and \mathbf{V} to be as small as possible. The global minimum of $g(\mathbf{W})$ is 0, which is achieved when \mathbf{U} and \mathbf{V} have the same Gram matrices, i.e., when \mathbf{W} belongs to

$$\mathcal{E} := \left\{ \mathbf{W} = \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix} : \mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V} = \mathbf{0} \right\}. \quad (5)$$

Informally, we can view (4) as finding a point from \mathcal{E} that also minimizes $f(\mathbf{U}\mathbf{V}^T)$. This is formally established in Theorem 2.

3.2. Main Results

Before presenting our main results, we lay out the necessary assumptions on the objective function $f(\mathbf{X})$. As is known, without any assumptions on the problem, even minimizing traditional quadratic objective functions is challenging. For this purpose, we focus on the model where $f(\mathbf{X})$ is $(2r, 4r)$ -restricted strongly convex and smooth, i.e., for any $n \times m$ matrices \mathbf{X}, \mathbf{G} with $\text{rank}(\mathbf{X}) \leq 2r$ and $\text{rank}(\mathbf{G}) \leq 4r$, the Hessian of $f(\mathbf{X})$ satisfies

$$\alpha \|\mathbf{G}\|_F^2 \leq [\nabla^2 f(\mathbf{X})](\mathbf{G}, \mathbf{G}) \leq \beta \|\mathbf{G}\|_F^2 \quad (6)$$

for some positive α and β . A similar assumption is also utilized in [23, Conditions 5.3 and 5.4]. With this assumption on $f(\mathbf{X})$, we now summarize our main results.

Our main argument is that the objective function $\rho(\mathbf{W})$ has no spurious local minima and satisfies the strict saddle property. This is equivalent to categorizing all the critical points into two types: 1) the global minima which correspond to the global solution of the original convex problem (1) and 2) strict saddles such that the Hessian matrix $\nabla^2 \rho(\mathbf{W})$ evaluated at these points has a strictly negative eigenvalue. We formally establish this in the following theorem, whose proof is omitted due to space limitations.

Theorem 2. For any $\mu > 0$, each critical point $\mathbf{W} = \begin{bmatrix} \mathbf{U} \\ \mathbf{V} \end{bmatrix}$ of $\rho(\mathbf{W})$ defined in (4) satisfies

$$\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V} = \mathbf{0}. \quad (7)$$

Furthermore, suppose the function $f(\mathbf{X})$ satisfies the $(2r, 4r)$ -restricted strong convexity and smoothness condition (6) with positive constants α and β satisfying $\frac{\beta}{\alpha} \leq 1.5$. Set $\mu \leq \frac{\alpha}{16}$ for the factored problem (4). Then $\rho(\mathbf{W})$ has no spurious local minimum, i.e., any local minimum of $\rho(\mathbf{W})$ is a global minimum corresponding to the global solution of the original convex problem (1): $\mathbf{U}\mathbf{V}^T = \mathbf{X}^*$. In addition, $\rho(\mathbf{W})$ obeys the strict saddle property that any critical point not being a local minimum is a strict saddle with

$$\lambda_{\min}(\nabla^2(\rho(\mathbf{W}))) \leq \begin{cases} -0.08\alpha\sigma_r(\mathbf{X}^*), & r = r^* \\ -0.05\alpha \cdot \min\{\sigma_{r^c}^2(\mathbf{W}), 2\sigma_{r^*}(\mathbf{X}^*)\}, & r > r^* \\ -0.1\alpha\sigma_{r^*}(\mathbf{X}^*), & r_c = 0, \end{cases} \quad (8)$$

where $r^c \leq r$ is the rank of \mathbf{W} , $\lambda_{\min}(\cdot)$ represents the smallest eigenvalue, and $\sigma_\ell(\cdot)$ denotes the ℓ -th largest singular value.

Remark 1. The above result implies that we can recover the rank- r^* global minimizer \mathbf{X}^* of (1) by many iterative algorithms (such as the trust region method [24] and stochastic gradient descent [18]) even from a random initialization. This is because 1) as guaranteed by Theorem 1, the strict saddle property ensures local search algorithms converge to a local minimum, and 2) there are no spurious local minima.

Remark 2. Since our main result only requires the $(2r, 4r)$ -restricted strong convexity and smoothness property (6), aside from low-rank matrix recovery [25], it can also be applied to many other low-rank matrix optimization problems [26] which do not necessarily involve quadratic loss functions. Typical examples include robust PCA [27], 1-bit matrix completion [28] and Poisson principal component analysis (PCA) [29].

Remark 3. Equation (7) shows that any critical point \mathbf{W} belongs to \mathcal{E} for the objective function in the factored problem (4) with any positive μ . This demonstrates the reason

for adding the regularizer $g(\mathbf{U}, \mathbf{V})$. Thus, any iterative optimization algorithm converging to some critical point of $\rho(\mathbf{W})$ results in a solution within \mathcal{E} .

Remark 4. For any critical point $\mathbf{W} \in \mathbb{R}^{(n+m) \times r}$ but not being a local minimum, the right hand side of (8) is strictly negative, implying \mathbf{W} is a strict saddle. We also note that Theorem 2 not only covers exact parameterization where $r = r^*$, but also includes over-parameterization where $r > r^*$.

Remark 5. The constants appearing in Theorem 2 are not optimized. We use $\mu \leq \frac{1}{16}\alpha$ simply to include $\mu = \frac{1}{16}$ which is utilized for the matrix sensing problem in [22, p.3]. If the ratio between the restricted strong convexity and smoothness constants $\frac{\beta}{\alpha} \leq 1.4$, then we can show $\rho(\mathbf{W})$ has no spurious local minima and obeys the strict saddle property for any $\mu \leq \frac{1}{4}\alpha$ (where $\mu = \frac{1}{4}$ is utilized for the matrix sensing problem in [13, p.3]). In all cases, a smaller μ yields a more negative constant in (8). This implies that when the restricted strong convexity constant α is not provided a priori, one can always choose a small μ to ensure the strict saddle property holds, and hence guarantee the global convergence of many iterative optimization algorithms. We finally note that the regularizer $g(\mathbf{W})$ is added mostly to avoid bad iterates \mathbf{W} . If we apply local search algorithms with a random initialization, then it is possible to drop the regularizer (i.e., set $\mu = 0$) for practical implementation; see the experiments in Section 4.1.

4. STYLIZED APPLICATION: MATRIX RECOVERY

We first consider the implication of Theorem 2 in the matrix sensing problem where

$$f(\mathbf{X}) = \frac{1}{2} \|\mathcal{A}(\mathbf{X} - \mathbf{X}^*)\|_2^2.$$

Here $\mathcal{A} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^p$ is a known measurement operator satisfying the following restricted isometry property.

Definition 3. (Restricted Isometry Property (RIP) [4, p.11]) The map $\mathcal{A} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^p$ satisfies the r -RIP with constant δ_r if

$$(1 - \delta_r) \|\mathbf{X}\|_F^2 \leq \|\mathcal{A}(\mathbf{X})\|^2 \leq (1 + \delta_r) \|\mathbf{X}\|_F^2 \quad (9)$$

holds for any $n \times m$ matrix \mathbf{X} with $\text{rank}(\mathbf{X}) \leq r$.

Note that in this case, the Hessian quadrature form $\nabla^2 f(\mathbf{X})[\mathbf{Y}, \mathbf{Y}]$ for any $n \times m$ matrices \mathbf{X} and \mathbf{Y} is given by

$$\nabla^2 f(\mathbf{X})[\mathbf{Y}, \mathbf{Y}] = \|\mathcal{A}(\mathbf{Y})\|^2.$$

If \mathcal{A} satisfies the $4r$ -restricted isometry property with constant δ_{4r} , then $f(\mathbf{X})$ satisfies the $(2r, 4r)$ -restricted strong convexity and smoothness condition (6) with constants $\alpha = 1 - \delta_{4r}$ and $\beta = 1 + \delta_{4r}$ since

$$(1 - \delta_{4r}) \|\mathbf{Y}\|_F^2 \leq \|\mathcal{A}(\mathbf{Y})\|^2 \leq (1 + \delta_{4r}) \|\mathbf{Y}\|_F^2$$

for any rank- $4r$ matrix \mathbf{Y} . Now applying Theorem 2, we can characterize the geometry for the following matrix sensing problem with the factorization approach:

$$\underset{\mathbf{U} \in \mathbb{R}^{n \times r}, \mathbf{V} \in \mathbb{R}^{m \times r}}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{A}(\mathbf{UV}^T - \mathbf{X}^*)\|_2^2 + g(\mathbf{U}, \mathbf{V}), \quad (10)$$

where $g(\mathbf{U}, \mathbf{V})$ is the added regularizer defined in (3).

Corollary 1. Suppose \mathcal{A} satisfies the $4r$ -RIP with constant $\delta_{4r} \leq \frac{1}{5}$, and set $\mu \leq \frac{1-\delta_{4r}}{16}$. Then the objective function in (10) has no spurious local minima and satisfies the strict saddle property.

This result follows directly from Theorem 2 by noting that $\frac{\beta}{\alpha} = \frac{1+\delta_{4r}}{1-\delta_{4r}} \leq 1.5$ if $\delta_{4r} \leq \frac{1}{5}$. We remark that Park et al. [13, Theorem 4.3] provided a similar geometric result for (10). Compared to their result which requires $\delta_{4r} \leq \frac{1}{100}$, our result has a much weaker requirement on the RIP of the measurement operator.

4.1. Experiments

In this section, we present some experiments to illustrate the performance of the factorization approach for matrix completion¹ where we want to recover a low-rank matrix \mathbf{X}^* from incomplete measurements $\{X_{ij}^*\}_{(i,j) \in \Omega}$, where $\Omega \subset [n] \times [m]$. We compare the performance of the matrix factorization approach with SVP [30], the convex approach, and singular value thresholding² (SVT) [31]. Let \mathcal{P}_Ω denote the projection onto the index set Ω . The convex approach (denoted by CVX) attempts to use the nuclear norm as a convex relaxation of the rankness and solves

$$\underset{\mathbf{X}}{\text{minimize}} \quad \|\mathbf{X}\|_*, \text{ subject to } \mathcal{P}_\Omega(\mathbf{X}) = \mathcal{P}_\Omega(\mathbf{X}^*). \quad (11)$$

In the first set of experiments, we set $n = m = 100$ and vary the rank r from 1 to 30. We generate a rank- r random matrix and randomly obtain p entries, i.e., $|\Omega| = p$. Figure 1 displays the phase transition for matrix factorization (with $\mu = 0$) solved by gradient descent with a random initialization, the SVP [30], the singular value thresholding (SVT) [31], and the convex approach. As can be seen, the matrix factorization approach has a similar phase transition to SVP, and is slightly better than SVT and the convex approach in terms of the number of measurements needed for successful recovery. Similar phase transitions for matrix factorization are also observed for different $\mu > 0$.

In the second set of experiments, we set $r = 5$ and $p = 3r(2n - r)$ (3 times the number of degrees of freedom within

a rank- r $n \times n$ matrix), and vary n from 40 to 5120. We compare the time needed for the four approaches in Figure 2; our matrix factorization approach is much faster than the other methods. The time savings for the matrix factorization approach comes from avoiding the SVD, which is needed both for SVT and SVP in each iteration. We also observe that the convex approach has highest computational complexity and is not scalable (which is the reason that we only present its time for n up to 640).

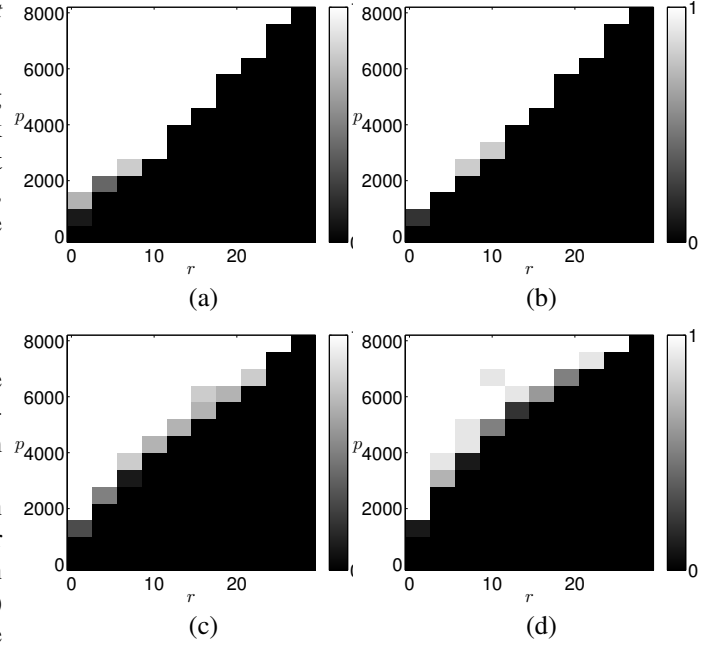


Fig. 1. Rate of success for matrix completion by (a) the matrix factorization approach with gradient descent; (b) SVP [30]; (c) solving the convex problem (11); (d) SVT [30]. 10 Monte Carlo trials are carried out and for each trial, we claim matrix recovery to be successful if the relative reconstruction error satisfies $\frac{\|\mathbf{X}^* - \hat{\mathbf{X}}\|_F}{\|\mathbf{X}^*\|_F} \leq 10^{-4}$.

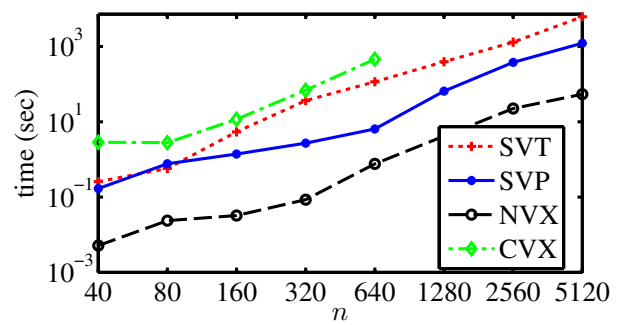


Fig. 2. Average computation time needed for different algorithms solving matrix completion.

¹Though \mathcal{P}_Ω does not satisfy the r -RIP (9) for all low-rank matrices \mathbf{X} , it satisfies the RIP when restricted to low-rank incoherent matrices [30, Theorem 4.2]. Thus, if the iterates of local search algorithms remain incoherent (which is experimentally observed), then Theorem 2 guarantees the global convergence of the matrix factorization approach with these algorithms.

²Software available at <http://svt.stanford.edu/>

5. REFERENCES

- [1] S. Aaronson, "The learnability of quantum states," in *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, vol. 463, pp. 3089–3114, The Royal Society, 2007.
- [2] T. Hong, H. Bai, S. Li, and Z. Zhu, "An efficient algorithm for designing projection matrix in compressive sensing based on alternating optimization," *Signal Processing*, vol. 125, pp. 9–20, 2016.
- [3] N. Srebro, J. Rennie, and T. S. Jaakkola, "Maximum-margin matrix factorization," in *Advances in Neural Information Processing Systems*, pp. 1329–1336, 2004.
- [4] B. Recht, M. Fazel, and P. A. Parrilo, "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization," *SIAM Review*, vol. 52, no. 3, pp. 471–501, 2010.
- [5] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Foundations of Computational Mathematics*, vol. 9, no. 6, pp. 717–772, 2009.
- [6] Z. Harchaoui, M. Douze, M. Paulin, M. Dudik, and J. Malick, "Large-scale image classification with trace-norm regularization," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3386–3393, IEEE, 2012.
- [7] M. A. Davenport and J. Romberg, "An overview of low-rank matrix recovery from incomplete observations," *IEEE Journal of Selected Topics in Signal Processing*, vol. 10, no. 4, pp. 608–622, 2016.
- [8] K. Mohan and M. Fazel, "Reweighted nuclear norm minimization with application to system identification," in *Proceedings of the 2010 American Control Conference*, pp. 2953–2959, IEEE, 2010.
- [9] S. Burer and R. D. Monteiro, "A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization," *Mathematical Programming*, vol. 95, no. 2, pp. 329–357, 2003.
- [10] S. Burer and R. D. Monteiro, "Local minima and convergence in low-rank semidefinite programming," *Mathematical Programming*, vol. 103, no. 3, pp. 427–444, 2005.
- [11] S. Bhojanapalli, B. Neyshabur, and N. Srebro, "Global optimality of local search for low rank matrix recovery," *arXiv preprint arXiv:1605.07221*, 2016.
- [12] R. Ge, J. D. Lee, and T. Ma, "Matrix completion has no spurious local minimum," *arXiv preprint arXiv:1605.07272*, 2016.
- [13] D. Park, A. Kyrillidis, C. Caramanis, and S. Sanghavi, "Non-square matrix sensing without spurious local minima via the Burer-Monteiro approach," *arXiv preprint arXiv:1609.03240*, 2016.
- [14] Q. Li and G. Tang, "The nonconvex geometry of low-rank matrix optimizations with general objective functions," *arXiv preprint arXiv:1611.03060*, 2016.
- [15] X. Li, Z. Wang, J. Lu, R. Arora, J. Haupt, H. Liu, and T. Zhao, "Symmetry, saddle points, and global geometry of nonconvex matrix factorization," *arXiv preprint arXiv:1612.09296*, 2016.
- [16] Q. Li, Z. Zhu, and G. Tang, "Geometry of factored nuclear norm regularization," *arXiv preprint arXiv:1704.01265*, 2017.
- [17] Z. Zhu, Q. Li, G. Tang, and M. B. Wakin, "The global optimization geometry of nonsymmetric matrix factorization and sensing," *arXiv preprint arXiv:1703.01256*, 2017.
- [18] R. Ge, F. Huang, C. Jin, and Y. Yuan, "Escaping from saddle points—online stochastic gradient for tensor decomposition," in *Proceedings of The 28th Conference on Learning Theory*, pp. 797–842, 2015.
- [19] A. R. Conn, N. I. Gould, and P. L. Toint, *Trust region methods*. SIAM, 2000.
- [20] J. Sun, Q. Qu, and J. Wright, "When are nonconvex problems not scary?," *arXiv preprint arXiv:1510.06096*, 2015.
- [21] J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht, "Gradient descent converges to minimizers," *University of California, Berkeley*, vol. 1050, p. 16, 2016.
- [22] S. Tu, R. Boczar, M. Soltanolkotabi, and B. Recht, "Low-rank solutions of linear matrix equations via Procrustes flow," *arXiv preprint arXiv:1507.03566*, 2015.
- [23] L. Wang, X. Zhang, and Q. Gu, "A unified computational and statistical framework for nonconvex low-rank matrix estimation," *arXiv preprint arXiv:1610.05275*, 2016.
- [24] J. Sun, Q. Qu, and J. Wright, "A geometric analysis of phase retrieval," *arXiv preprint arXiv:1602.06664*, 2016.
- [25] E. J. Candès and Y. Plan, "Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements," *IEEE Transactions on Information Theory*, vol. 57, no. 4, pp. 2342–2359, 2011.
- [26] M. Udell, C. Horn, R. Zadeh, and S. Boyd, "Generalized low rank models," *arXiv preprint arXiv:1410.0342*, 2014.
- [27] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?," *Journal of the ACM (JACM)*, vol. 58, no. 3, p. 11, 2011.
- [28] M. A. Davenport, Y. Plan, E. van den Berg, and M. Wotter, "1-bit matrix completion," *Information and Inference*, vol. 3, no. 3, pp. 189–223, 2014.
- [29] J. Salmon, Z. Harmany, C.-A. Deledalle, and R. Willett, "Poisson noise reduction with non-local PCA," *Journal of Mathematical Imaging and Vision*, vol. 48, no. 2, pp. 279–294, 2014.
- [30] P. Jain, R. Meka, and I. S. Dhillon, "Guaranteed rank minimization via singular value projection," in *Advances in Neural Information Processing Systems*, pp. 937–945, 2010.
- [31] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.