

ON THE ASYMPTOTIC LINEAR CONVERGENCE OF GRADIENT DESCENT FOR NON-SYMMETRIC MATRIX COMPLETION

Trung Vu^{1,2} and Raviv Raich¹

¹ School of EECS, Oregon State University, Corvallis, OR 97331-5501

² Department of CSEE, University of Maryland, Baltimore County, MD 21250-0002
{vutru, raich}@oregonstate.edu

ABSTRACT

This paper studies a factorization-based gradient descent approach for non-symmetric matrix completion. We introduce an objective that includes an orthogonality regularization for one of the factors. Additionally, we introduce a scaling term to ensure that the two factors are of equal magnitude to improve the convergence speed. For the proposed objective, we analyze the exact linear convergence rate of gradient descent via the asymptotically linear update equation for the error matrix. Our proposed result is the first closed-form expression of the exact linear rate. To illustrate the correctness and tightness of our analysis, we compare the empirical convergence rate against the analytical rate. Additional numerical experiments are done to verify the efficacy of the scaling approach.

1. INTRODUCTION

Matrix completion has numerous applications in machine learning and signal processing, namely system identification [1], collaborative filtering [2], and dimension reduction [3]. This problem can be described as follows. Let $\mathbf{M} \in \mathbb{R}^{n \times m}$ be a rank r matrix with $1 \leq r \leq \min(m, n)$, and $\Omega = \{(i, j) \mid M_{ij} \text{ is observed}\}$ be an index subset of cardinality s such that $s \leq mn$. We wish to recover the unobserved entries of \mathbf{M} by solving a rank minimization problem, subject to linear constraints consistent with the observations [4].

A body of work has been devoted to using first-order iterative methods for low-rank matrix completion that enforce the low-rank constraint using a factorization approach. By reparametrizing the $m \times n$ matrix as the product of two smaller matrices $\mathbf{M} = \mathbf{A}\mathbf{B}^\top$, for $\mathbf{A} \in \mathbb{R}^{n \times r}$ and $\mathbf{B} \in \mathbb{R}^{m \times r}$, the resulting equivalent problem is unconstrained and more computationally efficient to solve [5]. Although this problem is non-convex, recent progress shows that for such a problem any local minimum is also a global minimum [6, 7]. Thus, basic optimization algorithms such as gradient descent [6, 8, 9] and alternating minimization [10–13] can provably solve matrix completion under a specific sampling regime. Existing convergence analyses of algorithms for low-rank matrix completion often rely on *standard assumptions*, such as the incoherence of the underlying matrix \mathbf{M} and the uniform randomness of the sampling pattern Ω [4]. Under these assumptions and a sample complexity bound on the number of observed entries s , linear convergence to a global solution can be guaranteed (see [11] for alternating minimization, [9] for factorization-based gradient descent, and [14] for iterative hard thresholding), with an upper bound on the rate of convergence $\rho < 0.5$.

In this paper, we present a gradient descent (GD) approach for matrix completion that relies on the factorization of the matrix in the form of $\mathbf{A}\mathbf{B}^\top$. To ensure the linear convergence of gradient descent,

we incorporate orthogonality constraints (up to a scaling factor) on the left factor \mathbf{A} . Our focus is on the theoretical analysis of the asymptotic convergence rate. By exploiting the *local* structure of the problem, we characterize the exact linear rate of local convergence of the algorithm. The closed-form expression we obtained can be used to determine sufficient conditions that ensure local linear convergence. Moreover, since our expression is exact, one can identify conditions that are potentially less stringent than existing conditions. Similar to the previous work [15], we derive the asymptotically linear update equation for the error matrix and obtain the convergence rate in terms of this matrix. As a sanity check, we conduct numerical experiments to illustrate the correctness of the formula. We apply the algorithm to solve the matrix completion problem in a few scenarios and compare the empirical convergence rate against the analytical expression for the rate we derived. The agreement confirms our derivation. While other upper bounds for the rate of convergence of GD for matrix completion are available in the literature [6], the expression for the rate provided in this paper is the *first* to provide an exact prediction of the asymptotic rate of convergence of GD for non-symmetric matrix completion.

Notation Throughout the paper, we use the notations $\|\cdot\|_F$ and $\|\cdot\|_2$ to denote the Frobenius norm and the spectral norm of a matrix, respectively. Occasionally, $\|\cdot\|_2$ is used on a vector to denote the Euclidean norm. Boldfaced symbols are reserved for vectors and matrices, while the elements of a vector/matrix are unbold. In addition, \mathbf{I}_n denotes the $n \times n$ identity matrix, and $\mathbf{0}_{m \times n}$ denotes the $m \times n$ matrix of all zeros. We also use \otimes to denote the Kronecker product between two matrices. For a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, X_{ij} refers to the (i, j) element of \mathbf{X} . The spectral radius of \mathbf{X} is the largest absolute value of the eigenvalues of \mathbf{X} , denoted by $\rho(\mathbf{X})$. The notation $\text{vec}(\mathbf{X})$ denotes the vectorization of \mathbf{X} by stacking its columns on top of one another.

2. PROBLEM FORMULATION

In this paper, we consider a GD approach to solving the problem of low-rank matrix completion [6, 9]. In the (noiseless) low-rank matrix completion problem, the goal is to find a rank- r matrix \mathbf{X} such that

$$X_{ij} = M_{ij}, \quad (i, j) \in \Omega. \quad (1)$$

Here, \mathbf{X} and \mathbf{M} are $m \times n$ real-valued matrices and Ω is the sampling pattern that is a subset of $\{1, \dots, m\} \times \{1, \dots, n\}$. A factorization-based approach can be used to solve the aforementioned problem. Consider parameterizing \mathbf{X} as $\mathbf{X} = \mathbf{A}\mathbf{B}^\top$, where \mathbf{A} is $m \times r$ and \mathbf{B} is $n \times r$. This parameterization forces the rank of \mathbf{X} to be no more than r . Using this parametrization, one can

rephrase the low-rank matrix completion problem (in its noiseless case) as the problem of finding an $m \times r$ \mathbf{A} and an $n \times r$ \mathbf{B} such that $(\mathbf{A}\mathbf{B}^\top)_{ij} = M_{ij}$ for $(i, j) \in \Omega$. If we define the projection P_Ω such that

$$[P_\Omega(\mathbf{X})]_{ij} = \begin{cases} X_{ij} & (i, j) \in \Omega \\ 0 & (i, j) \notin \Omega \end{cases}, \quad (2)$$

then we can write the low-rank matrix completion problem as

$$P_\Omega(\mathbf{A}\mathbf{B}^\top - \mathbf{M}) = \mathbf{0}. \quad (3)$$

A simultaneous solution of all the equations in (3) is non-trivial. It is common to construct an objective function that is minimized when (3) holds. In [15], we analyzed the convergence of factorization-based gradient descent for such an objective function in the case of symmetric matrix completion. Adapting the convergence results for symmetric matrix completion in [15], one needs to perform some modification since the proposed matrix is not positive semi-definite and using symmetric matrix completion of the form $\mathbf{G} = \mathbf{X}\mathbf{X}^\top$ will not allow reducing the fitting error to zero. To remedy the problem, semidefinite lifting [16] in which $\mathbf{G} = [\mathbf{A}; \mathbf{B}][\mathbf{A}; \mathbf{B}]^\top = [\mathbf{A}\mathbf{A}^\top, \mathbf{M}; \mathbf{M}^\top, \mathbf{B}\mathbf{B}^\top]$ can be used, which in turn is compatible with symmetric matrix completion. However, the challenge is this method is not guaranteed to converge linearly due to ambiguity in the solution (\mathbf{A}, \mathbf{B}) . Moreover, a combination of semidefinite lifting and the analysis in [15] yields a rate of 1, i.e., no linear convergence guarantees. The following section presents a novel regularized objective for non-symmetric matrix completion and the details of the gradient descent approach for minimizing such an objective.

2.1. A Gradient Descent Approach

To solve this problem efficiently in a large-scale setting, we consider a gradient descent approach. Specifically, we consider the following objective

$$f(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \|P_\Omega(\mathbf{A}\mathbf{B}^\top - \mathbf{M})\|_F^2 + \frac{1}{4} \|\mathbf{A}^\top \mathbf{A} - c\mathbf{I}_r\|_F^2, \quad (4)$$

where¹

$$c = \sqrt{\frac{mn}{r|\Omega|}} \|P_\Omega(\mathbf{M})\|_F. \quad (5)$$

Note that the first term of the objective can be minimized to zero, by selecting \mathbf{A} and \mathbf{B} such that $P_\Omega(\mathbf{A}\mathbf{B}^\top - \mathbf{M}) = \mathbf{0}$, as in the low-rank matrix completion formulation in (3). The second term in the objective is used to alleviate the ambiguity of the solution. In particular, if the singular value decomposition of \mathbf{M} is given by $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ where \mathbf{U} is an $m \times r$ semi-orthogonal matrix, $\mathbf{\Sigma}$ is a non-negative diagonal matrix of dimension $r \times r$, and \mathbf{V} is an $n \times r$ semi-orthogonal matrix. Then, \mathbf{A} and \mathbf{B} that satisfy $\mathbf{A}\mathbf{B}^\top = \mathbf{M}$ are given by $\mathbf{A} = \mathbf{U}\mathbf{G}$ and $\mathbf{B} = \mathbf{V}\mathbf{\Sigma}(\mathbf{G}^{-1})^\top$ for any invertible $r \times r$ matrix \mathbf{G} . Minimizing the second term of the objective to zero implies a solution for \mathbf{A} of the form $\mathbf{A} = \sqrt{c}\mathbf{U}\mathbf{Q}$ where \mathbf{Q} is an orthogonal $r \times r$ matrix satisfying $\mathbf{Q}\mathbf{Q}^\top = \mathbf{Q}^\top\mathbf{Q} = \mathbf{I}_r$. We will show later, that this reduction of the ambiguity in the solution can allow us to establish a linear rate convergence. Note that the corresponding solution for \mathbf{B} is given by $\mathbf{B} = \frac{1}{\sqrt{c}}\mathbf{V}\mathbf{\Sigma}\mathbf{Q}$. Using the objective in (4), we can obtain the objective gradient with respect to

¹This choice of c is aimed to ensure that $\|\mathbf{A}\|_F \approx \|\mathbf{B}\|_F$. Due to space limitations, we omit the proof.

Algorithm 1 Factorization-based Gradient Descent

Input: $\mathbf{A}_0, \mathbf{B}_0, P_\Omega(\mathbf{M}), \eta$

Output: $\{\mathbf{A}_k, \mathbf{B}_k\}$

- 1: **for** $k = 0, 1, 2, \dots$ **do**
 - 2: $\mathbf{P}_k = P_\Omega(\mathbf{A}_k\mathbf{B}_k^\top - \mathbf{M})$
 - 3: $\mathbf{A}_{k+1} = \mathbf{A}_k - \eta(\mathbf{P}_k\mathbf{B}_k + \mathbf{A}_k(\mathbf{A}_k^\top\mathbf{A}_k - c\mathbf{I}_r))$
 - 4: $\mathbf{B}_{k+1} = \mathbf{B}_k - \eta\mathbf{P}_k^\top\mathbf{A}_k$
-

\mathbf{A} and \mathbf{B} as follows:

$$\begin{aligned} \frac{df}{d\mathbf{A}} &= P_\Omega(\mathbf{A}\mathbf{B}^\top - \mathbf{M})\mathbf{B} + \mathbf{A}(\mathbf{A}^\top\mathbf{A} - c\mathbf{I}_r) \\ \frac{df}{d\mathbf{B}} &= P_\Omega(\mathbf{A}\mathbf{B}^\top - \mathbf{M})^\top\mathbf{A}. \end{aligned} \quad (6)$$

Using the gradient, we can define the gradient descent iterations as:

$$\begin{aligned} \mathbf{A}_{k+1} &= \mathbf{A}_k - \eta[\mathbf{P}_k\mathbf{B}_k + \mathbf{A}_k(\mathbf{A}_k^\top\mathbf{A}_k - c\mathbf{I}_r)] \\ \mathbf{B}_{k+1} &= \mathbf{B}_k - \eta\mathbf{P}_k^\top\mathbf{A}_k, \end{aligned} \quad (7)$$

where $\mathbf{P}_k = P_\Omega(\mathbf{A}_k\mathbf{B}_k^\top - \mathbf{M})$ and η is the step-size parameter satisfying $\eta > 0$ (see Algorithm 1).

3. CONVERGENCE ANALYSIS

Due to the ambiguity in the solution for \mathbf{A} and \mathbf{B} , instead of considering convergence of \mathbf{A}_k and \mathbf{B}_k to \mathbf{A} and \mathbf{B} respectively, we consider the convergence of the product of matrices \mathbf{A}_k and \mathbf{B}_k with matrices \mathbf{A}_k^\top and \mathbf{B}_k^\top . In particular, we define the concatenation of \mathbf{A}_k and \mathbf{B}_k as $\mathbf{G}_k = [\mathbf{A}_k^\top\mathbf{B}_k^\top]^\top$ and similarly the concatenation of \mathbf{A} and \mathbf{B} as $\mathbf{G} = [\mathbf{A}^\top\mathbf{B}^\top]^\top$. We consider the convergence of $\mathbf{G}_k\mathbf{G}_k^\top$ to $\mathbf{G}\mathbf{G}^\top$. Note that in this approach, the limit $\mathbf{G}\mathbf{G}^\top$ is well-defined since $\mathbf{A}\mathbf{A}^\top, \mathbf{A}\mathbf{B}^\top$, and $\mathbf{B}\mathbf{B}^\top$ are all uniquely defined. Specifically, we have $\mathbf{A}\mathbf{A}^\top = c\mathbf{U}\mathbf{U}^\top$, $\mathbf{A}\mathbf{B}^\top = \mathbf{M}$, and $\mathbf{B}\mathbf{B}^\top = \frac{1}{c}\mathbf{V}\mathbf{\Sigma}^2\mathbf{V}^\top = \frac{1}{c}\mathbf{M}^\top\mathbf{M}$. Note that the ambiguous \mathbf{Q} in the solution for \mathbf{A} and \mathbf{B} cancels out in $\mathbf{A}\mathbf{A}^\top, \mathbf{A}\mathbf{B}^\top$, and $\mathbf{B}\mathbf{B}^\top$ and hence in $\mathbf{G}\mathbf{G}^\top$. Additionally, it can be shown that $\mathbf{G}_{k+1}\mathbf{G}_{k+1}^\top$ can be expressed as a function of $\mathbf{G}_k\mathbf{G}_k^\top$. Hence, a recursion of the form $\mathbf{G}_{k+1}\mathbf{G}_{k+1}^\top = f(\mathbf{G}_k\mathbf{G}_k^\top)$ and a fixed point approach can be used to facilitate the analysis. We start by stating the update equation on the elements of $\mathbf{G}_k\mathbf{G}_k^\top$ and proceed with the structural constraints on $\mathbf{G}_k\mathbf{G}_k^\top$.

3.1. Update Equations

The update equations are obtained by taking the products of the matrices in (7). Since our goal is to analyze the asymptotic behavior of the GD iterations, we focus our analysis on the leading terms. Specifically, if we denote $\mathbf{G}_k\mathbf{G}_k^\top - \mathbf{G}\mathbf{G}^\top$ by \mathbf{E}_k , then asymptotically (when $\mathbf{E}_k \rightarrow \mathbf{0}$) terms that are $\mathcal{O}(\|\mathbf{E}_k\|^2)$ can be neglected. Hence, we consider the following update equations and omit the $\mathcal{O}(\|\mathbf{E}_k\|^2)$ from the RHS. Due to space limitations, intermediate calculations are omitted and only key steps are provided.

$$\begin{aligned} \mathbf{A}_{k+1}\mathbf{A}_{k+1}^\top &= \mathbf{A}_k\mathbf{A}_k^\top - \eta[\mathbf{P}_k\mathbf{B}_k\mathbf{A}_k^\top + \mathbf{A}_k\mathbf{B}_k^\top\mathbf{P}_k^\top \\ &\quad + 2\mathbf{A}_k(\mathbf{A}_k^\top\mathbf{A}_k - c\mathbf{I}_r)\mathbf{A}_k^\top] \end{aligned} \quad (9)$$

$$\begin{aligned} \mathbf{A}_{k+1}\mathbf{B}_{k+1}^\top &= \mathbf{A}_k\mathbf{B}_k^\top - \eta[\mathbf{P}_k\mathbf{B}_k\mathbf{B}_k^\top + \mathbf{A}_k\mathbf{A}_k^\top\mathbf{P}_k \\ &\quad + \mathbf{A}_k(\mathbf{A}_k^\top\mathbf{A}_k - c\mathbf{I}_r)\mathbf{B}_k^\top] \end{aligned} \quad (10)$$

$$\mathbf{B}_{k+1}\mathbf{B}_{k+1}^\top = \mathbf{B}_k\mathbf{B}_k^\top - \eta[\mathbf{M}\mathbf{P}_k^\top\mathbf{A}_k\mathbf{B}_k^\top + \mathbf{B}_k\mathbf{A}_k^\top\mathbf{P}_k] \quad (11)$$

$$\mathbf{H} = \begin{bmatrix} 2c\mathbf{P}_A \otimes \mathbf{P}_A \\ \mathbf{I}_m \otimes \mathbf{M}^\top \\ \mathbf{M}^\top \otimes \mathbf{I}_m \\ \mathbf{0}_{n^2 \times m^2} \end{bmatrix} \begin{pmatrix} (\mathbf{I}_m \otimes \mathbf{M})\mathbf{TSS}^\top\mathbf{T}^\top \\ (\mathbf{I}_m \otimes \mathbf{B}\mathbf{B}^\top + c\mathbf{P}_A \otimes \mathbf{I}_n)\mathbf{TSS}^\top\mathbf{T}^\top - c\mathbf{P}_A^\perp \otimes \mathbf{I}_n \\ \mathbf{0}_{mn \times mn} \\ (\mathbf{M}^\top \otimes \mathbf{I}_n)\mathbf{TSS}^\top\mathbf{T}^\top \end{pmatrix} \begin{bmatrix} (\mathbf{M} \otimes \mathbf{I}_m)\mathbf{SS}^\top \\ \mathbf{0}_{mn \times mn} \\ (\mathbf{B}\mathbf{B}^\top \otimes \mathbf{I}_m) + c(\mathbf{I}_n \otimes \mathbf{P}_A)\mathbf{SS}^\top - c\mathbf{I}_n \otimes \mathbf{P}_A^\perp \\ (\mathbf{I}_n \otimes \mathbf{M}^\top)\mathbf{SS}^\top \end{bmatrix} \begin{bmatrix} \mathbf{0}_{m^2 \times n^2} \\ \mathbf{0}_{mn \times n^2} \\ \mathbf{0}_{mn \times n^2} \\ \mathbf{0}_{n^2 \times n^2} \end{bmatrix}. \quad (8)$$

Let $\mathbf{E}_{AA}^k = \mathbf{A}_k \mathbf{A}_k^\top - \mathbf{A}\mathbf{A}^\top$, $\mathbf{E}_{AB}^k = \mathbf{A}_k \mathbf{B}_k^\top - \mathbf{A}\mathbf{B}^\top$, $\mathbf{E}_{BA}^k = \mathbf{B}_k \mathbf{A}_k^\top - \mathbf{B}\mathbf{A}^\top$, and $\mathbf{E}_{BB}^k = \mathbf{B}_k \mathbf{B}_k^\top - \mathbf{B}\mathbf{B}^\top$. We can substitute the definition of the different \mathbf{E}_k terms into (9)-(11) and remove terms that are $\mathcal{O}(\|\mathbf{E}_k\|^2)$, simplify, and obtain update equations on the error terms:

$$\begin{aligned} \mathbf{E}_{AA}^{k+1} &= \mathbf{E}_{AA}^k - \eta[\mathbf{P}_k \mathbf{M}^\top + \mathbf{M} \mathbf{P}_k^\top + 2c\mathbf{P}_A \mathbf{E}_{AA}^k \mathbf{P}_A] \\ \mathbf{E}_{AB}^{k+1} &= \mathbf{E}_{AB}^k - \eta[\mathbf{P}_k \mathbf{B}\mathbf{B}^\top + c\mathbf{P}_A \mathbf{P}_k + \mathbf{E}_{AA}^k \mathbf{M} - c\mathbf{P}_A^\perp \mathbf{E}_{AB}^k] \\ \mathbf{E}_{BA}^{k+1} &= \mathbf{E}_{BA}^k - \eta[\mathbf{P}_k^\top \mathbf{M} + \mathbf{M}^\top \mathbf{P}_k]. \end{aligned} \quad (12)$$

Here, $\mathbf{P}_A = \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$ and $\mathbf{P}_A^\perp = \mathbf{I}_m - \mathbf{P}_A$. Using these updates, we can obtain an updated equation of the form

$$\mathbf{E}_{k+1} = f(\mathbf{E}_k) + \mathcal{O}(\|\mathbf{E}_k\|^2) \quad \text{where } \mathbf{E}_k = \begin{bmatrix} \mathbf{E}_{AA}^k & \mathbf{E}_{AB}^k \\ \mathbf{E}_{BA}^k & \mathbf{E}_{BB}^k \end{bmatrix}$$

and $f(\cdot)$ is a linear mapping $f: \mathcal{E} \rightarrow \mathcal{E}$, where

$$\mathcal{E} = \left\{ \mathbf{X} \mathbf{X}^\top - \begin{bmatrix} \mathbf{A}\mathbf{A}^\top & \mathbf{A}\mathbf{B}^\top \\ \mathbf{B}\mathbf{A}^\top & \mathbf{B}\mathbf{B}^\top \end{bmatrix} \mid \mathbf{X} \in \mathbb{R}^{(m+n) \times r} \right\}.$$

In other words, the space of all matrices can be written as the difference of two symmetric rank r $(m+n) \times (m+n)$ matrices. To obtain the linear mapping as a matrix, we consider using the vectorized version of \mathbf{E} . This involves the vectorized form of the error components, i.e., $\mathbf{e}_{aa}^k = \text{vec}(\mathbf{E}_{AA}^k)$, $\mathbf{e}_{ab}^k = \text{vec}(\mathbf{E}_{AB}^k)$, $\mathbf{e}_{ba}^k = \text{vec}(\mathbf{E}_{BA}^k)$, and $\mathbf{e}_{bb}^k = \text{vec}(\mathbf{E}_{BB}^k)$. Using this notation, we can write the update on the vectorized error terms by taking the vec operator on the corresponding matrix equations in (12). This yields:

$$\begin{aligned} \mathbf{e}_{aa}^{k+1} &= \mathbf{e}_{aa}^k - \eta[(\mathbf{M} \otimes \mathbf{I}_m)\mathbf{SS}^\top \mathbf{e}_{ab}^k \\ &\quad + (\mathbf{I}_m \otimes \mathbf{M})\mathbf{T}_{mn}\mathbf{SS}^\top\mathbf{T}_{mn}^\top \mathbf{e}_{ba}^k + 2c(\mathbf{P}_A \otimes \mathbf{P}_A)\mathbf{e}_{aa}^k] \\ \mathbf{e}_{ab}^{k+1} &= \mathbf{e}_{ab}^k - \eta[(\mathbf{B}\mathbf{B}^\top \otimes \mathbf{I}_m)\mathbf{SS}^\top \mathbf{e}_{ab}^k + c(\mathbf{I}_n \otimes \mathbf{P}_A)\mathbf{SS}^\top \mathbf{e}_{ab}^k \\ &\quad + (\mathbf{M}^\top \otimes \mathbf{I}_m)\mathbf{e}_{aa}^k - c(\mathbf{I}_n \otimes \mathbf{P}_A^\perp)\mathbf{e}_{ab}^k] \\ \mathbf{e}_{ba}^{k+1} &= \mathbf{e}_{ba}^k - \eta[(\mathbf{I}_m \otimes \mathbf{B}\mathbf{B}^\top)\mathbf{T}_{mn}\mathbf{SS}^\top\mathbf{T}_{mn}^\top \mathbf{e}_{ba}^k \\ &\quad + c(\mathbf{P}_A \otimes \mathbf{I}_n)\mathbf{T}_{mn}\mathbf{SS}^\top\mathbf{T}_{mn}^\top \mathbf{e}_{ba}^k + (\mathbf{I}_m \otimes \mathbf{M}^\top)\mathbf{e}_{aa}^k \\ &\quad - c(\mathbf{P}_A^\perp \otimes \mathbf{I})\mathbf{e}_{ba}^k] \\ \mathbf{e}_{bb}^{k+1} &= \mathbf{e}_{bb}^k - \eta[(\mathbf{M}^\top \otimes \mathbf{I}_m)\mathbf{T}_{mn}\mathbf{SS}^\top\mathbf{T}_{mn}^\top \mathbf{e}_{ab}^k \\ &\quad + (\mathbf{I}_m \otimes \mathbf{M}^\top)\mathbf{SS}^\top \mathbf{e}_{ab}^k]. \end{aligned} \quad (13)$$

Here \mathbf{T}_{mn} is the permutation matrix satisfying $\mathbf{T}_{mn} \text{vec}(\mathbf{X}) = \text{vec}(\mathbf{X}^\top)$ for all $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $\mathbf{S} \in \mathbb{R}^{mn \times s}$. In a matrix form,

$$\begin{bmatrix} \mathbf{e}_{aa}^{k+1} \\ \mathbf{e}_{ab}^{k+1} \\ \mathbf{e}_{ba}^{k+1} \\ \mathbf{e}_{bb}^{k+1} \end{bmatrix} = (\mathbf{I}_{(m+n)^2} - \eta\mathbf{H}) \begin{bmatrix} \mathbf{e}_{aa}^k \\ \mathbf{e}_{ba}^k \\ \mathbf{e}_{ab}^k \\ \mathbf{e}_{bb}^k \end{bmatrix}, \quad (14)$$

where \mathbf{H} is given in (8). We can write the vec operator of the error matrix \mathbf{E}_k in terms of the vectorized error terms as

$$\text{vec}(\mathbf{E}_k) = \text{vec} \left(\begin{bmatrix} \mathbf{E}_{AA}^k & \mathbf{E}_{AB}^k \\ \mathbf{E}_{BA}^k & \mathbf{E}_{BB}^k \end{bmatrix} \right) = \mathbf{Z} \begin{bmatrix} \mathbf{e}_{aa}^k \\ \mathbf{e}_{ba}^k \\ \mathbf{e}_{ab}^k \\ \mathbf{e}_{bb}^k \end{bmatrix}, \quad (15)$$

where \mathbf{Z} is a permutation matrix given by:

$$\mathbf{Z} = \text{diag} \left(\begin{bmatrix} \mathbf{I}_m \otimes \begin{bmatrix} \mathbf{I}_m \\ \mathbf{0}_{n \times m} \end{bmatrix}, \mathbf{I}_m \otimes \begin{bmatrix} \mathbf{0}_{m \times n} \\ \mathbf{I}_n \end{bmatrix} \right), \\ \left. \begin{bmatrix} \mathbf{I}_n \otimes \begin{bmatrix} \mathbf{I}_m \\ \mathbf{0}_{n \times m} \end{bmatrix}, \mathbf{I}_n \otimes \begin{bmatrix} \mathbf{0}_{m \times n} \\ \mathbf{I}_n \end{bmatrix} \right). \quad (16)$$

Using equations (14) and (15), we can now write the update equation for \mathbf{E}^k as follows:

$$\text{vec}(\mathbf{E}_{k+1}) = \mathbf{Z}(\mathbf{I}_{(m+n)^2} - \eta\mathbf{H})\mathbf{Z}^\top \text{vec}(\mathbf{E}_k). \quad (17)$$

As mentioned in [15], the integration of structural constraints is a necessary step in the analysis of the convergence behavior or the error matrix. We proceed by characterizing the structural constraints.

3.2. Integrating Structural Constraints

Recall the definition of \mathcal{E} the set of all error matrices. It can be shown that a symmetry constraint and a rank r constraint can be used to characterize the set \mathcal{E} . The two constraints can be captured as follows. For symmetry, we have $\mathbf{E} = (\mathbf{E} + \mathbf{E}^\top)/2$. To satisfy the rank constraint, we can consider the manifold of all rank r matrices, \mathcal{M}_r . The first-order term of the error matrix resides in the tangent space $T_{\mathcal{M}_r}(\mathbf{G}\mathbf{G}^\top)$ and hence it satisfies:

$$\mathbf{E} = \mathbf{E} - \mathbf{P}_G^\perp \mathbf{E} \mathbf{P}_G^\perp + \mathcal{O}(\|\mathbf{E}\|^2).$$

See [15] for details. Vectorizing the two requirements, we can write them as

$$\begin{aligned} \text{vec}(\mathbf{E}) &= \frac{1}{2}(\mathbf{I}_{(m+n)^2} + \mathbf{T}_{(m+n)^2})\text{vec}(\mathbf{E}) \\ \text{vec}(\mathbf{E}) &= (\mathbf{I}_{(m+n)^2} - \mathbf{P}_G^\perp \otimes \mathbf{P}_G^\perp)\text{vec}(\mathbf{E}), \end{aligned} \quad (18)$$

where \mathbf{T} is a commutation matrix such that $\mathbf{T} \text{vec}(\mathbf{E}) = \text{vec}(\mathbf{E}^\top)$. Letting $\mathbf{P}_1 = \frac{1}{2}(\mathbf{I}_{(m+n)^2} + \mathbf{T}_{(m+n)^2})$ and $\mathbf{P}_2 = (\mathbf{I}_{(m+n)^2} - \mathbf{P}_G^\perp \otimes \mathbf{P}_G^\perp)$, we can write the projection on the intersection of two requirements as

$$\mathbf{P} = \mathbf{P}_1 \mathbf{P}_2 = \mathbf{P}_2 \mathbf{P}_1. \quad (19)$$

Following [15], this product \mathbf{P} is also a projection. Hence, once can locally characterize the tangent space of the symmetric rank r $(m+n) \times (m+n)$ matrices at $\mathbf{G}\mathbf{G}^\top$ using $\text{vec}(\mathbf{E}) = \mathbf{P} \text{vec}(\mathbf{E})$. The structural constraints can be incorporated into the update equation in (17) as follows:

$$\text{vec}(\mathbf{E}_{k+1}) = \mathbf{P} \mathbf{Z}(\mathbf{I}_{(m+n)^2} - \eta\mathbf{H})\mathbf{Z}^\top \mathbf{P} \text{vec}(\mathbf{E}_k). \quad (20)$$

In other words, the structural constraints on \mathbf{E}_k are applied before and after the update. Since \mathbf{P} is a projection matrix it can be written as $\mathbf{P} = \mathbf{Q}\mathbf{Q}^\top$ where \mathbf{Q} is an $(m+n)^2 \times d$ orthogonal matrix and $d = r(2(m+n)+1-r)/2$ is the rank of \mathbf{P} that is also the dimension of the manifold of all $(m+n) \times (m+n)$ symmetric rank- r matrices. Using \mathbf{Q} , we can write an economy version of the update equation:

$$\mathbf{Q}^\top \text{vec}(\mathbf{E}_{k+1}) = \mathbf{Q}^\top \mathbf{Z}(\mathbf{I}_{(m+n)^2} - \eta\mathbf{H})\mathbf{Z}^\top \mathbf{Q} \mathbf{Q}^\top \text{vec}(\mathbf{E}_k). \quad (21)$$

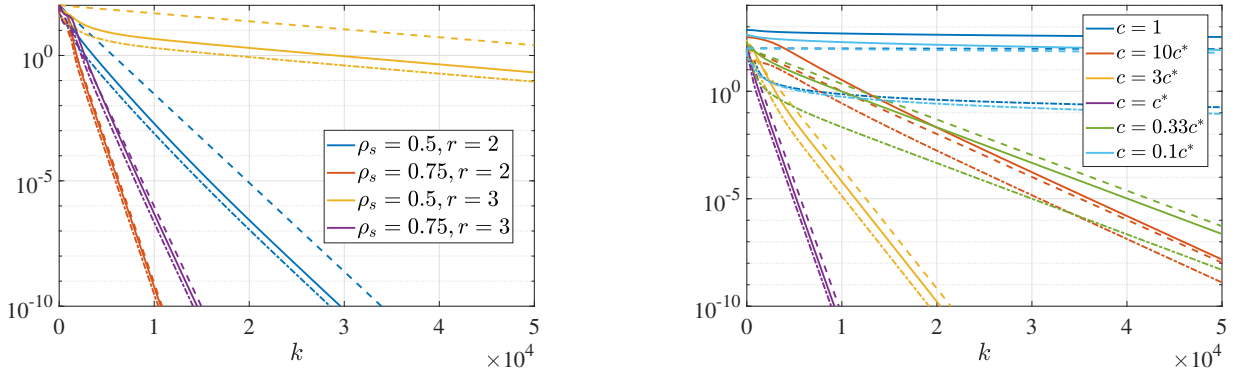


Fig. 1. (Left): a plot of the empirical error terms $\|\mathbf{E}_k\|$ (solid), $\|\mathbf{E}_{AB}^k\|$ (dash-dotted), and ρ^k (dashed) as a function of the number of iterations for the 4 settings described in Section 4 using a blue curve for $(\rho_s = 0.5, r = 2)$, red curves for $(0.75, 2)$, yellow curves for $(0.5, 3)$, and purple curves for $(0.75, 3)$. (Right): a plot of the empirical error terms $\|\mathbf{E}_k\|$ (solid), $\|\mathbf{E}_{AB}^k\|$ (dash-dotted), and ρ^k (dashed) as a function of the number of iterations for the setting of $n = 30$, $m = 20$, $r = 2$, and $\rho_s = 0.75$ for 6 different values of c in $\{1, 10c^*, 3c^*, c^*, c^*/3, c^*/10\}$ corresponding to the blue, red, yellow, purple, green, and cyan, respectively.

Note that $\|\mathbf{Q}^\top \text{vec}(\mathbf{E}^k)\| = \|\text{vec}(\mathbf{E}_k)\| + O(\|\text{vec}(\mathbf{E}_k)\|^2)$ and hence the asymptotic convergence rate of the error term $\mathbf{Q}^\top \text{vec}(\mathbf{E}_k)$ is the same the asymptotic convergence rate of the error term $\text{vec}(\mathbf{E}_k)$. Finally, the convergence rate can be obtained as the spectral radius of the matrix $\mathbf{Q}^\top \mathbf{Z}(\mathbf{I}_{(m+n)^2} - \eta \mathbf{H}) \mathbf{Z}^\top \mathbf{Q} = \mathbf{I}_d - \eta \hat{\mathbf{H}}$, where $\hat{\mathbf{H}} = \mathbf{Q}^\top \mathbf{Z} \mathbf{H} \mathbf{Z}^\top \mathbf{Q} \in \mathbb{R}^{d \times d}$.

3.3. Main Result

We conclude this section by formally stating our convergence rate analysis in the following theorem.

Theorem 1. *Let matrices \mathbf{P} , \mathbf{Z} , and \mathbf{H} be as in (19), (16), and (8), respectively. Additionally, the projection matrix \mathbf{P} can be decomposed as $\mathbf{P} = \mathbf{Q}\mathbf{Q}^\top$ where \mathbf{Q} is a $(m+n) \times d$ semi-orthogonal matrix. Finally, define matrix $\hat{\mathbf{H}} = \mathbf{Q}^\top \mathbf{Z} \mathbf{H} \mathbf{Z}^\top \mathbf{Q}$. If $\hat{\mathbf{H}}$ is non-singular, then Algorithm 1 produces a sequence of matrices $\mathbf{A}_k \mathbf{B}_k^\top$ converging to \mathbf{M} at an asymptotic linear rate $\rho(\mathbf{I}_d - \eta \hat{\mathbf{H}})$. Formally, there exists a neighborhood $\mathcal{N}(\mathbf{M})$ of \mathbf{M} such that for any $\mathbf{A}_0 \mathbf{B}_0^\top \in \mathcal{N}(\mathbf{M})$,*

$$\|\mathbf{A}_k \mathbf{B}_k^\top - \mathbf{M}\|_F \leq C \|\mathbf{A}_0 \mathbf{B}_0^\top\|_F \rho(\mathbf{I}_d - \eta \hat{\mathbf{H}})^k, \quad (22)$$

for some numerical constant $C > 0$.

4. NUMERICAL EXPERIMENTS

In this section, we conduct a numerical experiment to verify the validity of the theoretical analysis of the rate of convergence and to assess the provided choice of the constant c in the algorithm.

4.1. Theoretical Rate Verification

To assess the validity of our analysis, we perform the following experiment. We consider matrix completion for a 20×30 matrix ($m = 20$ and $n = 30$) and four settings for (ρ_s, r) , i.e., the pair (ratio of known entries, rank): $(0.5, 2)$, $(0.75, 2)$, $(0.5, 3)$, and $(0.75, 3)$. For each setting, we generate the factors \mathbf{A} and \mathbf{B} such that their entries are *i.i.d* following the standard normal distribution. We then select uniformly at random $\lceil \rho_s m n \rceil$ of the matrix entries to be the

known set of entries. For each of the four settings, we run Algorithm 1 with c as prescribed in (5) for 50,000 iteration. We use a small step size $\eta = 0.0005$ to make sure the linear convergence occurs in all settings.² We initialize \mathbf{A}_k and \mathbf{B}_k (at $k = 0$) such that their entries are *i.i.d* following the standard normal distribution. At each iteration, we compute $\|\mathbf{E}_{AB}^k\|$ and $\|\mathbf{E}_k\|$. In Fig. 1(Left), we plot for each of the aforementioned settings, $\|\mathbf{E}_k\|$ (solid), $\|\mathbf{E}_{AB}^k\|$ (dash-dotted), and ρ^k (dashed), where $\rho = \rho(\mathbf{I} - \eta \hat{\mathbf{H}})$ is the theoretical rate given by Theorem 1, as a function of the number of iterations k . We note that in each of the settings the rate at which $\|\mathbf{E}_{AB}^k\|$ and $\|\mathbf{E}_k\|$ decrease to zero matches the rate at which ρ^k decreases to zero.

4.2. The Impact of c

Here, we study the impact of the choice of c . We select the setting of $m = 20$, $n = 30$, $\rho_s = 0.5$, and $r = 2$. We generate a matrix using the aforementioned data generation process and apply the GD with the following values of c : $\{1, 10c^*, 3c^*, c^*, c^*/3, c^*/10\}$ where c^* is given in (5). In Fig. 1(Right), we plot for each value of c , the value of $\|\mathbf{E}_k\|$ (solid), $\|\mathbf{E}_{AB}^k\|$ (dash-dotted), and ρ^k (dashed) as functions of k . We observe that the choice of $c = c^*$ provides the fastest rate among the 6 values of c . In particular, the case of $c = 1$ highlights the challenge of using a GD with the same step size for the \mathbf{A} and \mathbf{B} update when the two matrices are significantly different in terms of their norms.

5. CONCLUSION

In this paper, we presented a variant of GD for matrix completion that is based on parameterizing the matrix as a product of two factors, wherein the first factor is assumed to be a semi-orthogonal times a constant. The GD algorithm provides a simultaneous update to the two factors by minimizing an objective consisting of a fitting term and a regularization term that ensures that the first factor is a semi-orthogonal matrix (up to a scaling factor). A detailed analysis was provided to establish conditions for linear convergence of the

²In each specific setting, one can choose an optimal step size by maximizing the rate $\rho(\mathbf{I}_d - \eta \hat{\mathbf{H}})$ w.r.t. $\eta > 0$.

algorithm and a formula for the asymptotic linear rate was provided. Additionally, a formula was provided to determine the scaling constant for the semi-orthogonal factor. Numerical experiments were conducted to verify the correctness of the linear rate expression and to illustrate the merit of selecting the scaling constant. In contrast, the proposed c^* yields identical norms for \mathbf{A} and \mathbf{B} and appears to produce the fastest rate of convergence.

6. REFERENCES

- [1] Zhang Liu and Lieven Vandenberghe, “Interior-point method for nuclear norm approximation with application to system identification,” *SIAM J. Matrix Anal. Appl.*, vol. 31, no. 3, pp. 1235–1256, 2009.
- [2] Jasson DM Rennie and Nathan Srebro, “Fast maximum margin matrix factorization for collaborative prediction,” in *Proc. Int. Conf. Mach. Learn.*, 2005, pp. 713–719.
- [3] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright, “Robust principal component analysis?,” *J. ACM*, vol. 58, no. 3, pp. 11, 2011.
- [4] Emmanuel J Candès and Benjamin Recht, “Exact matrix completion via convex optimization,” *Found. Comput. Math.*, vol. 9, no. 6, pp. 717, 2009.
- [5] Samuel Burer and Renato DC Monteiro, “Local minima and convergence in low-rank semidefinite programming,” *Math. Program.*, vol. 103, no. 3, pp. 427–444, 2005.
- [6] Ruoyu Sun and Zhi-Quan Luo, “Guaranteed matrix completion via non-convex factorization,” *IEEE Trans. Inf. Theory*, vol. 62, no. 11, pp. 6535–6579, 2016.
- [7] Rong Ge, Jason D Lee, and Tengyu Ma, “Matrix completion has no spurious local minimum,” in *Proc. Adv. Neural Inf. Process. Syst.*, 2016, pp. 2973–2981.
- [8] Yudong Chen and Martin J Wainwright, “Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees,” *arXiv preprint arXiv:1509.03025*, 2015.
- [9] Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen, “Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion,” in *Proc. Int. Conf. Mach. Learn.*, 2018, pp. 3345–3354.
- [10] Caihua Chen, Bingsheng He, and Xiaoming Yuan, “Matrix completion via an alternating direction method,” *IMA J. Numer. Anal.*, vol. 32, no. 1, pp. 227–245, 2012.
- [11] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi, “Low-rank matrix completion using alternating minimization,” in *Proc. Annu. ACM Symp. Theory Comput.*, 2013, pp. 665–674.
- [12] Moritz Hardt, “Understanding alternating minimization for matrix completion,” in *Proc. Annu. IEEE Symp. Found. Comput.*, 2014, pp. 651–660.
- [13] Moritz Hardt and Mary Wootters, “Fast matrix completion without the condition number,” in *Proc. Conf. Learn. Theory*, 2014, pp. 638–678.
- [14] Lijun Ding and Yudong Chen, “Leave-one-out approach for matrix completion: Primal and dual analysis,” *IEEE Trans. Inf. Theory*, vol. 66, no. 11, pp. 7274–7301, 2020.
- [15] Trung Vu and Raviv Raich, “Exact linear convergence rate analysis for low-rank symmetric matrix completion via gradient descent,” in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.* IEEE, 2021, pp. 3240–3244.
- [16] Qinqing Zheng and John Lafferty, “Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent,” *arXiv preprint arXiv:1605.07051*, 2016.