# Box Constrained Low-rank Matrix Approximation with Missing Values

Manami Tatsukawa[1] and Mirai Tanaka[2]

[1]*Department of Industrial Engineering and Economics, Tokyo Institute of Technology,*
*2-12-1 Ookayama, Meguro-ku, Tokyo 152-8550, Japan*
[2]*Department of Mathematical Analysis and Statistical Inference, The Institute of Statistical Mathematics,*
*10-3 Midori-cho, Tachikawa, Tokyo 190-8562, Japan*

Keywords:      Low-Rank Matrix Approximation, Missing Data, Matrix Completion with Noise, Principal Component Analysis with Missing Values, Collaborative Filtering, Block Coordinate Descent Method.

Abstract:        In this paper, we propose a new low-rank matrix approximation model for completing a matrix with missing values. Our proposed model contains a box constraint that arises from the context of collaborative filtering. Although it is unfortunately NP-hard to solve our model with high accuracy, we can construct a practical algorithm to obtain a stationary point. Our proposed algorithm is based on alternating minimization and converges to a stationary point under a mild assumption.

## 1 INTRODUCTION

### 1.1 Background

Low-rank matrix approximation is commonly used for feature extraction. This technique embeds high-dimensional data into a lower dimensional space, because relevant data in a high-dimensional space often lie in a lower dimensional space. Feature extraction enables us to identify potential features that may help to increase prediction accuracy when data are incomplete or contain some noise. We provide two examples of techniques for approximating a data matrix with missing values.

#### 1.1.1 PCA with Missing Value

Principal component analysis (PCA) is a classical technique used for extracting features. This technique embeds high-dimensional data into lower dimensional space. The components embedded into lower dimensional space are called principal components. To handle incompleteness of input data we often use nonlinear models; however, these models cause problems such as overfitting and bad locally optimal solutions. Tipping and Bishop (1999) introduced a probabilistic formulation of PCA. The probabilistic PCA is known to provide a good foundation for handling missing values (Ilin and Raiko, 2010). Probabilistic PCA is often solved by an expectation-maximization (EM) algorithm.

#### 1.1.2 Collaborative Filtering

Low-rank matrix approximation is utilized in recommendation systems such as those found on iTunes, Amazon, and Netflix. In these services, music, book, and movie recommendations are provided to users. Users rate items they have listened to, bought, or watched. Based on the ratings, items are recommended based on the user's preferences or the items' novelty to the user.

Let us consider the following situation: There are $m$ users and $n$ items. Every user rates some of items on a scale of one to five. One is the lowest score and it means that a user does not prefer the item. We represent this with an $m \times n$ matrix. Each row of the matrix represents each user and each column represents each item. We show an example of $3 \times 5$ matrix below:

$$
V = \begin{matrix} & \begin{matrix} a & b & c & d & e \end{matrix} \\ \begin{matrix} A \\ B \\ C \end{matrix} & \begin{pmatrix} 2 & 5 & ? & 4 & 1 \\ 1 & ? & 1 & 3 & 2 \\ ? & 1 & 4 & ? & 5 \end{pmatrix} \end{matrix}. \quad (1)
$$

Matrix $V$ represent that, for instance, user A rates item b five and does not rate item c, using the symbol "?". In what follows, we use the symbol "?" for representing missing values. In the example shown in Equation (1), users A and B seems to have similar trends. This means that an item that user A rates high may also be preferred by user B. For example, item b, which user A rates five, may also be rated high by

user B. If we complete the missing value of the matrix, we can predict how users prefer items that they have not evaluate yet. This technique of predicting is can be used for recommendation systems.

A technique used in recommendation systems is collaborative filtering (CF). CF algorithms have three main categories: memory-based, model-based, and hybrid (Su and Khoshgoftaar, 2009). Memory-based CF algorithms calculate similarities between users or items to predict users' preferences. Model-based CF algorithms learn a model in order to make predictions. Hybrid CF algorithms combine several CF techniques.

In memory-based CF algorithms, similarities between users or items are used to make predictions. As in measures of similarity, the vector cosine correlation and the Pearson correlation are often used. However, when many values are missing, it is difficult to compute similarities between users. In fact, the number of items might be greater than the number of users and each user may evaluate only a small number of items. For this reason, many items are evaluated by only a few users, while other users do not submit any evaluations.

To overcome the weakness in memory-based CF algorithms, model-based CF algorithms have been investigated. Model-based CF approaches use data mining or machine learning algorithms. One of the techniques of model-based CF is dimensionality reduction, such as PCA or singular value decomposition (SVD). As we mentioned above, high-dimensional data is thought to be expressed by lower dimensional data because related data lie in a lower dimensional space.

## 1.2 Related Work

Low-rank matrix completion and approximation have been studied. In this section, we briefly introduce some papers about these techniques.

Let us consider completing a matrix $V \in (\mathbb{R} \cup \{?\})^{m \times n}$ with missing values, where the symbol ? indicates that the corresponding value is missing. That is, $V_{ij} = ?$ indicates $V_{ij}$ is missing. Candés and Recht (2009) proposed rank minimization to complete matrix $V$. They considered the following problem:

$$
\begin{aligned}
& \text{minimize} \quad \text{rank}(\boldsymbol{X}) \\
& \text{subject to} \quad X_{ij} = V_{ij} \quad ((i,j) \in \Omega),
\end{aligned} \tag{2}
$$

where $\Omega$ is the set of indices of observed entries, *i.e.*, $\Omega = \{(i,j) : V_{ij} \in \mathbb{R}\}$. Problem (2) is NP-hard because it contains the $l_0$-norm minimization problem. This difficulty essentially arises from the nonconvexity and discontinuity of the rank function. Hence, they

introduced nuclear norm minimization as a convex relaxation of Problem (2). Nuclear norm minimization can be recast as a semidefinite optimization problem (SDP). There are many efficient algorithms and high-quality software packages available for solving SDP, including the interior-point method. However, the computation time for solving SDP is very sensitive to instance size and is unsuitable for solving large instances

Olsson and Oskarsson (2009); Gillis and Glineur (2011) studied the following problem to complete $V$:

$$
\begin{aligned}
& \text{minimize} \quad \sum_{i=1}^{m} \sum_{j=1}^{n} W_{ij}(X_{ij} - V_{ij})^2 \\
& \text{subject to} \quad \text{rank}(\boldsymbol{X}) \le r,
\end{aligned} \tag{3}
$$

where decision variable $\boldsymbol{X} \in \mathbb{R}^{m \times n}$ is a completed matrix of $V$ and $\boldsymbol{W} \in \{0,1\}^{m \times n}$ is a given weight matrix corresponding to an observation, *i.e.*, $W_{ij} = 1$ for $V_{ij} \in \mathbb{R}$ and otherwise $W_{ij} = 0$. Using $\Omega$, we obtain the following equivalent formulation of Problem (3):

$$
\begin{aligned}
& \text{minimize} \quad \sum_{(i,j) \in \Omega} (X_{ij} - V_{ij})^2 \\
& \text{subject to} \quad \text{rank}(\boldsymbol{X}) \le r.
\end{aligned} \tag{4}
$$

Our formulation is similar to this one and this model is helpful to understand our model. Olsson and Oskarsson (2009) proposed a heuristic based on an approximated continuous (but nonconvex) formulation of Problems (3). Gillis and Glineur (2011) proved the NP-hardness of Problem (3), and equivalently, Problem (4).

On the other hand, the low-rank matrix approximation problem is easily solved when no values are missing. In fact, when $\Omega$ is an entire set of indices, Problem (4) is equivalent to the following problem:

$$
\begin{aligned}
& \text{minimize} \quad \|\boldsymbol{X} - \boldsymbol{V}\|_{\mathrm{F}}^2 \\
& \text{subject to} \quad \text{rank}(\boldsymbol{X}) \le r,
\end{aligned}
$$

where $\| \cdot \|_{\mathrm{F}}$ denotes the Frobenius norm defined by

$$
\|\boldsymbol{A}\|_{\mathrm{F}} = \sqrt{\sum_{i=1}^{m} \sum_{j=1}^{n} A_{ij}^2}
$$

for $\boldsymbol{A} \in \mathbb{R}^{m \times n}$. This problem is nonconvex; however, a global optimal solution is obtained by the truncated SVD of $V$. Specifically, it is well known that an optimal solution of this problem can be written as $\sum_{l=1}^{r} \sigma_l \boldsymbol{p}_l \boldsymbol{q}_l^{\top}$ (Trefethen and Bau, 1997, Theorem 5.9), where $\sigma_l$, $\boldsymbol{p}_l$, and $\boldsymbol{q}_l$ respectively represent the $l$-th largest singular value and corresponding singular vectors of $V$. The computation of all singular values and the vectors of $V$ is expensive. More specifically, it requires a computation time in $O(\min\{m^2 n, m n^2\})$, which can be too heavy for a

large instance. Instead, the $r$ largest singular values and corresponding singular vectors can be quickly computed by an iterative method if $r$ is small.

When we use the truncated SVD for a matrix containing missing values, we need to complete the matrix. One method for doing this involves completing the input matrix with the average of the non-missing entries in the same rows or columns (Sarwar et al., 2000). Specifically, their method comprises the following steps:

1. Temporarily fill in the missing values of incomplete input matrix $\boldsymbol{V}$ with the average of the non-missing entries in the same columns. We call the completed matrix $\hat{\boldsymbol{V}}$. That is, $\hat{V}_{ij} = (1/|\{i' : (i',j) \in \Omega\}|)\sum_{i':(i',j)\in\Omega} V_{i'j}$ for $(i,j) \notin \Omega$.

2. Compute row average vector $\boldsymbol{\mu}$. The $i$-th element $\mu_i$ of $\boldsymbol{\mu}$ is the average of the $i$-th row of $\boldsymbol{V}$. That is, $\mu_i = (1/|\{j : (i,j) \in \Omega\}|)\sum_{j:(i,j)\in\Omega} V_{ij}$.

3. Compute the best rank $r$ approximation matrix $\boldsymbol{V}_r$ of $\hat{\boldsymbol{V}} - \boldsymbol{\mu}\boldsymbol{1}^\top$ by using the truncated SVD, where $\boldsymbol{1}$ is a vector of all ones.

4. Return $\boldsymbol{V}_r + \boldsymbol{\mu}\boldsymbol{1}^\top$ as a low-rank approximation of input matrix $\boldsymbol{V}$.

## 1.3 Our Contribution and Structure of This Paper

The remainder of this paper is organized as follows: Section 2 proposes a new model for low-rank matrix approximation contains a box constraint that arises from the context of CF. The rank minimization models was previously proposed by Candés and Recht (2009) and Olsson and Oskarsson (2009). Their model can be recast as an SDP. There are many software packages available for solving SDP, but SDP is unsuitable for solving large instances. Our model includes Problem (4), which is proved NP-hard by Gillis and Glineur (2011). However, we can solve our model by truncated SVD. Moreover, our model is more suitable for CF than previous models owing to the box constraint. We need not our model recast as an SDP and suitable for large instances. Section 3 proves the NP-hardness of our proposed model. Section 4 proposes an algorithm for solving our proposed model and proves the convergence of the algorithm. Section 5 reports preliminary numerical results to assess our proposed model and algorithm. Section 6 presents concluding remarks.

## 2 MODEL

In this section, we propose a new model for low-rank matrix approximation that is easier to calculate. Here, matrix $\boldsymbol{X} \in \mathbb{R}^{m \times n}$ denotes a completed matrix of matrix $\boldsymbol{V} \in (\mathbb{R} \cup \{?\})^{m \times n}$ with missing values. We consider minimizing the difference between $\boldsymbol{V}$ and $\boldsymbol{X}$. This idea is described in previous studies.

We consider imposing a box constraint $\boldsymbol{L} \leq \boldsymbol{X} \leq \boldsymbol{U}$, where $\boldsymbol{L} \in (\mathbb{R} \cup \{-\infty\})^{m \times n}$ and $\boldsymbol{U} \in (\mathbb{R} \cup \{+\infty\})^{m \times n}$ satisfy $\boldsymbol{L} \leq \boldsymbol{U}$ and the inequalities are entrywise. For example, $\boldsymbol{L} \leq \boldsymbol{X}$ means $L_{ij} \leq X_{ij}$ for all $i,j$. The box constraint is important when applying low-rank matrix approximation with missing values to CF, because users evaluate items in some range. For example, an Amazon user evaluates an item by assigning it one to five stars. We model such an evaluation value range as the box constraint. Matrices $\boldsymbol{L}, \boldsymbol{U}$ represent the lower and upper bounds of an evaluation value range, respectively. This constraint is helpful in making exact or approximate predictions and in eliminating outliers.

To prevent a case in which both the rank constraint and the box constraint are not simultaneously fulfilled, we consider variables that fulfil each constraint separately. For this reason, we introduce another variable, $\boldsymbol{Y} \in \mathbb{R}^{m \times n}$, and impose the box constraint on $\boldsymbol{Y}$. We then add the squared Frobenius norm $\|\boldsymbol{X} - \boldsymbol{Y}\|_{\mathrm{F}}^2$ of the difference between $\boldsymbol{X}$ and $\boldsymbol{Y}$ to the objective function as a penalty.

Summarizing the previous argument, we formulate low-rank matrix approximation with missing values as the following optimization problem:

$$
\begin{aligned}
\text{minimize} \quad & \|\boldsymbol{X} - \boldsymbol{Y}\|_{\mathrm{F}}^2 + \lambda \sum_{(i,j)\in\Omega} (Y_{ij} - V_{ij})^2 \\
\text{subject to} \quad & \mathrm{rank}(\boldsymbol{X}) \leq r, \\
& \boldsymbol{L} \leq \boldsymbol{Y} \leq \boldsymbol{U},
\end{aligned}
\tag{5}
$$

where $\lambda$ is a parameter to determine the weight of the two objectives. When we use large $\lambda$, the second part in the objective function is emphasized, so that we expect to obtain $\boldsymbol{Y}$ close to $\boldsymbol{V}$ permitting small violation of the low-rank constraint. When we use small $\lambda$, the first part in the objective function is emphasized, so that we expect to obtain $\boldsymbol{Y}$ satisfying the two constraints simultaneously permitting small difference from $\boldsymbol{V}$.

## 3 HARDNESS

In this section, we show the NP-hardness of Problem (5). Specifically, we prove the following result:

**Theorem 1.** *When* $V \in ([0,1] \cup \{?\})^{m \times n}$ *and* $r = 1$*, it is NP-hard to find an approximate solution to Problem* (5) *with an objective function accuracy of less than* $2^{-12}(mn)^{-7}$.

To prove this theorem, we employ the following theorem:

**Theorem 2** (Gillis and Glineur (2011, Theorem 1.2))**.** *When* $V \in ([0,1] \cup \{?\})^{m \times n}$ *and* $r = 1$*, it is NP-hard to find an approximate solution to Problem* (4) *with an objective function accuracy of less than* $2^{-12}(mn)^{-7}$.

*Proof of Theorem 1.* In Problem (5), we set $L_{ij} = U_{ij} = V_{ij}$ for $(i,j) \in \Omega$; otherwise, $L_{ij} = -\infty$ and $U_{ij} = +\infty$. Then, for $(i,j) \in \Omega$, $Y_{ij}$ is fixed to $V_{ij}$ and the second part of the objective function is removed. Thus, the resulting objective function can be written as $\sum_{(i,j) \in \Omega}(X_{ij} - V_{ij})^2 + \sum_{(i,j) \notin \Omega}(X_{ij} - Y_{ij})^2$. The latter part of this expression is also removed because $Y_{ij}$ for $(i,j) \notin \Omega$ is unconstrained and is able to coincide with $X_{ij}$. As a result, we obtain Problem (4). In this procedure, we have reduced Problem (4) to Problem (5). This reduction is clearly in polynomial time. Hence, the hardness for Problem (4) also holds for Problem (5). $\square$

*Remark* 1. Theorem 1 is easily generalized to any $r$ because Theorem 2 is generalized to any $r$ (Gillis and Glineur, 2011, Remark 3).

# 4 ALGORITHM

We propose an alternating minimization algorithm for solving Problem (5). In this section, we use the following extended-real-valued functions:

$$f_0(X,Y) = \|X - Y\|_F^2 + \lambda \sum_{(i,j) \in \Omega}(Y_{ij} - V_{ij})^2,$$

$$f_1(X) = \iota(\text{rank}(X) \le r),$$
$$f_2(Y) = \iota(L \le Y \le U),$$
$$f(X,Y) = f_0(X,Y) + f_1(X) + f_2(Y),$$

where $\iota$ is the indicator function. That is, $f_1(X) = 0$ if $\text{rank}(X) \le r$; otherwise $f_1(X) = +\infty$; and $f_2(Y) = 0$ if $L \le Y \le U$; otherwise $f_2(Y) = +\infty$. Note that the minimization of $f(X,Y)$ is equivalent to Problem (5). We consider the alternating minimization of $f(X,Y)$ as shown in Algorithm 1. Each iteration of this algorithm is easily computed as we see below.

In the update of $X$, we set $X^{(k+1)}$ to the best rank-$r$ approximation of $Y^{(k)}$. That is, $X^{(k+1)}$ is an optimal solution of the following subproblem:

$$\begin{aligned} \text{minimize} \quad & \|X - Y^{(k)}\|_F^2 \\ \text{subject to} \quad & \text{rank}(X) \le r. \end{aligned}$$

---

**Algorithm 1:** Alternating minimization algorithm for solving Problem (5).

---

Take initial guess $(X^{(0)}, Y^{(0)}) \in \text{dom} f_1 \times \text{dom} f_2$.
**for** $k = 0, 1, 2, \dots$ until convergence:
    Update $X^{(k+1)} = \text{argmin}_X f(X, Y^{(k)})$.
    Update $Y^{(k+1)} = \text{argmin}_Y f(X^{(k+1)}, Y)$.

---

Although this subproblem is a nonconvex optimization problem, the optimal solution is easily computed by the truncated SVD of $Y^{(k)}$.

The update of $Y$ is also easily computed. In fact, we solve the following subproblem to update $Y$:

$$\begin{aligned} \text{minimize} \quad & \|X^{(k+1)} - Y\|_F^2 + \lambda \sum_{(i,j) \in \Omega}(Y_{ij} - V_{ij})^2 \\ \text{subject to} \quad & L \le Y \le U. \end{aligned}$$

(6)

Note that this subproblem is separable. The separated subproblem for each entry reduces to the minimization of a univariate convex quadratic function over a closed interval. Specifically, we solve

$$\begin{aligned} \text{minimize} \quad & (1+\lambda)Y_{ij}^2 - 2(X_{ij}^{(k+1)} + \lambda V_{ij})Y_{ij} \\ \text{subject to} \quad & L_{ij} \le Y_{ij} \le U_{ij} \end{aligned}$$

for each $(i,j) \in \Omega$ and

$$\begin{aligned} \text{minimize} \quad & Y_{ij}^2 - 2X_{ij}^{(k+1)}Y_{ij} \\ \text{subject to} \quad & L_{ij} \le Y_{ij} \le U_{ij} \end{aligned}$$

for each $(i,j) \notin \Omega$. These have the following closed-form solution:

$$Y_{ij}^{(k+1)} = \begin{cases} L_{ij} & (A_{ij} \le L_{ij}), \\ A_{ij} & (L_{ij} < A_{ij} < U_{ij}), \\ U_{ij} & (U_{ij} \le A_{ij}), \end{cases}$$

where

$$A_{ij} = \begin{cases} \dfrac{X_{ij}^{(k+1)} + \lambda V_{ij}}{1 + \lambda} & ((i,j) \in \Omega), \\ X_{ij}^{(k+1)} & ((i,j) \notin \Omega). \end{cases}$$

Hence, we can solve Subproblem (6) by thresholding, which requires a computation time in $O(mn)$.

A sequence generated by Algorithm 1 converges to a stationary point of $f$ under a mild assumption. Here, we say point $(\bar{X}, \bar{Y})$ is a stationary point of $f$ in the sense of Tseng (2001) if $(\bar{X}, \bar{Y}) \in \text{dom} f = \{(X,Y) : f(X,Y) < +\infty\}$ and

$$f'(\bar{X}, \bar{Y}; \Delta X, \Delta Y) \ge 0 \ (\forall(\Delta X, \Delta Y)),$$

where $f'(\bar{X}, \bar{Y}; \Delta X, \Delta Y)$ is the lower directional derivative of $f$ at $(\bar{X}, \bar{Y})$ in the direction $(\Delta X, \Delta Y)$, *i.e.,*

$$f'(\bar{X}, \bar{Y}; \Delta X, \Delta Y)$$

$$= \liminf_{\varepsilon \downarrow 0} \frac{f(\bar{X} + \varepsilon \Delta X, \bar{Y} + \varepsilon \Delta Y) - f(\bar{X}, \bar{Y})}{\varepsilon}.$$

Note that this definition works even if $f$ is nonsmooth. Specifically, the following theorem holds.

**Theorem 3.** *Let $\{(X^{(k)}, Y^{(k)})\}$ be a sequence generated by Algorithm 1 and assume that level set $L = \{(X, Y) : f(X, Y) \leq f(X^{(0)}, Y^{(0)})\}$ is bounded. Then, $\{(X^{(k)}, Y^{(k)})\}$ has at least one cluster point. In addition, every cluster point is a stationary point of $f$.*

To prove this theorem, the following technical lemma is required.

**Lemma 1.** *Effective domain $\operatorname{dom} f$ of $f$ is closed.*

*Proof.* Clearly, $\operatorname{dom} f = \operatorname{dom} f_1 \times \operatorname{dom} f_2$ and $\operatorname{dom} f_2$ is closed. Thus, we only need to show the closedness of $\operatorname{dom} f_1$. Take $X \notin \operatorname{dom} f_1$ arbitrarily. Then, the $r$-th largest singular value $\sigma_r(X)$ of $X$ is positive. We arbitrarily take $E$ such that $\|E\|_2 < \sigma_r(X)$. From Golub and van Loan (2013, Corollary 8.6.2), we can easily prove that $\sigma_l(X + E) \geq \sigma_l(X) - \|E\|_2 > 0$ for $l = 1, \ldots, r$, so that $X + E \notin \operatorname{dom} f_1$. This indicates the closedness of $\operatorname{dom} f_1$. $\square$

Here, we provide a proof of Theorem 3. The following proof is essentially based on the discussion in Tseng (2001, Sections 3 and 4).

*Proof of Theorem 3.* From the optimality in each update, sequence $\{(X^{(k)}, Y^{(k)})\}$ is contained in $L \subset \operatorname{dom} f$, so that $\{(X^{(k)}, Y^{(k)})\}$ is bounded. Hence, $\{(X^{(k)}, Y^{(k)})\}$ has at least one cluster point on $\operatorname{dom} f$ because $\operatorname{dom} f$ is closed from Lemma 1.

Let $\{(X^{(k_j)}, Y^{(k_j)})\}$ be a subsequence of $\{(X^{(k)}, Y^{(k)})\}$ that converges to a cluster point $(\bar{X}, \bar{Y})$. To show that $(\bar{X}, \bar{Y})$ is a stationary point, we only have to prove $f'(\bar{X}, \bar{Y}; \Delta X, \Delta Y) \geq 0$ for any $(\Delta X, \Delta Y)$ because $(\bar{X}, \bar{Y}) \in \operatorname{dom} f$. From the optimality in each update, we obtain

$$f(X^{(k_j+1)}, Y^{(k_j+1)}) \leq f(X^{(k_j+1)}, Y^{(k_j)})$$
$$\leq f(X, Y^{(k_j)}) \ (\forall X).$$

Taking the limit as $j$ tends to infinity, we obtain

$$f(\bar{X}, \bar{Y}) \leq f(X, \bar{Y}) \ (\forall X) \qquad (7)$$

because $f$ is continuous on $\operatorname{dom} f$. In addition,

$$f(X^{(k_j)}, Y^{(k_j)}) \leq f(X^{(k_j)}, Y) \ (\forall Y)$$

holds. Taking the limit as $j$ tends to infinity, we obtain

$$f(\bar{X}, \bar{Y}) \leq f(\bar{X}, Y) \ (\forall Y). \qquad (8)$$

Because $f_0$ is differentiable, the following relationship holds for any $(\Delta X, \Delta Y)$:

$$f'(\bar{X}, \bar{Y}; \Delta X, \Delta Y)$$

$$= \langle \nabla f_0(\bar{X}, \bar{Y}), (\Delta X, \Delta Y) \rangle$$
$$+ \liminf_{\varepsilon \downarrow 0} \left( \frac{f_1(\bar{X} + \varepsilon \Delta X) - f_1(\bar{X})}{\varepsilon} \right.$$
$$\left. + \frac{f_2(\bar{Y} + \varepsilon \Delta Y) - f_2(\bar{Y})}{\varepsilon} \right)$$
$$\geq \langle \nabla_X f_0(\bar{X}, \bar{Y}), \Delta X \rangle + \langle \nabla_Y f_0(\bar{X}, \bar{Y}), \Delta Y \rangle$$
$$+ \liminf_{\varepsilon \downarrow 0} \frac{f_1(\bar{X} + \varepsilon \Delta X) - f_1(\bar{X})}{\varepsilon}$$
$$+ \liminf_{\varepsilon \downarrow 0} \frac{f_2(\bar{Y} + \varepsilon \Delta Y) - f_2(\bar{Y})}{\varepsilon}$$
$$= \langle \nabla_X f_0(\bar{X}, \bar{Y}), \Delta X \rangle + \langle \nabla_Y f_0(\bar{X}, \bar{Y}), \Delta Y \rangle$$
$$+ f_1'(\bar{X}; \Delta X) + f_2'(\bar{Y}; \Delta Y)$$
$$= \liminf_{\varepsilon \downarrow 0} \frac{f(\bar{X} + \varepsilon \Delta X, \bar{Y}) - f(\bar{X}, \bar{Y})}{\varepsilon}$$
$$+ \liminf_{\varepsilon \downarrow 0} \frac{f(\bar{X}, \bar{Y} + \varepsilon \Delta Y) - f(\bar{X}, \bar{Y})}{\varepsilon}$$
$$\geq 0,$$

where the last inequality holds because of inequalities (7) and (8). $\square$

*Remark 2.* The boundedness of level set $L = \{(X, Y) : f(X, Y) \leq f(X^{(0)}, Y^{(0)})\}$ holds if, for example, $L_{ij} > -\infty$ and $U_{ij} < +\infty$ for all $i, j$.

*Remark 3.* We can extend our proposed algorithm to a closed convex constraint on $Y$, instead of the box constraint. If we can efficiently solve the subproblem to update $Y$, the algorithm works well. In practice, the subproblem can be solved faster than the subproblem to update $X$ with the truncated SVD on $Y$.

We proved the convergence to a stationary point above instead of an optimal solution. A set of stationary points contains every local optimal solution and of course the global optimal solution. Thus, in practice, we run our algorithm from multiple initial points and select the best stationary point provided by our algorithm.

# 5 PRELIMINARY EXPERIMENTS

We show preliminary numerical results using synthetic dataset to investigate a convergence rate of our algorithm and an effect given by differences of initial points. We executed all experiments on a macOS Sierra 10.12.6 with an Intel Core m3, a 1.1 GHz clock speed, and 8 GB of physical memory. We implemented our algorithms in MATLAB (R2017b).

We generated matrix $A \in \mathbb{R}^{20 \times 100}$ such that $\operatorname{rank}(A) = 10$ and $0.5 \leq A_{ij} \leq 5.5$ for all $(i, j)$ by multiplying two matrices $B \in \mathbb{R}^{20 \times 9}$ and $C \in \mathbb{R}^{9 \times 100}$

Figure 1: Comparison of iteration numbers.



Figure 2: Comparison of objective values.

and some scaling, where all entries of $B$ and $C$ are identically sampled from the standard normal distribution. Then, we round all entries of $A$, *i.e.*, $\bar{A}_{ij} \in \{1, 2, \ldots, 5\}$, and randomly missed 80% of them. We used the incomplete matrix as $V$. In our algorithm, we set parameters to $r = 10$ for computing truncated SVD and $\lambda = 1$ for the weight parameter in the objective function.

In this experiment, we compared four methods of generating initial points:

**SKKR** We applied the method proposed by Sarwar et al. (2000) to given $V$ and used the output as an initial point of our algorithm;

**perturb + SKKR** We applied the method proposed by Sarwar et al. (2000) to perturbed $V$'s and used the outputs as initial points of our algorithm;

**low rank rand** We generated initial points of our algorithm in the same way of generating $A$;

**rand** We generated $20 \times 100$ matrices whose entries were identically sampled from the standard normal distribution and used it as initial points.

In each method of perturb + SKKR, low rank rand, and rand, we generated 10 initial points and obtained 10 different solutions.

The differences of iteration numbers necessary to converge to stationary points are shown in Figure 1. The horizontal axis is iteration number and the vertical axis is objective value. From Figure 1, we can see that if we run our algorithm from initial points generated by perturb + SKKR and low rank rand, it converged faster than from an initial point provided by SKKR. On the other hand, it took much more time to converge if it started from initial points generated by rand.

The objective values at the resulting stationary points are shown in Figure 2. The vertical axis is objective value. The black horizontal line around

$2 \times 10^{-5}$ depicts the objective value of the stationary point provided by our algorithm starting from an initial point generated by SKKR. From Figure 2, we can see that all solutions provided by our algorithm starting from initial points generated by perturb + SKKR and low rank rand are better than that by SKKR. Some solutions resulting from rand is also smaller than that from SKKR.

## 6 CONCLUSION

In this paper, we proposed a low-rank matrix approximation model for completing a matrix with missing values. Our proposed model utilizes not only a rank constraint but also a box constraint. Owing to the box constraint, this model shows promise for use in recommendation systems. In addition, we proposed an alternating minimization algorithm for solving our proposed model, and proved that a sequence generated by our proposed algorithm converges to a stationary point under a mild assumption. Our numerical results are preliminary, however, we showed that our proposed algorithm converges quickly and provides better solution than the existing method proposed by Sarwar et al. (2000) in a specific case. The performance of our algorithm depends on an initial guess. Hence, we should try multiple initial guesses and take the best one. In our experiments, we fixed parameters $r$ and $\lambda$. However, we should employ cross validation to decide appropriate values for such parameters. We left extensive numerical experiments as our future work.

## ACKNOWLEDGEMENTS

## REFERENCES

Candés, E. J. and Recht, B. (2009). Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772.

Gillis, N. and Glineur, F. (2011). Low-rank matrix approximation with weights or missing data is NP-hard. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1149–1165.

Golub, G. H. and van Loan, C. F. (2013). *Matrix Computations*. The Johns Hopkins University Press, fourth edition.

Ilin, A. and Raiko, T. (2010). Practical approaches to principal component analysis in the presence of missing values. *Journal of Machine Learning Research*, 11(Jul):1957–2000.

Olsson, C. and Oskarsson, M. (2009). A convex approach to low rank matrix approximation with missing data. In Salberg, A.-B., Hardeberg, J. Y., and Jenssen, R., editors, *Proceedings of the 16th Scandinavian Conference on Image Analysis (SCIA '09)*, pages 301–309.

Sarwar, B., Karypis, G., Konstan, J., and Riedl, J. (2000). Application of dimensionality reduction in recommender system—a case study. In Kohavi, R., Masand, B., Spiliopoulou, M., and Srivastava, J., editors, *Proceedings of the ACM WEBKDD 2000 Workshop*.

Su, X. and Khoshgoftaar, T. M. (2009). A survey of collaborative filtering techniques. *Advances in Artificial Intelligence*, 2009(421425):1–19.

Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. *Journal of the Royal Statistical Society. Series B*, 61(3):611–622.

Trefethen, L. N. and Bau, D. (1997). *Numerical Linear Algebra*. SIAM.

Tseng, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, 109(3):475–494.