

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

# Interpretable Matrix Completion: A Discrete Optimization Approach

Dimitris Bertsimas

Sloan School of Management and Operations Research Center, Massachusetts Institute of Technology, Cambridge, MA 02139,  
dbertsim@mit.edu

Michael Lingzhi Li

Operations Research Center, Massachusetts Institute of Technology, Cambridge, MA 02139, mlli@mit.edu

We consider the problem of matrix completion on an  $n \times m$  matrix. We introduce the problem of *Interpretable Matrix Completion* that aims to provide meaningful insights for the low-rank matrix using side information. We show that the problem can be reformulated as a binary convex optimization problem. We design OptComplete, based on a novel concept of stochastic cutting planes to enable efficient scaling of the algorithm up to matrices of sizes  $n = 10^6$  and  $m = 10^6$ . We report experiments on both synthetic and real-world datasets that show that OptComplete has favorable scaling behavior and accuracy when compared with state-of-the-art methods for other types of matrix completion, while providing insight on the factors that affect the matrix.

*Key words:* Matrix Completion, Mixed-Integer Optimization, Stochastic Approximation

---

## 1. Introduction

Low-rank matrix completion has attracted much attention after the successful application in the Netflix Competition. It is now widely utilized in far-reaching areas such as computer vision (Candes and Plan (2010)), signal processing (Ji et al. (2010)), and control theory (Boyd et al. (1994)) to generate a completed matrix from partially observed entries.

The classical low-rank matrix completion problem considers the following problem: Given a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  with entries only partially known (denote  $\Omega \subset \{1, \dots, n\} \times \{1, \dots, m\}$  as the set of known entries), we aim to recover a matrix  $\mathbf{X} \in \mathbb{R}^{n \times m}$  of rank  $k$  that minimizes a certain distance metric between  $\mathbf{X}$  and  $\mathbf{A}$  on the known entries of  $\mathbf{A}$ :

$$\min_{\mathbf{X}} \frac{1}{nm} \sum_{(i,j) \in \Omega} \|X_{ij} - A_{ij}\| \quad \text{subject to} \quad \text{Rank}(\mathbf{X}) = k,$$

where we normalized the objective so that it is  $O(1)$ . The rank  $k$  constraint on  $\mathbf{X}$  can be equivalently formulated as the existence of two matrices  $\mathbf{U} \in \mathbb{R}^{n \times k}$ ,  $\mathbf{V} \in \mathbb{R}^{m \times k}$  such that  $\mathbf{X} = \mathbf{UV}^T$ . Therefore, the problem can be restated as:

$$\min_{\mathbf{U}} \min_{\mathbf{V}} \frac{1}{nm} \sum_{(i,j) \in \Omega} \|X_{ij} - A_{ij}\| \quad \text{subject to} \quad \mathbf{X} = \mathbf{UV}^T. \quad (1)$$

In many applications for matrix completion, it is customary for each row of the data to represent an individual and each column a product or an item of interest, and  $A_{ij}$  being the response data of individual  $i$  on item  $j$ . Therefore, the matrices  $\mathbf{U}$  and  $\mathbf{V}$  are commonly interpreted as the “user matrix” and “product matrix” respectively.

Let us denote each row of  $\mathbf{U}$  as  $\mathbf{u}_i$  and each column as  $\mathbf{u}^i$  (similarly for  $\mathbf{V}$ ). Then,  $\mathbf{u}^i$  ( $\mathbf{v}^i$ ) represents a “latent feature” for users (products), and in total there are  $k$  latent features for users (products). The goal of matrix completion is thus to discover such latent features of the users and the products, so that the dot product of such features on user  $i$  and product  $j$ ,  $X_{ij} = \mathbf{u}_i \cdot \mathbf{v}_j$ , is the intended response of user  $i$  on product  $j$ .

While this interpretation is intuitive, it does not offer insight on what the latent features of users and products mean. Inductive Matrix Completion, first considered in Dhillon et al. (2013), aims to rectify such problem by asserting that each of the  $k$  latent features is a linear combination of  $p > k$  known features. We focus on the *one-sided information* case, where only one of the user/product matrix is subject to such constraint. This case is more relevant as features related to users are often fragmented and increasingly constrained by data privacy regulations, while features about

products are easy to obtain. We would also have a short discussion later on why the two-sided information case is not interesting under the context of this paper.

Without loss of generality, we would assume the information is on the product matrix and denote the known feature matrix with  $p$  features as  $\mathbf{B} \in \mathbb{R}^{m \times p}$ . As an example, if the items are movies, then  $\mathbf{b}^j$  represents feature  $j$  for a movie (actors, budget, running time, etc), and each product feature  $\mathbf{v}^j$  needs to be linear combination of such features. Mathematically, this translates to the constraint:

$$\mathbf{V} = \mathbf{B}\mathbf{S}$$

where  $\mathbf{S} \in \mathbb{R}^{p \times k}$ . Therefore, the inductive version of the problem in (1) can be written as:

$$\min_{\mathbf{U}} \min_{\mathbf{S}} \frac{1}{nm} \sum_{(i,j) \in \Omega} \|X_{ij} - A_{ij}\| \quad \text{subject to} \quad \mathbf{X} = \mathbf{U}\mathbf{S}^T\mathbf{B}^T. \quad (2)$$

Although the inductive version of the problem adds more interpretability to the product latent features, they are still far from fully interpretable. For example, if we take the items to be movies, and features to be running time, budget, box office, and number of top 100 actors, then the generated features could look like:

$$5.6 \times \text{running time} - 0.00067 \times \text{budget} + 12 \times \# \text{ of Top 100 actors},$$

$$0.25 \times \text{box office} - 5 \times \# \text{ of Top 100 actors}.$$

These features, although a linear combination of interpretable features, are not very interpretable itself due to the involvement of multiple factors with different units. Therefore, it cannot significantly help decision makers to understand “what is important” about the product. Furthermore, the appearance of the same factor in multiple features with different signs (as shown above) further complicates any attempt at understanding the result.

Therefore, we argue that instead of supposing the  $k$  features are linear combinations of the  $p$  known features, we should assume that the  $k$  features are *selected* from the  $p$  known features. This

formulation alleviates the two previous problems mentioned: it guarantees the latent features to be interpretable (as long as the original features are), and it prevents any duplicating features in the selected  $k$  latent features. We denote this *Interpretable Matrix Completion*. We use the term interpretable, as opposed to inductive, to highlight that our approach, like sparse linear regression, gives actionable insights on what are the important features of matrix  $\mathbf{A}$ . We note that Interpretable Matrix Completion is considerably harder than inductive or the classical matrix completion problem as it is a discrete problem and selecting  $k$  out of  $p$  factors is exponential in complexity.

In this paper, we show that the Interpretable Matrix Completion problem can be written as a mixed integer convex optimization problem. Inspired by Bertsimas and van Parys (2020) for sparse linear regression, we reformulate the interpretable matrix completion problem as a binary convex optimization problem. Then we introduce a new algorithm OptComplete, based on stochastic cutting planes, to enable scalability for matrices of sizes on the order of  $(n, m) = (10^6, 10^6)$ . In addition, we provide empirical evidence on both synthetic and real-world data that OptComplete is able to match or exceed current state-of-the-art methods for inductive and general matrix completion on both speed and accuracy, despite OptComplete solving a more difficult problem.

Specifically, our contributions in this paper are as follows:

1. We introduce the interpretable matrix completion problem, and reformulate it as a binary convex optimization problem that can be solved using cutting planes methods.
2. We propose a new novel approach to cutting planes by introducing stochastic cutting planes. We prove that the new algorithm converges to an optimal solution of the interpretable matrix completion problem with exponentially vanishing failure probability.
3. We present computational results on both synthetic and real datasets that show that the algorithm matches or outperforms current state-of-the-art methods in terms of both scalability and accuracy.

The structure of the paper is as follows. In Section 2, we introduce the binary convex reformulation of the low-rank interpretable matrix completion problem, and how it can be solved through a cutting plane algorithm, which we denote CutPlanes. In Section 3, we introduce OptComplete, a stochastic cutting planes method designed to scale the CutPlanes algorithm in Section 2, and show that it recovers the optimal solution of CutPlanes with exponentially vanishing failure probability. In Section 4, we report on computational experiments with synthetic data that compare OptComplete to Inductive Matrix Completion (IMC) introduced in Natarajan and Dhillon (2014) and SoftImpute-ALS (SIALS) by Hastie et al. (2015), two state-of-the-art matrix completion algorithms for inductive and general completion. We also compare OptComplete to CutPlanes to demonstrate the 20x to 60x speedup of the stochastic algorithm. In Section 5, we report computational experiments on the Netflix Prize dataset. In Section 6 we provide our conclusions.

## Literature

Matrix completion has been applied successfully to many tasks, including recommender systems Koren et al. (2009), social network analysis Chiang et al. (2014) and clustering Chen et al. (2014b). After Candès and Tao (2010) proved a theoretical guarantee for the retrieval of the exact matrix under the nuclear norm convex relaxation, a lot of methods have focused on the nuclear norm problem (see Mazumder et al. (2010), Beck and Teboulle (2009), Jain et al. (2010), and Tanner and Wei (2013) for examples). Alternative methods include alternating projections by Recht and Ré (2013) and Grassmann manifold optimization by Keshavan et al. (2009). There has also been work where the uniform distributional assumptions required by the theoretical guarantees are violated, such as Negahban and Wainwright (2012) and Chen et al. (2014a).

Interest in inductive matrix completion intensified after Xu et al. (2013) showed that given predictive side information, one only needs  $O(\log n)$  samples to retrieve the full matrix. Thus, most of this work (see Xu et al. (2013), Jain and Dhillon (2013), Farhat et al. (2013),

Natarajan and Dhillon (2014)) have focused on the case in which the side information is assumed to be perfectly predictive so that the theoretical bound of  $O(\log n)$  sample complexity Xu et al. (2013) can be achieved. Chiang et al. (2015) explored the case in which the side information is corrupted with noise, while Shah et al. (2017) and Si et al. (2016) incorporated nonlinear combination of factors into the side information. Surprisingly, as pointed out by a recent article Nazarov et al. (2018), there is a considerable lack of effort to introduce sparsity/interpretability into inductive matrix completion, with Lu et al. (2016), Soni et al. (2016) and Nazarov et al. (2018) being among the only works that attempt to do so. Our work differs from the previous attempts in that previous attempts mainly focus on choosing latent features which are *sparse* linear combinations of the given features. In contrast *interpretable matrix completion* is aimed to *select* exactly  $k$  features from the known features.

## 2. Interpretable Matrix Completion

In this section, we present the mathematical formulation of Interpretable Matrix Completion and how it can be reformulated as a binary convex problem that is based on Bertsimas and van Parys (2020). We show how this naturally leads to a cutting plane algorithm, and discuss its computational complexity. We also discuss the two-sided information case, and how that reduces to the sparse regression problem.

### 2.1. Binary Convex Reformulation of Interpretable Matrix Completion

The (one-sided) interpretable matrix completion problem can be written as a mixed binary optimization problem:

$$\min_U \min_{\mathbf{s} \in S_k^p} \frac{1}{nm} \sum_{(i,j) \in \Omega} \|X_{ij} - A_{ij}\| \quad \text{subject to} \quad \mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{B}^T,$$

where  $\mathbf{S} = \text{Diag}\{s_1, \dots, s_p\} \in \mathbb{R}^{p \times p}$  and:

$$S_k^p = \left\{ \mathbf{s} = (s_1, \dots, s_p)^T \in \{0, 1\}^p : \sum_{i=1}^p s_i = k \right\}.$$

We note that given that  $\sum_{i=1}^p s_i = k$ , the rank of matrix  $\mathbf{X}$  is indeed  $k$ . We further note that the coefficients of  $\mathbf{S}$  can be taken to be binary without loss of generality, since if they are not and  $\mathbf{S} = \text{Diag}(1/d_1, \dots, 1/d_p)$ , then by applying the transformation:

$$\mathbf{U} \rightarrow \mathbf{UD} \quad \mathbf{S} \rightarrow \mathbf{SD}^{-1} \quad (3)$$

for  $\mathbf{D} = \text{Diag}(d_1, \dots, d_p)$ , results in an equivalent problem with the coefficients of  $\mathbf{S}$  being binary.

For this paper, we consider the squared norm, and for robustness purposes (see Bertsimas and van Parys (2020) and Bertsimas and Copenhaver (2018)), we add a Tikhonov regularization term to the original problem. Specifically, the (one-sided) interpretable matrix completion problem with regularization we address is

$$\min_{\mathbf{U}} \min_{\mathbf{s} \in S_k^p} \frac{1}{nm} \left( \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \frac{1}{\gamma} \|\mathbf{U}\|_2^2 \right) \quad \text{subject to} \quad \mathbf{X} = \mathbf{USB}^T. \quad (4)$$

In this section, we show how that problem (4) can be reformulated as a binary convex optimization problem, and can be solved to optimality using a cutting plane algorithm. The main theorem and proof is presented below:

**Theorem 1** *Problem (4) can be reformulated as a binary convex optimization problem:*

$$\min_{\mathbf{s} \in S_k^p} c(\mathbf{s}) = \frac{1}{nm} \sum_{i=1}^n \bar{\mathbf{a}}_i \left( \mathbf{I}_m + \gamma \mathbf{W}_i \left( \sum_{j=1}^p s_j \mathbf{K}_j \right) \mathbf{W}_i \right)^{-1} \bar{\mathbf{a}}_i^T,$$

where  $\mathbf{W}_1, \dots, \mathbf{W}_n \in \mathbb{R}^{m \times m}$  are diagonal matrices:

$$(\mathbf{W}_i)_{jj} = \begin{cases} 1, & (i, j) \in \Omega, \\ 0, & (i, j) \notin \Omega, \end{cases}$$

$\bar{\mathbf{a}}_i = \mathbf{W}_i \mathbf{a}_i$ ,  $i = 1, \dots, n$ , where  $\mathbf{a}_i \in \mathbb{R}^{1 \times m}$  is the  $i$ th row of  $\mathbf{A}$  with unknown entries taken to be 0, and  $\mathbf{K}_j = \mathbf{b}^j (\mathbf{b}^j)^T \in \mathbb{R}^{m \times m}$ ,  $j = 1, \dots, p$  with  $\mathbf{b}^j \in \mathbb{R}^{m \times 1}$  the  $j$ th column of  $\mathbf{B}$ .

*Proof:* With the diagonal matrices  $\mathbf{W}_i$  defined above, we can rewrite the sum in (4) over known entries of  $\mathbf{A}$ ,  $\sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2$ , as a sum over the rows of  $\mathbf{A}$ :

$$\sum_{i=1}^n \|(\mathbf{x}_i - \mathbf{a}_i) \mathbf{W}_i\|_2^2,$$

where  $\mathbf{x}_i \in \mathbb{R}^{1 \times m}$  is the  $i$ th row of  $\mathbf{X}$ . Using  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{B}^T$ , then  $\mathbf{x}_i = \mathbf{u}_i\mathbf{S}\mathbf{B}^T$  where  $\mathbf{u}_i \in \mathbb{R}^{1 \times m}$  is the  $i$ th row of  $\mathbf{U}$ . Moreover,

$$\|\mathbf{U}\|_2^2 = \sum_{i=1}^n \|\mathbf{u}_i\|_2^2.$$

Then, Problem (4) becomes:

$$\min_{\mathbf{s} \in S_k^p} \min_{\mathbf{U}} \frac{1}{nm} \left( \sum_{i=1}^n \left( \|(\mathbf{u}_i\mathbf{S}\mathbf{B}^T - \mathbf{a}_i)\mathbf{W}_i\|_2^2 + \frac{1}{\gamma} \|\mathbf{u}_i\|_2^2 \right) \right).$$

We then notice that within the sum  $\sum_{i=1}^n$  each row of  $\mathbf{U}$  can be optimized separately, leading to:

$$\min_{\mathbf{s} \in S_k^p} \frac{1}{nm} \left( \sum_{i=1}^n \min_{\mathbf{u}_i} \left( \|(\mathbf{u}_i\mathbf{S}\mathbf{B}^T - \mathbf{a}_i)\mathbf{W}_i\|_2^2 + \frac{1}{\gamma} \|\mathbf{u}_i\|_2^2 \right) \right). \quad (5)$$

The inner optimization problem  $\min_{\mathbf{u}_i} \|(\mathbf{u}_i\mathbf{S}\mathbf{B}^T - \mathbf{a}_i)\mathbf{W}_i\|_2^2 + \frac{1}{\gamma} \|\mathbf{u}_i\|_2^2$  can be solved in closed form given  $\mathbf{S}$ , as it is a weighted linear regression problem with Tikhonov regularization, see Bertsimas and van Parys (2020). The closed form solution is:

$$\min_{\mathbf{u}_i} \|(\mathbf{u}_i\mathbf{S}\mathbf{B}^T - \mathbf{a}_i)\mathbf{W}_i\|_2^2 + \frac{1}{\gamma} \|\mathbf{u}_i\|_2^2 = \bar{\mathbf{a}}_i (\mathbf{I}_m + \gamma \mathbf{W}_i \mathbf{B} \mathbf{S} \mathbf{B}^T \mathbf{W}_i)^{-1} \bar{\mathbf{a}}_i^T. \quad (6)$$

So Problem (5) can be simplified to:

$$\min_{\mathbf{s} \in S_k^p} \frac{1}{nm} \left( \sum_{i=1}^n \bar{\mathbf{a}}_i (\mathbf{I}_m + \gamma \mathbf{W}_i \mathbf{B} \mathbf{S} \mathbf{B}^T \mathbf{W}_i)^{-1} \bar{\mathbf{a}}_i^T \right).$$

Finally, we notice that

$$\mathbf{B}\mathbf{S}\mathbf{B}^T = \sum_{j=1}^p s_j \mathbf{b}^j (\mathbf{b}^j)^T = \sum_{j=1}^p s_j \mathbf{K}_j,$$

and we obtain the required expression. Since  $\mathbf{K}_j$  are positive semi-definite, and the inverse of positive semi-definite matrices is a convex function, the entire function is convex in  $\mathbf{s}$ .  $\square$

With Theorem 1, our original problem can now be restated as:

$$\min_{\mathbf{s} \in S_k^p} c(\mathbf{s}) = \frac{1}{nm} \sum_{i=1}^n \bar{\mathbf{a}}_i \left( \mathbf{I}_m + \gamma \mathbf{W}_i \left( \sum_{j=1}^p s_j \mathbf{K}_j \right) \mathbf{W}_i \right)^{-1} \bar{\mathbf{a}}_i^T. \quad (7)$$

This can be solved utilizing the cutting plane algorithm first introduced by Duran and Grossmann (1986), summarized as Algorithm 1.



---

**Algorithm 1** Cutting-plane algorithm for matrix completion with side information.

---

```

1: procedure CUTPLANES( $\mathbf{A}, \mathbf{B}$ )           # masked matrix  $\mathbf{A}$ , and feature matrix  $\mathbf{B}$ 
2:    $t \leftarrow 1$ 
3:    $\mathbf{s}_1 \leftarrow$  warm start           # Heuristic Warm Start
4:    $\eta \leftarrow 0$                      # Initialize feasible solution variable
5:   while  $\eta_t < c(\mathbf{s}_t)$  do         # While the current solution is not optimal
6:      $\mathbf{s}_{t+1}, \eta_{t+1} \leftarrow \arg \min_{\mathbf{s} \in S_k^p, \eta > 0} \eta$  s.t.  $\eta \geq c(\mathbf{s}_i) + \nabla c(\mathbf{s}_i)^T (\mathbf{s} - \mathbf{s}_i) \quad \forall i \in [t]$ 
7:      $t \leftarrow t + 1$ 
8:   end while
9:    $\mathbf{s} \leftarrow \mathbf{s}_t$ 
10:   $i \leftarrow 1$ 
11:  for  $i < n$  do                       # Fill each row  $\mathbf{x}_i$  of final output matrix  $\mathbf{X}$ 
12:     $\mathbf{x}_i \leftarrow \mathbf{B}^s ((\mathbf{B}^s)^T \mathbf{W}_i \mathbf{B}^s)^{-1} (\mathbf{B}^s)^T \bar{\mathbf{a}}_i^T$  #  $\mathbf{B}^s$  is submatrix of  $\mathbf{B}$  with  $\mathbf{s}$  columns
13:  end for
14:  return  $\mathbf{X}$                            # Return the filled matrix  $\mathbf{X}$ 
15: end procedure

```

---

The cutting plane algorithm, at iteration  $t$ , adds a linear approximation of  $c(\mathbf{s})$  at the current feasible solution  $\mathbf{s}_t$  to the set of constraints:

$$\eta \geq c(\mathbf{s}_t) + \nabla c(\mathbf{s}_t)^T (\mathbf{s} - \mathbf{s}_t), \quad (8)$$

and we solve the mixed-integer linear programming problem:

$$\begin{aligned} \min_{\mathbf{s} \in S_k^p, \eta \geq 0} \quad & \eta \\ & \eta \geq c(\mathbf{s}_t) + \nabla c(\mathbf{s}_t)^T (\mathbf{s} - \mathbf{s}_t), \quad i \in [t] \end{aligned}$$

to obtain  $\mathbf{s}_{t+1}, \eta_{t+1}$ . We see that  $\eta_{t+1}$  is exactly the minimum value of the current approximation of  $c(\mathbf{s})$ ,  $c_t(\mathbf{s})$ , defined below:

$$\eta_{t+1} = \min_{\mathbf{s}} \max_{i \in [t]} c(\mathbf{s}_t) + \nabla c(\mathbf{s}_t)^T (\mathbf{s} - \mathbf{s}_t) = \min_{\mathbf{s}} c_t(\mathbf{s}).$$

Since  $c(\mathbf{s})$  is convex, the piecewise linear approximation  $c_t(\mathbf{s})$  is an outer approximation ( $c_t(\mathbf{s}) \leq c(\mathbf{s}) \forall \mathbf{s}$ ), so  $\eta_t \leq c(\mathbf{s}_t) \forall t$ . As the algorithm progresses, the set of linear approximations form an increasingly better approximation of  $c(\mathbf{s})$ , and  $\eta_t$  increases with  $t$ . The algorithm terminates once  $\eta_t$  does not further increase, as it implies the linear approximation shares the same minimum value as the true function  $c(\mathbf{s})$ , which is the desired value.

Once the optimal solution  $\mathbf{s}^*$  is reached, we can obtain the optimal  $\mathbf{U}$  using the closed form solution in (6) and recover  $\mathbf{X}$ . In the next section, we discuss how this algorithm can be implemented in the context of  $c(\mathbf{s})$  in (7) and derive its computational complexity.

## 2.2. Implementation and Computational Complexity of CutPlanes

The computational complexity of the cutting plane comes from calculating  $c(\mathbf{s})$  and its derivative  $\nabla c(\mathbf{s})$ . We first introduce the notations  $\alpha_i(\mathbf{s}) \in \mathbb{R}$  and  $\gamma_i(\mathbf{s}) \in \mathbb{R}^{m \times 1}$ .

$$\alpha_i(\mathbf{s}) = \frac{1}{m} \bar{\mathbf{a}}_i^T \gamma_i(\mathbf{s}) = \frac{1}{m} \bar{\mathbf{a}}_i^T \left[ \left( \mathbf{I}_m + \gamma \mathbf{W}_i \left( \sum_{j=1}^p s_j \mathbf{K}_j \right) \mathbf{W}_i \right)^{-1} \bar{\mathbf{a}}_i^T \right], \quad i = 1, \dots, n. \quad (9)$$

Then, the function  $c(\mathbf{s})$  in (7) can be expressed as

$$c(\mathbf{s}) = \frac{1}{n} \sum_{i=1}^n \alpha_i(\mathbf{s}) = \frac{1}{n} \sum_{i=1}^n \frac{\bar{\mathbf{a}}_i^T \gamma_i(\mathbf{s})}{m}. \quad (10)$$

To calculate the derivative  $\nabla c(\mathbf{s})$ , it is easier to utilize the expression in Theorem 1 and then utilize the chain rule. After some manipulations, we obtain

$$\nabla c(\mathbf{s}) = \frac{1}{n} \sum_{i=1}^n -\frac{\gamma(\mathbf{B}^T \mathbf{W}_i \gamma_i(\mathbf{s}))^2}{m}. \quad (11)$$

Therefore, we would focus on calculating  $\gamma_i(\mathbf{s})$ . First, by the Matrix Inversion Lemma (Woodbury (1949)) we have

$$\begin{aligned} \gamma_i(\mathbf{s}) &= \left( \mathbf{I}_m + \gamma \mathbf{W}_i \left( \sum_{j=1}^p s_j \mathbf{K}_j \right) \mathbf{W}_i \right)^{-1} \bar{\mathbf{a}}_i^T \\ &= \left( \mathbf{I}_m - \mathbf{V} \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}^T \mathbf{W}_i \mathbf{V} \right)^{-1} \mathbf{V}^T \right) \bar{\mathbf{a}}_i^T \\ &= \left( \bar{\mathbf{a}}_i^T - \mathbf{V} \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}^T \mathbf{W}_i \mathbf{V} \right)^{-1} \mathbf{V}^T \bar{\mathbf{a}}_i^T \right), \end{aligned} \quad (12)$$

where  $\mathbf{V} \in \mathbb{R}^{m \times k}$  is the feature matrix formed by the  $k$  columns of  $\mathbf{B}$  such that  $s_j = 1$ , and we have suppressed the dependency of  $\mathbf{V}$  on  $\mathbf{s}$  for notation ease. Note that in order to compute  $\gamma_i(\mathbf{s})$  using Eq. (9) we need to invert an  $m \times m$  matrix, while from Eq. (12) we need to invert a  $k \times k$  matrix  $\frac{I_k}{\gamma} + \mathbf{V}^T \mathbf{W}_i \mathbf{V}$ , which only requires  $O(k^3)$  calculations. Furthermore, note that calculating  $\bar{\mathbf{a}}_i \bar{\mathbf{a}}_i^T$  only requires  $|\Omega_i|$  multiplications where  $\Omega_i$  is the number of known entries in row  $i$  of  $\mathbf{A}$  as we do not need to multiply on the unknown entries. Similarly, we can compute  $\mathbf{V}^T \bar{\mathbf{a}}_i^T$  in  $|\Omega_i|k$  multiplications, and  $\mathbf{V}^T \mathbf{W}_i \mathbf{V}$  in  $|\Omega_i|k^2$  multiplications.

Therefore, we can compute  $\gamma_i(\mathbf{s})$  in floating point complexity of  $O(|\Omega_i|k^2 + k^3)$ . Then to calculate  $\bar{\mathbf{a}}_i \gamma_i(\mathbf{s})$  in (10) and  $-\gamma(\mathbf{B}^T \mathbf{W}_i \gamma_i(\mathbf{s}))^2$  (11) only requires  $O(|\Omega_i|)$  and  $O(|\Omega_i|p)$  calculations respectively. Thus, the total complexity of generating a full cutting plane is:

$$\sum_{i=1}^n O(|\Omega_i|p + |\Omega_i|k^2 + k^3) = O(|\Omega|(p + k^2) + nk^3). \quad (13)$$

### 2.3. Two-sided Information Case

In this section, we briefly discuss the matrix completion problem under the two-sided information case, and how it reduces to the problem of sparse linear regression. The two sided interpretable matrix completion problem with Tikhonov regularization can be stated as follows:

$$\min_{\mathbf{L}} \frac{1}{nm} \left( \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \frac{1}{\gamma} \|\mathbf{L}\|_2^2 \right) \quad \text{subject to} \quad \mathbf{X} = \mathbf{U} \mathbf{L} \mathbf{B}^T \quad \|\mathbf{L}\|_0 = k, \quad (14)$$

where  $\mathbf{U} \in \mathbb{R}^{n \times p_1}$  is a known matrix of  $p_1$  features of each row,  $\mathbf{B} \in \mathbb{R}^{m \times p_2}$  is a known matrix of  $p_2$  features of each column, and  $\mathbf{L} \in \mathbb{R}^{p_1 \times p_2}$  is a sparse matrix that has  $k$  nonzero entries, ensuring that  $\text{Rank}(\mathbf{X}) \leq k$ . We note that in Eq. (14) we restrict the support of matrix  $\mathbf{L}$  to be  $k$ , rather than forcing the entries of  $\mathbf{L}$  to be binary. This is because unlike in the one-sided case, both  $\mathbf{U}$  and  $\mathbf{B}$  are known, so we cannot apply the scaling transformation in (3).

We denote by  $\mathbf{u}^i \in \mathbb{R}^{n \times 1}$  the  $i$ th column of  $\mathbf{U}$  and  $\mathbf{b}^j \in \mathbb{R}^{m \times 1}$  the  $j$ th column of  $\mathbf{B}$ . We introduce the matrices  $\mathbf{W}_i$  as in Theorem 1. Using  $\mathbf{X} = \mathbf{U} \mathbf{L} \mathbf{B}$ , we can write

$$X_{ij} = \sum_{q=1}^{p_1} \sum_{\ell=1}^{p_2} L_{q,\ell} D_{ij}^{q,\ell},$$

where  $D_{ij}^{q,\ell} = (\mathbf{u}^q(\mathbf{b}^\ell)^T)_{ij}$  is the  $(i, j)$ th entry of the matrix formed by multiplying  $q$ th column of  $\mathbf{U}$  with  $\ell$ th column of  $\mathbf{B}$ . Then, Problem (14) becomes:

$$\min_{\mathbf{L}} \frac{1}{nm} \left( \sum_{(i,j) \in \Omega} \left( \sum_{q=1}^{p_1} \sum_{\ell=1}^{p_2} L_{q,\ell} D_{ij}^{q,\ell} - A_{ij} \right)^2 + \frac{1}{\gamma} \|\mathbf{L}\|_2^2 \right) \quad \text{subject to} \quad \|\mathbf{L}\|_0 = k. \quad (15)$$

As every  $\mathbf{D}$  matrix is known, this becomes a sparse regression problem where there are  $p_1 p_2$  features to choose from (the  $\mathbf{D}$  matrices), there are  $|\Omega|$  samples (the  $\mathbf{A}$  matrix), the sparsity requirement is  $k$ , the regression coefficients are  $\mathbf{L}$ , and we have Tikhonov regularization. Vectorizing  $\mathbf{D}$ ,  $\mathbf{L}$ , and  $\mathbf{A}$  reduces the problem back to the familiar form of sparse linear regression, that can be solved by the algorithm developed in Bertsimas and van Parys (2020) at scale.

### 3. OptComplete: The Stochastic Cutting Plane Speedup

In this section, we introduce OptComplete, a stochastic version of the cutting plane algorithm introduced in Section 2. We present theoretical results to show that the stochastic algorithm recovers the true optimal solution of the original algorithm with high probability without distributional assumptions. We also include a discussion on the dependence of such probability with various factors and its favorable theoretical computational complexity.

#### 3.1. Introduction of OptComplete

In the previous section, we showed that through careful evaluation, we can calculate a full cutting plane in  $O(|\Omega|(p+k^2) + nk^3)$  calculations. However, in very high dimensions where  $|\Omega|, n, m$  are extremely large, the cost of generating the full cutting plane is still prohibitive. Thus, we consider generating approximations of the cutting plane that would enable the algorithm to scale for high values for  $n$  and  $m$ . Specifically, consider the cutting plane function in (10), reproduced below:

$$c(\mathbf{s}) = \frac{1}{n} \sum_{i=1}^n \alpha_i(\mathbf{s}),$$

where:

$$\alpha_i(\mathbf{s}) = \frac{1}{m} \left( \bar{\mathbf{a}}_i \bar{\mathbf{a}}_i^T - \bar{\mathbf{a}}_i \mathbf{V} \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}^T \mathbf{W}_i \mathbf{V} \right)^{-1} \mathbf{V}^T \bar{\mathbf{a}}_i^T \right).$$

We approximate the inner term  $\alpha_i(\mathbf{s})$  by choosing  $1 \leq f < m$  samples from  $\{1, \dots, m\}$  *without replacement*, with the set denoted  $F$ . Then we formulate the submatrix  $\mathbf{V}_F$  with such selected rows, and similarly with  $\bar{\mathbf{a}}_{Fi}$ . Then we calculate the approximation:

$$\alpha_i(\mathbf{s}) \approx \alpha_i^F(\mathbf{s}) = \frac{1}{f} \left( \bar{\mathbf{a}}_{Fi} \bar{\mathbf{a}}_{Fi}^T - \bar{\mathbf{a}}_{Fi} \mathbf{V}_F \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}_F^T \mathbf{W}_i \mathbf{V}_F \right)^{-1} \mathbf{V}_F^T \bar{\mathbf{a}}_{Fi}^T \right).$$

Then we choose  $1 \leq g < n$  samples from  $\{1, \dots, n\}$  without replacement, with the set denoted  $G$ .

We can then calculate an approximation of  $c(\mathbf{s})$  using the approximated  $\alpha_i^F(\mathbf{s})$ :

$$c(\mathbf{s}) \approx \tilde{c}_G^F(\mathbf{s}) = \frac{1}{r} \sum_{i \in G} \alpha_i^F(\mathbf{s}),$$

where the set  $F$  is chosen independently for every row  $i \in G$ . Then the derivative of  $\tilde{c}_G^F(\mathbf{s})$  is:

$$\nabla \tilde{c}_G^F(\mathbf{s}) = \frac{1}{r} \sum_{i \in G} \frac{(\mathbf{B}^T \mathbf{W}_i \boldsymbol{\gamma}_i^F(\mathbf{s}))^2}{f}.$$

Using such approximations, we can derive a stochastic cutting-plane algorithm, which we call OptComplete presented as Algorithm 2.

For this algorithm to work, we need the approximation  $\tilde{c}_G^F(\mathbf{s})$  and its derivative to be close to the nominal values. Furthermore, the approximated cutting planes should not cutoff the true solution. In the next section, we show that OptComplete enjoys such properties with high probability. In Section 3.3, we discuss how to select the size of  $f$  and  $g$ .

### 3.2. Main Theoretical Results

We would first show that the inner approximation is close to the true term with high probability:

**Theorem 2** *Let  $\mathbf{A}$  be a partially known matrix,  $\mathbf{B}$  a known feature matrix, and  $\mathbf{W}_i$  as defined in Theorem 1. Let  $F$  be a random sample of size  $f$  from the set  $\{1, \dots, m\}$ , chosen without replacement. With probability at least  $1 - \epsilon$ , we have*

$$\begin{aligned} |\alpha_i(\mathbf{s}) - \alpha_i^F(\mathbf{s})| &\leq \sqrt{\frac{Mk \log(\frac{k}{\epsilon})}{f}}, \quad \forall i \in \{1, \dots, m\}, \quad \forall \mathbf{s} \in S_k^p, \\ \left\| \frac{(\mathbf{B}^T \mathbf{W}_i \boldsymbol{\gamma}_i(\mathbf{s}))^2}{m} - \frac{(\mathbf{B}^T \mathbf{W}_i \boldsymbol{\gamma}_i^F(\mathbf{s}))^2}{f} \right\|_2 &\leq \sqrt{\frac{M'(p+k) \log(\frac{k}{\epsilon})}{f}}, \quad \forall i \in \{1, \dots, m\}, \quad \forall \mathbf{s} \in S_k^p, \end{aligned}$$

where  $M, M'$  are absolute constants.

---

**Algorithm 2** Stochastic Cutting-plane algorithm for matrix completion with side information.
 

---

```

1: procedure OPTCOMPLETE( $\mathbf{A}, \mathbf{B}$ )           # masked matrix  $\mathbf{A}$ , and feature matrix  $\mathbf{B}$ 
2:    $t \leftarrow 1$ 
3:    $\mathbf{s}_1 \leftarrow$  random initialization
4:    $\eta \leftarrow 0$                            # Initialize feasible solution variable
5:   while  $\eta_t < c(\mathbf{s}_t)$  do                 # While the current solution is not optimal
6:      $G \leftarrow$   $|g|$ -sized random sample of  $\{1, \dots, n\}$  with replacement
7:     for  $i \in G$  do                           # Generate  $F$  for each row in random sample
8:        $F_i \leftarrow$   $|f|$ -sized random sample of  $\{1, \dots, m\}$  with replacement
9:     end for
10:     $\mathbf{s}_{t+1}, \eta_{t+1} \leftarrow \arg \min_{\mathbf{s} \in S_k^p, \eta > 0} \eta$  s.t.  $\eta \geq \tilde{c}_G^F(\mathbf{s}_i) + \nabla \tilde{c}_G^F(\mathbf{s}_i)^T (\mathbf{s} - \mathbf{s}_i) \quad \forall i \in [t]$ 
11:     $t \leftarrow t + 1$ 
12:  end while
13:   $\mathbf{s} \leftarrow \mathbf{s}_t$ 
14:   $i \leftarrow 1$ 
15:  for  $i < n$  do                               # Fill each row  $\mathbf{x}_i$  of final output matrix  $\mathbf{X}$ 
16:     $\mathbf{x}_i \leftarrow \mathbf{B}^s ((\mathbf{B}^s)^T \mathbf{W}_i \mathbf{B}^s)^{-1} (\mathbf{B}^s)^T \bar{\mathbf{a}}_i^T$  #  $\mathbf{B}^s$  is submatrix of  $\mathbf{B}$  with  $\mathbf{s}$  columns
17:  end for
18:  return  $\mathbf{X}$                                    # Return the filled matrix  $\mathbf{X}$ 
19: end procedure

```

---

We see that, without assumptions on the data, the inner approximation for both the value and the derivative follows a bound with  $O\left(\sqrt{\frac{(p+k)}{f}}\right)$  terms with very high probability. Furthermore, inverting the statements give that, for all  $i \in \{1, \dots, m\}$  and all  $\mathbf{s} \in S_k^p$ :

$$\mathbb{P}\left(|\alpha_i(\mathbf{s}) - \alpha_i^F(\mathbf{s})| \geq \delta\right) \leq k \exp\left(-\frac{f\delta^2}{Mk}\right),$$

$$\mathbb{P}\left(\left\|\frac{(\mathbf{B}^T \mathbf{W}_i \gamma_i(\mathbf{s}))^2}{m} - \frac{(\mathbf{B}^T \mathbf{W}_i \gamma_i^F(\mathbf{s}))^2}{f}\right\|_2 \geq \delta\right) \leq k \exp\left(-\frac{f\delta^2}{M'(p+k)}\right).$$

So the failure probability drops off exponentially with increasing bound  $\delta$ , reflecting a Gaussian tail structure for  $\alpha_i(\mathbf{s})$  and  $\mathbf{B}^T \mathbf{W}_i \boldsymbol{\gamma}_i^F(\mathbf{s})$ . The proof is contained in Appendix A.

Using this result, we are able to prove a tight deviation bound for the approximated cost function  $c_G^F(\mathbf{s})$  and  $\nabla c_G^F(\mathbf{s})$ :

**Theorem 3** *Let  $\mathbf{A}$  be a partially known matrix,  $\mathbf{B}$  a known feature matrix, and  $\mathbf{W}_i$  as defined in Theorem 1. Let  $G$  be a random sample of size  $g$  from  $\{1, \dots, n\}$  chosen without replacement. Then for each  $i \in G$ , we let  $F_i$  be a random sample of size  $f$  from the set  $\{1, \dots, m\}$ , all chosen without replacement. We have, with probability at least  $1 - \epsilon$ :*

$$\begin{aligned} |\tilde{c}_G^F(\mathbf{s}) - c(\mathbf{s})| &\leq \sqrt{\frac{Ak \log\left(\frac{k}{\epsilon}\right)}{g}}, & \forall \mathbf{s} \in S_k^p, \\ \|\nabla \tilde{c}_G^F(\mathbf{s}) - \nabla c(\mathbf{s})\|_2 &\leq \sqrt{\frac{B(p+k) \log\left(\frac{k}{\epsilon}\right)}{g}}, & \forall \mathbf{s} \in S_k^p, \end{aligned}$$

where  $A, B$  are absolute constants.

Similar to the inner approximations,  $c_G^F(\mathbf{s})$  and  $\nabla c_G^F(\mathbf{s})$  has Gaussian tails. Furthermore, the scaling here only depends on  $g$  and not  $f$ : This shows that the error of the inner approximation is dominated by the outer sampling of the rows  $G$ . The proof is contained in Appendix B.

Then, using this result, we are able to prove our main result for OptComplete. We would first introduce a new definition:

**Definition 1** *The **convexity parameter**  $a$  of the cost function  $c(\mathbf{s})$  is defined as the largest positive number for which the the following statement is true:*

$$c(\mathbf{s}) \geq c(\mathbf{s}_0) + \nabla c(\mathbf{s}_0)^T (\mathbf{s} - \mathbf{s}_0) + \frac{a^2}{2} (\mathbf{s} - \mathbf{s}_0)^T (\mathbf{s} - \mathbf{s}_0) \quad \forall \mathbf{s}, \mathbf{s}_0 \in S_k^p \quad \forall i \quad (16)$$

We have the following proposition which shows that unless the cost function is degenerate (i.e. different sets of  $k$  features doesn't change the solution), we always have a positive convexity parameter:

**Proposition 1** *Assume that there does not exist  $\mathbf{s}_1, \mathbf{s}_0 \in S_k^p$  such that  $c(\mathbf{s}_1) = c(\mathbf{s}_0)$ . Then  $a > 0$ .*

The proof is contained in Appendix C. Now, we state our main theorem for OptComplete.

**Theorem 4** *For the matrix completion problem (4), let  $\mathbf{B} \in \mathbb{R}^{m \times p}$  be a known feature matrix,  $\mathbf{A} \in \mathbb{R}^{n \times m}$  a matrix with entries partially known, and OptComplete as defined in Algorithm 2. Assume that Problem (4) is feasible. Then, OptComplete terminates in a finite number of steps  $C$ , and finds an optimal solution of (4) with probability at least  $1 - kC \exp\left(-\frac{Da^4g}{(p+k)}\right)$  where  $D$  is an absolute constant independent of  $C, f, g, k, p$ , and  $a$  is the convexity parameter of the functions  $\tilde{\alpha}_i^s(\mathbf{s})$ .*

The proof is contained in Appendix D. This theorem shows that as long as the original problem is feasible, OptComplete is able to find the optimal solution of the original binary convex problem with exponentially vanishing failure probability that scales as  $O\left(\exp\left(\frac{-g}{(p+k)}\right)\right)$ . The theorem requires no assumptions on the data, and thus applies generally. We again note that the bound does not depend on  $f$  and only on  $g$ : we would discuss how this would inform our selection of the size of  $f$  and  $g$  in the next section.

### 3.3. Sampling Size and Computational Complexity

To select an appropriate  $f$  and  $g$ , we first note that Candès and Tao (2010) showed that to complete a square  $N \times N$  matrix of rank  $k$ , we need at least  $O(kN \log N)$  elements. Assume an average known rate of  $\alpha = \frac{|\Omega|}{mn}$  in the original matrix  $\mathbf{A}$ , the expected number of known elements under a sampling of  $f$  and  $g$  is  $\alpha fg$ . Using  $N^2 = mn$ , we need that:

$$\alpha fg \geq c \cdot k \sqrt{nm} \log(\sqrt{nm}) \quad (17)$$

for some constant  $c$ . Theorem 4 showed that the bound on failure probability scales with  $O\left(\exp\left(\frac{-g}{(p+k)}\right)\right)$ , and thus we cannot have  $g$  too small. Using (13), the complexity of the cutting plane with  $f$  and  $g$  samples are:

$$O(\alpha fg(p+k^2) + gk^3). \quad (18)$$



Therefore, if we fix the expected known elements ( $\alpha fg$ ) constant, it is more advantageous to select a smaller  $g$ , as  $g$  scales with  $k^3$ . Thus, we set:

$$f = \min\left(\frac{ck\sqrt{mn}\log(\sqrt{mn})}{\alpha \min(g_0, n)}, m\right), \quad g = \min(g_0, n), \quad (19)$$

Experimentally, we found  $g_0 = 100$ ,  $c = 1$  to generate good results (the results were similar for  $\frac{1}{2} \leq c \leq 2$ ). Therefore, by (18), the approximated cutting plane has a computational complexity of:

$$O(k\sqrt{mn}\log(\sqrt{mn})(p+k^2)). \quad (20)$$

This scales in a square root fashion in  $n$  and  $m$ , rather than linearly in  $n$  and  $m$  for the full cutting plane. This allows OptComplete to enjoy a considerable speedup compared to CutPlanes, as demonstrated in Section 4.

#### 4. Synthetic Data Experiments

We assume that the matrix  $\mathbf{A} = \mathbf{UV} + \mathbf{E}$ , where  $\mathbf{U} \in \mathbb{R}^{n \times k}$ ,  $\mathbf{V} \in \mathbb{R}^{k \times m}$ , and  $\mathbf{E}$  is an error matrix with individual elements sampled from  $N(0, 0.01)$ . We sample the elements of  $\mathbf{U}$  and  $\mathbf{V}$  from a uniform distribution of  $[0, 1]$ , and then randomly select a fraction  $\mu = 1 - \alpha$  to be missing. We formulate the feature matrix  $\mathbf{B}$  by combining  $\mathbf{V} \in \mathbb{R}^{k \times m}$  with a confounding matrix  $\mathbf{Z} \in \mathbb{R}^{(p-k) \times m}$  that contains unnecessary factors sampled similarly from the Uniform  $[0, 1]$  distribution. We run OptComplete on a server with 16 CPU cores, using Gurobi 8.1.0. For each combination  $(m, n, p, k, \mu)$ , we ran 10 tests and report the median value for every statistic.

We report the following statistics with  $\mathbf{s}^*$  being the ground-truth factor vector, and  $\bar{\mathbf{s}}$  the estimated factor vector.

- $n, m$  - the dimensions of  $\mathbf{A}$ .
- $p$  - the number of features in the feature matrix.
- $k$  - the true number of features.
- $\mu$  - The fraction of missing entries in  $\mathbf{A}$ .
- $T$  - the total time taken for the algorithm.

- MAPE - the Mean Absolute Percentage Error (MAPE) for the retrieved matrix  $\hat{\mathbf{A}}$ :

$$\text{MAPE} = \frac{1}{|\mathcal{S}|} \sum_{(i,j) \in \mathcal{S}} \frac{|\hat{A}_{ij} - A_{ij}|}{|A_{ij}|},$$

where  $\mathcal{S} = \Omega^c$  is the set of missing data in  $\mathbf{A}$ .

Since the concept of Interpretable Matrix Completion is new, there is a lack of directly comparable algorithms in the literature. Thus, in lieu, we compare OptComplete to state-of-the-art solvers for Inductive Matrix Completion and general matrix completion, which are:

- IMC by Natarajan and Dhillon (2014) - This algorithm is a well-accepted benchmark for testing Inductive Matrix Completion algorithms.
- SoftImpute-ALS (SIALS) by Hastie et al. (2015) - This is widely recognized as a state-of-the-art matrix completion method without feature information. It has among the best scaling behavior across all classes of matrix completion algorithms as it utilizes fast alternating least squares to achieve scalability.

We use the best existing implementations of IMC (Matlab 2018b) and SIALS (R 3.4.4, package softImpute) with parallelization on the same server.

We further compare our algorithm to CutPlanes, the original cutting plane algorithm developed in Section 2. It is known that for general mixed-integer convex problems, the cutting plane algorithm has the best overall performance (see e.g. Lubin et al. (2016) for details), and thus CutPlanes represent a good baseline of comparison for OptComplete.

We randomly selected 20% of those elements masked to serve as a validation set. The regularization parameter  $\gamma$  of OptComplete, the rank parameter of IMC and the penalization parameter  $\lambda$  of IMC and SIALS are selected using the validation set. The results are separated into sections below. The first five sections modify one single variable out of  $n, m, p, k, \mu$  to investigate OptComplete's scalability, where the leftmost column indicates the variable modified. The last section compares the four algorithms scalability for a variety of parameters that reflect more realistic scenarios.

	$n$	$m$	$p$	$k$	$\mu\%$	OptComplete		CutPlanes		IMC		SIALS	
						$T$	MAPE	$T$	MAPE	$T$	MAPE	$T$	MAPE
$\mu$	100	100	15	5	20%	1.7s	0.1%	6.0s	0.1%	0.03s	0.01%	0.02s	0.3%
	100	100	15	5	50%	0.9s	0.02%	4.5s	0.02%	0.07s	0.5%	0.03s	0.9%
	100	100	15	5	80%	0.6s	0.03%	2.5s	0.03%	0.09s	1.3%	0.06s	5.6%
	100	100	15	5	95%	0.2s	0.04%	1.2s	0.04%	0.12s	12.1%	0.12s	7.4%
$n$	100	100	15	5	50%	0.9s	0.02%	4.5s	0.02%	0.07s	0.5%	0.03s	0.9%
	$10^3$	100	15	5	50%	3.1s	0.01%	72.5s	0.01%	0.6s	0.4%	0.1s	0.2%
	$10^4$	100	15	5	50%	9.5s	0.004%	957s	0.004%	4.5s	0.3%	6.5s	0.5%
	$10^5$	100	15	5	50%	18.0s	0.003%	10856s	0.003%	32.7s	0.1%	38s	3.0%
$m$	100	100	15	5	50%	0.9s	0.02%	4.5s	0.02%	0.07s	0.5%	0.03s	0.9%
	100	$10^3$	15	5	50%	0.7s	0.01%	18.6s	0.01%	0.8s	0.3%	0.1s	0.5%
	100	$10^4$	15	5	50%	1.2s	0.004%	68.5s	0.004%	6.2s	0.2%	0.8s	0.3%
	100	$10^5$	15	5	50%	3.0s	0.002%	259s	0.002%	56.2s	0.1%	12.7s	0.8%
$p$	100	100	15	5	50%	0.9s	0.02%	4.5s	0.02%	0.07s	0.5%	0.03s	0.9%
	100	100	50	5	50%	2.0s	0.02%	18.0s	0.02%	0.3s	0.6%	0.03s	0.9%
	100	100	200	5	50%	12.1s	0.02%	95.9s	0.02%	1.9s	0.8%	0.03s	0.9%
	100	100	$10^3$	5	50%	90.3s	0.02%	680s	0.02%	10.4s	1.0%	0.03s	0.9%
$k$	100	100	50	5	50%	2.0s	0.02%	18.0s	0.02%	0.3s	0.5%	0.03s	0.9%
	100	100	50	10	50%	20.7s	0.06%	130s	0.06%	0.20s	1.2%	0.1s	0.8%
	100	100	50	20	50%	240s	0.07%	1584s	0.07%	0.35s	2.1%	0.21s	1.0%
	100	100	50	30	50%	980s	0.09%	8461s	0.09%	0.5s	3.3%	0.43s	2.8%
	100	100	15	5	95%	0.2s	0.04%	1.2s	0.04%	0.12s	12.1%	0.12s	7.4%
	$10^3$	$10^3$	50	5	95%	1.4s	0.006%	3.5s	0.006%	4.6s	4.7%	2.8s	12.5%
	$10^4$	$10^3$	100	5	95%	5.7s	0.002%	35.2s	0.002%	18s	2.5%	20.7s	12.6%
	$10^5$	$10^3$	200	10	95%	52s	0.001%	1520s	0.001%	295s	1.7%	420s	4.6%
	$10^5$	$10^4$	200	10	95%	98s	0.001%	5769s	0.001%	1750s	0.5%	4042s	4.1%
	$10^6$	$10^4$	200	10	95%	480s	0.001%	N/A	N/A	13750s	0.3%	25094s	2.5%
	$10^6$	$10^5$	200	10	95%	680s	0.001%	N/A	N/A	N/A	N/A	N/A	N/A
	$10^6$	$10^6$	200	10	95%	1415s	0.001%	N/A	N/A	N/A	N/A	N/A	N/A

**Table 1** Comparison of OptComplete, IMC and SIALS on synthetic data. *N/A* means the algorithm did not complete running in 20 hours, corresponding to 72000 seconds.

Overall, we see that OptComplete achieves near-exact retrieval on all datasets evaluated, and successfully recovers the factors in the ground truth. The solutions (and its error) also matches with that of CutPlanes, the standard cutting plane algorithm. The non-zero MAPE is due to the random noise added resulting in slightly perturbed coefficients.

For the realistic and large data sizes in the last panel, we see that OptComplete not only achieves near-exact retrieval, it does so while requiring considerably less time than IMC and SIALS at the same time. For  $m, n$  on the scale of  $n = 10^6$  and  $m = 10^4$ , OptComplete is over 20 times faster than IMC and over 40 times faster than SIALS. At the scale of  $n = 10^6$  and  $m = 10^5$ , IMC and SIALS did not finish running within 20 hours, while OptComplete completed in just under 12 minutes. We also see that OptComplete achieves very significant speedups compared to the standard cutting plane algorithm - up to 60x at the scale of  $n = 10^5$  and  $m = 10^4$ .

We analyze the scaling of OptComplete as a function of:

1.  $\mu$  - The algorithm is able to retrieve the exact factors used even with 95% of missing data. Furthermore, the running time decreased with increasing missing entries, consistent with the fact that is computational complexity scales with  $|\Omega|$ .
2.  $n$  - The algorithm has good scalability in  $n$ , reflecting its  $O(\sqrt{n} \log(n))$  type complexity. This allows the algorithm to support matrices with  $n$  in the  $10^6$  range. Its scaling behavior is superior to both IMC and SIALS.
3.  $m$  - The algorithm scales exceptionally well in  $m$ . We observe that empirically the algorithm runtime seems to grow much slower than the theoretical  $O(\sqrt{m} \log(m))$  dependence. A closer examination reveals that as  $m$  increases, the number of cutting planes generated by Gurobi is decreasing. Qualitatively, this can be explained by a larger  $m$  giving the algorithm more signal to find which  $k$  features are the correct ones out of the  $p$  ones. We note that such behavior is also exhibited by CutPlanes, as it roughly scales as  $O(\sqrt{m})$  rather than the  $O(m)$  as expected. We see that IMC and SIALS scales as  $O(m)$ .
4.  $p$  - The algorithm scales relatively well in  $p$ , which reflects the performance of the Gurobi solver. We empirically observe that Gurobi is generating roughly  $O(1) - O(p)$  cutting planes. Thus, as each cutting plane is  $O(p)$ , we expect  $O(p) - O(p^2)$  dependence, as is observed here. We note that OptComplete achieves similar scaling behavior as IMC in  $p$ . Note here the SIALS algorithm does not utilize feature information and thus a change in  $p$  does not affect the algorithm's run speed.

5.  $k$  - The algorithm does not scale very well in  $k$ . We empirically observe that Gurobi solver is roughly generating  $O(k)$  cutting planes and each cutting plane has cubic dependence on  $k$ . It appears that SIALS and IMC almost have a linear scaling behavior. However, in most applications, such as recommendation systems or low-rank retrieval,  $k$  is usually kept very low ( $k \leq 30$ ), so this is not a particular concern.

## 5. Real-World Experiments

In this section, we report on the performance of OptComplete on the Netflix Prize dataset (Bennett et al. 2007). This dataset was released in a competition to predict ratings of customers on unseen movies, given over 10 million ratings scattered across 500,000 people and 16,000 movies. Thus, when presented in a matrix  $\mathbf{A}$  where  $A_{ij}$  represents the rating of individual  $i$  on movie  $j$ , the goal is to complete the matrix  $\mathbf{A}$  under a low-rank assumption.

The feature matrix  $\mathbf{B}$  of OptComplete is constructed using data from the TMDB Database, and covers 59 features that measure geography, popularity, top actors/actresses, box office, runtime, genre and more. The full list of 59 features is contained in Appendix E.

For this experiment, we included movies where all 59 features are available, and people who had at least 5 ratings present. This gives a matrix of 471,268 people and 14,538 movies. The slight reduction of size from the original data is due to the lack of features for about 2,000 niche movies. To observe the scalability of OptComplete, we created five data sets:

1. Base -  $\mathbf{A}_1$  has dimensions  $3,923 \times 103$ .
2. Small -  $\mathbf{A}_2$  has dimensions  $18,227 \times 323$ .
3. Medium -  $\mathbf{A}_3$  has dimensions  $96,601 \times 788$ .
4. Large -  $\mathbf{A}_4$  has dimensions  $471,268 \times 1760$ .
5. Full -  $\mathbf{A}$  has dimensions  $471,268 \times 14,538$ .

These sizes are constructed such that the total number of elements in  $\mathbf{A}$  in the successive sizes are approximately different by approximately an order of magnitude.

For each individual matrix, we uniformly randomly withhold 20% of the ratings as a test set  $\mathcal{S}$ , and use the remaining 80% of ratings to impute a complete matrix  $\hat{\mathbf{A}}$  - we perform cross-validation on the appropriate hyperparameters. Then, we report MAPE.

For comparison, we again use IMC and SIALS. We set the maximum rank of SIALS to be  $k$  - the rank optimized for in OptComplete. The results are listed below:

$n$	$m$	$p$	$k$	$\mu\%$	OptComplete		IMC		SIALS	
					$T$	MAPE	$T$	MAPE	$T$	MAPE
3,923	103	59	5	92.6%	6.0s	29.4%	0.6s	34.2%	0.3s	31.2%
18,227	323	59	5	94.8%	12.2s	21.8%	5.2s	29.1%	4.1s	24.1%
96,601	788	59	5	94.2%	25.5s	20.9%	38.1s	28.7%	30.4s	21.3%
471,268	1,760	59	5	93.6%	102s	18.8%	460s	24.6%	430s	19.8%
471,268	14,538	59	5	94.1%	170s	15.7%	3921s	21.5%	5300s	16.7%

**Table 2** Comparison of methods on Netflix data for  $k = 5$ .

$n$	$m$	$p$	$k$	$\mu\%$	OptComplete		IMC		SIALS	
					$T$	MAPE	$T$	MAPE	$T$	MAPE
3,923	103	59	10	92.6%	11.0s	30.4%	1.4s	36.7%	0.8s	35.8%
18,227	323	59	10	94.8%	20.3s	24.0%	12.5s	32.5%	7.0s	28.9%
96,601	788	59	10	94.2%	45.9s	22.3%	84.2s	29.6%	50.7s	22.8%
471,268	1,760	59	10	93.6%	260s	20.7%	1022s	24.8%	870s	20.7%
471,268	14,538	59	10	94.1%	380s	19.6%	8704s	23.1%	10240s	20.0%

**Table 3** Comparison of methods on Netflix data for  $k = 10$ .

We can see that OptComplete outperforms both IMC and SIALS in accuracy across the datasets under different  $k$ ; furthermore in the two largest datasets OptComplete ran 10x to 20x faster than IMC and SIALS. Here we see that an increase from  $k = 5$  to  $k = 10$  actually decreased out-of-sample performance as additional factors are actually not very helpful in predictive customer tastes. The decline for OptComplete and IMC were especially higher due to the fact that the possible factors are fixed and thus an increase in the number of factors caused some non-predictive factors to be included.

For the  $k = 5$  case, OptComplete identified the following as the top factors that influences an individual's rating:

- IMDB Rating
- Genre: Drama
- Released within last 10 years
- Number of Top 100 Actors
- Produced in US

These factors provide an intuitive explanation of the individual ratings of each customer in terms of a small number of factors, while exceeding the high predictive accuracy of SIALS.

## 6. Conclusions

We have presented OptComplete, a scalable algorithm to retrieve a low-rank matrix in the presence of side information. Compared with state of the art algorithms for matrix completion, OptComplete exceeds current benchmarks on both scalability and accuracy and provides insight on the factors that affect the ratings.

# Appendices

## A. Proof of Theorem 2

We first note that since  $S_k^p$  is a finite set, we only need to prove the result for a particular  $\mathbf{s} \in S_k^p$ , and it would apply for all  $\mathbf{s}$ . Therefore, we would assume  $\mathbf{s}$  is fixed below. For simplicity, we would only demonstrate the proof for  $\alpha_i(\mathbf{s}) = \bar{\mathbf{a}}_i \gamma_i(\mathbf{s})$ , as the one for  $(\mathbf{B}\mathbf{W}_i \gamma_i(\mathbf{s}))^2$  follows in the same exact fashion. Furthermore, since we are only focusing on one particular  $i \in \{1, \dots, n\}$ , we would drop all  $i$  subscripts below for ease of notation. The quantities of interest are therefore:

$$\begin{aligned}\alpha(\mathbf{s}) &= \frac{1}{m} \left( \bar{\mathbf{a}}\bar{\mathbf{a}}^T - \bar{\mathbf{a}}\mathbf{V} \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}^T \mathbf{W} \mathbf{V} \right)^{-1} \mathbf{V}^T \bar{\mathbf{a}}^T \right) \\ \alpha^F(\mathbf{s}) &= \frac{1}{f} \left( \bar{\mathbf{a}}_F \bar{\mathbf{a}}_F^T - \bar{\mathbf{a}}_F \mathbf{V}_F \left( \frac{\mathbf{I}_k}{\gamma} + \mathbf{V}_F^T \mathbf{W} \mathbf{V}_F \right)^{-1} \mathbf{V}_F^T \bar{\mathbf{a}}_F^T \right).\end{aligned}$$

First, let  $\bar{\mathbf{V}} = \mathbf{W}\mathbf{V}$ , and let us consider a reduced QR factorization of  $\bar{\mathbf{V}} = \mathbf{Q}\mathbf{R}$  where  $\mathbf{Q} \in \mathbb{R}^{m \times k}$  has orthogonal columns such that  $\mathbf{Q}^T \mathbf{Q} = m \cdot \mathbf{I}_k$ , and  $\mathbf{R}_i \in \mathbb{R}^{k \times k}$ . Note such definition implies  $\|\mathbf{R}\| = O(1)$ . Then, we would rewrite the terms as follows:

$$\begin{aligned}\alpha(\mathbf{s}) &= \frac{\bar{\mathbf{a}}\bar{\mathbf{a}}^T}{m} - \frac{\bar{\mathbf{a}}\mathbf{V}}{m} \left( \frac{\mathbf{I}_k}{m\gamma} + \frac{\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R}}{m} \right)^{-1} \frac{\mathbf{V}^T \bar{\mathbf{a}}^T}{m}, \\ \alpha^F(\mathbf{s}) &= \frac{\bar{\mathbf{a}}_F \bar{\mathbf{a}}_F^T}{f} - \frac{\bar{\mathbf{a}}_F \mathbf{V}_F}{f} \left( \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T \mathbf{Q}_F^T \mathbf{Q}_F \mathbf{R}}{f} \right)^{-1} \frac{\mathbf{V}_F^T \bar{\mathbf{a}}_F^T}{f}.\end{aligned}$$

We note that

$$\begin{aligned}\bar{\mathbf{a}}\bar{\mathbf{a}}^T &= \sum_{i=1}^m \bar{a}_i^2, & \bar{\mathbf{a}}\mathbf{V} &= \sum_{i=1}^m \bar{a}_i \mathbf{v}_i, \\ \bar{\mathbf{a}}_F \bar{\mathbf{a}}_F^T &= \sum_{i \in F} \bar{a}_i^2, & \bar{\mathbf{a}}_F \mathbf{V}_F &= \sum_{i \in F} \bar{a}_i \mathbf{v}_i.\end{aligned}$$

Therefore, if we treat  $\bar{a}_1^2, \dots, \bar{a}_m^2$  as a finite population, then  $\bar{\mathbf{a}}_F \bar{\mathbf{a}}_F^T$  is a random sample of  $f$  points drawn without replacement from that set, and similarly for  $\bar{\mathbf{a}}_F \mathbf{V}_F$ . Therefore, we can then utilize Hoeffding's inequality to bound the deviation of these terms, as reproduced below:

**Proposition 2 (Hoeffding's Inequality)** *Let  $\mathcal{X} = (x_1, \dots, x_n)$  be a finite population of  $N$  points and  $X_1, \dots, X_n$  be a random sample drawn without replacement from  $\mathbf{X}$ . Let*

$$a = \min_{1 \leq i \leq n} x_i \quad \text{and} \quad b = \max_{1 \leq i \leq n} x_i.$$



Then, for all  $\epsilon > 0$ , we have

$$\mathbb{P} \left( \left| \frac{\sum_{i=1}^n X_i}{n} - \mu \right| \geq \epsilon \right) \leq 2 \exp \left( -\frac{2n\epsilon^2}{(b-a)^2} \right). \quad (\text{A1})$$

For a proof, see for example Boucheron et al. (2013). Then, applying Proposition 2 to  $\bar{\mathbf{a}}_f \bar{\mathbf{a}}_f^T$ ,  $\bar{\mathbf{a}}_f \mathbf{V}_f$ , and inverting the inequality, we have

$$\mathbb{P} \left( \left| \frac{\bar{\mathbf{a}}_F \bar{\mathbf{a}}_F^T}{f} - \frac{\bar{\mathbf{a}} \bar{\mathbf{a}}^T}{m} \right| \leq \sqrt{\frac{A \log(\frac{1}{\epsilon})}{f}} \right) \geq 1 - \epsilon, \quad (\text{A2})$$

$$\mathbb{P} \left( \left\| \frac{\bar{\mathbf{a}}_F \mathbf{V}_F}{f} - \frac{\bar{\mathbf{a}} \mathbf{V}}{m} \right\| \leq \sqrt{\frac{Bk \log(\frac{k}{\epsilon})}{f}} \right) \geq 1 - \epsilon, \quad (\text{A3})$$

where  $A, B$  are constants independent of  $k, f, m, \epsilon$ .

Now we would show that  $\frac{\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R}}{m}$  is close to  $\frac{\mathbf{R}^T \mathbf{Q}_F^T \mathbf{Q}_F \mathbf{R}}{f}$ :

**Lemma 1**

$$\mathbb{P} \left( \left\| \frac{\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R}}{m} - \frac{\mathbf{R}^T \mathbf{Q}_F^T \mathbf{Q}_F \mathbf{R}}{f} \right\| \leq \sqrt{\frac{Ck \log(\frac{k}{\epsilon})}{f}} \right) \geq 1 - \epsilon. \quad (\text{A4})$$

To prove this, we would first introduce a matrix analog of the well-known Chernoff bound, the proof of which can be found in Tropp (2012):

**Lemma 2** *Let  $\mathcal{X} \in \mathbb{R}^{k \times k}$  be a finite set of positive-semidefinite matrices, and suppose that*

$$\max_{\mathbf{X} \in \mathcal{X}} \lambda_{\max}(\mathbf{X}) \leq D,$$

where  $\lambda_{\min}/\lambda_{\max}$  is the minimum/maximum eigenvalue function. Sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_\ell\}$  uniformly at random without replacement. Compute:

$$\mu_{\min} := \ell \cdot \lambda_{\min}(\mathbb{E} \mathbf{X}_1) \quad \mu_{\max} := \ell \cdot \lambda_{\max}(\mathbb{E} \mathbf{X}_1).$$

Then,

$$\mathbb{P} \left\{ \lambda_{\min} \left( \sum_j \mathbf{X}_j \right) \leq (1 - \delta) \mu_{\min} \right\} \leq k \cdot \exp \left( \frac{-\delta^2 \mu_{\min}}{4D} \right), \quad \text{for } \delta \in [0, 1),$$

$$\mathbb{P} \left\{ \lambda_{\max} \left( \sum_j \mathbf{X}_j \right) \leq (1 + \delta) \mu_{\max} \right\} \leq k \cdot \exp \left( \frac{-\delta^2 \mu_{\max}}{4D} \right), \quad \text{for } \delta \geq 0.$$

Using this lemma, we would proceed with the proof of Lemma 1.

*Proof of Lemma 1:* First, we note that

$$\begin{aligned}\mathbf{Q}^T \mathbf{Q} &= \sum_{i=1}^m \mathbf{q}_i^T \mathbf{q}_i, \\ \mathbf{Q}_F^T \mathbf{Q}_F &= \sum_{i \in F} \mathbf{q}_i^T \mathbf{q}_i,\end{aligned}$$

where  $\mathbf{q}_i^T \mathbf{q}_i \in \mathbb{R}^{k \times k}$  rank-one positive semi-definite matrices. Therefore, we can take  $\mathbf{Q}_F^T \mathbf{Q}_F$  as a random sample of size  $f$  from the set  $\mathcal{X} = \{\mathbf{q}_i^T \mathbf{q}_i\}_{i=1, \dots, m}$ , which satisfies the conditions in Lemma 2 with  $D = O(k)$ . Furthermore, with  $\mathcal{X}$ , we observe that we have  $\mathbb{E} \mathbf{X}_1 = \frac{\mathbf{Q}^T \mathbf{Q}}{m} = \mathbf{I}_k$ , so we have

$$\lambda_{\min}(\mathbb{E} \mathbf{X}_1) = \lambda_{\max}(\mathbb{E} \mathbf{X}_1) = 1.$$

Therefore, we apply Lemma 2 to  $\mathbf{Q}_F^T \mathbf{Q}_F$  and obtain

$$\begin{aligned}\mathbb{P} \left\{ \lambda_{\min}(\mathbf{Q}_F^T \mathbf{Q}_F) \leq (1 - \delta)f \right\} &\leq k \cdot \exp\left(\frac{-\delta^2 f}{kD'}\right), \\ \mathbb{P} \left\{ \lambda_{\max}(\mathbf{Q}_F^T \mathbf{Q}_F) \geq (1 + \delta)f \right\} &\leq k \cdot \exp\left(\frac{-\delta^2 f}{kD'}\right),\end{aligned}$$

where we set  $D = \frac{kD'}{4}$  with  $D' = O(1)$ . Some rearrangement gives:

$$\mathbb{P} \left\{ \lambda_{\min} \left( \frac{\mathbf{Q}_F^T \mathbf{Q}_F}{f} \right) \geq 1 - \sqrt{\frac{kD' \log\left(\frac{2k}{\epsilon}\right)}{f}} \text{ and } \lambda_{\max} \left( \frac{\mathbf{Q}_F^T \mathbf{Q}_F}{f} \right) \leq 1 + \sqrt{\frac{kD' \log\left(\frac{2k}{\epsilon}\right)}{f}} \right\} \geq 1 - \epsilon. \quad (\text{A5})$$

Now since  $\frac{\mathbf{Q}^T \mathbf{Q}}{m} = \mathbf{I}_k$ , we have

$$\lambda_{\min} \left( \frac{\mathbf{Q}^T \mathbf{Q}}{m} \right) = \lambda_{\max} \left( \frac{\mathbf{Q}^T \mathbf{Q}}{m} \right) = 1 \quad (\text{A6})$$

Combining equation (A6) and (A5) gives

$$\mathbb{P} \left\{ \left\| \frac{\mathbf{Q}_F^T \mathbf{Q}_F}{f} - \frac{\mathbf{Q}^T \mathbf{Q}}{m} \right\| \leq \sqrt{\frac{kD' \log\left(\frac{2k}{\epsilon}\right)}{f}} \right\} \geq 1 - \epsilon. \quad (\text{A7})$$

Then, we have

$$\mathbb{P} \left\{ \left\| \frac{\mathbf{R}^T \mathbf{Q}_F^T \mathbf{Q}_F \mathbf{R}}{f} - \frac{\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R}}{m} \right\| \leq \|\mathbf{R}\|^2 \sqrt{\frac{kD' \log\left(\frac{2k}{\epsilon}\right)}{f}} \right\} \geq 1 - \epsilon. \quad (\text{A8})$$

Taking  $C = D' \|\mathbf{R}\|^4 \log(2)$  gives the required result. (Note  $\|\mathbf{R}\| = O(1)$  as we setup the QR decomposition to have  $\mathbf{Q}^T \mathbf{Q} = O(m)$ ).  $\square$

With equations (A2), (A3), and (A4), we are now ready to bound  $\alpha(\mathbf{s})$  and  $\alpha^F(\mathbf{s})$ . We first introduce another lemma from matrix perturbation theory (for proof, see e.g. Stewart (1990)).

**Lemma 3** *Let  $\mathbf{A}, \mathbf{B}$  be invertible matrices and let  $\mathbf{B} = \mathbf{A} + \Delta$ . Then, we have*

$$\|\mathbf{A}^{-1} - \mathbf{B}^{-1}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{B}^{-1}\| \|\Delta\|. \quad (\text{A9})$$

Then, we have

$$\begin{aligned} \|\alpha(\mathbf{s}) - \alpha^F(\mathbf{s})\| &\leq \left| \frac{\bar{\mathbf{a}}\bar{\mathbf{a}}^T}{m} - \frac{\bar{\mathbf{a}}_F\bar{\mathbf{a}}_F^T}{f} \right| \\ &\quad + \left| \frac{\bar{\mathbf{a}}\mathbf{V}}{m} \left( \frac{\mathbf{I}_k}{m\gamma} + \frac{\mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R}}{m} \right)^{-1} \frac{\mathbf{V}^T\bar{\mathbf{a}}^T}{m} - \frac{\bar{\mathbf{a}}_F\mathbf{V}_F}{f} \left( \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T\mathbf{Q}_F^T\mathbf{Q}_F\mathbf{R}}{f} \right)^{-1} \frac{\mathbf{V}_F^T\bar{\mathbf{a}}_F^T}{f} \right|. \end{aligned}$$

Using (A2) and triangle inequality, we have

$$\begin{aligned} &\leq \sqrt{\frac{A \log(\frac{1}{\epsilon})}{f}} + \left| \left( \frac{\bar{\mathbf{a}}\mathbf{V}}{m} - \frac{\bar{\mathbf{a}}_F\mathbf{V}_F}{f} \right) \left( \frac{\mathbf{I}_k}{m\gamma} + \frac{\mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R}}{m} \right)^{-1} \frac{\mathbf{V}^T\bar{\mathbf{a}}^T}{m} \right| \\ &\quad + \left| \frac{\bar{\mathbf{a}}_F\mathbf{V}_F}{f} \left( \left( \frac{\mathbf{I}_k}{m\gamma} + \frac{\mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R}}{m} \right)^{-1} - \left( \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T\mathbf{Q}_F^T\mathbf{Q}_F\mathbf{R}}{f} \right)^{-1} \right) \frac{\mathbf{V}^T\bar{\mathbf{a}}^T}{m} \right| \\ &\quad + \left| \frac{\bar{\mathbf{a}}_F\mathbf{V}_F}{f} \left( \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T\mathbf{Q}_F^T\mathbf{Q}_F\mathbf{R}}{f} \right)^{-1} \left( \frac{\mathbf{V}^T\bar{\mathbf{a}}^T}{m} - \frac{\mathbf{V}_F^T\bar{\mathbf{a}}_F^T}{f} \right) \right|. \end{aligned}$$

Using (A3) and Lemma 3, we have

$$\begin{aligned} &\leq \sqrt{\frac{A \log(\frac{1}{\epsilon})}{f}} + \sqrt{\frac{B'k \log(\frac{k}{\epsilon})}{f}} \\ &\quad + D' \left| \frac{\bar{\mathbf{a}}_F\mathbf{V}_F}{f} \left( \frac{\mathbf{I}_k}{m\gamma} - \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R}}{m} - \frac{\mathbf{R}^T\mathbf{Q}_F^T\mathbf{Q}_F\mathbf{R}}{f} \right) \frac{\mathbf{V}^T\bar{\mathbf{a}}^T}{m} \right| \\ &\quad + \sqrt{\frac{B''k \log(\frac{k}{\epsilon})}{f}}, \end{aligned}$$

for some  $O(1)$  constant  $D'$ . Note that  $\|\frac{\mathbf{I}_k}{m\gamma} - \frac{\mathbf{I}_k}{f\gamma}\| \leq \frac{M}{f}$ , so using (A4), we obtain

$$\mathbb{P} \left( \left\| \frac{\mathbf{I}_k}{m\gamma} - \frac{\mathbf{I}_k}{f\gamma} + \frac{\mathbf{R}^T\mathbf{Q}^T\mathbf{Q}\mathbf{R}}{m} - \frac{\mathbf{R}^T\mathbf{Q}_F^T\mathbf{Q}_F\mathbf{R}}{f} \right\| \leq \sqrt{\frac{C'k \log(\frac{k}{\epsilon})}{f}} \right) \geq 1 - \epsilon, \quad (\text{A10})$$

for some new  $O(1)$  constant  $C'$ . Then, using (A10), we have

$$\leq \sqrt{\frac{A \log(\frac{1}{\epsilon})}{f}} + \sqrt{\frac{B'k \log(\frac{k}{\epsilon})}{f}} + \sqrt{\frac{C''k \log(\frac{k}{\epsilon})}{f}} + \sqrt{\frac{B''k \log(\frac{k}{\epsilon})}{f}}.$$

Since  $k \geq 1$ , we have

$$\leq \sqrt{\frac{Mk \log\left(\frac{k}{\epsilon}\right)}{f}},$$

For some  $O(1)$  constant  $M$  with probability at least  $1 - 4\epsilon$ . Therefore, taking  $\epsilon' = \frac{\epsilon}{4}$  and absorbing the extra constant into  $M$  gives the required result.

## B. Proof of Theorem 3

We first prove a Hoeffding-type bound as follows

**Proposition 3** *Let  $X_1, \dots, X_n$  be independent (but not necessarily identically distributed) random variables which satisfy*

$$\mathbb{P}(|X_i - a_i| \geq t) \leq \exp\left(\frac{-t^2}{\sigma_i^2}\right).$$

*Then, we have*

$$\mathbb{P}\left(\left|\frac{\sum_{i=1}^n X_i - a_i}{n}\right| \geq t\right) \leq 4 \exp\left(-\frac{n^2 t^2}{4 \sum_{i=1}^n \sigma_i^2}\right). \quad (\text{A11})$$

We note that such statement differs from Hoeffding's inequality as we do not require  $\mathbb{E}[X_i] = a_i$ .

The proof is as follows.

*Proof of Proposition 3:* We first introduce a lemma known as Chernoff's bounding method (proven in Chernoff (1952)):

**Lemma 4** *Let  $Z$  be a random variable on  $\mathbb{R}$ . Then for all  $t > 0$ , we have*

$$\mathbb{P}(Z \geq t) \leq \inf_{s>0} [e^{-st} M_Z(s)],$$

where  $M_Z(s)$  is the moment generating function of  $Z$ .

Let  $Y_i = X_i - a_i$ . Then, we have

$$\mathbb{P}\left(\sum_{i=1}^n Y_i \geq t\right) \leq \min_{s>0} e^{-st} \prod_{i=1}^n \mathbb{E}[e^{sY_i}]. \quad (\text{A12})$$

We aim to bound  $\mathbb{E}[e^{sY_i}]$ , subject to  $\mathbb{P}(|Y_i| \geq t) \leq \exp\left(\frac{-t^2}{\sigma_i^2}\right)$ . Since  $e^{st}$  is an increasing function of  $t$ ,  $\mathbb{E}[e^{sY_i}]$  is maximized when  $\mathbb{P}(Y_i \geq t) = \exp\left(\frac{-t^2}{\sigma_i^2}\right)$ , which results in a probability distribution function of  $f_{Y_i}(y) = 2\frac{t}{\sigma_i^2} \exp\left(\frac{-t^2}{\sigma_i^2}\right)$  for  $Y_i$ . Then, we have

$$\begin{aligned} \mathbb{E}[e^{sY_i}] &\leq \int_0^\infty 2\frac{t}{\sigma_i^2} \exp\left(st - \frac{t^2}{\sigma_i^2}\right) dt \\ &= \int_0^\infty 2\frac{t}{\sigma_i^2} \exp\left(-\left(\frac{t}{\sigma_i} - \frac{s\sigma_i}{2}\right)^2 + \frac{s^2\sigma_i^2}{4}\right) dt \\ &= 2 \exp\left(\frac{s^2\sigma_i^2}{4}\right) \int_0^\infty t \exp\left(-\left(t - \frac{s\sigma_i}{2}\right)^2\right) dt \\ &= \exp\left(\frac{s^2\sigma_i^2}{4}\right) \left(\exp\left(-\frac{s^2\sigma_i^2}{4}\right) + \sqrt{\pi} \frac{s\sigma_i}{2} \left(\operatorname{erf}\left(\frac{s\sigma_i}{2}\right) + 1\right)\right) \end{aligned}$$

where erf is the standard error function. As the error function is upper bounded by 1, the last expression is less than or equal to:

$$\begin{aligned} &\leq 1 + \sqrt{\pi} s \sigma_i \exp\left(\frac{s^2 \sigma_i^2}{4}\right) \\ &\leq 2 \exp\left(s^2 \sigma_i^2\right). \end{aligned}$$

We then substitute this result into (A12) and obtain

$$\mathbb{P}\left(\sum_{i=1}^n Y_i \geq t\right) \leq \min_{s>0} 2 \exp\left(-st + s^2 \sum_{i=1}^n \sigma_i^2\right)$$

Note that this is minimized at  $s = \frac{t}{2\sqrt{\sum_{i=1}^n \sigma_i^2}}$ , so we have

$$\leq 2 \exp\left(-\frac{t^2}{4 \sum_{i=1}^n \sigma_i^2}\right).$$

Therefore,

$$\mathbb{P}\left(\frac{\sum_{i=1}^n Y_i}{n} \geq t\right) \leq 2 \exp\left(-\frac{n^2 t^2}{4 \sum_{i=1}^n \sigma_i^2}\right).$$

By applying the previous derivation to  $-Y_1, \dots, -Y_n$ , we obtain

$$\mathbb{P}\left(\frac{\sum_{i=1}^n Y_i}{n} \leq -t\right) \leq 2 \exp\left(-\frac{n^2 t^2}{4 \sum_{i=1}^n \sigma_i^2}\right).$$

Combining the two results completes the proof.  $\square$

Using Proposition 3, we bound the difference between  $\tilde{c}_G^F(\mathbf{s}_t)$  and  $c(\mathbf{s}_t)$ . We have

$$\begin{aligned} |\tilde{c}_G^F(\mathbf{s}_t) - c(\mathbf{s}_t)| &= \left| \frac{1}{n} \sum_{i=1}^n \alpha_i(\mathbf{s}) - \frac{1}{g} \sum_{i \in G} \alpha_i^F(\mathbf{s}) \right| \\ &\leq \left| \frac{1}{n} \sum_{i=1}^n \alpha_i(\mathbf{s}) - \frac{1}{g} \sum_{i \in G} \alpha_i(\mathbf{s}) \right| + \left| \frac{1}{g} \sum_{i \in G} \alpha_i(\mathbf{s}) - \frac{1}{g} \sum_{i \in G} \alpha_i^F(\mathbf{s}) \right|. \end{aligned}$$

The first term can be seen as the tail bound for a random sample of size  $g$  chosen without replacement from the finite set  $\{\alpha_i(\mathbf{s})\}_{i=1, \dots, n}$ . Thus, we can apply Hoeffding's theorem in Proposition 2, and obtain that with probability at least  $1 - \epsilon$

$$\left| \frac{1}{n} \sum_{i=1}^n \alpha_i(\mathbf{s}) - \frac{1}{g} \sum_{i \in G} \alpha_i(\mathbf{s}) \right| \leq \sqrt{\frac{M \log\left(\frac{1}{\epsilon}\right)}{g}}. \quad (\text{A13})$$

Substituting (A13) into the expression above shows that with probability at least  $1 - \epsilon$  we have:

$$\leq \sqrt{\frac{M \log\left(\frac{1}{\epsilon}\right)}{g}} + \left| \frac{1}{g} \sum_{i \in G} \alpha_i(\mathbf{s}) - \frac{1}{g} \sum_{i \in G} \alpha_i^F(\mathbf{s}) \right|.$$

Note that, for any fixed set  $G$ , for  $i, j \in G$  with  $i \neq j$ , we have that  $\alpha_i^F(\mathbf{s})$  and  $\alpha_j^F(\mathbf{s})$  are independent (as we construct  $F$  separately for every  $i$ ). Furthermore, Theorem 2 can be inverted to read

$$\mathbb{P}(|\alpha_i(\mathbf{s}) - \alpha_i^F(\mathbf{s})| \geq t) \leq \exp\left(-\frac{ft^2}{Mk \log(k)}\right).$$

Therefore,  $\alpha_i^F(\mathbf{s})$  satisfies the conditions of Proposition 3 with  $X_i = \alpha_i^F(\mathbf{s})$  and parameters  $\sigma_i^2 = \frac{Mk \log(k)}{f}$ ,  $a_i = \alpha_i(\mathbf{s})$ . Then, applying Proposition 3 to the second term, we have

$$\leq \sqrt{\frac{M \log\left(\frac{1}{\epsilon}\right)}{g}} + \sqrt{\frac{M'k \log\left(\frac{k}{\epsilon}\right)}{fg}},$$

with probability  $1 - 2\epsilon$ . As  $k, f \geq 1$ , a loose bound is therefore

$$\leq \sqrt{\frac{M''k \log\left(\frac{k}{\epsilon}\right)}{g}},$$

with probability  $1 - 2\epsilon$ . Therefore, taking  $\epsilon' = 2\epsilon$  and  $A = M'' \log(2)$  we have

$$\mathbb{P}\left(|\tilde{c}_G^F(\mathbf{s}_t) - c(\mathbf{s}_t)| \leq \sqrt{\frac{Ak \log\left(\frac{k}{\epsilon'}\right)}{g}}\right) \geq 1 - \epsilon'. \quad (\text{A14})$$

Through a similar derivation, we have

$$\mathbb{P}\left(\|\nabla \tilde{c}_G^F(\mathbf{s}_t) - \nabla c(\mathbf{s}_t)\|_2 \leq \sqrt{\frac{B(p+k) \log\left(\frac{k}{\epsilon'}\right)}{g}}\right) \geq 1 - \epsilon'. \quad (\text{A15})$$

## C. Proof of Proposition 1

We would prove the contrapositive. Assume that  $a = 0$ . We would calculate the second derivative of the cost function expression in 7. First, define  $\mathbf{T}_{ij} = \frac{\mathbf{I}}{\gamma} + \mathbf{W}_i \mathbf{K}_j \mathbf{W}_i$ . Then our cost function can be rewritten as:

$$c(\mathbf{s}) = \frac{1}{nm} \sum_{i=1}^n \bar{\mathbf{a}}_i \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T$$

Then the second derivative can be easily calculated as:

$$\frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} = 2 \sum_{i=1}^n \bar{\mathbf{a}}_i \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \mathbf{T}_{ik} \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \mathbf{T}_{il} \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T$$

Thus, for any vector  $\mathbf{t} \in \mathbb{R}^p$ , we have:

$$\mathbf{t}^T \frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} \mathbf{t} = 2 \sum_{i=1}^n \bar{\mathbf{a}}_i \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \left( \sum_{k=1}^p t_k \mathbf{T}_{ik} \right) \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \left( \sum_{l=1}^p t_l \mathbf{T}_{il} \right) \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T$$

Now define  $\mathbf{v}_i = \left( \sum_{l=1}^p t_l \mathbf{T}_{il} \right) \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T$ , then our expression becomes:

$$\mathbf{t}^T \frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} \mathbf{t} = 2 \sum_{i=1}^n \mathbf{v}_i^T \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \mathbf{v}_i$$

Now since  $a = 0$ , the Hessian cannot be positive definite for all  $\mathbf{s}$  (as if it was, then  $a > 0$ ). Then there exist  $\mathbf{s}^0, \mathbf{s}^1 \in S_k^p$  such that for  $\mathbf{t} = \mathbf{s}^0 - \mathbf{s}^1$ , we have:

$$\mathbf{t}^T \frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} \Big|_{\mathbf{s}=\mathbf{s}^0} \mathbf{t} = 2 \sum_{i=1}^n \mathbf{v}_i^T \left( \sum_{j=1}^p s_j^0 \mathbf{T}_{ij} \right)^{-1} \mathbf{v}_i = 0$$

and  $\mathbf{t}$  Note that for any  $\gamma > 0$ ,  $\mathbf{T}_{ij}$  is positive definite, so  $\sum_{j=1}^p s_j \mathbf{T}_{ij}$  is positive definite for any  $\mathbf{s} \in S_k^p$  and all  $i$ . Therefore if  $a = 0$ , we must have  $\mathbf{v}_i = 0$  for all  $i$ , which means that:

$$\left( \sum_{l=1}^p t_l \mathbf{T}_{il} \right) \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T = \mathbf{0}$$

for all  $i$ . Now note that we have:

$$c(\mathbf{s}^1) = c(\mathbf{s}^0) + \nabla c(\mathbf{s}^0)^T \mathbf{t} + \mathbf{t}^T \frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} \Big|_{\mathbf{s}=\mathbf{s}^0} \mathbf{t} + \dots$$

We have that:

$$\nabla c(\mathbf{s}^0)^T \mathbf{t} = \sum_{i=1}^n \left( \sum_{l=1}^p t_l \mathbf{T}_{il} \right) \left( \sum_{j=1}^p s_j \mathbf{T}_{ij} \right)^{-1} \bar{\mathbf{a}}_i^T = 0$$



And similarly, we have  $\mathbf{t}^T \frac{\partial c(\mathbf{s})}{\partial s_k \partial s_l} \Big|_{\mathbf{s}=\mathbf{s}^0} \mathbf{t} = 0$  and all higher derivatives being 0.

Thus, we have:

$$c(\mathbf{s}^1) = c(\mathbf{s}^0)$$

Therefore, if we have  $a = 0$ , then we must have equality on the cost function for two feasible solutions. Thus, if the cost function is not degenerate for feasible solutions, then  $a > 0$ .

## D. Proof of Theorem 4

In this section, we provide the proof of Theorem 4. We first note that OptComplete is a specific implementation of the outer approximation algorithm, which Fletcher and Leyffer (1994) have shown to always terminate in finite number of steps  $C$ . Thus, given that we assumed the problem is feasible, for OptComplete to not return an optimal solution, it would have to cut it off during the course of its execution. Let  $\mathbf{s}^*$  be an optimal solution for Problem (4). Let  $\mathbf{s}_t$  be an optimal solution at the  $t$ -th iteration of OptComplete,  $t \in [C]$ . The cutting plane constraint for OptComplete at the point of an optimal solution  $\mathbf{s}^*$  is

$$\eta_t \geq \tilde{c}_G^F(\mathbf{s}_t) + \nabla \tilde{c}_G^F(\mathbf{s}_t)^T (\mathbf{s}^* - \mathbf{s}_t).$$

If  $c(\mathbf{s}^*) < \eta_t$ , then  $\mathbf{s}^*$  will be cut off, and OptComplete will not find  $\mathbf{s}^*$ . Applying the definition of the convexity parameter (16) and letting  $\|\mathbf{s}^* - \mathbf{s}_t\| = \theta_t$  (noting that  $\theta_t \geq 1$ ) we obtain

$$c(\mathbf{s}^*) \geq c(\mathbf{s}_t) + \nabla c(\mathbf{s}_t)^T (\mathbf{s}^* - \mathbf{s}_t) + \frac{\theta_t^2 a^2}{2}. \quad (\text{A16})$$

Therefore, if

$$c(\mathbf{s}_t) + \nabla c(\mathbf{s}_t)^T (\mathbf{s}^* - \mathbf{s}_t) + \frac{\theta_t^2 a^2}{2} \leq c(\mathbf{s}^*) < \tilde{c}_r(\mathbf{s}_t) + \nabla \tilde{c}_r(\mathbf{s}_t)^T (\mathbf{s}^* - \mathbf{s}_t),$$

or equivalently if

$$\zeta_t := [\tilde{c}_G^F(\mathbf{s}_t) - c(\mathbf{s}_t)] + [\nabla \tilde{c}_G^F(\mathbf{s}_t) - \nabla c(\mathbf{s}_t)]^T (\mathbf{s}^* - \mathbf{s}_t) > \frac{\theta_t^2 a^2}{2}, \quad (\text{A17})$$

then OptComplete will not find  $\mathbf{s}^*$ . Therefore, for OptComplete to succeed,  $\zeta_t$  should satisfy  $\zeta_t \leq \theta_t^2 a^2 / 2$  for all  $t \in [C]$ .

Let  $S =$  the event that OptComplete succeeds in finding  $\mathbf{s}^*$ . Then,

$$\mathbb{P}(S) \geq \mathbb{P} \left( \zeta_t \leq \frac{\theta_t^2 a^2}{2} \quad t \in [C] \right).$$

Since at each step of OptComplete we randomly sample  $r$  new rows and  $s$  new columns, the events  $\zeta_t \leq \frac{\theta_t^2 a^2}{2}$  are independent for different  $t \in [C]$ , and hence

$$\mathbb{P}(S) \geq \prod_{t=1}^C \left( 1 - \mathbb{P} \left( \zeta_t > \frac{\theta_t^2 a^2}{2} \right) \right). \quad (\text{A18})$$

Therefore, we focus on calculating  $\mathbb{P}(\zeta_t \geq \theta_t^2 a^2 / 2)$ . Since

$$\zeta_t := [\tilde{c}_G^F(\mathbf{s}_t) - c(\mathbf{s}_t)] + [\nabla \tilde{c}_G^F(\mathbf{s}_t) - \nabla c(\mathbf{s}_t)]^T (\mathbf{s}^* - \mathbf{s}_t),$$

using Theorem 3, we can thus provide the following bound for deviation of  $\zeta_t$ , for some constant  $C$ :

$$\mathbb{P} \left( \zeta_t \leq \theta_t \sqrt{\frac{C(p+k) \log \left( \frac{k}{\epsilon} \right)}{g}} \right) \geq 1 - \epsilon, \quad (\text{A19})$$

where  $\theta_t = \|\mathbf{s}^* - \mathbf{s}_t\|$ . Then we can invert this to calculate  $\mathbb{P}(\zeta_t \geq \theta_t^2 a^2 / 2)$

$$\mathbb{P}(\zeta_t \geq \theta_t^2 a^2 / 2) \leq k \exp \left( -\frac{a^4 g}{4C(p+k)\theta_t^2} \right) \leq k \exp \left( -\frac{Da^4 g}{(p+k)} \right), \quad (\text{A20})$$

taking  $D = \frac{1}{4C}$ , and noting that  $\theta_t \geq 1$ . Then, we substitute (A20) into (A18) to obtain

$$\begin{aligned} \mathbb{P}(S) &\geq \left( 1 - k \exp \left( -\frac{Da^4 g}{(p+k)} \right) \right)^C \\ &\geq 1 - kC \exp \left( -\frac{Da^4 g}{(p+k)} \right). \end{aligned}$$

completing the proof.

## E. List of Features Used in the Netflix Problem

- 24 Indicator Variables for Genres: Action, Adventure, Animation, Biography, Comedy, Crime, Documentary, Drama, Family, Fantasy, Film Noir, History, Horror, Music, Musical, Mystery, Romance, Sci-Fi, Short, Sport, Superhero, Thriller, War, Western
- 5 Indicator Variables for Release Date: Within last 10 years, Between 10-20 years, Between 20-30 years, Between 30-40 years, Between 40-50 Years
- 6 Indicator Variables for Top Actors/Actresses defined by their Influence Score at time of release: Top 100 Actors, Top 100 Actresses, Top 250 Actors, Top 250 Actresses, Top 1000 Actors, Top 1000 Actresses
  - IMDB Rating
  - Number of Reviews
  - Total Production Budget
  - Total Runtime
  - Total Box Office Revenue
  - Indicator Variable for whether it is US produced
- 11 Indicator Variables for Month of Year Released (January removed to prevent multicollinearity)
  - Number of Original Music Score
  - Number of Male Actors
  - Number of Female Factors
- 3 Indicator Variables for Film Language: English, French, Japanese
- Constant

## References

- Beck A, Teboulle M (2009) A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* 2(1):183–202.
- Bennett J, Lanning S, et al. (2007) The netflix prize. *KDD Cup and Workshop 2007* (Citeseer).
- Bertsimas D, Copenhaver MS (2018) Characterization of the equivalence of robustification and regularization in linear, median, and matrix regression. *European Journal of Operations Research* 270:931–942.
- Bertsimas D, van Parys B (2020) Sparse high-dimensional regression: Exact scalable algorithms and phase transitions. *Annals of Statistics* 48(1):300–323.
- Boucheron S, Lugosi G, Massart P (2013) *Concentration inequalities: A nonasymptotic theory of independence* (Oxford university press).
- Boyd S, El Ghaoui L, Feron E, Balakrishnan V (1994) *Linear matrix inequalities in system and control theory*, volume 15 (SIAM).
- Candes EJ, Plan Y (2010) Matrix completion with noise. *Proceedings of the IEEE* 98(6):925–936.
- Candès EJ, Tao T (2010) The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory* 56(5):2053–2080.
- Chen Y, Bhojanapalli S, Sanghavi S, Ward R (2014a) Coherent matrix completion. *International Conference on Machine Learning*, 674–682.
- Chen Y, Jalali A, Sanghavi S, Xu H (2014b) Clustering partially observed graphs via convex optimization. *The Journal of Machine Learning Research* 15(1):2213–2238.
- Chernoff H (1952) A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics* 23(4):493–507.
- Chiang KY, Hsieh CJ, Dhillon IS (2015) Matrix completion with noisy side information. *Advances in Neural Information Processing Systems*, 3447–3455.
- Chiang KY, Hsieh CJ, Natarajan N, Dhillon IS, Tewari A (2014) Prediction and clustering in signed networks: a local to global perspective. *The Journal of Machine Learning Research* 15(1):1177–1213.
- Dhillon P, Lu Y, Foster DP, Ungar L (2013) New subsampling algorithms for fast least squares regression. *Advances in Neural Information Processing Systems*, 360–368.

- Duran MA, Grossmann IE (1986) An outer-approximation algorithm for a class of mixed-integer nonlinear programs. *Mathematical Programming* 36(3):307–339.
- Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, Warren RM, Streicher EM, Calver A, Sloutsky A, et al. (2013) Genomic analysis identifies targets of convergent positive selection in drug-resistant mycobacterium tuberculosis. *Nature Genetics* 45(10):1183.
- Fletcher R, Leyffer S (1994) Solving mixed integer nonlinear programs by outer approximation. *Mathematical programming* 66(1-3):327–349.
- Hastie T, Mazumder R, Lee JD, Zadeh R (2015) Matrix completion and low-rank svd via fast alternating least squares. *The Journal of Machine Learning Research* 16(1):3367–3402.
- Jain P, Dhillon IS (2013) Provable inductive matrix completion. *arXiv preprint arXiv:1306.0626* .
- Jain P, Meka R, Dhillon IS (2010) Guaranteed rank minimization via singular value projection. *Advances in Neural Information Processing Systems*, 937–945.
- Ji H, Liu C, Shen Z, Xu Y (2010) Robust video denoising using low rank matrix completion. *Computer Society Conference on Computer Vision and Pattern Recognition (IEEE)*.
- Keshavan RH, Oh S, Montanari A (2009) Matrix completion from a few entries. *Information Theory, 2009. ISIT 2009. IEEE International Symposium on*, 324–328 (IEEE).
- Koren Y, Bell R, Volinsky C (2009) Matrix factorization techniques for recommender systems. *Computer* 30–37.
- Lu J, Liang G, Sun J, Bi J (2016) A sparse interactive model for matrix completion with side information. *Advances in Neural Information Processing Systems*, 4071–4079.
- Lubin M, Yamangil E, Bent R, Vielma JP (2016) Extended formulations in mixed-integer convex programming. *International Conference on Integer Programming and Combinatorial Optimization*, 102–113 (Springer).
- Mazumder R, Hastie T, Tibshirani R (2010) Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research* 11(Aug):2287–2322.
- Natarajan N, Dhillon IS (2014) Inductive matrix completion for predicting gene–disease associations. *Bioinformatics* 30(12):160–168.

- Nazarov I, Shirokikh B, Burkina M, Fedonin G, Panov M (2018) Sparse group inductive matrix completion. *arXiv preprint arXiv:1804.10653* .
- Negahban S, Wainwright MJ (2012) Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research* 13(May):1665–1697.
- Recht B, Ré C (2013) Parallel stochastic gradient algorithms for large-scale matrix completion. *Mathematical Programming Computation* 5(2):201–226.
- Shah V, Rao N, Ding W (2017) Matrix factorization with side and higher order information. *Stat* 1050:4.
- Si S, Chiang KY, Hsieh CJ, Rao N, Dhillon IS (2016) Goal-directed inductive matrix completion. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1165–1174 (ACM).
- Soni A, Chevalier T, Jain S (2016) Noisy inductive matrix completion under sparse factor models. *arXiv preprint arXiv:1609.03958* .
- Stewart GW (1990) Matrix perturbation theory .
- Tanner J, Wei K (2013) Normalized iterative hard thresholding for matrix completion. *SIAM Journal on Scientific Computing* 35(5):S104–S125.
- Tropp JA (2012) User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics* 12(4):389–434.
- Woodbury MA (1949) The stability of out-input matrices. *Chicago, Ill* .
- Xu M, Jin R, Zhou ZH (2013) Speedup matrix completion with side information: Application to multi-label learning. *Advances in Neural Information Processing Systems*, 2301–2309.