

Predictive Low Rank Matrix Learning under Partial Observations: Mixed-Projection ADMM

Dimitris Bertsimas

*Massachusetts Institute of Technology
Cambridge, MA 02139, USA*

DBERTSIM@MIT.EDU

Nicholas A. G. Johnson

*Massachusetts Institute of Technology
Cambridge, MA 02139, USA*

NAGJ@MIT.EDU

Editor: TBD

Abstract

We study the problem of learning a partially observed matrix under the low rank assumption in the presence of fully observed side information that depends linearly on the true underlying matrix. This problem consists of an important generalization of the Matrix Completion problem, a central problem in Statistics, Operations Research and Machine Learning, that arises in applications such as recommendation systems, signal processing, system identification and image denoising. We formalize this problem as an optimization problem with an objective that balances the strength of the fit of the reconstruction to the observed entries with the ability of the reconstruction to be predictive of the side information. We derive a mixed-projection reformulation of the resulting optimization problem and present a strong semidefinite cone relaxation. We design an efficient, scalable alternating direction method of multipliers algorithm that produces high quality feasible solutions to the problem of interest. Our numerical results demonstrate that in the small rank regime ($k \leq 15$), our algorithm outputs solutions that achieve on average 79% lower objective value and 90.1% lower ℓ_2 reconstruction error than the solutions returned by the experiment-wise best performing benchmark method. The runtime of our algorithm is competitive with and often superior to that of the benchmark methods. Our algorithm is able to solve problems with $n = 10000$ rows and $m = 10000$ columns in less than a minute.

Keywords: Matrix Completion; Rank; Mixed-Projection; ADMM;

1. Introduction

In many real world applications, we are faced with the problem of recovering a (often large) matrix from a (often small) subset of its entries. This problem, known as Matrix Completion (MC), has gained significant attention due to its broad range of applications in areas such as signal processing (Candes and Plan, 2010), system identification (Liu, 2019) and image denoising (Ji et al., 2010). The fundamental task in MC is to accurately reconstruct the missing entries of a matrix given a limited number of observed entries. The challenge is particularly pronounced when the number of observed entries is small relatively to the dimension of the matrix, yet this is the common scenario in practice.

One of the most prominent uses of MC is in recommendation systems, where the goal is to predict user preferences for items (e.g., movies, products) based on a partially observed user-

item rating matrix (Ramlatchan et al., 2018). The Netflix Prize competition highlighted the potential of MC techniques, where the objective was to predict missing ratings in a user-movie matrix to improve recommendation accuracy (Bell and Koren, 2007). The success of such systems hinges on the assumption that the underlying rating matrix is low rank, meaning that the preferences of users can be well-approximated by a small number of factors. Indeed, it has been well studied that many real world datasets are low rank (Udell and Townsend, 2019).

In many practical applications, in addition to a collection of observed matrix entries we additionally have access to auxiliary side information that can be leveraged when performing the reconstruction. For example, in a recommendation system, side information might consist of social network data or item attributes. The vast majority of existing approaches to MC in the presence of side information incorporate the side information by making additional structural restrictions on the reconstructed matrix beyond the usual low rank assumption (see, for example, Xu et al. (2013); Chiang et al. (2015); Bertsimas and Li (2020)). In this work, we take an alternate approach by assuming that the side information can be well modelled as a linear function of the underlying full matrix. In this setting, the side information can be thought of as labels for a regression problem where the unobserved matrix consists of the regression features. This assumption is in keeping with ideas from the predictive low rank kernel learning literature (Bach and Jordan, 2005) (note however that low rank kernel learning assumes a fully observed input matrix).

Formally, let $\Omega \subseteq [n] \times [m]$ denote a collection of revealed entries of a partially observed matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$, let $\mathbf{Y} \in \mathbb{R}^{n \times d}$ denote a matrix of side information and let k denote a specified target rank. We consider the problem given by

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times m}, \boldsymbol{\alpha} \in \mathbb{R}^{m \times d}} \quad & \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \|\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}\|_F^2 + \gamma \|\mathbf{X}\|_* \\ \text{s.t.} \quad & \text{rank}(\mathbf{X}) \leq k, \end{aligned} \tag{1}$$

where $\lambda, \gamma > 0$ are hyperparameters that in practice can either take a default value or can be cross-validated by minimizing a validation metric (Owen and Perry, 2009) to obtain strong out-of-sample performance (Bousquet and Elisseeff, 2002). We assume that the ground truth matrix \mathbf{A} has low rank and that the side information can be well approximated as $\mathbf{Y} = \mathbf{A}\boldsymbol{\alpha} + \mathbf{N}$ for some weighting matrix $\boldsymbol{\alpha}$ and noise matrix \mathbf{N} . The first term in the objective function of (1) measures how well the observed entries of the unknown matrix are fit by the estimated matrix \mathbf{X} , the second term of the objective function measures how well the side information \mathbf{Y} can be represented as a linear function of the estimated matrix \mathbf{X} and the final term of the objective is a regularization term. To the best of our knowledge, Problem (1) has not previously been directly studied despite its very natural motivation.

1.1 Contribution and Structure

In this paper, we tackle (1) by developing novel mixed-projection optimization techniques (Bertsimas et al., 2022). We show that solving (1) is equivalent to solving an appropriately defined robust optimization problem. We develop an exact reformulation of (1) by combining a parametrization of the \mathbf{X} decision variable as the product of two low rank factors with the introduction of a projection matrix to model the column space of \mathbf{X} . We derive a

semidefinite cone convex relaxation for our mixed-projection reformulation and we present an efficient, scalable alternating direction method of multipliers (ADMM) algorithm that produces high quality feasible solutions to (1). Our numerical results show that across all experiments in the small rank regime ($k \leq 15$), our algorithm outputs solutions that achieve on average 79% lower objective value in (1) and 90.1% lower ℓ_2 reconstruction error than the solutions returned by the experiment-wise best performing benchmark method. For the 5 experiments with $k > 15$, the only benchmark that returns a solution with superior quality than that returned by our algorithm takes on average 3 times as long to execute. The runtime of our algorithm is competitive with and often superior to that of the benchmark methods. Our algorithm is able to solve problems with $n = 10000$ rows and $m = 10000$ columns in less than a minute.

The rest of the paper is laid out as follows. In Section 2, we review previous work that is closely related to (1). In Section 3, we study (1) under a robust optimization lens and investigate formulating (1) as a two stage optimization problem where the inner stage is a regression problem that can be solved in closed form. We formulate (1) as a mixed-projection optimization problem in Section 4 and present a natural convex relaxation. In Section 5, we present and rigorously study our ADMM algorithm. Finally, in Section 6 we investigate the performance of our algorithm against benchmark methods on synthetic data.

Notation: We let nonbold face characters such as b denote scalars, lowercase bold faced characters such as \mathbf{x} denote vectors, uppercase bold faced characters such as \mathbf{X} denote matrices, and calligraphic uppercase characters such as \mathcal{Z} denote sets. We let $[n]$ denote the set of running indices $\{1, \dots, n\}$. We let $\mathbf{0}_n$ denote an n -dimensional vector of all 0's, $\mathbf{0}_{n \times m}$ denote an $n \times m$ -dimensional matrix of all 0's, and \mathbf{I}_n denote the $n \times n$ identity matrix. We let \mathcal{S}^n denote the cone of $n \times n$ symmetric matrices and \mathcal{S}_+^n denote the cone of $n \times n$ positive semidefinite matrices

2. Literature Review

In this section, we review a handful of notable approaches from the literature that have been employed to solve MC and to solve general low rank optimization problems. As an exhaustive literature review of MC methods is outside of the scope of this paper, we focus our review on a handful of well studied approaches which we will employ as benchmark methods in this work. We additionally give an overview of the ADMM algorithmic framework which is of central relevance to this work. For a more detailed review of the MC literature, we refer the reader to Ramlatchan et al. (2018) and Nguyen et al. (2019).

2.1 Matrix Completion Methods

2.1.1 ITERATIVE-SVD

Iterative-SVD is an expectation maximization style algorithm (Dempster et al., 1977) that generates a solution to the MC problem by iteratively computing a singular value decomposition (SVD) of the current iterate and estimating the missing values by performing a regression against the low rank factors returned by SVD (Troyanskaya et al., 2001). This is one of a handful of methods in the literature that leverage the SVD as their primary al-

gorithmic workhorse (Billsus et al., 1998; Sarwar et al., 2000). Concretely, given a partially observed matrix $\{X_{ij}\}_{(i,j)\in\Omega}$ where $\Omega \subseteq [n] \times [m]$ and a target rank $k \in \mathbf{N}_+$, Iterative-SVD proceeds as follows:

1. Initialize the iteration count $t \leftarrow 0$ and initialize missing entries of X_{ij} , $(i, j) \notin \Omega$ with the row average $X_{ij} = \frac{\sum_{l:(i,l)\in\Omega} X_{il}}{|\{(i,l)\in\Omega\}|}$.
2. Compute a rank k SVD $\mathbf{X}_t = \mathbf{U}_t \boldsymbol{\Sigma}_t \mathbf{V}_t^T$ of the current iterate where $\mathbf{U}_t \in \mathbb{R}^{n \times k}$, $\boldsymbol{\Sigma}_t \in \mathbb{R}^{k \times k}$, $\mathbf{V}_t \in \mathbb{R}^{m \times k}$.
3. For each $(i, j) \notin \Omega$, estimate the missing value $(\mathbf{X}_{t+1})_{ij}$ by regressing all other entries in row i against all except the j^{th} row of \mathbf{V}_t . Concretely, letting $\tilde{\mathbf{x}} = (\mathbf{X}_t)_{i, \star \setminus j} \in \mathbb{R}^{m-1}$ denote the column vector consisting of the i^{th} row of \mathbf{X}_t excluding the j^{th} entry, letting $\tilde{\mathbf{V}} = (\mathbf{V}_t)_{\star \setminus j, \star} \in \mathbb{R}^{(m-1) \times k}$ denote the matrix formed by eliminating the j^{th} row from \mathbf{V}_t and letting $\hat{\mathbf{v}} = (\mathbf{V}_t)_{j, \star} \in \mathbb{R}^k$ denote the column vector consisting of the j^{th} row of \mathbf{V}_t , we set $(\mathbf{X}_{t+1})_{ij} = \hat{\mathbf{v}}^T (\tilde{\mathbf{V}}^T \tilde{\mathbf{V}})^{-1} \tilde{\mathbf{V}}^T \tilde{\mathbf{x}}$.
4. Terminate if the total change between \mathbf{X}_t and \mathbf{X}_{t+1} is less than 0.01. Otherwise, increment t and return to Step 2.

2.1.2 SOFT-IMPUTE

Soft-Impute is a convex relaxation inspired algorithm that leverages the nuclear norm as a low rank inducing regularizer (Mazumder et al., 2010). This approach is one of a broad class of methods that tackle MC from a nuclear norm minimization lens (Fazel, 2002; Candès and Tao, 2010; Candès and Recht, 2012). Seeking a reconstruction with minimum nuclear norm is typically motivated by the observation that the nuclear norm ball given by $\mathcal{B} = \{\mathbf{X} \in \mathbb{R}^{n \times n} : \|\mathbf{X}\|_* \leq k\}$ is the convex hull of the nonconvex set $\mathcal{X} = \{\mathbf{X} \in \mathbb{R}^{n \times n} : \text{rank}(\mathbf{X}) \leq k, \|\mathbf{X}\|_\sigma \leq 1\}$, where $\|\cdot\|_\sigma$ denotes the spectral norm. Moreover, several conditions have been established under which nuclear norm minimization methods are guaranteed to return the ground truth matrix (Candès and Tao, 2010; Candès and Recht, 2012) though such conditions tend to be strong and hard to verify in practice. Soft-Impute proceeds by iteratively replacing the missing elements of the matrix with those obtained from a soft thresholded low rank singular value decomposition. Accordingly, similarly to Iterative-SVD, Soft-Impute relies on the computation of a low rank SVD as the primary algorithmic workhorse. The approach relies on the result that for an arbitrary matrix \mathbf{X} , the solution of the problem $\min_{\mathbf{Z}} \frac{1}{2} \|\mathbf{X} - \mathbf{Z}\|_F^2 + \lambda \|\mathbf{Z}\|_*$ is given by $\hat{\mathbf{Z}} = S_\lambda(\mathbf{X})$ where $S_\lambda(\cdot)$ denotes the soft-thresholding operation (Donoho et al., 1995). Explicitly, Soft-Impute proceeds as follows for a given regularization parameter $\lambda > 0$ and termination threshold $\epsilon > 0$:

1. Initialize the iteration count $t \leftarrow 0$ and initialize $\mathbf{Z}_t = \mathbf{0}_{n \times m}$.
2. Compute $\mathbf{Z}_{t+1} = S_\lambda(P_\Omega(\mathbf{X}) + P_\Omega^\perp(\mathbf{Z}_t))$ where $P_\Omega(\cdot)$ denotes the operation that projects onto the revealed entries of \mathbf{X} while $P_\Omega^\perp(\cdot)$ denotes the operation that projects onto the missing entries of \mathbf{X} .
3. Terminate if $\frac{\|\mathbf{Z}_t - \mathbf{Z}_{t+1}\|_F^2}{\|\mathbf{Z}_t\|_F^2}$. Otherwise, increment t and return to Step 2.

2.1.3 FAST-IMPUTE

Fast-Impute is a projected gradient descent approach to MC that has desirable global convergence properties (Bertsimas and Li, 2020). Fast-Impute belongs to the broad class of methods that solve MC by factorizing the target matrix as $\mathbf{X} = \mathbf{U}\mathbf{V}^T$ where $\mathbf{U} \in \mathbb{R}^{n \times k}$, $\mathbf{V} \in \mathbb{R}^{m \times k}$ and performing some variant of gradient descent (or alternating minimization) on the matrices \mathbf{U} and \mathbf{V} (Koren et al., 2009; Jain and Netrapalli, 2015; Zheng and Lafferty, 2016; Jin et al., 2016). We note that we leverage this common factorization in the approach to (1) presented in this work. Gradient descent based methods have shown great success. Despite the non-convexity of the factorization, it has been shown that in many cases gradient descent and its variants will nevertheless converge to a globally optimal solution (Chen and Wainwright, 2015; Ge et al., 2015; Sun and Luo, 2016; Ma et al., 2018; Bertsimas and Li, 2020). Fast-Impute takes the approach of expressing \mathbf{U} as a closed form function of \mathbf{V} after performing the factorization and directly performs projected gradient descent updates on \mathbf{V} with classic Nesterov acceleration (Nesterov, 1983). Moreover, to enhance scalability of their method, Bertsimas and Li (2020) design a stochastic gradient extension of Fast-Impute that estimates the gradient at each update step by only considering a sub sample of the rows and columns of the target matrix.

2.2 Low Rank Optimization Methods

2.2.1 SCALEDGD

ScaledGD is a highly performant method to obtain strong solutions to low rank matrix estimation problems that take the following form:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times m}} f(\mathbf{X}) = \frac{1}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2 \quad \text{s.t. } \text{rank}(\mathbf{X}) \leq k,$$

where $\mathcal{A}(\cdot) : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^l$ models some measurement process and we have $\mathbf{y} \in \mathbb{R}^l$ (Tong et al., 2021). ScaledGD proceeds by factorizing the target matrix as $\mathbf{X} = \mathbf{U}\mathbf{V}^T$ and iteratively performing gradient updates on the low rank factors \mathbf{U}, \mathbf{V} after preconditioning the gradients with an adaptive matrix that is efficient to compute. Doing so yields a linear convergence rate that is notably independent of the condition number of the low rank matrix. In so doing, ScaledGD combines the desirable convergence rate of alternating minimization with the desirable low per-iteration cost of gradient descent. Explicitly, letting $\mathcal{L}(\mathbf{U}, \mathbf{V}) = f(\mathbf{U}\mathbf{V}^T)$, ScaledGD updates the low rank factors as:

$$\begin{aligned} \mathbf{U}_{t+1} &\leftarrow \mathbf{U}_t - \eta \nabla_{\mathbf{U}} \mathcal{L}(\mathbf{U}_t, \mathbf{V}_t) (\mathbf{V}_t^T \mathbf{V}_t)^{-1}, \\ \mathbf{V}_{t+1} &\leftarrow \mathbf{V}_t - \eta \nabla_{\mathbf{V}} \mathcal{L}(\mathbf{U}_t, \mathbf{V}_t) (\mathbf{U}_t^T \mathbf{U}_t)^{-1}, \end{aligned}$$

where $\eta > 0$ denotes the step size.

2.2.2 MIXED-PROJECTION CONIC OPTIMIZATION

Mixed-projection conic optimization is a recently proposed modelling and algorithmic framework designed to tackle a broad class of matrix optimization problems (Bertsimas et al., 2022, 2023c). Specifically, this approach considers problems that have the following form:

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times m}} \quad & \langle \mathbf{C}, \mathbf{X} \rangle + \lambda \cdot \text{rank}(\mathbf{X}) + \Omega(\mathbf{X}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{X} = \mathbf{B}, \text{rank}(\mathbf{X}) \leq k, \mathbf{X} \in \mathcal{K}, \end{aligned} \tag{2}$$

where $\mathbf{C} \in \mathbb{R}^{n \times m}$ is a cost matrix, $\lambda > 0$, $k \in \mathbb{N}_+$, $\mathbf{A} \in \mathbb{R}^{l \times n}$, $\mathbf{B} \in \mathbb{R}^{l \times m}$, \mathcal{K} denotes a proper cone in the sense of Boyd et al. (2004) and $\Omega(\cdot)$ is a frobenius norm regularization function or a spectral norm regularization function of the input matrix. The main workhorse of mixed-projection conic optimization is the use of a projection matrix to cleverly model the rank terms in (2). This can be viewed as the matrix generalization of using binary variables to model the sparsity of a vector in mixed-integer optimization. Bertsimas et al. (2022) show that for an arbitrary matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$, we have

$$\text{rank}(\mathbf{X}) \leq k \iff \exists \mathbf{P} \in \mathbb{S}^n : \mathbf{P}^2 = \mathbf{P}, \text{Tr}(\mathbf{P}) \leq k, \mathbf{X} = \mathbf{P}\mathbf{X}.$$

Introducing projection matrices allows the rank functions to be eliminated from (2) at the expense of introducing non convex quadratic equality constraints. From here, most existing works that leverage mixed-projection conic optimization have either focused on obtaining strong semidefinite based convex relaxations (Bertsimas et al., 2022, 2023c) or have focused on obtaining certifiably optimal solutions for small and moderately sized problem instances (Bertsimas et al., 2023a,b). In this work, we leverage the mixed-projection framework to scalably obtain high quality solutions to large problem instances.

2.3 Alternating Direction Method of Multipliers

Alternating direction method of multipliers (ADMM) is an algorithm that was originally designed to solve linearly constrained convex optimization problems of the form

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m} f(\mathbf{x}) + g(\mathbf{z}) \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} = \mathbf{c}, \tag{3}$$

where we have $\mathbf{A} \in \mathbb{R}^{l \times n}$, $\mathbf{B} \in \mathbb{R}^{l \times m}$, $\mathbf{c} \in \mathbb{R}^l$ and the functions f and g are assumed to be convex (Boyd et al., 2011). The main benefit of ADMM is that it can combine the decomposition benefits of dual ascent with the desirable convergence properties of the method of multipliers. Letting $\mathbf{y} \in \mathbb{R}^l$ denote the dual variable, the augmented lagrangian of (3) is given by

$$\mathcal{L}_\rho(\mathbf{x}, \mathbf{z}, \mathbf{y}) = f(\mathbf{x}) + g(\mathbf{z}) + \mathbf{y}^T(\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}) + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{z} - \mathbf{c}\|_2^2,$$

where $\rho > 0$ is the augmented lagrangian parameter. ADMM then proceeds by iteratively updating the primal variable \mathbf{x} , updating the primal variable \mathbf{z} and taking a gradient ascent step on the dual variable \mathbf{y} . Explicitly, ADMM consists of the following updates:

1. $\mathbf{x}_{t+1} \leftarrow \text{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{z}_t, \mathbf{y}_t)$,
2. $\mathbf{z}_{t+1} \leftarrow \text{argmin}_{\mathbf{z}} \mathcal{L}(\mathbf{x}_{t+1}, \mathbf{z}, \mathbf{y}_t)$,
3. $\mathbf{y}_{t+1} \leftarrow \mathbf{y}_t + \rho(\mathbf{A}\mathbf{x}_{t+1} + \mathbf{B}\mathbf{z}_{t+1} - \mathbf{c})$.

Under very mild regularity conditions on f, g and \mathcal{L} , it is well known that ADMM is guaranteed to produce a sequence of primal iterates that converges to the optimal value of (3) and a sequence of dual iterates that converge to the optimal dual variable (note that there is no guarantee of primal variable convergence) (Boyd et al., 2011). Importantly, although ADMM was originally designed for linearly constrained convex optimization, it has often been applied to non convex optimization problems and yielded empirically strong results (Xu et al., 2016). This observation has motivated work to explore the theoretical convergence behavior of ADMM and its variants on specific classes of non convex optimization problems (Guo et al., 2017; Wang et al., 2019; Wang and Zhao, 2021).

3. Formulation Properties

In this section, we rigorously investigate certain key features of (1). Specifically, we establish an equivalence between (1) and an appropriately defined robust optimization problem. Moreover, we illustrate that (1) can be reduced to an optimization problem over only \mathbf{X} and establish that the resulting objective function is not convex, not concave and non-smooth. Finally, we study how efficient evaluations of the reduced problem objective function can be performed.

3.1 Equivalence Between Regularization and Robustness

Real-world datasets frequently contain inaccuracies and missing values, which hinder the ability of machine learning models to generalize effectively to new data when these inconsistencies are not appropriately modelled. Consequently, robustness is a crucial quality for machine learning models, both in theory and application (Xu et al., 2009; Bertsimas and den Hertog, 2020). In this section, we show that our regularized problem (1) can be viewed as a robust optimization (RO) problem. This finding justifies the inclusion of the nuclear norm regularization term in (1) and is in a similar flavor as known results from the robust optimization literature in the case of vector (Bertsimas and Copenhaver, 2018) and matrix (Bertsimas et al., 2023a) problems.

Proposition 1 *Problem (1) is equivalent to the following robust optimization problem:*

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times m}, \boldsymbol{\alpha} \in \mathbb{R}^{m \times d}} \max_{\boldsymbol{\Delta} \in \mathcal{U}} \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \|\mathbf{Y} - \mathbf{X}\boldsymbol{\alpha}\|_F^2 + \langle \mathbf{X}, \boldsymbol{\Delta} \rangle \\ \text{s.t.} \quad \text{rank}(\mathbf{X}) \leq k_0, \end{aligned} \tag{4}$$

where $\mathcal{U} = \{\boldsymbol{\Delta} \in \mathbb{R}^{n \times m} : \|\boldsymbol{\Delta}\|_\sigma \leq \gamma\}$

Proof To establish this result, it suffices to argue that $\max_{\boldsymbol{\Delta} \in \mathcal{U}} \langle \mathbf{X}, \boldsymbol{\Delta} \rangle = \gamma \|\mathbf{X}\|_\star$. This equivalence follows immediately from the fact that the nuclear norm is dual to the spectral norm. So as to keep this manuscript self contained, we present a proof of this equivalence below.

Consider any matrix $\bar{\boldsymbol{\Delta}} \in \mathbb{R}^{n \times m}$ such that $\|\bar{\boldsymbol{\Delta}}\|_\sigma \leq \gamma$. Let $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$ be a singular value decomposition of \mathbf{X} where we let $r = \text{rank}(\mathbf{X})$ and we have $\mathbf{U} \in \mathbb{R}^{n \times r}, \boldsymbol{\Sigma} \in \mathbb{R}^{r \times r}, \mathbf{V} \in$

$\mathbb{R}^{m \times r}$. We have

$$\begin{aligned} \langle \mathbf{X}, \bar{\Delta} \rangle &= \text{Tr}(\bar{\Delta}^T \mathbf{U} \Sigma \mathbf{V}^T) = \text{Tr}(\mathbf{V}^T \bar{\Delta}^T \mathbf{U} \Sigma) = \langle \mathbf{U}^T \bar{\Delta} \mathbf{V}, \Sigma \rangle = \sum_{i=1}^r \Sigma_{ii} (\mathbf{U}^T \bar{\Delta} \mathbf{V})_{ii} \\ &= \sum_{i=1}^r \Sigma_{ii} U_i^T \bar{\Delta} V_i \leq \sum_{i=1}^r \Sigma_{ii} \sigma_1(\bar{\Delta}) \leq \gamma \sum_{i=1}^r \Sigma_{ii} = \gamma \|\mathbf{X}\|_* \end{aligned}$$

where we have used the fact that Σ is a diagonal matrix and the columns of \mathbf{U} and \mathbf{V} have unit length. Thus, we have shown that $\gamma \|\mathbf{X}\|_*$ is an upper bound for $\max_{\Delta \in \mathcal{U}} \langle \mathbf{X}, \Delta \rangle$. To show that the upper bound is always achieved, consider the matrix $\tilde{\Delta} = \gamma \mathbf{U} \mathbf{V}^T \in \mathbb{R}^{n \times m}$ where \mathbf{U} and \mathbf{V} are taken from a singular value decomposition of \mathbf{X} . Observe that

$$\|\tilde{\Delta}\|_\sigma = \gamma \|\mathbf{U} \mathbf{V}^T\|_\sigma \leq \gamma \implies \tilde{\Delta} \in \mathcal{U}.$$

We conclude by noting that

$$\langle \mathbf{X}, \tilde{\Delta} \rangle = \text{Tr}(\mathbf{V} \Sigma \mathbf{U}^T \gamma \mathbf{U} \mathbf{V}^T) = \gamma \text{Tr}(\mathbf{V}^T \mathbf{V} \Sigma \mathbf{U}^T \mathbf{U}) = \gamma \text{Tr}(\mathbf{I} \Sigma \mathbf{I}) = \gamma \|\mathbf{X}\|_*.$$

■

Proposition 1 implies that solving the nuclear norm regularized (1) is equivalent to solving an unregularized robust optimization problem that protects against adversarial perturbations that are bounded in spectral norm. This result is not surprising given the duality of norms, yet is nevertheless insightful.

3.2 A Partial Minimization

Let $g(\mathbf{X}, \alpha)$ denote the objective function of (1). Note that $g(\mathbf{X}, \alpha)$ is bi-convex in (\mathbf{X}, α) but is not jointly convex due to the product $\mathbf{X}\alpha$. Observe that we can simplify (1) by performing a partial minimization in α . For any \mathbf{X} , the problem in α requires finding the unconstrained minimum of a convex quadratic function. The gradient of g with respect to α is given by $\nabla_\alpha g(\mathbf{X}, \alpha) = 2\lambda \mathbf{X}^T (\mathbf{X}\alpha - \mathbf{Y})$. Setting $\nabla_\alpha g(\mathbf{X}, \alpha)$ to 0 yields $\alpha^* = (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y}$ as a minimizer of g over α . Letting $f(\mathbf{X})$ correspond to the partially minimized objective function of (1), we have

$$\begin{aligned} f(\mathbf{X}) &= \min_{\alpha} g(\mathbf{X}, \alpha) = \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \|(\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T) \mathbf{Y}\|_F^2 + \gamma \|\mathbf{X}\|_* \\ &= \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T (\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T) \mathbf{Y}) + \gamma \|\mathbf{X}\|_* \end{aligned}$$

We note that α^* corresponds to the well studied ordinary least squares solution. When $\mathbf{X}^T \mathbf{X}$ has full rank, α is the unique minimizer of g . If $\mathbf{X}^T \mathbf{X}$ is rank deficient, α^* corresponds to the minimizer with minimum norm.

Though we have simplified the objective function of (1), $f(\mathbf{X})$ is not a particularly well behaved function. We formalize this statement in Proposition (2).

Proposition 2 *The function $f(\mathbf{X})$ is in general neither convex nor concave and is non-smooth.*

Proof To illustrate that $f(\mathbf{X})$ is in general neither convex nor concave, suppose that $\Omega = \emptyset$, $n = 2$ and $m = d = \lambda = \gamma = 1$. In this setting, we have $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{2 \times 1}$. Assuming that $\mathbf{x} \neq \mathbf{0}_2$, we can write the objective function as

$$\begin{aligned} f(\mathbf{x}) &= \text{Tr}(\mathbf{y}^T (\mathbf{I}_2 - \mathbf{x}(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T) \mathbf{y}) + \|\mathbf{x}\|_* \\ &= \mathbf{y}^T \mathbf{y} - \frac{\mathbf{y}^T \mathbf{x} \mathbf{x}^T \mathbf{y}}{\mathbf{x}^T \mathbf{x}} + \|\mathbf{x}\|_2 \\ &= \mathbf{y}^T \mathbf{y} - \frac{(\mathbf{y}^T \mathbf{x})^2}{\mathbf{x}^T \mathbf{x}} + \sqrt{\mathbf{x}^T \mathbf{x}}. \end{aligned}$$

For $\mathbf{x} = \mathbf{0}_2$, the objective value $f(\mathbf{0}_2)$ is equal to $\mathbf{y}^T \mathbf{y}$. Let $\mathbf{y} = \mathbf{1}_2$ and consider the line in \mathbf{R}^2 defined by $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^2 : x_2 = x_1 + 1\}$. The restriction of $f(\mathbf{x})$ to the line defined by \mathcal{X} is a univariate function given by

$$f_{\mathcal{X}}(t) = 2 - \frac{(2t + 1)^2}{2t^2 + 2t + 1} + \sqrt{2t^2 + 2t + 1}$$

where $t \in \mathbb{R}$ is a dummy variable. Observe that we have $f_{\mathcal{X}}(-1) = f_{\mathcal{X}}(0) = 2$, $f_{\mathcal{X}}(-0.5) = 2 + \frac{\sqrt{2}}{2}$ and $f_{\mathcal{X}}(-4) = f_{\mathcal{X}}(3) = 5.04$. Thus, the point $(-0.5, f_{\mathcal{X}}(0.5))$ lies above the chord connecting $(-1, f_{\mathcal{X}}(-1))$ and $(0, f_{\mathcal{X}}(0))$, so $f_{\mathcal{X}}(t)$ is not a convex function. Moreover, the point $(-1, f_{\mathcal{X}}(-1))$ lies below the chord connecting $(-4, f_{\mathcal{X}}(-4))$ and $(-0.5, f_{\mathcal{X}}(-0.5))$, so $f_{\mathcal{X}}(t)$ is not a concave function. Since a function is convex (respectively concave) if and only if its restriction to every line is convex (respectively concave), we have established that $f(\mathbf{X})$ is neither convex nor concave since $f_{\mathcal{X}}(t)$ is neither convex nor concave. To conclude the proof of Proposition 2, note that the non-smooth property of $f(\mathbf{X})$ follows immediately from the non-smooth property of the nuclear norm function. \blacksquare

Although the above closed form partial minimization in $\boldsymbol{\alpha}$ eliminates $m \times d$ variables from (1), this comes at the expense of introducing a $m \times m$ matrix pseudo-inverse term into the objective function which can be computationally expensive to evaluate. Efficient evaluation of an objective function is crucial in many optimization problems to quickly measure solution quality. A plethora of modern optimization techniques require iterative objective function evaluations. As a result, the computational cost of evaluating an objective function can quickly become the bottleneck of an algorithm's complexity. Directly evaluating $f(\mathbf{X})$ naively requires $O(|\Omega|)$ operations for the first term, $O(m^2n + m^3 + n^2d)$ for the second term (forming the matrix $\mathbf{X}^T \mathbf{X}$ is $O(m^2n)$, taking the pseudo-inverse is $O(m^3)$, computing the products involving \mathbf{Y} is $O(n^2d)$) and requires $O(mn \min(m, n))$ for the third term (the nuclear norm can be evaluated by computing a singular value decomposition of \mathbf{X}). We observe that computing the second term of $f(\mathbf{X})$ involving the pseudo-inverse dominates the complexity calculation. Indeed, the overall complexity of evaluating $f(\mathbf{X})$ naively is $O(m^2n + m^3 + n^2d)$.

Fortunately, it is possible to make evaluations of $f(\mathbf{X})$ without explicitly forming the product $\mathbf{X}^T \mathbf{X}$ or taking a pseudo-inverse. Proposition 3 illustrates that it suffices (in terms of computational complexity) to take a singular value decomposition of \mathbf{X} . Moreover, a large class of optimization algorithms require only function evaluations for feasible solutions. If we consider only those values of \mathbf{X} that are feasible to (1), it is sufficient (in terms of

computational complexity) to take a rank k truncated singular value decomposition of \mathbf{X} to make functions evaluations of $f(\mathbf{X})$.

Proposition 3 *The function $f(\mathbf{X})$ can equivalently be written as*

$$f(\mathbf{X}) = \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T (\mathbf{I}_n - \mathbf{U}\mathbf{U}^T) \mathbf{Y}) + \gamma \|\mathbf{X}\|_*,$$

where $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ is a singular value decomposition of \mathbf{X} where we let $r = \text{rank}(\mathbf{X})$ and we have $\mathbf{U} \in \mathbb{R}^{n \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$, $\mathbf{V} \in \mathbb{R}^{m \times r}$.

Proof To establish the result, it suffices to show that

$$\text{Tr}(\mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y}) = \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{U}^T \mathbf{Y}).$$

Let $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be a singular value decomposition of \mathbf{X} where $r = \text{rank}(\mathbf{X})$ and $\mathbf{U} \in \mathbb{R}^{n \times r}$, $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$, $\mathbf{V} \in \mathbb{R}^{m \times r}$. Observe that

$$\begin{aligned} \text{Tr}(\mathbf{Y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T \mathbf{Y}) &= \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T (\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^\dagger \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{Y}) \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T (\mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T)^\dagger \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{Y}) \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{V}\mathbf{\Sigma}^{-2} \mathbf{V}^T \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{Y}) \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^{-2} \mathbf{\Sigma}\mathbf{U}^T \mathbf{Y}) \\ &= \text{Tr}(\mathbf{Y}^T \mathbf{U}\mathbf{U}^T \mathbf{Y}), \end{aligned}$$

where we have repeatedly invoked the property that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_r$. ■

In light of Proposition 3, evaluating $f(\mathbf{X})$ for feasible solutions still requires $O(|\Omega|)$ operations for the first term, but the second term can be evaluated using $O(kn(m+d))$ operations (performing a truncated singular value decomposition is $O(knm)$ and computing the products involving \mathbf{Y} is $O(knd)$) and the third term can be evaluated using $O(knm)$ operations (by performing a truncated singular value decomposition) for an overall complexity of $O(kn(m+d))$. This is significantly less expensive than the $O(m^2n + m^3 + n^2d)$ complexity of naive direct evaluation of $f(\mathbf{X})$ introduced previously.

4. An Exact Mixed-Projection Formulation

In this section, we reformulate (1) as a mixed-projection optimization problem and further reduce the dimension of the resulting problem in a commonly studied manner by parameterizing \mathbf{X} as the matrix product of two low dimensional matrices. Thereafter, we illustrate how to employ the matrix generalization of the perspective relaxation (Günlük and Linderoth, 2012; Bertsimas et al., 2022, 2023c,a) to construct a convex relaxation of (1).

We first note that given the result of Section 3.2, we can rewrite (1) as an optimization problem only over \mathbf{X} as follows:

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times m}} \quad & \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T (\mathbf{I}_n - \mathbf{X} (\mathbf{X}^T \mathbf{X})^\dagger \mathbf{X}^T) \mathbf{Y}) + \gamma \|\mathbf{X}\|_* \\ \text{s.t.} \quad & \text{rank}(\mathbf{X}) \leq k. \end{aligned} \tag{5}$$

Observe that the matrix $\mathbf{X}(\mathbf{X}^T\mathbf{X})^\dagger\mathbf{X}^T$ is the linear transformation that projects vectors onto the subspace spanned by the columns of the matrix \mathbf{X} . Drawing on ideas presented in Bertsimas et al. (2022, 2023c,a), we introduce an orthogonal projection matrix $\mathbf{P} \in \mathcal{P}_k$ to model the column space of \mathbf{X} where $\mathcal{P}_\eta = \{\mathbf{P} \in \mathcal{S}^n : \mathbf{P}^2 = \mathbf{P}, \text{tr}(\mathbf{P}) \leq \eta\}$ for $\eta \geq 0$. We can express the desired relationship between \mathbf{P} and \mathbf{X} as $\mathbf{X} = \mathbf{P}\mathbf{X}$ since projecting a matrix onto its own column space leaves the matrix unchanged. This gives the following reformulation of (1):

$$\begin{aligned} \min_{\mathbf{X} \in \mathbb{R}^{n \times m}, \mathbf{P} \in \mathbb{R}^{n \times n}} \quad & \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Y}) + \gamma \|\mathbf{X}\|_* \\ \text{s.t.} \quad & (\mathbf{I}_n - \mathbf{P})\mathbf{X} = \mathbf{0}_{n \times m}, \mathbf{P} \in \mathcal{P}_{\min(k, \text{rank}(\mathbf{X}))}. \end{aligned} \quad (6)$$

Observe that the matrix pseudo-inverse term has been eliminated from the objective function, however we have introduced the bilinear constraint $\mathbf{X} = \mathbf{P}\mathbf{X}$ which is non convex in the optimization variables as well as the non convex constraint $\mathbf{P} \in \mathcal{P}_{\min(k, \text{rank}(\mathbf{X}))}$. We now have the following result:

Proposition 4 *Problem (6) is a valid reformulation of (5).*

Proof We show that given a feasible solution to (6), we can construct a feasible solution to (5) that achieves the same objective value and vice versa.

Consider an arbitrary feasible solution $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ to (6). Since $\bar{\mathbf{P}}\bar{\mathbf{X}} = \bar{\mathbf{X}}$ and $\bar{\mathbf{P}} \in \mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{X}}))}$, we have $\text{rank}(\bar{\mathbf{X}}) \leq k$. We claim that $\bar{\mathbf{X}}$ achieves the same objective value in (5) as $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ achieves in (6). To show this, it suffices to illustrate that for all $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ feasible to (6) we have $H(\bar{\mathbf{X}}) := \bar{\mathbf{X}}(\bar{\mathbf{X}}^T\bar{\mathbf{X}})^\dagger\bar{\mathbf{X}}^T = \bar{\mathbf{P}}$. The matrix $\bar{\mathbf{P}}$ is an orthogonal projection matrix since it is symmetric and satisfies $\bar{\mathbf{P}}^2 = \bar{\mathbf{P}}$. Moreover, since $\text{rank}(\bar{\mathbf{P}}) = \text{rank}(\bar{\mathbf{X}})$ and $\bar{\mathbf{P}}\bar{\mathbf{X}} = \bar{\mathbf{X}}$ we know that $\bar{\mathbf{P}}$ is an orthogonal projection onto the subspace spanned by the columns of $\bar{\mathbf{X}}$. Similarly, it can easily be verified that $H(\bar{\mathbf{X}})$ is symmetric and satisfies $H(\bar{\mathbf{X}})^2 = H(\bar{\mathbf{X}})$, $\text{rank}(H(\bar{\mathbf{X}})) = \text{rank}(\bar{\mathbf{X}})$ and $H(\bar{\mathbf{X}})\bar{\mathbf{X}} = \bar{\mathbf{X}}$. Thus, $H(\bar{\mathbf{X}})$ is also an orthogonal projection matrix onto the subspace spanned by the columns of $\bar{\mathbf{X}}$. To conclude, we invoke the property that given a subspace $\mathcal{V} \subset \mathbb{R}^n$ the orthogonal projection onto \mathcal{V} is uniquely defined. To see this, suppose \mathbf{P}_1 and \mathbf{P}_2 are two orthogonal projections onto \mathcal{V} . Let $l = \dim(\mathcal{V})$. Let $\{\mathbf{e}_i\}_{i=1}^l$ be an orthogonal basis for \mathcal{V} and let $\{\mathbf{e}_i\}_{l+1}^n$ be an orthogonal basis for \mathcal{V}^\perp . Since \mathbf{P}_1 is an orthogonal projection onto \mathcal{V} , we have $\mathbf{P}_1\mathbf{e}_i = \mathbf{e}_i$ for all $1 \leq i \leq l$ and $\mathbf{P}_1\mathbf{e}_i = \mathbf{0}_n$ for all $l+1 \leq i \leq n$. However, the same must hold for \mathbf{P}_2 which implies that $\mathbf{P}_1 = \mathbf{P}_2$.

Consider an arbitrary feasible solution $\bar{\mathbf{X}}$ to (5). Let $r = \text{rank}(\bar{\mathbf{X}})$ and $\bar{\mathbf{X}} = \bar{\mathbf{U}}\bar{\mathbf{\Sigma}}\bar{\mathbf{V}}^T$ be a singular value decomposition of $\bar{\mathbf{X}}$ where we have $\bar{\mathbf{U}} \in \mathbb{R}^{n \times r}$, $\bar{\mathbf{\Sigma}} \in \mathbb{R}^{r \times r}$, $\bar{\mathbf{V}} \in \mathbb{R}^{m \times r}$. Define $\bar{\mathbf{P}} = \bar{\mathbf{U}}\bar{\mathbf{U}}^T$. By construction, we have $\bar{\mathbf{P}} \in \mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{X}}))}$ since $r \leq k$. Moreover, it is easy to verify that

$$\bar{\mathbf{P}}\bar{\mathbf{X}} = \bar{\mathbf{U}}\bar{\mathbf{U}}^T\bar{\mathbf{U}}\bar{\mathbf{\Sigma}}\bar{\mathbf{V}}^T = \bar{\mathbf{U}}\bar{\mathbf{\Sigma}}\bar{\mathbf{V}}^T = \bar{\mathbf{X}},$$

where we have used the property $\bar{\mathbf{U}}^T\bar{\mathbf{U}} = \mathbf{I}_r$. Finally, Proposition 3 immediately implies that $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ achieves the same objective in (6) as $\bar{\mathbf{X}}$ achieves in (5). This completes the proof. \blacksquare

Optimizing explicitly over the space of $n \times m$ matrices can rapidly become prohibitively

costly in terms of runtime and memory requirements. Accordingly, we adopt the common approach of factorizing $\mathbf{X} \in \mathbb{R}^{n \times m}$ as $\mathbf{U}\mathbf{V}^T$ for $\mathbf{U} \in \mathbb{R}^{n \times k}$, $\mathbf{V} \in \mathbb{R}^{m \times k}$. This leads to the following formulation:

$$\begin{aligned} & \min_{\substack{\mathbf{U} \in \mathbb{R}^{n \times k}, \mathbf{V} \in \mathbb{R}^{m \times k}, \\ \mathbf{P} \in \mathbb{R}^{n \times n}}} \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Y}) + \frac{\gamma}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\ \text{s.t.} \quad & (\mathbf{I}_n - \mathbf{P})\mathbf{U} = \mathbf{0}_{n \times k}, \mathbf{P} \in \mathcal{P}_{\min(k, \text{rank}(\mathbf{U}\mathbf{V}^T))}. \end{aligned} \tag{7}$$

Notice that we have replaced $n \times m$ optimization variables with $k \times (n + m)$ optimization variables, an often significant dimension reduction in practice. Attentive readers may object that though this is true, we have introduced n^2 decision variables through the introduction of the projection matrix variable \mathbf{P} which nullifies any savings introduced through the factorization of \mathbf{X} . Note, however, that it is possible to factor any feasible projection matrix as $\mathbf{P} = \mathbf{M}\mathbf{M}^T$ for some $\mathbf{M} \in \mathbb{R}^{n \times k}$. In Section 5, we leverage this fact so that the presence of the projection matrix incurs a cost of $n \times k$ additional variables rather than n^2 variables. We have the following result:

Proposition 5 *Problem (7) is a valid reformulation of (6).*

Proof We show that given a feasible solution to (7), we can construct a feasible solution to (6) that achieves the same or lesser objective value and vice versa.

Consider an arbitrary feasible solution $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ to (7). Let $\bar{\mathbf{X}} = \bar{\mathbf{U}}\bar{\mathbf{V}}^T$. We will show that $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ is feasible to (6) and achieves the same or lesser objective as $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ does in (7). Feasibility of (7) implies that $\bar{\mathbf{P}} \in \mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{U}}\bar{\mathbf{V}}^T))} = \mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{X}}))}$ and also that

$$(\mathbf{I}_n - \bar{\mathbf{P}})\bar{\mathbf{X}} = (\mathbf{I}_n - \bar{\mathbf{P}})\bar{\mathbf{U}}\bar{\mathbf{V}}^T = \mathbf{0}_{n \times k}\bar{\mathbf{V}}^T = \mathbf{0}_{n \times m},$$

thus the solution $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ is certainly feasible for (6). To see that $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ achieves the same or lesser objective value, it suffices to argue that $\|\bar{\mathbf{X}}\|_* \leq \frac{1}{2}(\|\bar{\mathbf{U}}\|_F^2 + \|\bar{\mathbf{V}}\|_F^2)$. This follows immediately from the following lemma established by Mazumder et al. (2010) (see Appendix A.5 in their paper for a proof):

Lemma 6 *For any matrix \mathbf{Z} , the following holds:*

$$\|\mathbf{Z}\|_* = \min_{\mathbf{U}, \mathbf{V}: \mathbf{Z} = \mathbf{U}\mathbf{V}^T} \frac{1}{2} (\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2).$$

If $\text{rank}(\mathbf{Z}) = k \leq \min(m, n)$, then the minimum above is attained at a factor decomposition $\mathbf{U}_{n \times k}\mathbf{V}_{m \times k}^T$. Letting $\mathbf{Z}_{n \times m} = \mathbf{L}_{n \times k}\mathbf{\Sigma}_{k \times k}\mathbf{R}_{m \times k}^T$ denote a singular value decomposition of \mathbf{Z} , the minimum above is attained at $\mathbf{U}_{n \times k} = \mathbf{L}_{n \times k}\mathbf{\Sigma}_{k \times k}^{\frac{1}{2}}$, $\mathbf{V}_{m \times k} = \mathbf{R}_{m \times k}\mathbf{\Sigma}_{k \times k}^{\frac{1}{2}}$.

Consider now an arbitrary feasible solution $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ to (6). Let $\bar{\mathbf{X}} = \mathbf{L}\mathbf{\Sigma}\mathbf{R}^T$ be a singular value decomposition of $\bar{\mathbf{X}}$ where $\mathbf{L} \in \mathbb{R}^{n \times k}$, $\mathbf{\Sigma} \in \mathbb{R}^{k \times k}$, $\mathbf{R} \in \mathbb{R}^{m \times k}$ and define $\bar{\mathbf{U}} = \mathbf{L}\mathbf{\Sigma}^{\frac{1}{2}}$, $\bar{\mathbf{V}} = \mathbf{R}\mathbf{\Sigma}^{\frac{1}{2}}$. Feasibility of $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ in (6) implies that $\bar{\mathbf{P}} \in \mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{X}}))} =$

$\mathcal{P}_{\min(k, \text{rank}(\bar{\mathbf{U}}\bar{\mathbf{V}}^T))}$. Moreover, since the columns of \mathbf{L} form an orthogonal basis for the column space of $\bar{\mathbf{X}}$, the condition $(\mathbf{I}_n - \bar{\mathbf{P}})\bar{\mathbf{X}} = \mathbf{0}_{n \times m}$ implies that

$$(\mathbf{I}_n - \bar{\mathbf{P}})\bar{\mathbf{U}} = (\mathbf{I}_n - \bar{\mathbf{P}})\mathbf{L}\boldsymbol{\Sigma}^{\frac{1}{2}} = \mathbf{0}_{n \times k}\boldsymbol{\Sigma}^{\frac{1}{2}} = \mathbf{0}_{n \times k}.$$

Thus, the solution $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ is feasible to (7). Moreover, by Lemma 6 we have $\frac{1}{2}(\|\bar{\mathbf{U}}\|_F^2 + \|\bar{\mathbf{V}}\|_F^2) = \|\bar{\mathbf{X}}\|_*$ so $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ achieves the same objective in (7) as $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ does in (6). This completes the proof. ■

In the remainder of the paper, we will relax the constraint $\mathbf{P} \in \mathcal{P}_{\min(k, \text{rank}(\mathbf{U}\mathbf{V}^T))}$ to $\mathbf{P} \in \mathcal{P}_k$ and develop a scalable algorithm to obtain high quality feasible solutions. Explicitly, we consider the problem given by:

$$\begin{aligned} & \min_{\substack{\mathbf{U} \in \mathbb{R}^{n \times k}, \mathbf{V} \in \mathbb{R}^{m \times k}, \\ \mathbf{P} \in \mathbb{R}^{n \times n}}} \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Y}) + \frac{\gamma}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) \\ \text{s.t.} \quad & (\mathbf{I}_n - \mathbf{P})\mathbf{U} = \mathbf{0}_{n \times k}, \mathbf{P} \in \mathcal{P}_k. \end{aligned} \tag{8}$$

It is straightforward to see that the optimal value of (8) is no greater than the optimal value of (7). Unfortunately, the converse does not necessarily hold. To see why the optimal value of (8) can be strictly less than that of (7) in certain pathological cases, suppose we had $k = n = m$, $\Omega = \emptyset$. In this setting, letting $\bar{\mathbf{P}} = \mathbf{I}_n$, $\bar{\mathbf{U}} = \mathbf{0}_{n \times k}$ and $\bar{\mathbf{V}} = \mathbf{0}_{m \times k}$, the solution $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ would be feasible to (8) and achieve an objective value of 0. However the optimal value of (7) would be strictly greater than 0 in this setting as long as $\mathbf{Y} \neq \mathbf{0}$. Although (8) is a relaxation of (1), we will see in Section 6 that the solutions we will obtain to (8) will be high quality solutions for (1), the main problem of interest.

4.1 A Positive Semidefinite Cone Relaxation

Convex relaxations are useful in non convex optimization primarily for two reasons. Firstly, given the objective value achieved by an arbitrary feasible solution, strong convex relaxations can be used to upperbound the worst case suboptimality of said solution. Secondly, convex relaxations can often be used as building blocks for global optimization procedures. In this section, we present a natural convex relaxation of (8) that leverages the matrix generalization of the perspective relaxation (Günlük and Linderoth, 2012; Bertsimas et al., 2022, 2023c,a).

Rather than working directly with (8), consider the equivalent formulation (6) with $\mathcal{P}_{\min(k, \text{rank}(\mathbf{X}))}$ replaced by \mathcal{P}_k . Before proceeding, we will assume knowledge of an upper bound $M \in \mathbb{R}_+$ on the spectral norm of an optimal \mathbf{X} to (6). Tighter bounds M are desirable as they will lead to stronger convex relaxations of (6). We note that it is always possible to specify such an upper bound M without prior knowledge of an optimal solution to (6). To see this, note that setting $\mathbf{X} = \mathbf{0}_{n \times m}$ in (6) produces an objective value of $\sum_{(i,j) \in \Omega} A_{ij}^2 + \lambda \|\mathbf{Y}\|_F^2$. Thus, any \mathbf{X} such that $\gamma \|\mathbf{X}\|_* > \sum_{(i,j) \in \Omega} A_{ij}^2 + \lambda \|\mathbf{Y}\|_F^2$ cannot possibly be optimal to (6). Finally, since the nuclear norm is an upper bound on the

spectral norm of a matrix, we must have

$$\|\mathbf{X}\|_\sigma \leq \frac{\sum_{(i,j) \in \Omega} A_{ij}^2 + \lambda \|\mathbf{Y}\|_F^2}{\gamma},$$

for any matrix \mathbf{X} that is optimal to (6). We can therefore take $M = \frac{\sum_{(i,j) \in \Omega} A_{ij}^2 + \lambda \|\mathbf{Y}\|_F^2}{\gamma}$.

Notice that the non convexity in (6) is captured entirely by the bilinear constraint $(\mathbf{I}_n - \mathbf{P})\mathbf{X} = \mathbf{0}_{n \times m}$ and the quadratic constraint $\mathbf{P}^2 = \mathbf{P}$. In keeping with the approach presented in Bertsimas et al. (2023c,a), we leverage the matrix perspective to convexify the bilinear term and solve over the convex hull of the set \mathcal{P}_k . Recalling that the nuclear norm is semidefinite representable, we have the following formulation:

$$\begin{aligned} & \min_{\substack{\mathbf{P}, \mathbf{W}_1 \in \mathbb{R}^{n \times n}, \\ \mathbf{X} \in \mathbb{R}^{n \times m}, \mathbf{W}_2 \in \mathbb{R}^{m \times m}}} \sum_{(i,j) \in \Omega} (X_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T (\mathbf{I}_n - \mathbf{P}) \mathbf{Y}) + \frac{\gamma}{2} (\text{Tr}(\mathbf{W}_1) + \text{Tr}(\mathbf{W}_2)) \\ \text{s.t.} \quad & \mathbf{I}_n \succeq \mathbf{P} \succeq \mathbf{0}, \text{Tr}(\mathbf{P}) \leq k, \\ & \begin{pmatrix} M\mathbf{P} & \mathbf{X} \\ \mathbf{X}^T & M\mathbf{I}_m \end{pmatrix} \succeq \mathbf{0}, \begin{pmatrix} \mathbf{W}_1 & \mathbf{X} \\ \mathbf{X}^T & \mathbf{W}_2 \end{pmatrix} \succeq \mathbf{0}. \end{aligned} \tag{9}$$

We now have the following result:

Proposition 7 *Problem (9) is a valid convex relaxation of (8).*

Proof Problem (9) is clearly a convex optimization problem. We will show that the optimal value of (9) is a lower bound on the optimal value of (8) by showing that given any optimal solution to (8), we can construct a feasible solution to (9) that achieves the same objective value.

Consider any optimal solution $(\bar{\mathbf{U}}, \bar{\mathbf{V}}, \bar{\mathbf{P}})$ to (8). From the proof of Proposition 5, we know that the solution $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ where $\bar{\mathbf{X}} = \bar{\mathbf{U}}\bar{\mathbf{V}}^T$ is feasible to (6) (where we replace the constraint $\mathbf{P} \in \mathcal{P}_{\min(k, \text{rank}(\mathbf{UV}^T))}$ with $\mathbf{P} \in \mathcal{P}_k$) and must also be optimal. Let $\bar{\mathbf{X}} = \mathbf{L}\boldsymbol{\Sigma}\mathbf{R}^T$ be a singular value decomposition of $\bar{\mathbf{X}}$ with $\mathbf{L} \in \mathbb{R}^{n \times k}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{k \times k}$ and $\mathbf{R} \in \mathbb{R}^{m \times k}$. Let $\bar{\mathbf{W}}_1 = \mathbf{L}\boldsymbol{\Sigma}\mathbf{L}^T$ and $\bar{\mathbf{W}}_2 = \mathbf{R}\boldsymbol{\Sigma}\mathbf{R}^T$. We claim that $(\bar{\mathbf{X}}, \bar{\mathbf{P}}, \bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ is feasible to (9) and achieves the same objective value as $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ does in (6).

From the feasibility of $\bar{\mathbf{P}}$ in (8), we know that $\bar{\mathbf{P}} \in \mathcal{P}_k$ which implies $\mathbf{I}_n \succeq \bar{\mathbf{P}} \succeq \mathbf{0}$ and $\text{Tr}(\bar{\mathbf{P}}) \leq k$. By the generalized Schur complement lemma (see Boyd et al. (1994), Equation 2.41), we know that

$$\begin{pmatrix} M\bar{\mathbf{P}} & \bar{\mathbf{X}} \\ \bar{\mathbf{X}}^T & M\mathbf{I}_m \end{pmatrix} \succeq \mathbf{0} \iff M\mathbf{I}_m \succeq \mathbf{0}, \text{ and } M\mathbf{I}_m - \bar{\mathbf{X}}^T(M\bar{\mathbf{P}})^\dagger \bar{\mathbf{X}} \succeq \mathbf{0}.$$

We trivially have $M\mathbf{I}_m \succeq \mathbf{0}$. To see that the second condition holds, note that since $\bar{\mathbf{P}}$ is a projection matrix and $\bar{\mathbf{P}}\bar{\mathbf{X}} = \bar{\mathbf{X}}$, we have $\bar{\mathbf{X}}^T(M\bar{\mathbf{P}})^\dagger \bar{\mathbf{X}} = \frac{1}{M} \bar{\mathbf{X}}^T \bar{\mathbf{P}} \bar{\mathbf{X}} = \frac{1}{M} \bar{\mathbf{X}}^T \bar{\mathbf{X}}$. Furthermore, since $\bar{\mathbf{X}}$ is optimal to (6), we have $\|\bar{\mathbf{X}}\|_\sigma \leq M$ by assumption. Thus, we have

$$\|\bar{\mathbf{X}}\|_\sigma \leq M \implies \|\bar{\mathbf{X}}^T \bar{\mathbf{X}}\|_\sigma \leq M^2 \implies M^2 \mathbf{I}_m \succeq \bar{\mathbf{X}}^T \bar{\mathbf{X}} \implies M\mathbf{I}_m \succeq \frac{1}{M} \bar{\mathbf{X}}^T \bar{\mathbf{X}}.$$

Finally, observe that

$$\begin{pmatrix} \bar{\mathbf{W}}_1 & \bar{\mathbf{X}} \\ \bar{\mathbf{X}}^T & \bar{\mathbf{W}}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{L}\boldsymbol{\Sigma}\mathbf{L}^T & \mathbf{L}\boldsymbol{\Sigma}\mathbf{R}^T \\ \mathbf{R}\boldsymbol{\Sigma}\mathbf{L}^T & \mathbf{R}\boldsymbol{\Sigma}\mathbf{R}^T \end{pmatrix} = \begin{pmatrix} \mathbf{L} \\ \mathbf{R} \end{pmatrix} \boldsymbol{\Sigma} \begin{pmatrix} \mathbf{L} \\ \mathbf{R} \end{pmatrix}^T.$$

Since $\boldsymbol{\Sigma}$ is a diagonal matrix with non negative entries, the matrix $\begin{pmatrix} \bar{\mathbf{W}}_1 & \bar{\mathbf{X}} \\ \bar{\mathbf{X}}^T & \bar{\mathbf{W}}_2 \end{pmatrix}$ is certainly positive semidefinite. Thus we have shown that $(\bar{\mathbf{X}}, \bar{\mathbf{P}}, \bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ is indeed feasible to (9). To conclude the proof, we note that

$$\begin{aligned} \frac{\gamma}{2}(\text{Tr}(\bar{\mathbf{W}}_1) + \text{Tr}(\bar{\mathbf{W}}_2)) &= \frac{\gamma}{2}(\text{Tr}(\mathbf{L}\boldsymbol{\Sigma}\mathbf{L}^T) + \text{Tr}(\mathbf{R}\boldsymbol{\Sigma}\mathbf{R}^T)) = \frac{\gamma}{2}(\text{Tr}(\mathbf{L}^T\mathbf{L}\boldsymbol{\Sigma}) + \text{Tr}(\mathbf{R}^T\mathbf{R}\boldsymbol{\Sigma})) \\ &= \frac{\gamma}{2}(\text{Tr}(\boldsymbol{\Sigma}) + \text{Tr}(\boldsymbol{\Sigma})) = \gamma\|\bar{\mathbf{X}}\|_*, \end{aligned}$$

thus $(\bar{\mathbf{X}}, \bar{\mathbf{P}}, \bar{\mathbf{W}}_1, \bar{\mathbf{W}}_2)$ achieves the same objective value in (9) as $(\bar{\mathbf{X}}, \bar{\mathbf{P}})$ achieves in (6). ■

In general, an optimal solution to (9) will have $\mathbf{P} \notin \mathcal{P}_k$. We briefly note that to obtain a stronger convex relaxation, one could leverage eigenvector disjunctions (Bertsimas et al., 2023b; Saxena et al., 2010) to iteratively cut off solutions to (9) with $\mathbf{P} \notin \mathcal{P}_k$ and form increasingly tighter disjunctive approximations to the set \mathcal{P}_k .

5. Mixed-Projection ADMM

In this section, we present an alternating direction method of multipliers (ADMM) algorithm that is scalable and obtains high quality solutions for (8). Rather than forming the augmented lagrangian directly for (8), we first modify our problem formulation by introducing a dummy variable $\mathbf{Z} \in \mathbb{R}^{n \times k}$ that is an identical copy of \mathbf{U} . Additionally, rather than directly enforcing the constraint $\mathbf{P} \in \mathcal{P}_k$, we introduce an indicator function penalty $\mathbb{I}_{\mathcal{P}_k}(\mathbf{P})$ into the objective function where $\mathbb{I}_{\mathcal{X}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{X}$, otherwise $\mathbb{I}_{\mathcal{X}}(\mathbf{x}) = \infty$. Explicitly, we consider the following problem:

$$\begin{aligned} \min_{\substack{\mathbf{U}, \mathbf{Z} \in \mathbb{R}^{n \times k}, \\ \mathbf{V} \in \mathbb{R}^{m \times k}, \\ \mathbf{P} \in \mathbb{R}^{n \times n}}} \quad & \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Y}) + \frac{\gamma}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) + \mathbb{I}_{\mathcal{P}_k}(\mathbf{P}) \\ \text{s.t.} \quad & (\mathbf{I}_n - \mathbf{P})\mathbf{U} = \mathbf{0}_{n \times k}, \mathbf{U} - \mathbf{Z} = \mathbf{0}_{n \times k}. \end{aligned} \tag{10}$$

It is trivial to see that (10) is equivalent to (8). We will see in this section that working with formulation (10) leads to an ADMM algorithm with favorable decomposition properties. Introducing dual variables $\boldsymbol{\Phi}, \boldsymbol{\Psi} \in \mathbb{R}^{n \times k}$ for the constraints $(\mathbf{I}_n - \mathbf{P})\mathbf{U} = \mathbf{0}_{n \times k}$ and $\mathbf{U} - \mathbf{Z} = \mathbf{0}_{n \times k}$ respectively, the augmented lagrangian \mathcal{L} for (10) is given by:

$$\begin{aligned} \mathcal{L}(\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}, \boldsymbol{\Phi}, \boldsymbol{\Psi}) &= \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \lambda \text{Tr}(\mathbf{Y}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Y}) \\ &+ \frac{\gamma}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2) + \mathbb{I}_{\mathcal{P}_k}(\mathbf{P}) + \text{Tr}(\boldsymbol{\Phi}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Z}) \\ &+ \text{Tr}(\boldsymbol{\Psi}^T(\mathbf{Z} - \mathbf{U})) + \frac{\rho_1}{2}\|(\mathbf{I}_n - \mathbf{P})\mathbf{Z}\|_F^2 + \frac{\rho_2}{2}\|\mathbf{Z} - \mathbf{U}\|_F^2, \end{aligned} \tag{11}$$

where $\rho_1, \rho_2 > 0$ are non-negative penalty parameters. In what follows, we show that performing a partial minimization of the augmented lagrangian (11) over each of the primal variables $\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}$ yields a subproblem that can be solved efficiently. We present each subproblem and investigate the complexity of computing the corresponding subproblem solutions.

5.1 Subproblem in \mathbf{U}

First, suppose we fix variables $\mathbf{V}, \mathbf{P}, \mathbf{Z}, \Phi, \Psi$ and seek to minimize $\mathcal{L}(\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}, \Phi, \Psi)$ over \mathbf{U} . Eliminating terms that do not depend on \mathbf{U} , the resulting subproblem is given by

$$\min_{\mathbf{U} \in \mathbb{R}^{n \times k}} \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \frac{\gamma}{2} \|\mathbf{U}\|_F^2 - \text{Tr}(\Psi^T \mathbf{U}) + \frac{\rho_2}{2} \|\mathbf{Z} - \mathbf{U}\|_F^2. \quad (12)$$

We now have the following result:

Proposition 8 *The optimal solution $\bar{\mathbf{U}}$ for (12) is given by*

$$\bar{U}_{i,\star} = [2\mathbf{V}^T \mathbf{W}_i \mathbf{V} + (\gamma + \rho_2) \mathbf{I}_k]^{-1} [2\mathbf{V}^T \mathbf{W}_i A_{i,\star} + \Psi_{i,\star} + \rho_2 Z_{i,\star}], \quad (13)$$

for each $i \in \{1, \dots, n\}$ where each $\mathbf{W}_i \in \mathbb{R}^{m \times m}$ is a diagonal matrix satisfying $(\mathbf{W}_i)_{jj} = 1$ if $(i, j) \in \Omega$, otherwise $(\mathbf{W}_i)_{jj} = 0$. Here, the column vectors $\bar{U}_{i,\star} \in \mathbb{R}^k, A_{i,\star} \in \mathbb{R}^m, \Psi_{i,\star} \in \mathbb{R}^k, Z_{i,\star} \in \mathbb{R}^k$ denote the i^{th} row of the matrices $\bar{\mathbf{U}}, \mathbf{A}, \Psi, \mathbf{Z}$ respectively, where the unknown entries of \mathbf{A} are taken to be 0.

Proof Let $f(\mathbf{U})$ denote the objective function of (12). With $\{\mathbf{W}_i\}_{i=1}^n$ defined as in Proposition 8, observe that we can write $f(\mathbf{U})$ as

$$\begin{aligned} f(\mathbf{U}) &= \sum_{i=1}^n \|\mathbf{W}_i(\mathbf{V}\mathbf{U}_{i,\star} - A_{i,\star})\|_2^2 + \frac{\gamma}{2} \sum_{i=1}^n \|U_{i,\star}\|_2^2 - \sum_{i=1}^n \Psi_{i,\star}^T U_{i,\star} + \frac{\rho_2}{2} \sum_{i=1}^n \|Z_{i,\star} - U_{i,\star}\|_2^2 \\ &= \sum_{i=1}^n \left[\|\mathbf{W}_i(\mathbf{V}\mathbf{U}_{i,\star} - A_{i,\star})\|_2^2 + \frac{\gamma}{2} \|U_{i,\star}\|_2^2 - \Psi_{i,\star}^T U_{i,\star} + \frac{\rho_2}{2} \|Z_{i,\star} - U_{i,\star}\|_2^2 \right] = \sum_{i=1}^n g_i(\mathbf{U}), \end{aligned}$$

where we define $g_i(\mathbf{U}) = \|\mathbf{W}_i(\mathbf{V}\mathbf{U}_{i,\star} - A_{i,\star})\|_2^2 + \frac{\gamma}{2} \|U_{i,\star}\|_2^2 - \Psi_{i,\star}^T U_{i,\star} + \frac{\rho_2}{2} \|Z_{i,\star} - U_{i,\star}\|_2^2$. Thus, we have shown that $f(\mathbf{U})$ is separable over the rows of the matrix \mathbf{U} . Each function $g_i(\mathbf{U})$ is a (strongly) convex quadratic. Thus, we can minimize $g_i(\mathbf{U})$ by setting its gradient to 0. For any fixed row $i \in \{1, \dots, n\}$, we can differentiate and obtain

$$\nabla_{U_{i,\star}} g_i(\mathbf{U}) = 2\mathbf{V}^T \mathbf{W}_i(\mathbf{V}\mathbf{U}_{i,\star} - A_{i,\star}) + \gamma U_{i,\star} - \Psi_{i,\star} - \rho_2(Z_{i,\star} - U_{i,\star}).$$

By equating the gradient $\nabla_{U_{i,\star}} g_i(\mathbf{U})$ to 0 and rearranging, we obtain that the optimal vector $\bar{U}_{i,\star}$ is given by (13). This completes the proof. \blacksquare

Observe that since the matrix $\mathbf{V}^T \mathbf{W}_i \mathbf{V}$ is positive semidefinite and $\gamma + \rho_2 > 0$, the matrix inverse $[2\mathbf{V}^T \mathbf{W}_i \mathbf{V} + (\gamma + \rho_2) \mathbf{I}_k]^{-1}$ is well defined for all $i \in \{1, \dots, n\}$. Computing the optimal solution to (12) requires computing n different $k \times k$ matrix inverses (where in general $k \ll \min\{m, n\}$). Computing a single $k \times k$ matrix inverse requires $O(k^3)$ time

and forming the matrix product $\mathbf{V}^T \mathbf{W}_i \mathbf{V}$ requires $O(k^2 m)$ time for a given i . Thus, the complexity of computing the optimal solution for a single column is $O(k^3 + k^2 m)$. Notice that each column of $\bar{\mathbf{U}}$ can be computed independently of the other columns. We leverage this observation by developing a multi-threaded implementation of the algorithm presented in this section. Letting w denote the number of compute threads available, computing the optimal solution $\bar{\mathbf{U}}$ of (12) requires $O\left(\frac{k^3 n + k^2 m n}{\min\{w, n\}}\right)$ time (the term $\min\{w, n\}$ in the denominator reflects that fact that increasing the number of available compute threads beyond the number of columns of $\bar{\mathbf{U}}$ does not yield additional reduction in compute complexity).

5.2 Subproblem in \mathbf{V}

Now, suppose we fix variables $\mathbf{U}, \mathbf{P}, \mathbf{Z}, \Phi, \Psi$ and seek to minimize $\mathcal{L}(\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}, \Phi, \Psi)$ over \mathbf{V} . Eliminating terms that do not depend on \mathbf{V} , the resulting subproblem is given by

$$\min_{\mathbf{V} \in \mathbb{R}^{m \times k}} \sum_{(i,j) \in \Omega} ((\mathbf{U}\mathbf{V}^T)_{ij} - A_{ij})^2 + \frac{\gamma}{2} \|\mathbf{V}\|_F^2. \quad (14)$$

We now have the following result:

Proposition 9 *The optimal solution $\bar{\mathbf{V}}$ for (14) is given by*

$$\bar{V}_{j,\star} = [2\mathbf{U}^T \mathbf{W}_j \mathbf{U} + \gamma \mathbf{I}_k]^{-1} [2\mathbf{U}^T \mathbf{W}_j A_{\star,j}], \quad (15)$$

for each $j \in \{1, \dots, m\}$ where each $\mathbf{W}_j \in \mathbb{R}^{n \times n}$ is a diagonal matrix satisfying $(\mathbf{W}_j)_{ii} = 1$ if $(i, j) \in \Omega$, otherwise $(\mathbf{W}_j)_{ii} = 0$. Here, the column vector $\bar{V}_{j,\star} \in \mathbb{R}^k$ denotes the j^{th} row of $\bar{\mathbf{V}}$ while the column vector $A_{\star,j} \in \mathbb{R}^n$ denotes the j^{th} column of \mathbf{A} where the unknown entries of \mathbf{A} are taken to be 0.

Proof This proof follows the proof of Proposition 8. Let $f(\mathbf{V})$ denote the objective function of (14). With $\{\mathbf{W}_j\}_{j=1}^m$ defined as in Proposition 9, observe that we can write $f(\mathbf{V})$ as

$$\begin{aligned} f(\mathbf{V}) &= \sum_{j=1}^m \|\mathbf{W}_j(\mathbf{U}V_{j,\star} - A_{\star,j})\|_2^2 + \frac{\gamma}{2} \sum_{j=1}^m \|V_{j,\star}\|_2^2 \\ &= \sum_{j=1}^m \left[\|\mathbf{W}_j(\mathbf{U}V_{j,\star} - A_{\star,j})\|_2^2 + \frac{\gamma}{2} \|V_{j,\star}\|_2^2 \right] = \sum_{j=1}^m g_j(\mathbf{V}), \end{aligned}$$

where we define $g_j(\mathbf{V}) = \|\mathbf{W}_j(\mathbf{U}V_{j,\star} - A_{\star,j})\|_2^2 + \frac{\gamma}{2} \|V_{j,\star}\|_2^2$. Thus, we have shown that $f(\mathbf{V})$ is separable over the rows of the matrix \mathbf{V} . Each function $g_j(\mathbf{V})$ is a (strongly) convex quadratic. Thus, we can minimize $g_j(\mathbf{V})$ by setting its gradient to 0. For any fixed row $j \in \{1, \dots, m\}$, we can differentiate and obtain

$$\nabla_{\mathbf{V}_{j,\star}} g_j(\mathbf{V}) = 2\mathbf{U}^T \mathbf{W}_j(\mathbf{U}V_{j,\star} - A_{\star,j}) + \gamma V_{j,\star}.$$

By equating the gradient $\nabla_{\mathbf{V}_{j,\star}} g_j(\mathbf{V})$ to 0 and rearranging, we obtain that the optimal vector $\bar{U}_{i,\star}$ is given by (15). This completes the proof. \blacksquare

Observe that since the matrix $\mathbf{U}^T \mathbf{W}_j \mathbf{U}$ is positive semidefinite and $\gamma > 0$, the matrix inverse $[2\mathbf{U}^T \mathbf{W}_j \mathbf{U} + \gamma \mathbf{I}_k]^{-1}$ is well defined for all $j \in \{1, \dots, m\}$. Computing the optimal solution to (14) requires computing m different $k \times k$ matrix inverses. Forming the matrix product $\mathbf{U}^T \mathbf{W}_j \mathbf{U}$ requires $O(k^2 n)$ time for a given j . Thus, the complexity of computing the optimal solution for a single column is $O(k^3 + k^2 n)$. Notice that, similarly to the solution of (12), each column of $\bar{\mathbf{V}}$ can be computed independently of the other columns. As before, we leverage this observation in our multi-threaded implementation of the algorithm presented in this section. Letting w denote the number of compute threads available, computing the optimal solution $\bar{\mathbf{V}}$ of (14) requires $O\left(\frac{k^3 m + k^2 m n}{\min\{w, m\}}\right)$ time.

The optimal solution $\bar{\mathbf{V}}$ to (14) reveals that the frobenius norm regularization term on \mathbf{V} in (8) (which emerges from the nuclear norm regularization term on \mathbf{X} in (1)) has computational benefits. Indeed, if we had $\gamma = 0$, it is possible that the matrix $\mathbf{U}^T \mathbf{W}_j \mathbf{U}$ be singular at certain iterates of our ADMM algorithm, in which case the corresponding matrix inverse would be undefined. This observation is in keeping with several recent works in the statistics, machine learning and operations research literatures where the presence of a regularization penalty in the objective function yields improved out of sample performance as well as benefits in computational tractability (see for example Bertsimas et al. (2020, 2021, 2023c,a); Bertsimas and Johnson (2024)).

5.3 Subproblem in \mathbf{P}

Now, suppose we fix variables $\mathbf{U}, \mathbf{V}, \mathbf{Z}, \Phi, \Psi$ and seek to minimize $\mathcal{L}(\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}, \Phi, \Psi)$ over \mathbf{P} . Eliminating terms that do not depend on \mathbf{P} , the resulting subproblem is given by

$$\begin{aligned} \min_{\mathbf{P} \in \mathbb{S}_+^n} \quad & -\lambda \text{Tr}(\mathbf{Y}^T \mathbf{P} \mathbf{Y}) - \text{Tr}(\Phi^T \mathbf{P} \mathbf{Z}) + \frac{\rho_1}{2} \|(\mathbf{I}_n - \mathbf{P}) \mathbf{Z}\|_F^2 \\ \text{s.t.} \quad & \mathbf{P}^2 = \mathbf{P}, \text{Tr}(\mathbf{P}) \leq k. \end{aligned} \tag{16}$$

We now have the following result:

Proposition 10 *Let $\mathbf{M} \Sigma \mathbf{M}^T$ be a rank k truncated singular value decomposition for the matrix given by:*

$$\left(\lambda \mathbf{Y} \mathbf{Y}^T + \frac{\rho_1}{2} \mathbf{Z} \mathbf{Z}^T + \frac{1}{2} (\Phi \mathbf{Z}^T + \mathbf{Z} \Phi^T) \right),$$

where $\Sigma \in \mathbb{R}^{k \times k}$, $\mathbf{M} \in \mathbb{R}^{n \times k}$, $\mathbf{M}^T \mathbf{M} = \mathbf{I}_k$. The optimal solution $\bar{\mathbf{P}}$ for (16) is given by $\bar{\mathbf{P}} = \mathbf{M} \mathbf{M}^T$.

Proof Let $f(\mathbf{P})$ denote the objective function of (16). Observe that for any \mathbf{P} that is feasible to (16), we can write $f(\mathbf{P})$ as:

$$\begin{aligned} f(\mathbf{P}) &= -\lambda \text{Tr}(\mathbf{Y}^T \mathbf{P} \mathbf{Y}) - \text{Tr}(\Phi^T \mathbf{P} \mathbf{Z}) + \frac{\rho_1}{2} \|(\mathbf{I}_n - \mathbf{P}) \mathbf{Z}\|_F^2 \\ &= -\lambda \text{Tr}(\mathbf{Y} \mathbf{Y}^T \mathbf{P}) - \text{Tr}(\mathbf{Z} \Phi^T \mathbf{P}) + \frac{\rho_1}{2} \text{Tr}(\mathbf{Z} \mathbf{Z}^T (\mathbf{I}_n - \mathbf{P})) \\ &= \frac{\rho_1}{2} \text{Tr}(\mathbf{Z} \mathbf{Z}^T) - \langle \lambda \mathbf{Y} \mathbf{Y}^T + \mathbf{Z} \Phi^T + \frac{\rho_1}{2} \mathbf{Z} \mathbf{Z}^T, \mathbf{P} \rangle. \end{aligned}$$

Thus, it is immediately clear that a solution will be optimal to (16) if and only if it is optimal to the problem given by:

$$\begin{aligned} & \max_{\mathbf{P} \in \mathbb{S}_+^n} \langle \mathbf{C}, \mathbf{P} \rangle \\ & \text{s.t.} \quad \mathbf{P}^2 = \mathbf{P}, \text{Tr}(\mathbf{P}) \leq k, \end{aligned} \quad (17)$$

where we define the matrix $\mathbf{C} \in \mathbb{R}^{n \times n}$ as $\mathbf{C} = \lambda \mathbf{Y} \mathbf{Y}^T + \mathbf{Z} \Phi^T + \frac{\rho_1}{2} \mathbf{Z} \mathbf{Z}^T$. Let $\bar{\mathbf{C}} = \frac{1}{2}(\mathbf{C} + \mathbf{C}^T)$ denote the symmetric part of \mathbf{C} . Observe that for any symmetric matrix \mathbf{P} , we can consider $\langle \bar{\mathbf{C}}, \mathbf{P} \rangle$ in place of $\langle \mathbf{C}, \mathbf{P} \rangle$ since we have

$$\begin{aligned} \langle \mathbf{C}, \mathbf{P} \rangle &= \sum_{i=1}^n \sum_{j=1}^n P_{ij} C_{ij} = \sum_{i=1}^n P_{ii} C_{ii} + \sum_{i=1}^{n-1} \sum_{j=i+1}^n P_{ij} (C_{ij} + C_{ji}) \\ &= \sum_{i=1}^n P_{ii} \bar{C}_{ii} + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n P_{ij} \bar{C}_{ij} = \sum_{i=1}^n \sum_{j=1}^n P_{ij} \bar{C}_{ij} = \langle \bar{\mathbf{C}}, \mathbf{P} \rangle. \end{aligned}$$

Let $\bar{\mathbf{C}} = \mathbf{M} \Sigma \mathbf{M}^T$ be a full singular value decomposition of $\bar{\mathbf{C}}$ with $\mathbf{M}, \Sigma \in \mathbb{R}^{n \times n}$, $\mathbf{M}^T \mathbf{M} = \mathbf{M} \mathbf{M}^T = \mathbf{I}_n$. The matrix Σ is the diagonal matrix of (ordered) singular values of $\bar{\mathbf{C}}$ and we let σ_i denote the i^{th} singular value. Any feasible matrix \mathbf{P} to (17) can be written as $\mathbf{P} = \mathbf{L} \mathbf{L}^T$ where $\mathbf{L} \in \mathbb{R}^{n \times k}$, $\mathbf{L}^T \mathbf{L} = \mathbf{I}_k$. Thus, for any \mathbf{P} feasible to (17) we express the objective value as:

$$\langle \bar{\mathbf{C}}, \mathbf{P} \rangle = \text{Tr}(\mathbf{M} \Sigma \mathbf{M}^T \mathbf{L} \mathbf{L}^T) = \text{Tr}(\Sigma (\mathbf{M}^T \mathbf{L} \mathbf{L}^T \mathbf{M})) = \sum_{i=1}^n \sigma_i \|(\mathbf{M}^T \mathbf{L})_{i,\star}\|_2^2.$$

Let $\mathbf{N} = \mathbf{M}^T \mathbf{L} \in \mathbb{R}^{n \times k}$. Note that we have $\mathbf{N}^T \mathbf{N} = \mathbf{L}^T \mathbf{M} \mathbf{M}^T \mathbf{L} = \mathbf{L}^T \mathbf{L} = \mathbf{I}_k$, which implies that the columns of \mathbf{N} are orthonormal. This immediately implies that we have $N_{i,\star}^T N_{i,\star} = \|(\mathbf{M}^T \mathbf{L})_{i,\star}\|_2^2 \leq 1$. Moreover, we have

$$\sum_{i=1}^n \|(\mathbf{M}^T \mathbf{L})_{i,\star}\|_2^2 = \sum_{i=1}^n N_{i,\star}^T N_{i,\star} = \sum_{i=1}^n \sum_{j=1}^k N_{ij}^2 = \sum_{j=1}^k \sum_{i=1}^n N_{ij}^2 = \sum_{j=1}^k 1 = k.$$

We can therefore upper bound the optimal objective value of (17) as

$$\langle \bar{\mathbf{C}}, \mathbf{P} \rangle = \sum_{i=1}^n \sigma_i \|(\mathbf{M}^T \mathbf{L})_{i,\star}\|_2^2 \leq \sum_{i=1}^k \sigma_i.$$

To conclude the proof, notice that by taking $\bar{\mathbf{P}} = \bar{\mathbf{M}} \bar{\mathbf{M}}^T$ where $\bar{\mathbf{M}} \in \mathbb{R}^{n \times k}$ is the matrix that consists of the first k columns of \mathbf{M} we can achieve the upper bound on (17):

$$\langle \bar{\mathbf{C}}, \bar{\mathbf{P}} \rangle = \text{Tr}(\mathbf{M} \Sigma \mathbf{M}^T \bar{\mathbf{M}} \bar{\mathbf{M}}^T) = \text{Tr}(\bar{\mathbf{M}}^T \mathbf{M} \Sigma \mathbf{M}^T \bar{\mathbf{M}}) = \sum_{i=1}^k \sigma_i. \quad \blacksquare$$

To compute the optimal solution of (16), we need to compute a rank k singular value decomposition of the matrix $\bar{\mathbf{C}} = (\lambda\mathbf{Y}\mathbf{Y}^T + \frac{\rho_1}{2}\mathbf{Z}\mathbf{Z}^T + \frac{1}{2}(\mathbf{\Phi}\mathbf{Z}^T + \mathbf{Z}\mathbf{\Phi}^T))$ which requires $O(kn^2)$ time since $\bar{\mathbf{C}} \in \mathbb{R}^{n \times n}$. Moreover, explicitly forming the matrix $\bar{\mathbf{C}}$ in memory from its constituent matrices $\mathbf{Y}, \mathbf{Z}, \mathbf{\Phi}$ requires $O(n^2(d+k))$ operations. Thus, naively computing the optimal solution to (16) has complexity $O(n^2(d+k))$ where the bottleneck operation from a complexity standpoint is explicitly forming the matrix $\bar{\mathbf{C}}$.

Fortunately, it is possible to compute the optimal solution to (16) more efficiently. Observe that we can equivalently express the matrix $\bar{\mathbf{C}}$ as $\bar{\mathbf{C}} = \mathbf{F}_1\mathbf{F}_2^T$ where $\mathbf{F}_1, \mathbf{F}_2 \in \mathbb{R}^{n \times (d+3k)}$ are defined as

$$\begin{aligned}\mathbf{F}_1 &= \left(\sqrt{\lambda}\mathbf{Y} \quad \sqrt{\frac{\rho_1}{2}}\mathbf{Z} \quad \sqrt{\frac{1}{2}}\mathbf{\Phi} \quad \sqrt{\frac{1}{2}}\mathbf{Z} \right), \\ \mathbf{F}_2 &= \left(\sqrt{\lambda}\mathbf{Y} \quad \sqrt{\frac{\rho_1}{2}}\mathbf{Z} \quad \sqrt{\frac{1}{2}}\mathbf{Z} \quad \sqrt{\frac{1}{2}}\mathbf{\Phi} \right).\end{aligned}$$

Computing a truncated singular value decomposition requires only computing repeated matrix vector products. Therefore, rather than explicitly forming the matrix $\bar{\mathbf{C}}$ in memory at a cost of $O(n^2(d+k))$ operations, in our implementation we design a custom matrix class where matrix vector products between $\bar{\mathbf{C}}$ and arbitrary vectors $\mathbf{x} \in \mathbb{R}^n$ are computed by first evaluating the matrix vector product $\boldsymbol{\nu} = \mathbf{F}_2^T\mathbf{x}$ and subsequently evaluating the matrix vector product $\bar{\mathbf{C}}\mathbf{x} = \mathbf{F}_1\boldsymbol{\nu}$. In so doing, we can evaluate matrix vector products $\bar{\mathbf{C}}\mathbf{x}$ in $O(n(d+k))$ time rather than $O(n^2)$ time (in general, we will have $d+k \ll n$). Computing a truncated singular value decomposition of $\bar{\mathbf{C}}$ with this methodology of evaluating matrix vector products requires only $O(k^2n + knd)$ operations. Thus, our custom matrix class implementation avoids needing to explicitly form $\bar{\mathbf{C}}$ in memory and allows the optimal solution to (16) to be computed in $O(k^2n + knd)$ time.

5.4 Subproblem in \mathbf{Z}

Now, suppose we fix variables $\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{\Phi}, \mathbf{\Psi}$ and seek to minimize $\mathcal{L}(\mathbf{U}, \mathbf{V}, \mathbf{P}, \mathbf{Z}, \mathbf{\Phi}, \mathbf{\Psi})$ over \mathbf{Z} . Eliminating terms that do not depend on \mathbf{Z} , the resulting subproblem is given by

$$\min_{\mathbf{Z} \in \mathbb{R}^{n \times k}} \text{Tr}(\mathbf{\Phi}^T(\mathbf{I}_n - \mathbf{P})\mathbf{Z}) + \text{Tr}(\mathbf{\Psi}^T\mathbf{Z}) + \frac{\rho_1}{2}\|(\mathbf{I}_n - \mathbf{P})\mathbf{Z}\|_F^2 + \frac{\rho_2}{2}\|\mathbf{Z} - \mathbf{U}\|_F^2. \quad (18)$$

We now have the following result:

Proposition 11 *The optimal solution $\bar{\mathbf{Z}}$ for (18) is given by*

$$\begin{aligned}\bar{\mathbf{Z}} &= \frac{1}{\rho_1 + \rho_2} \left(\mathbf{I}_n + \frac{\rho_1}{\rho_2}\mathbf{P} \right) \left(\rho_2\mathbf{U} - (\mathbf{I}_n - \mathbf{P})\mathbf{\Phi} - \mathbf{\Psi} \right) \\ &= \frac{1}{\rho_1 + \rho_2} \left(\rho_2\mathbf{U} - \mathbf{\Phi} + \mathbf{P}\mathbf{\Phi} - \mathbf{\Psi} + \rho_1\mathbf{P}\mathbf{U} - \frac{\rho_1}{\rho_2}\mathbf{P}\mathbf{\Psi} \right).\end{aligned} \quad (19)$$

Proof Let $f(\mathbf{Z})$ denote the objective function of (18). The function $f(\mathbf{Z})$ is a convex quadratic, thus it can be minimized by setting its gradient to 0. Differentiating $f(\mathbf{Z})$, we obtain:

$$\nabla_{\mathbf{Z}} f(\mathbf{Z}) = (\mathbf{I}_n - \mathbf{P})^T\mathbf{\Phi} + \mathbf{\Psi} + \rho_1(\mathbf{I}_n - \mathbf{P})^T(\mathbf{I}_n - \mathbf{P})\mathbf{Z} + \rho_2(\mathbf{Z} - \mathbf{U}).$$

Moreover, for any matrix \mathbf{P} for which the augmented lagrangian (11) takes finite value, we will have $\mathbf{P} \in \mathcal{P}_k$ which implies that $\mathbf{P}^T = \mathbf{P}$ and $\mathbf{P}^2 = \mathbf{P}$. We can therefore simplify $\nabla_{\mathbf{Z}} f(\mathbf{Z})$ as:

$$\nabla_{\mathbf{Z}} f(\mathbf{Z}) = (\mathbf{I}_n - \mathbf{P})\Phi + \Psi + \rho_1(\mathbf{I}_n - \mathbf{P})\mathbf{Z} + \rho_2(\mathbf{Z} - \mathbf{U}).$$

By equating the gradient $\nabla_{\mathbf{Z}} f(\mathbf{Z})$ to 0 and rearranging, we obtain that the optimal matrix $\bar{\mathbf{Z}}$ is given by:

$$\bar{\mathbf{Z}} = \left(\rho_1(\mathbf{I}_n - \mathbf{P}) + \rho_2\mathbf{I}_n \right)^{-1} \left(\rho_2\mathbf{U} - (\mathbf{I}_n - \mathbf{P})\Phi - \Psi \right)$$

To conclude the proof, it remains to show that $\left(\rho_1(\mathbf{I}_n - \mathbf{P}) + \rho_2\mathbf{I}_n \right)^{-1} = \frac{1}{\rho_1 + \rho_2} \left(\mathbf{I}_n + \frac{\rho_1}{\rho_2}\mathbf{P} \right)$. Let $\mathbf{P} = \mathbf{M}\mathbf{M}^T$ where $\mathbf{M} \in \mathbb{R}^{n \times k}$, $\mathbf{M}^T\mathbf{M} = \mathbf{I}_k$. Such a matrix \mathbf{M} is guaranteed to exist for any $\mathbf{P} \in \mathcal{P}_k$. We have

$$\begin{aligned} \rho_1(\mathbf{I}_n - \mathbf{P}) + \rho_2\mathbf{I}_n &= \rho_1(\mathbf{I}_n - \mathbf{M}\mathbf{M}^T) + \rho_2\mathbf{I}_n \\ &= (\rho_1 + \rho_2)\mathbf{I}_n + \mathbf{M}(-\rho_1\mathbf{I}_n)\mathbf{M}^T \\ &= \frac{1}{\rho_1 + \rho_2}\mathbf{I}_n - \frac{1}{(\rho_1 + \rho_2)^2}\mathbf{M} \left(\frac{1}{\rho_1 + \rho_2}\mathbf{M}^T\mathbf{M} - \frac{1}{\rho_1}\mathbf{I}_k \right)^{-1} \mathbf{M}^T \\ &= \frac{1}{\rho_1 + \rho_2}\mathbf{I}_n - \frac{1}{(\rho_1 + \rho_2)^2}\mathbf{M} \left(\frac{-\rho_1(\rho_1 + \rho_2)}{\rho_2}\mathbf{I}_k \right) \mathbf{M}^T \\ &= \frac{1}{\rho_1 + \rho_2} \left(\mathbf{I}_n + \frac{\rho_1}{\rho_2}\mathbf{P} \right), \end{aligned}$$

where the third equality follows from the Woodbury matrix inversion lemma (see Petersen et al. (2008), Section 3.2.2). As a sanity check, one can verify that the product of $\left(\rho_1(\mathbf{I}_n - \mathbf{P}) + \rho_2\mathbf{I}_n \right)$ and $\frac{1}{\rho_1 + \rho_2} \left(\mathbf{I}_n + \frac{\rho_1}{\rho_2}\mathbf{P} \right)$ is indeed the n dimensional identity matrix. ■

Evaluating the optimal solution to (18) requires only matrix-matrix multiplications. Computing the products of $\mathbf{P}\Phi$, $\mathbf{P}\mathbf{U}$, $\mathbf{P}\Psi$ in the definition of $\bar{\mathbf{Z}}$ from (19) requires $O(kn^2)$ operations. Thus, the naive cost of forming $\bar{\mathbf{Z}}$ is $O(kn^2)$. However, notice that if we had a factored representation of the matrix \mathbf{P} as $\mathbf{P} = \mathbf{M}\mathbf{M}^T$ with $\mathbf{M} \in \mathbb{R}^{n \times k}$, for any matrix $\mathbf{R} \in \mathbb{R}^{n \times k}$ we could compute matrix-matrix products $\mathbf{P}\mathbf{R}$ by first computing $\mathbf{S} = \mathbf{M}^T\mathbf{R}$ and thereafter computing $\mathbf{P}\mathbf{R} = \mathbf{M}\mathbf{S}$ for a total complexity of $O(k^2n)$. One might object that this ignores the time required to compute such a matrix \mathbf{M} . However, observe that in computing a matrix \mathbf{P} that is optimal to (16), we in fact must already generate such a matrix \mathbf{M} (see proposition 10). In fact, in our implementation we never explicitly form a $n \times n$ matrix \mathbf{P} as it suffices to only store a copy of its low rank factorization matrix \mathbf{M} . Thus, the optimal solution to (18) can be evaluated in $O(k^2n)$ time.

5.5 An ADMM Algorithm

Having illustrated that the partial minimization of the lagrangian (11) across each of the primal variables (Problems (12), (14), (16), (18)) can be solved efficiently, we can now present the overall approach Algorithm 1.

We initialize primal iterates $\mathbf{U}_0 = \mathbf{Z}_0 = \mathbf{L}\Sigma^{\frac{1}{2}}$, $\mathbf{P}_0 = \mathbf{L}\mathbf{L}^T$, $\mathbf{V}_0 = \mathbf{R}\Sigma^{\frac{1}{2}}$ where $\mathbf{L}\Sigma\mathbf{R}$ denotes a rank k truncated singular value decomposition of \mathbf{A} (the missing entries of \mathbf{A}

Algorithm 1: Mixed-Projection ADMM

Data: $n, m, k \in \mathbb{Z}^+, \Omega \subset [n] \times [m], \{A_{ij}\}_{(i,j) \in \Omega}, \lambda, \gamma \in \mathbb{R}^+$. Tolerance parameter $\epsilon > 0$. Maximum iteration parameter $T \in \mathbb{Z}^+$

Result: $(\bar{U}, \bar{V}, \bar{P})$ that is feasible to (8).

$(\mathbf{U}_0, \mathbf{P}_0, \mathbf{V}_0, \mathbf{Z}_0) \leftarrow (\mathbf{L}\Sigma^{\frac{1}{2}}, \mathbf{L}\mathbf{L}^T, \mathbf{R}\Sigma^{\frac{1}{2}}, \mathbf{L}\Sigma^{\frac{1}{2}})$ where $\mathbf{L}\Sigma\mathbf{R}$ is a rank k truncated SVD of \mathbf{A} and missing entries are filled in with 0;

$(\Phi_0, \Psi_0) \leftarrow (\mathbf{1}_{n \times k}, \mathbf{1}_{n \times k});$

$t \leftarrow 0;$

while $t < T$ and $\max\{\|(\mathbf{I}_n - \mathbf{P}_t)\mathbf{Z}_t\|_F^2, \|\mathbf{Z}_t - \mathbf{U}_t\|_F^2\} > \epsilon$ **do**

$(\mathbf{U}_{t+1}, \mathbf{P}_{t+1}) \leftarrow \operatorname{argmin}_{\mathbf{U}, \mathbf{P}} \mathcal{L}(\mathbf{U}, \mathbf{V}_t, \mathbf{P}, \mathbf{Z}_t, \Phi_t, \Psi_t);$

$(\mathbf{V}_{t+1}, \mathbf{Z}_{t+1}) \leftarrow \operatorname{argmin}_{\mathbf{V}, \mathbf{Z}} \mathcal{L}(\mathbf{U}_{t+1}, \mathbf{V}, \mathbf{P}_{t+1}, \mathbf{Z}, \Phi_t, \Psi_t);$

$\Phi_{t+1} \leftarrow \Phi_t + \rho_1(\mathbf{I} - \mathbf{P}_{t+1})\mathbf{Z}_{t+1};$

$\Psi_{t+1} \leftarrow \Psi_t + \rho_2(\mathbf{Z}_{t+1} - \mathbf{U}_{t+1});$

$t \leftarrow t + 1;$

end

return $(\mathbf{U}_t, \mathbf{V}_t, \mathbf{P}_t)$

are filled in with 0s) and we initialize dual iterates $\Phi_0 = \Psi_0 = \mathbf{1}_{n \times k}$. Observe that the subproblems (12) and (16) can be solved simultaneously. Similarly, the subproblems (14) and (18) can be solved simultaneously. At each iteration of Algorithm 1, we first update the iterates $\mathbf{U}_{t+1}, \mathbf{P}_{t+1}$ by solving problems (12) and (16) with $(\mathbf{V}_t, \mathbf{Z}_t, \Phi_t, \Psi_t)$ fixed. Next, we update the iterates $\mathbf{V}_{t+1}, \mathbf{Z}_{t+1}$ by solving problems (14) and (18) with $(\mathbf{U}_{t+1}, \mathbf{P}_{t+1}, \Phi_t, \Psi_t)$ fixed. Finally, we update the dual iterates Φ, Ψ by taking a gradient ascent step. The gradients of the augmented lagrangian (11) with respect to Φ and Ψ are given by the primal residuals $(\mathbf{I}_n - \mathbf{P}_{t+1})\mathbf{Z}_{t+1}$ and $\mathbf{Z}_{t+1} - \mathbf{U}_{t+1}$ respectively. We use ρ_1 and ρ_2 respectively as the step size. We proceed until the squared norm of each primal residual is below a numerical tolerance parameter ϵ or until we reach an input maximum number of iterations T . We know have the following result:

Proposition 12 *Assume that the number of compute threads w is less than $\min\{n, m\}$. The per iteration complexity of Algorithm 1 is $O(k^2n + knd + \frac{k^3(n+m)+k^2nm}{w})$.*

Proof The result follows from the complexity analysis of problems (12), (14), (16) and (18). ■

6. Computational Results

We evaluate the performance of Algorithm 1 implemented in Julia 1.7.3. Throughout, we fix $\rho_1 = \rho_2 = 10$, set the maximum number of iterations $T = 20$ and set the number of compute threads $w = 24$. Note that given the novelty of Problem (1), there are no pre-existing specialized methods to benchmark against. Accordingly, we compare the performance of Algorithm 1 against well studied methods for the very closely related MC problem as well as a highly performant generic method for low rank matrix optimization

problems. The MC methods we consider are Fast-Impute (Bertsimas and Li, 2020), Soft-Impute (Mazumder et al., 2010) and Iterative-SVD (Troyanskaya et al., 2001) which we introduced formally in Section 2.1. We utilize the implementation of Fast-Impute made publicly available by (Bertsimas and Li, 2020) while we use the implementation of Soft-Impute and Iterative-SVD from the python package fancyimpute 0.7.0 (Rubinsteyn and Feldman, 2016). The matrix optimization method we consider is ScaledGD (scaled gradient descent) (Tong et al., 2021) which we introduced formally in Section 2.2.1 and implement ourselves. All experiments were performed using synthetic data on MIT’s Supercloud Cluster (Reuther et al., 2018), which hosts Intel Xeon Platinum 8260 processors. To bridge the gap between theory and practice, we have made our code freely available on GitHub at github.com/NicholasJohnson2020/LearningLowRankMatrices.

To evaluate the performance of Algorithm 1, Fast-Impute, Soft-Impute, Iterative-SVD and ScaledGD on synthetic data, we consider the objective value achieved by a returned solution in (1), the ℓ_2 reconstruction error between a returned solution and the ground truth, the numerical rank of a returned solution and the execution time of each algorithm. Explicitly, let $\hat{\mathbf{X}} \in \mathbb{R}^{n \times m}$ denote the solution returned by a given method (where we define $\hat{\mathbf{X}} = \hat{\mathbf{U}}\hat{\mathbf{V}}^T$ if the method outputs low rank factors $\hat{\mathbf{U}}, \hat{\mathbf{V}}$) and let $\mathbf{A}^{true} \in \mathbb{R}^{n \times m}$ denote the ground truth matrix. We define the the ℓ_2 reconstruction error of $\hat{\mathbf{X}}$ as

$$ERR_{\ell_2}(\hat{\mathbf{X}}) = \frac{\|\hat{\mathbf{X}} - \mathbf{A}^{true}\|_F^2}{\|\mathbf{A}^{true}\|_F^2}.$$

We compute the numerical rank of $\hat{\mathbf{X}}$ by calling the default rank function from the Julia LinearAlgebra package. We aim to answer the following questions:

1. How does the performance of Algorithm 1 compare to existing methods such as Fast-Impute, Soft-Impute, Iterative-SVD and ScaledGD on synthetic data?
2. How is the performance of Algorithm 1 affected by the number of rows n , the number of columns m , the dimension of the side information d and the underlying rank k of the ground truth?
3. Empirically, which subproblem solution update is the computational bottleneck of Algorithm 1?

6.1 Synthetic Data Generation

To generate synthetic data, we specify a number of rows $n \in \mathbb{Z}_+$, a number of columns $m \in \mathbb{Z}_+$, a desired rank $k \in \mathbb{Z}_+$ with $k < \min\{n, m\}$, the dimension of the side information $d \in \mathbb{Z}_+$, a fraction of missing values $\alpha \in (0, 1)$ and a noise parameter $\sigma \in \mathbb{R}_+$ that controls the signal to noise ratio. We sample matrices $\mathbf{U} \in \mathbb{R}^{n \times k}$, $\mathbf{V} \in \mathbb{R}^{m \times k}$, $\boldsymbol{\beta} \in \mathbb{R}^{m \times d}$ by drawing each entry $U_{ij}, V_{ij}, \beta_{ij}$ independently from the uniform distribution on the interval $[0, 1]$. Furthermore, we sample a noise matrix $\mathbf{N} \in \mathbb{R}^{n \times d}$ by drawing each entry N_{ij} independently from the univariate normal distribution with mean 0 and variance σ^2 . We let $\mathbf{A} = \mathbf{UV}^T$ and we let $\mathbf{Y} = \mathbf{A}\boldsymbol{\beta} + \mathbf{N}$. Lastly, we sample $\lfloor \alpha \cdot n \cdot m \rfloor$ indices uniformly at random from the collection $\mathcal{I} = \{(i, j) : 1 \leq i \leq n, 1 \leq j \leq m\}$ to be the set of missing indices, which we denote by Γ . The set of revealed entries can then be defined as $\Omega = \mathcal{I} \setminus \Gamma$. We

fix $\alpha = 0.9, \sigma = 2$ throughout our experiments and report numerical results for various different combinations of (n, m, d, k) .

6.2 Sensitivity to Row Dimension

We present a comparison of Algorithm 1 with ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD as we vary the number of rows n . In these experiments, we fixed $m = 100, k = 5$, and $d = 150$ across all trials. We varied $n \in \{100, 200, 400, 800, 1000, 2000, 5000, 10000\}$ and we performed 20 trials for each value of n . For ScaledGD, we set the step size to be $\eta = \frac{1}{10\sigma_1(\mathbf{A})}$ where $\sigma_1(\mathbf{A})$ denotes the largest singular value of the input matrix \mathbf{A} where we fill the unobserved entries with the value 0. Letting $f(\mathbf{U}_t, \mathbf{V}_t)$ denote the objective value achieved after iteration t of ScaledGD, we terminate ScaledGD when either $t > 1000$ or $\frac{f(\mathbf{U}_{t-1}, \mathbf{V}_{t-1}) - f(\mathbf{U}_t, \mathbf{V}_t)}{f(\mathbf{U}_{t-1}, \mathbf{V}_{t-1})} < 10^{-3}$. In words, we terminate ScaledGD after 1000 iterations or after the relative objective value improvement between two iterations is less than 0.1%.

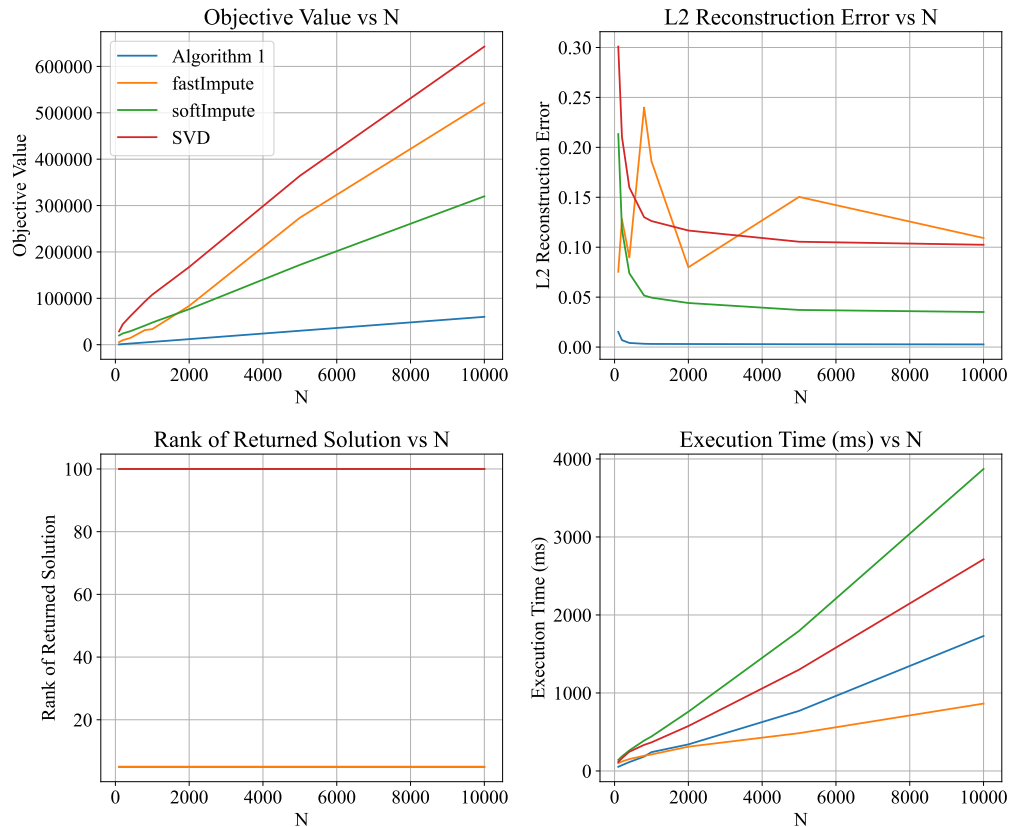


Figure 1: Objective value (top left), ℓ_2 reconstruction error (top right), fitted rank (bottom left) and execution time (bottom right) versus n with $m = 100, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

We report the objective value, ℓ_2 reconstruction error, fitted rank and execution time for Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Figure 1. We additionally report the objective value, reconstruction error and execution time for ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Tables 1, 2 and 3 of Appendix A. In Figure 2, we plot the average cumulative time spent solving subproblems (12), (14), (16), (18) during the execution of Algorithm 1 versus n . Our main findings from this set of experiments are:

1. Algorithm 1 systematically produces higher quality solutions than ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD (see Table 1), sometimes achieving an objective value that is an order of magnitude superior than the next best method. On average, Algorithm 1 outputs a solution whose objective value is 86% lesser than the objective value achieved by the best performing alternative method (Fast-Impute). We remind the reader that Fast-Impute, Soft-Impute and Iterative-SVD are methods designed for the generic MC problem and are not custom built to solve (1) so it should not come as a surprise that Algorithm 1 significantly outperforms these 3 methods in terms of objective value. ScaledGD however has explicit knowledge of the objective function of (1) along with its gradient, yet surprisingly produces the weakest average objective value across these experiments. We note that we use the default hyperparameters for ScaledGD recommended by the authors of this method (Tong et al., 2021). We observe that the objective value achieved by all methods increases linearly as the number of rows n increases.
2. In terms of ℓ_2 reconstruction error, Algorithm 1 again systematically produces solutions that are of higher quality than ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD (see Table 2), often achieving an error that is an order of magnitude superior than the next best method. On average, Algorithm 1 outputs a solution whose ℓ_2 reconstruction error is 92% lesser than the reconstruction error achieved by the best performing alternative method (Soft-Impute in all but one parameter configuration). This is especially noteworthy since Algorithm 1 is not designed explicitly with reconstruction error minimization as the objective, unlike Fast-Impute and Soft-Impute, and suggests that the side information \mathbf{Y} is instrumental in recovering high quality low rank estimates of the partially observed data matrix.
3. We observe that the fitted rank of the solutions returned by Algorithm 1, ScaledGD and Fast-Impute always matched the specified target rank as would be expected, but surprisingly the solutions returned by Soft-Impute and Iterative-SVD were always of full rank despite the fact that these methods were provided with the target rank explicitly. This is potentially due to a numerical issues in the computation of the rank due to presence of extremely small singular values.
4. The runtime of Algorithm 1 is competitive with that of the other methods. The runtime of Algorithm 1 is less than that of Soft-Impute and Iterative-SVD but greater than that of Fast-Impute. For experiments with $n \leq 2000$, Table 3 illustrates that ScaledGD was the method with the fastest execution time (however as previously mentioned the returned solutions were of low quality). The runtime of Algorithm 1, Fast-Impute, Soft-Impute and iterate SVD appear to grow linearly with n .

5. Figure 2 illustrates that the computation of the solution for (12) is the computational bottleneck in the execution of Algorithm 1 in this set of experiments, followed next by the computation of the solution for (16). Empirically, we observe that the solution time of (12), (14), (16) and (18) appear to scale linearly with the number of rows n . This observation is consistent with the computational complexities derived for each subproblem of Algorithm 1 in Section 5.

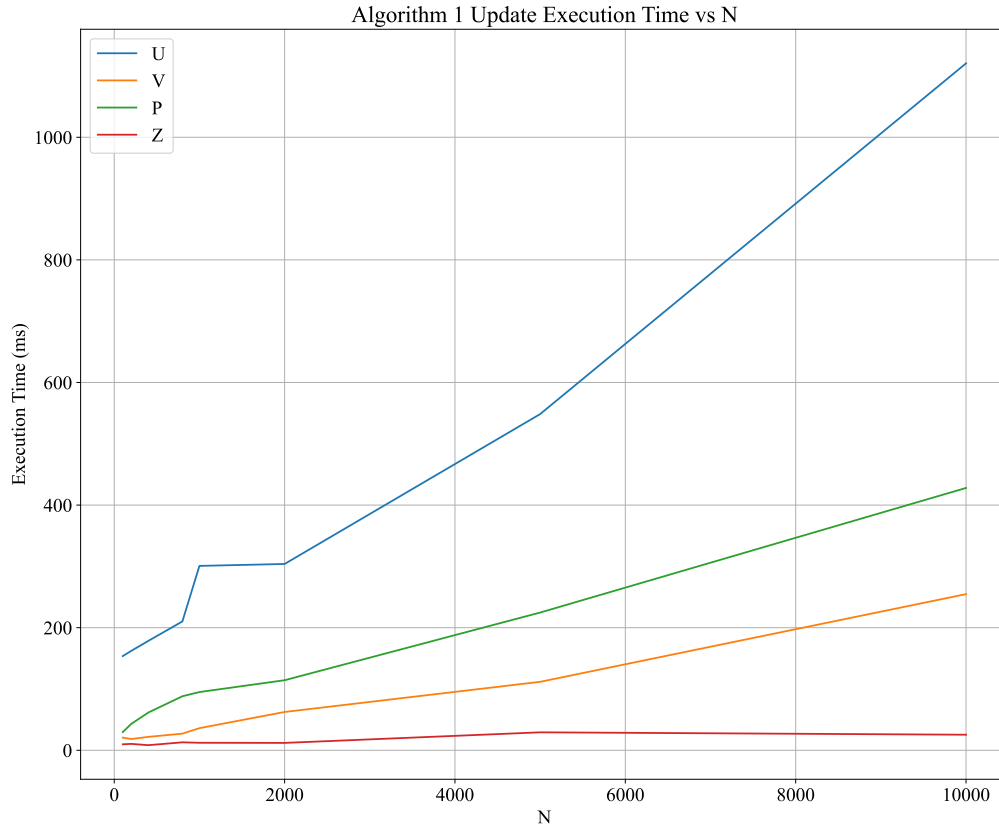


Figure 2: Cumulative time spent solving each subproblem of Algorithm 1 versus n with $m = 100$, $k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

6.3 Sensitivity to Column Dimension

Here, we present a comparison of Algorithm 1 with ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD as we vary the number of columns m . We fixed $n = 1000$, $k = 5$, and $d = 150$ across all trials. We varied $m \in \{100, 200, 400, 800, 1000, 2000, 5000, 10000\}$ and we performed 20 trials for each value of m .

We report the objective value, ℓ_2 reconstruction error, fitted rank and execution time for Algorithm 1, Fast-Impute and Soft-Impute in Figure 3. We additionally report the

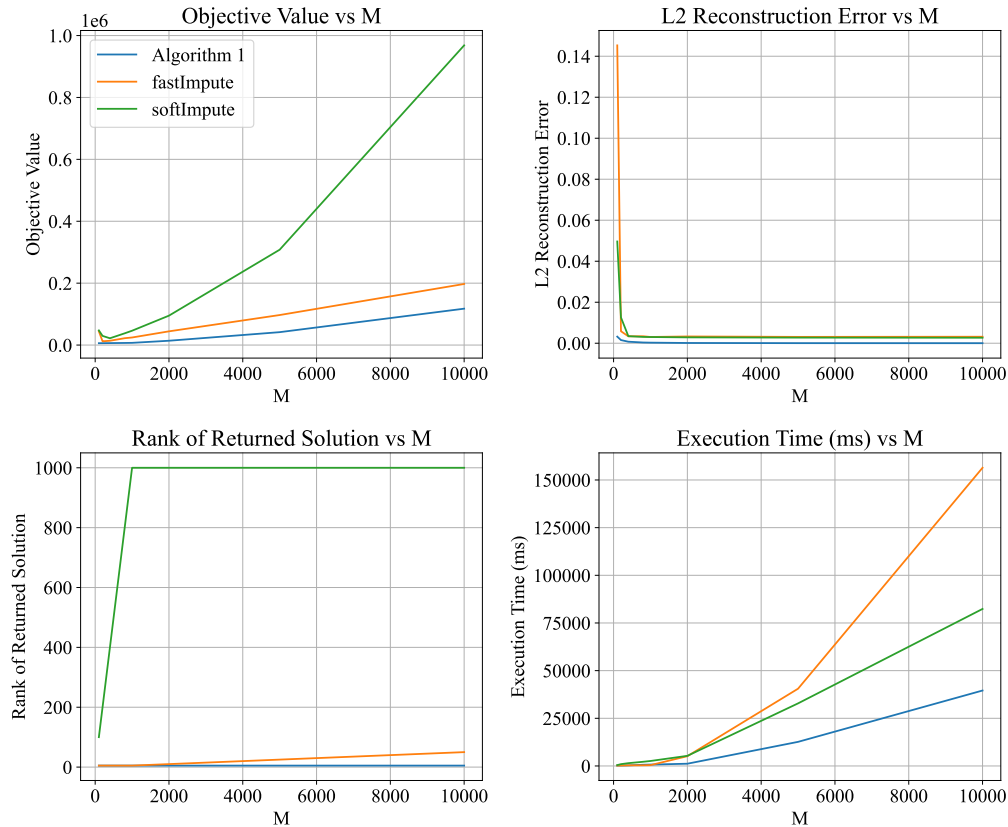


Figure 3: Objective value (top left), ℓ_2 reconstruction error (top right), fitted rank (bottom left) and execution time (bottom right) versus m with $n = 1000, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

objective value, reconstruction error and execution time for ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Tables 4, 5 and 6 of Appendix A. In Figure 4, we plot the average cumulative time spent solving subproblems (12), (14), (16), (18) during the execution of Algorithm 1 versus m . Our main findings from this set of experiments are as follows:

1. Here again, Algorithm 1 systematically produces higher quality solutions than ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD (see Table 4). On average, Algorithm 1 outputs a solution whose objective value is 62% lesser than the objective value achieved by the best performing alternative method (Fast-Impute). Here again, ScaledGD produces the weakest average objective value across these experiments. We observe that the objective value achieved by each method appears to increase super-linearly as the number of columns m increases.
2. In terms of ℓ_2 reconstruction error, Algorithm 1 again systematically produces solutions that are of higher quality than ScaledGD, Fast-Impute, Soft-Impute and

Iterative-SVD (see Table 5), often achieving an error that is an order of magnitude superior than the next best method. On average, Algorithm 1 outputs a solution whose ℓ_2 reconstruction error is 90% lesser than the reconstruction error achieved by the best performing alternative method (Soft-Impute in all but one parameter configuration).

3. The fitted rank of the solutions returned by Algorithm 1, ScaledGD and Fast-Impute always matched the specified target rank, but the solutions returned by Soft-Impute and Iterative-SVD were always of full rank despite the fact that these methods were provided with the target rank explicitly.
4. The runtime of Algorithm 1 exhibits the most favorable scaling behavior among the methods tested in these experiments. For instances with $m \geq 2000$, Table 6 shows that Algorithm 1 had the fastest runtime. For instances with $m < 2000$, ScaledGD had the fastest execution time but produced low quality solutions. The runtime of all methods tested grow super-linearly with m .
5. Figure 4 illustrates that the computation of the solution for (12) and (14) are the computational bottlenecks in the execution of Algorithm 1 in this set of experiments while the computation of the solution for (18) and (16) appear to be a constant function of m . This observation is consistent with the complexity analysis performed for each subproblem of Algorithm 1 in Section 5. Indeed, this analysis indicated that solve times for (18) and (16) are independent of m while the solve times for (12) and (14) scale linearly with m when the number of threads w satisfies $w < m$.

6.4 Sensitivity to Side Information Dimension

We present a comparison of Algorithm 1 with ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD as we vary the dimension of the side information d . In these experiments, we fixed $n = 1000$, $m = 100$ and $k = 5$ across all trials. We varied $d \in \{10, 50, 100, 150, 200, 250, 500, 1000\}$ and we performed 20 trials for each value of d .

We report the objective value, ℓ_2 reconstruction error, fitted rank and execution time for Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Figure 5. We additionally report the objective value, reconstruction error and execution time for ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Tables 7, 8 and 9 of Appendix A. In Figure 6, we plot the average cumulative time spent solving subproblems (12), (14), (16), (18) during the execution of Algorithm 1 versus d . Our main findings from this set of experiments are:

1. Just as in Sections 6.2 and 6.3, Algorithm 1 systematically produces higher quality solutions than ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD (see Table 7). On average, Algorithm 1 outputs a solution whose objective value is 85% lesser than the objective value achieved by the best performing alternative method (Fast-Impute). ScaledGD produces the weakest average objective value across these experiments. The objective value achieved by each method appears to increase linearly as the dimension d of the side information increases.

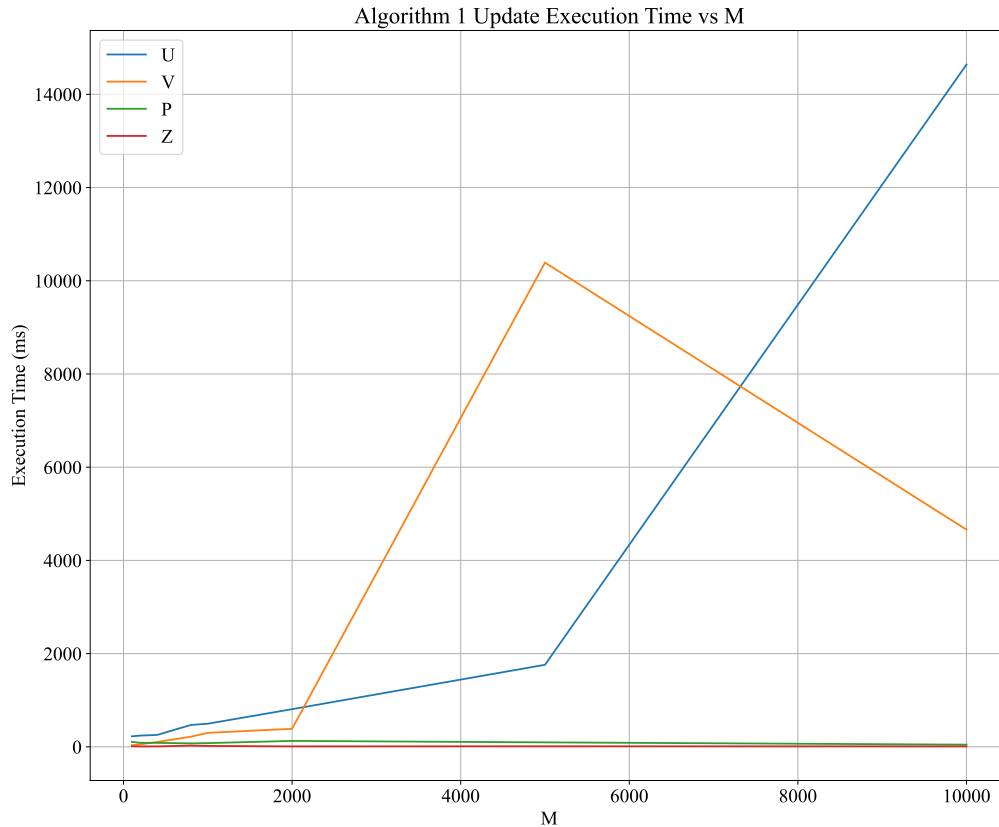


Figure 4: Cumulative time spent solving each subproblem of Algorithm 1 versus m with $n = 1000, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

2. In terms of ℓ_2 reconstruction error, just as in Sections 6.2 and 6.3 Algorithm 1 produces solutions that are of higher quality than ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD (see Table 8), often achieving an error that is an order of magnitude superior than the next best method. On average, Algorithm 1 outputs a solution whose ℓ_2 reconstruction error is 93% lesser than the reconstruction error achieved by the best performing alternative method (Soft-Impute). The performance of Algorithm 1 improves as d increases, consistent with the intuition that recovering the partially observed matrix \mathbf{A} becomes easier as more side information becomes available.
3. The runtime of Algorithm 1 is competitive with that of the other methods. The runtime of Algorithm 1 is less than that of Soft-Impute and Iterative-SVD but greater than that of Fast-Impute. Table 9 illustrates that ScaledGD was the fastest performing method, however its solutions were of the lowest quality. The runtime of Algorithm 1 and ScaledGD grows with d while Fast-Impute, Soft-Impute and iterate SVD are

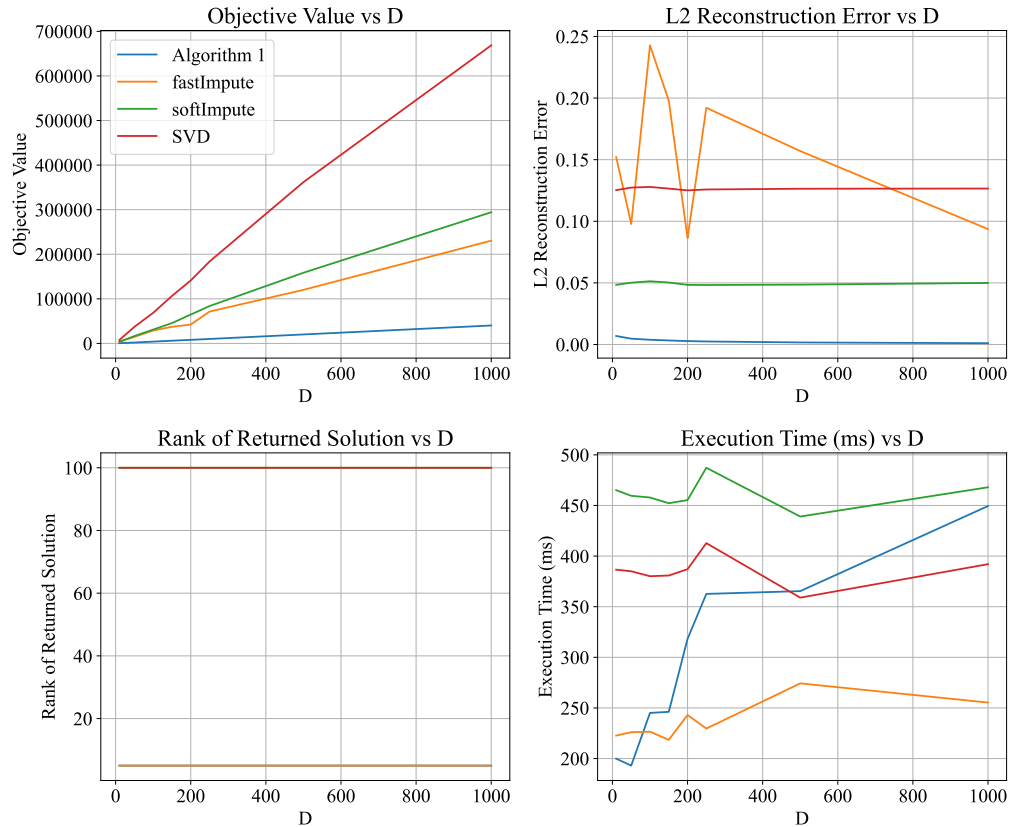


Figure 5: Objective value (top left), ℓ_2 reconstruction error (top right), fitted rank (bottom left) and execution time (bottom right) versus d with $n = 1000, m = 100$ and $k = 5$. Averaged over 20 trials for each parameter configuration.

constant with d which should be expected as these methods do not act on the side information matrix \mathbf{Y} .

- Figure 6 illustrates that the computation of the solution for (16) is the computational bottleneck in the execution of Algorithm 1 in this set of experiments, followed next by the computation of the solution to (12). The solution times for (12), (14) and (18) appear constant as a function of d . This is consistent with the complexity analysis from Section 5 which found that the solve time for (16) is linear in d while the solve time for the 3 other subproblems are independent of d .

6.5 Sensitivity to Target Rank

We present a comparison of Algorithm 1 with ScaledGD, Fast-Impute, Soft-Impute and Iterative-SVD as we vary the rank of the underlying matrix k . In these experiments, we fixed

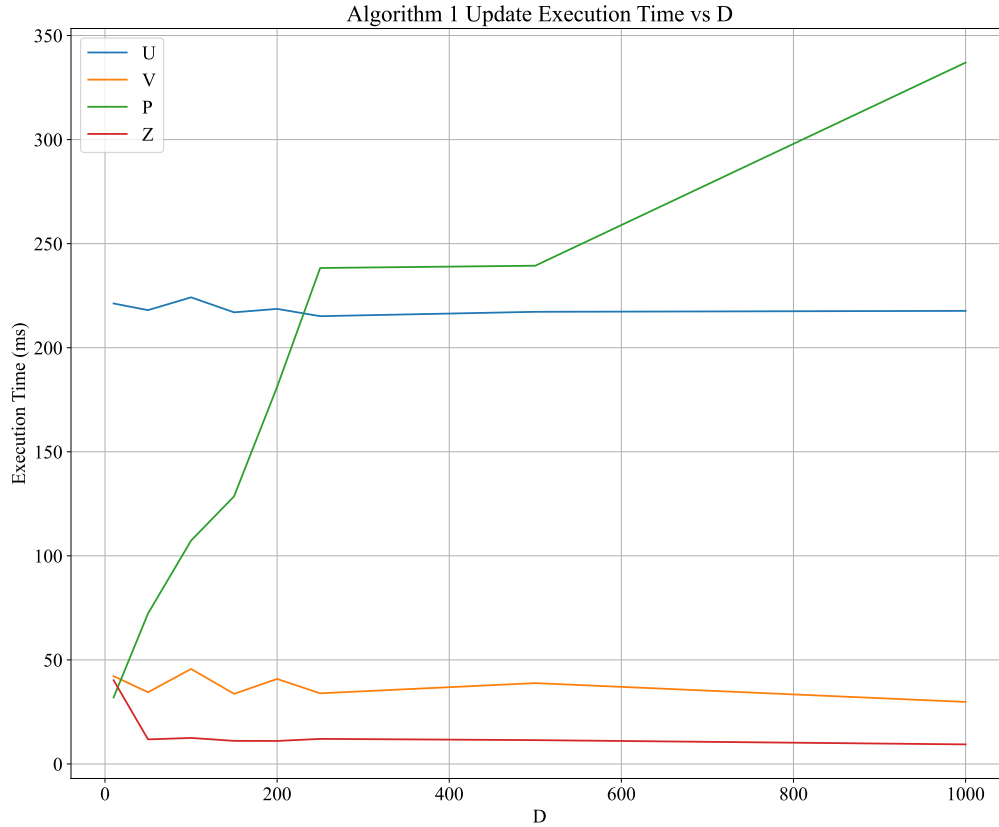


Figure 6: Cumulative time spent solving each subproblem of Algorithm 1 versus d with $n = 1000, m = 100$ and $k = 5$. Averaged over 20 trials for each parameter configuration.

$n = 1000, m = 100$ and $d = 150$ across all trials. We varied $k \in \{5, 10, 15, 20, 25, 30, 35, 40\}$ and we performed 20 trials for each value of d .

We report the objective value, ℓ_2 reconstruction error, fitted rank and execution time for Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Figure 7. We additionally report the objective value, reconstruction error and execution time for ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and Iterative-SVD in Tables 10, 11 and 12 of Appendix A. In Figure 8, we plot the average cumulative time spent solving subproblems (12), (14), (16), (18) during the execution of Algorithm 1 versus k . Our main findings from this set of experiments are as follows:

1. Unlike in Sections 6.2, 6.3 and 6.4, Algorithm 1 only produced higher quality solutions than all benchmark methods in 3 out of 8 of the tested parameter configurations where $k \leq 15$ (see Table 10). Fast-Impute was the best performing method in 3 configurations and Soft-Impute was best in the remaining 2 configurations. ScaledGD produces the weakest average objective value across these experiments.

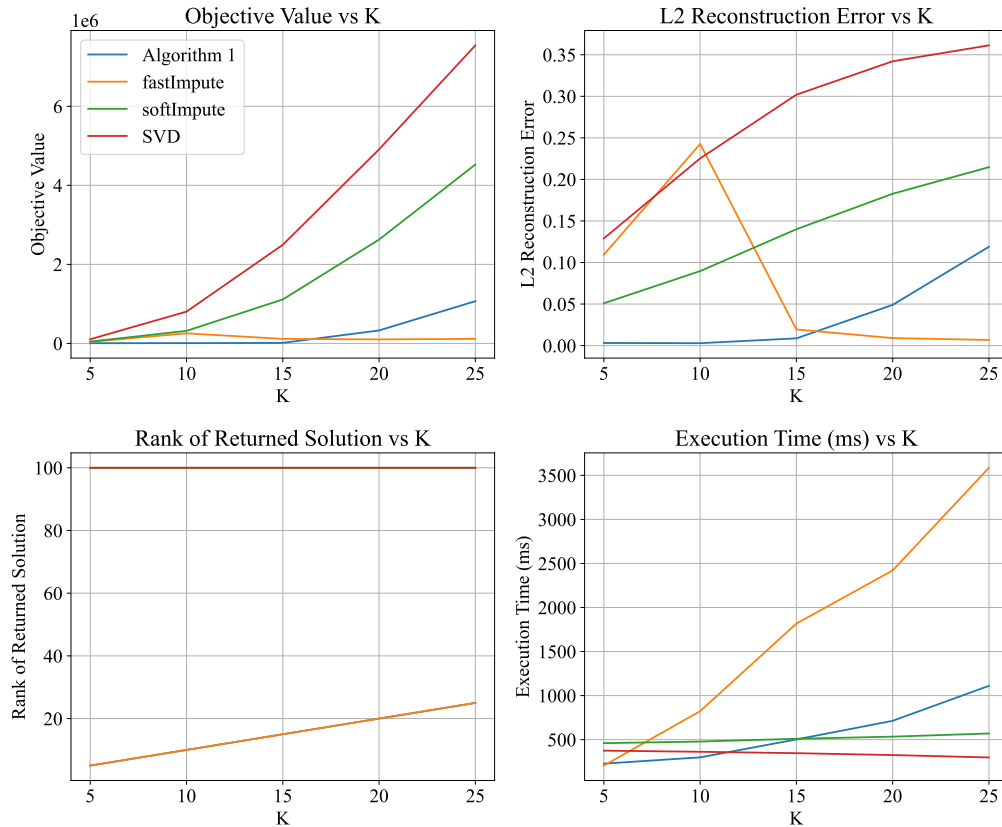


Figure 7: Objective value (top left), ℓ_2 reconstruction error (top right), fitted rank (bottom left) and execution time (bottom right) versus k with $n = 1000, m = 100$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

2. In terms of ℓ_2 reconstruction error, Algorithm 1 again produced higher quality solutions than all benchmark methods in 3 out of 8 of the tested parameter configurations where $k \leq 15$ (see Table 11). Fast-Impute produced solutions achieving the lowest error in the other 5 parameter configurations.
3. The runtime of Algorithm 1 is competitive with that of the other methods. Table 12 illustrates that ScaledGD was the fastest performing method, however its solutions were of the lowest quality. The runtime of Algorithm 1 is most competitive with Soft-Impute and Iterative-SVD for small values of k . Though Fast-Impute is the best performing method in terms of objective in 3 out of 8 configurations and the best in terms of ℓ_2 error in 5 out of 8 configurations, it takes on average 3 times as long as Algorithm 1 to execute.
4. Figure 8 illustrates that the computation of the solution for (12) is the computational bottleneck in the execution of Algorithm 1 in this set of experiments, followed next by the computation of the solution to (14) and (18).

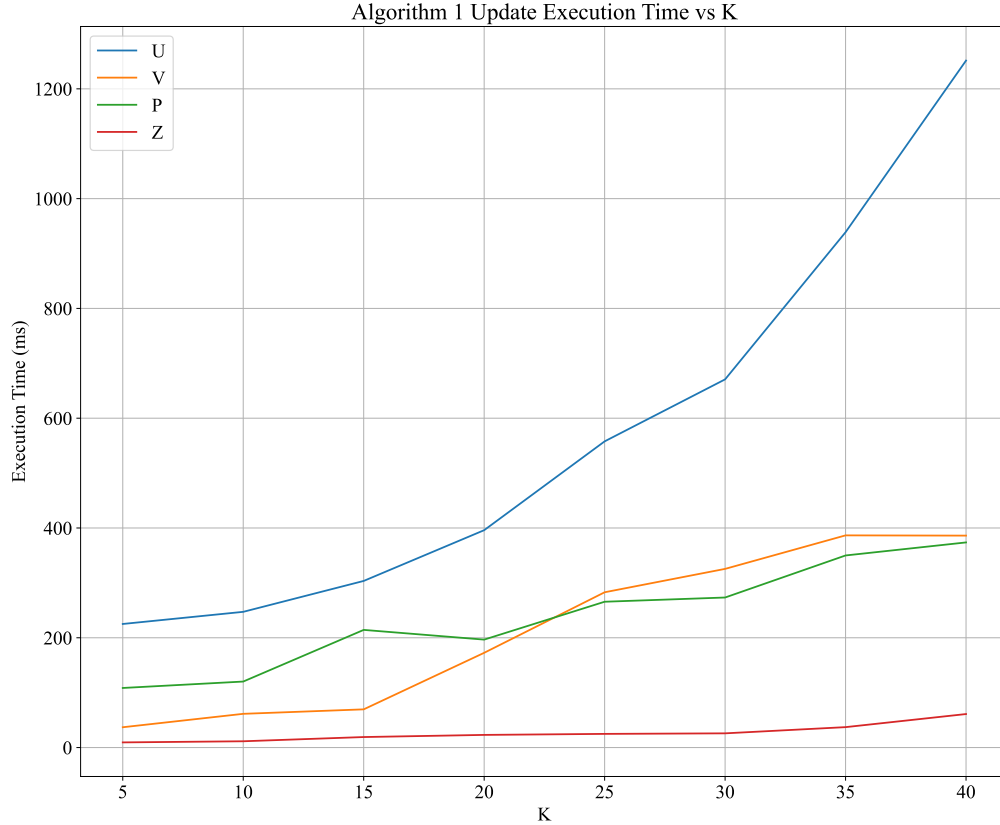


Figure 8: Cumulative time spent solving each subproblem of Algorithm 1 versus k with $n = 1000, m = 100$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

6.6 Summary of Findings

We now summarize our findings from our numerical experiments. In Sections 6.2-6.5, we see that across all experiments using synthetic data and target rank $k \leq 15$, Algorithm 1 produces solutions that achieve on average 79% lower objective value and 90.1% lower ℓ_2 reconstruction error than the solutions returned by the experiment-wise best performing benchmark method. In the regime where $k > 15$, we see in Section 6.5 that Fast-Impute outperforms Algorithm 1. We see that the execution time of Algorithm 1 is competitive with and often notably faster than the benchmark methods on synthetic data. Importantly, in the regime $k > 15$, although Fast-Impute returns higher quality solutions than Algorithm 1, the former has an execution time that is on average 3 times as long as our method. Our computational results are consistent with the complexity analysis performed in Section 5 for Problems (12), (14), (16) and (18). We observe that solution time for (16) becomes the bottleneck as the target rank k scales, otherwise the solution time for (12) is the bottleneck.

7. Conclusion

In this paper, we introduced Problem (1) which seeks to reconstruct a partially observed matrix that is predictive of fully observed side information. We illustrate that (1) has a natural interpretation as a robust optimization problem and can be reformulated as a mixed-projection optimization problem. We derive a semidefinite cone relaxation (9) to (1) and we present Algorithm 1, a mixed-projection alternating direction method of multipliers algorithm that obtains scalable, high quality solutions to (1). We rigorously benchmark the performance of Algorithm 1 on synthetic data against benchmark methods Fast-Impute, Soft-Impute, Iterative-SVD and ScaledGD. We find that across all experiments with $k \leq 15$, Algorithm 1 outputs solutions that achieve on average 79% lower objective value in (1) and 90.1% lower ℓ_2 reconstruction error than the solutions returned by the experiment-wise best performing benchmark method. For the 5 experiments with $k > 15$, Fast-Impute returns superior quality solutions than Algorithm 1, however the former takes on average 3 times as long as Algorithm 1 to execute. The runtime of Algorithm 1 is competitive with and often superior to that of the benchmark methods. Algorithm 1 is able to solve problems with $n = 10000$ rows and $m = 10000$ columns in less than a minute. Future work could expand the mixed-projection ADMM framework introduced in this work to incorporate positive semidefinite constraints and general linear constraints. Additionally, future work could empirically investigate the strength of the semidefinite relaxation (9) and could explore how to leverage this lower bound to certify globally optimal solutions.

References

- Francis R Bach and Michael I Jordan. Predictive low-rank decomposition for kernel methods. In *Proceedings of the 22nd international conference on Machine learning*, pages 33–40, 2005.
- Robert M Bell and Yehuda Koren. Lessons from the netflix prize challenge. *Acm Sigkdd Explorations Newsletter*, 9(2):75–79, 2007.
- Dimitris Bertsimas and Martin S Copenhaver. Characterization of the equivalence of robustification and regularization in linear and matrix regression. *European Journal of Operational Research*, 270(3):931–942, 2018.
- Dimitris Bertsimas and Dick den Hertog. *Robust and adaptive optimization*. Dynamic Ideas LLC, 2020.
- Dimitris Bertsimas and Nicholas A. G. Johnson. Compressed sensing: A discrete optimization approach. *Machine Learning*, 2024. URL <https://doi.org/10.1007/s10994-024-06577-0>.
- Dimitris Bertsimas and Michael Lingzhi Li. Fast exact matrix completion: A unified optimization framework for matrix completion. *Journal of Machine Learning Research*, 21(231):1–43, 2020.
- Dimitris Bertsimas, Jean Pauphilet, and Bart Van Parys. Rejoinder: Sparse regression: scalable algorithms and empirical performance. 2020.
- Dimitris Bertsimas, Ryan Cory-Wright, and Jean Pauphilet. A unified approach to mixed-integer optimization problems with logical constraints. *SIAM Journal on Optimization*, 31(3):2340–2367, 2021.
- Dimitris Bertsimas, Ryan Cory-Wright, and Jean Pauphilet. Mixed-projection conic optimization: A new paradigm for modeling rank constraints. *Operations Research*, 70(6):3321–3344, 2022.
- Dimitris Bertsimas, Ryan Cory-Wright, and Nicholas A. G. Johnson. Sparse plus low rank matrix decomposition: A discrete optimization approach. *Journal of Machine Learning Research*, 24(267):1–51, 2023a. URL <http://jmlr.org/papers/v24/21-1130.html>.
- Dimitris Bertsimas, Ryan Cory-Wright, Sean Lo, and Jean Pauphilet. Optimal low-rank matrix completion: Semidefinite relaxations and eigenvector disjunctions. *arXiv preprint arXiv:2305.12292*, 2023b.
- Dimitris Bertsimas, Ryan Cory-Wright, and Jean Pauphilet. A new perspective on low-rank optimization. *Mathematical Programming, articles in advance*, pages 1–46, 2023c.
- Daniel Billsus, Michael J Pazzani, et al. Learning collaborative information filters. In *Icml*, volume 98, pages 46–54, 1998.
- Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.

- Stephen Boyd, Laurent El Ghaoui, Eric Feron, and Venkataramanan Balakrishnan. *Linear matrix inequalities in system and control theory*. SIAM, 1994.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, USA, 2004.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- Emmanuel Candes and Benjamin Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
- Emmanuel J Candes and Yaniv Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE transactions on information theory*, 56(5):2053–2080, 2010.
- Yudong Chen and Martin J Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *arXiv preprint arXiv:1509.03025*, 2015.
- Kai-Yang Chiang, Cho-Jui Hsieh, and Inderjit S Dhillon. Matrix completion with noisy side information. *Advances in neural information processing systems*, 28, 2015.
- Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1):1–22, 1977.
- David L Donoho, Iain M Johnstone, Gérard Kerkycharian, and Dominique Picard. Wavelet shrinkage: asymptopia? *Journal of the Royal Statistical Society: Series B (Methodological)*, 57(2):301–337, 1995.
- Maryam Fazel. *Matrix rank minimization with applications*. PhD thesis, PhD thesis, Stanford University, 2002.
- Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on learning theory*, pages 797–842. PMLR, 2015.
- Oktay Günlük and Jeff Linderoth. Perspective reformulation and applications. In *Mixed Integer Nonlinear Programming*, pages 61–89. Springer, 2012.
- Ke Guo, Deren Han, David ZW Wang, and Tingting Wu. Convergence of admm for multi-block nonconvex separable optimization models. *Frontiers of Mathematics in China*, 12: 1139–1162, 2017.
- Prateek Jain and Praneeth Netrapalli. Fast exact matrix completion with finite samples. In *Conference on Learning Theory*, pages 1007–1034. PMLR, 2015.

- Hui Ji, Chaoqiang Liu, Zuowei Shen, and Yuhong Xu. Robust video denoising using low rank matrix completion. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 1791–1798. IEEE, 2010.
- Chi Jin, Sham M Kakade, and Praneeth Netrapalli. Provable efficient online matrix completion via non-convex stochastic gradient descent. *Advances in Neural Information Processing Systems*, 29, 2016.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.
- Qi Liu. Power network system identification and recovery based on the matrix completion. In *Journal of Physics: Conference Series*, volume 1237, page 032059. IOP Publishing, 2019.
- Cong Ma, Kaizheng Wang, Yuejie Chi, and Yuxin Chen. Implicit regularization in non-convex statistical estimation: Gradient descent converges linearly for phase retrieval and matrix completion. In *International Conference on Machine Learning*, pages 3345–3354. PMLR, 2018.
- Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research*, 11(80):2287–2322, 2010. URL <http://jmlr.org/papers/v11/mazumder10a.html>.
- Yurii Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. In *Dokl akad nauk Ssr*, volume 269, page 543, 1983.
- Luong Trung Nguyen, Junhan Kim, and Byonghyo Shim. Low-rank matrix completion: A contemporary survey. *IEEE Access*, 7:94215–94237, 2019.
- Art B. Owen and Patrick O. Perry. Bi-cross-validation of the SVD and the nonnegative matrix factorization. *The Annals of Applied Statistics*, 3(2):564 – 594, 2009.
- Kaare Brandt Petersen, Michael Syskind Pedersen, et al. The matrix cookbook. *Technical University of Denmark*, 7(15):510, 2008.
- Andy Ramlatchan, Mengyun Yang, Quan Liu, Min Li, Jianxin Wang, and Yaohang Li. A survey of matrix completion methods for recommendation systems. *Big Data Mining and Analytics*, 1(4):308–323, 2018.
- Albert Reuther, Jeremy Kepner, Chansup Byun, Siddharth Samsi, William Arcand, David Bestor, Bill Bergeron, Vijay Gadepally, Michael Houle, Matthew Hubbell, Michael Jones, Anna Klein, Lauren Milechin, Julia Mullen, Andrew Prout, Antonio Rosa, Charles Yee, and Peter Michaleas. Interactive supercomputing on 40,000 cores for machine learning and data analysis. In *2018 IEEE High Performance extreme Computing Conference (HPEC)*, pages 1–6. IEEE, 2018.
- Alex Rubinsteyn and Sergey Feldman. fancyimpute: An imputation library for python, 2016. URL <https://github.com/iskandr/fancyimpute>.

- Badrul Sarwar, George Karypis, Joseph Konstan, and John T Riedl. Application of dimensionality reduction in recommender system—a case study. 2000.
- Anureet Saxena, Pierre Bonami, and Jon Lee. Convex relaxations of non-convex mixed integer quadratically constrained programs: extended formulations. *Mathematical programming*, 124(1):383–411, 2010.
- Ruoyu Sun and Zhi-Quan Luo. Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579, 2016.
- Tian Tong, Cong Ma, and Yuejie Chi. Accelerating ill-conditioned low-rank matrix estimation via scaled gradient descent. *Journal of Machine Learning Research*, 22(150):1–63, 2021.
- Olga Troyanskaya, Michael Cantor, Gavin Sherlock, Pat Brown, Trevor Hastie, Robert Tibshirani, David Botstein, and Russ B Altman. Missing value estimation methods for dna microarrays. *Bioinformatics*, 17(6):520–525, 2001.
- Madeleine Udell and Alex Townsend. Why are big data matrices approximately low rank? *SIAM Journal on Mathematics of Data Science*, 1(1):144–160, 2019.
- Junxiang Wang and Liang Zhao. Nonconvex generalization of alternating direction method of multipliers for nonlinear equality constrained problems. *Results in Control and Optimization*, 2:100009, 2021.
- Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of admm in nonconvex nonsmooth optimization. *Journal of Scientific Computing*, 78:29–63, 2019.
- Huan Xu, Constantine Caramanis, and Shie Mannor. Robustness and regularization of support vector machines. *Journal of Machine Learning Research*, 10(7), 2009.
- Miao Xu, Rong Jin, and Zhi-Hua Zhou. Speedup matrix completion with side information: Application to multi-label learning. *Advances in neural information processing systems*, 26, 2013.
- Zheng Xu, Soham De, Mario Figueiredo, Christoph Studer, and Tom Goldstein. An empirical study of admm for nonconvex problems. *arXiv preprint arXiv:1612.03349*, 2016.
- Qinqing Zheng and John Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. *arXiv preprint arXiv:1605.07051*, 2016.

Appendix A. Supplemental Computational Results

Table 1: Comparison of the objective value of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus n with $m = 100, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

N	Objective				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	249262.99	655.46	5530.48	19893.08	28677.73
200	306738.68	1280.82	9756.14	24251.93	44054.89
400	417643.27	2483.49	14321.46	28932.85	61112.62
800	421032.49	4813.05	31520.96	41179.38	93119.34
1000	522586.08	6010.34	33557.75	47701.68	107851.28
2000	563033.20	11975.48	83669.84	76566.52	167458.68
5000	1226489.68	30060.61	273747.04	172093.66	364065.64
10000	1973665.62	60082.27	521189.88	319915.98	642759.84

Table 2: Comparison of the reconstruction error of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus n with $m = 100, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

N	ℓ_2 Reconstruction Error				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	100.22460	0.01520	0.07540	0.21330	0.30090
200	58.37210	0.00695	0.12800	0.11770	0.21050
400	33.92500	0.00412	0.08980	0.07390	0.15980
800	14.97890	0.00328	0.23990	0.05160	0.13010
1000	12.66500	0.00312	0.18600	0.04950	0.12620
2000	5.54420	0.00304	0.07990	0.04410	0.11670
5000	2.47260	0.00282	0.15040	0.03720	0.10550
10000	1.32070	0.00267	0.10920	0.03510	0.10240

Table 3: Comparison of the execution time of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus n with $m = 100, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

N	Execution Time (ms)				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	10.84	53.47	99.95	141.42	115.26
200	41.11	73.84	121.05	187.26	163.11
400	54.95	113.89	152.16	262.95	246.42
800	68.00	184.11	195.05	389.63	334.16
1000	44.11	241.16	211.63	442.05	366.53
2000	124.47	340.84	311.32	759.58	575.79
5000	813.53	770.11	484.58	1795.58	1298.63
10000	18828.21	1730.00	863.84	3871.53	2713.26

Table 4: Comparison of the objective value of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus m with $n = 1000, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

M	Objective				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	530097.31	6014.93	44337.08	47334.20	103403.04
200	2483913.82	6131.52	12159.08	29560.77	114448.43
400	14226534.31	6361.90	13875.31	21942.31	90652.40
800	99356634.18	6800.63	22152.43	37924.99	87895.33
1000	105451997.06	7126.60	24435.40	46327.69	128499.51
2000	591164404.77	13964.62	44333.51	95044.30	815807.16
5000	4002087935.12	41679.23	96985.27	308044.93	11294104.83
10000	9826251365.01	117558.76	197362.37	968255.00	60913874.17

Table 5: Comparison of the reconstruction error of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus m with $n = 1000, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

ℓ_2 Reconstruction Error					
M	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	13.68740	0.00322	0.14530	0.04960	0.12560
200	40.58900	0.00154	0.00590	0.01260	0.06640
400	127.71450	0.00075	0.00340	0.00340	0.02240
800	508.24550	0.00036	0.00340	0.00310	0.00460
1000	443.26620	0.00029	0.00310	0.00300	0.00350
2000	1292.61610	0.00012	0.00330	0.00290	0.00300
5000	3658.18810	0.00004	0.00310	0.00270	0.00540
10000	4559.27360	0.00002	0.00320	0.00270	0.00650

Table 6: Comparison of the execution time of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus m with $n = 1000, k = 5$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

Execution Time (ms)					
M	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
100	44.16	222.21	193.16	444.89	365.74
200	54.32	237.37	223.32	953.32	760.11
400	80.95	311.32	315.47	1466.32	1511.79
800	121.37	637.53	360.53	2198.21	2564.00
1000	154.89	728.58	434.47	2611.58	3009.21
2000	4652.11	1181.47	5127.37	5308.16	6062.89
5000	28587.11	12645.16	40526.16	32824.79	35015.21
10000	108255.05	39569.37	156399.42	82361.37	86762.84

Table 7: Comparison of the objective value of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus d with $n = 1000, m = 100$ and $k = 5$. Averaged over 20 trials for each parameter configuration.

D	Objective				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
10	11691.78	475.17	4222.64	3367.41	7710.45
50	96771.27	2067.20	14565.83	16511.97	37203.79
100	229740.41	4033.46	28967.06	31000.81	68634.31
150	532018.26	6057.23	37244.64	45581.92	106504.32
200	734648.30	7994.51	42289.45	64771.47	141252.05
250	1195065.81	9984.20	71433.91	83752.90	183720.00
500	4165782.61	20094.24	120114.14	158458.64	361740.57
1000	13578263.54	40191.29	230467.93	294273.13	668550.27

Table 8: Comparison of the reconstruction error of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus d with $n = 1000, m = 100$ and $k = 5$. Averaged over 20 trials for each parameter configuration.

D	ℓ_2 Reconstruction Error				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
10	0.60780	0.00690	0.15210	0.04840	0.12530
50	1.41600	0.00471	0.09780	0.05020	0.12730
100	5.03590	0.00382	0.24280	0.05130	0.12790
150	14.00330	0.00326	0.19790	0.05030	0.12650
200	22.26870	0.00276	0.08640	0.04840	0.12510
250	39.41630	0.00245	0.19210	0.04830	0.12580
500	177.68800	0.00165	0.15700	0.04860	0.12640
1000	679.11770	0.00104	0.09360	0.04990	0.12660

Table 9: Comparison of the execution time of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus d with $n = 1000, m = 100$ and $k = 5$. Averaged over 20 trials for each parameter configuration.

D	Execution Time (ms)				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
10	84.00	199.89	222.68	465.11	386.53
50	80.79	193.05	226.00	459.53	385.00
100	108.63	245.11	226.53	457.84	380.11
150	113.47	246.16	218.47	452.21	380.79
200	117.32	318.26	243.05	455.32	387.05
250	152.79	362.63	229.63	487.21	412.79
500	176.21	365.37	274.26	439.05	358.95
1000	138.74	449.42	255.32	467.95	392.00

Table 10: Comparison of the objective value of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus k with $n = 1000, m = 100$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

K	Objective				
	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
5	514330.09	6021.76	41376.50	45858.20	106390.58
10	1892278.99	7393.65	255228.60	318398.98	805396.86
15	5213393.44	14104.62	115383.43	1112396.97	2495972.46
20	10196279.89	328671.00	101812.52	2628073.04	4910386.51
25	16816442.74	1069103.04	116005.02	4526388.13	7541300.54
30	-	34567679.83	10695127.17	6864577.33	10634436.81
35	39536651.09	187701091.79	144464.25	9424715.13	14192827.60
40	-	723504611.03	191276652.97	12529277.05	18290215.22

Table 11: Comparison of the reconstruction error of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus k with $n = 1000, m = 100$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

ℓ_2 Reconstruction Error					
K	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
5	12.95500	0.00314	0.10940	0.05090	0.12900
10	5.41720	0.00288	0.24270	0.08970	0.22520
15	3.45210	0.00871	0.01940	0.14010	0.30200
20	2.34590	0.04900	0.00903	0.18270	0.34220
25	1.73610	0.11890	0.00678	0.21480	0.36130
30	-	0.20030	0.00562	0.23810	0.37240
35	1.16780	0.17330	0.00488	0.25260	0.37740
40	-	0.22320	0.00452	0.26550	0.38020

Table 12: Comparison of the execution time of ScaledGD, Algorithm 1, Fast-Impute, Soft-Impute and SVD versus k with $n = 1000, m = 100$ and $d = 150$. Averaged over 20 trials for each parameter configuration.

Execution Time (ms)					
K	ScaledGD	Algorithm 1	Fast-Impute	Soft-Impute	SVD
5	55.79	227.47	205.16	460.11	374.79
10	80.79	298.63	823.21	477.95	361.95
15	107.89	502.21	1817.16	509.68	346.05
20	95.53	713.42	2420.89	533.79	324.79
25	111.53	1110.89	3586.05	569.68	297.58
30	-	1353.95	4435.63	591.21	280.37
35	107.21	1822.16	6212.63	640.05	271.00
40	-	2281.95	8168.68	645.11	263.37