



HAL
open science

Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames

Geneviève Robin, Hoi-To Wai, Julie Josse, Olga Klopp, Éric Moulines

► **To cite this version:**

Geneviève Robin, Hoi-To Wai, Julie Josse, Olga Klopp, Éric Moulines. Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames. 32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Dec 2018, Montréal, Canada. hal-01959188v1

HAL Id: hal-01959188

<https://hal.science/hal-01959188v1>

Submitted on 19 Dec 2018 (v1), last revised 5 Apr 2019 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Low-rank Interaction with Sparse Additive Effects Model for Large Data Frames

Geneviève Robin

Centre de Mathématiques Appliquées
École Polytechnique, XPOP, INRIA
91120 Palaiseau, France
genevieve.robin@polytechnique.edu

Hoi-To Wai

Department of SE&EM
The Chinese University of Hong Kong
Shatin, Hong Kong
htwai@se.cuhk.edu.hk

Julie Josse

Centre de Mathématiques Appliquées
École Polytechnique, XPOP, INRIA
91120 Palaiseau, France
julie.josse@polytechnique.edu

Olga Klopp

ESSEC Business School
CREST, ENSAE
95021 Cergy, France
klopp@essec.edu

Éric Moulines

Centre de Mathématiques Appliquées
École Polytechnique, XPOP, INRIA
91120 Palaiseau, France
eric.moulines@polytechnique.edu

Abstract

Many applications of machine learning involve the analysis of large data frames – matrices collecting heterogeneous measurements (binary, numerical, counts, etc.) across samples – with missing values. Low-rank models, as studied by Udell et al. [30], are popular in this framework for tasks such as visualization, clustering and missing value imputation. Yet, available methods with statistical guarantees and efficient optimization do not allow explicit modeling of main additive effects such as row and column, or covariate effects. In this paper, we introduce a *low-rank interaction and sparse additive effects* (LORIS) model which combines matrix regression on a dictionary and low-rank design, to estimate main effects and interactions simultaneously. We provide statistical guarantees in the form of upper bounds on the estimation error of both components. Then, we introduce a *mixed coordinate gradient descent* (MCGD) method which provably converges sub-linearly to an optimal solution and is computationally efficient for large scale data sets. We show on simulated and survey data that the method has a clear advantage over current practices, which consist in dealing separately with additive effects in a preprocessing step.

1 Introduction

Recently, a lot of effort has been devoted towards the efficient analysis of large data frames, a term coined by Udell et al. [30]. A data frame is a large table of heterogeneous data (binary, numerical, counts) with missing entries, where each row represents an example and each column a feature. In order to analyze them, a powerful technique is to use *low-rank models* that embed rows and columns

of data frames into low-dimensional spaces [18, 28, 30], enabling effective data analytics such as clustering, visualization and missing value imputation; see also [22] and the references therein.

Characterizing additive effects of side information – such as covariates, row or column effects – *simultaneously* with low rank interactions is an important extension to plain low-rank models. For example, in data frames obtained from recommender systems, user information and item characteristics are known to influence the ratings in addition to interactions between users and items [9]. These modifications to the low rank model have been advocated in the statistics literature, but they have been implemented only for small data frames [1].

In the large-scale low-rank matrix estimation literature, available methods either do not take additive effects into account [8, 24, 30, 26, 10], or only handle the numerical data [15, 14]. As a common heuristics for preprocessing, prior work such as [24, 30] remove the row and column means and apply some normalization of the row and column variance. We show in numerical experiments this apparently benign operation is not appropriate for large and heterogenous data frames, and can cause severe impairments in the analysis.

The present work investigates a generalization of previous contributions in the analysis of data frames. Our contributions can be summarized as follows.

Contributions We present a new framework that is *statistically* and *computationally* efficient for analyzing large and incomplete heterogeneous data frames.

- We describe in Section 2 the *low-rank interaction with sparse additive effects* (LORIS) model, which combines matrix regression on a dictionary with low rank approximation. We propose a convex doubly penalized quasi-maximum likelihood approach, where the rank constraint is relaxed with a nuclear norm penalty, to estimate the regression coefficients and the low rank component simultaneously. We establish non-asymptotic upper bounds on the estimation errors.
- We propose in Section 3 a Mixed Coordinate Gradient Descent (MCGD) method to solve efficiently the LORIS estimation problem. It uses a mixed update strategy including a proximal update for the sparse component and a conditional gradient (CG) for the low-rank component. We show that the MCGD method converges to an ϵ -optimal solution in $\mathcal{O}(1/\epsilon)$ iterations. We also outline an extension to efficient distributed implementation.
- We demonstrate in Section 4 the efficacy of our method both in terms of estimation and imputation quality on simulated and survey data examples.

Related work Our statistical model and analysis are related to prior work on *low-rank plus sparse matrix decomposition* [32, 5, 6, 16, 21]; these papers provide statistical results for a particular case where the loss function is quadratic and the sparse component is entry-wise sparse. In comparison, the originality of the present work is two-fold. First, the sparsity pattern of the main effects is not restricted to entry-wise sparsity. Second, the data fitting term is not quadratic, but a heterogeneous exponential family quasi log-likelihood. This new framework enables us to tackle many more data sets combining heterogeneous data, main effects and interactions.

For the algorithmic development, our proposed method is related to the prior work such as [25, 29, 7, 14, 33, 17, 27, 11, 23, 4, 12]. These are based on various first-order optimization methods and shall be reviewed in detail in Section 3. Among others, the MCGD method is mostly related to the recent FW-T method by Mu et al. [27] that uses a mixed update rule to tackle a similar estimation problem. There are two differences: first, FW-T is focused on a quadratic loss which is a special case of the statistical estimation problem that we analyze; second, the per-iteration complexity of MCGD is lower as the update rules are simpler. Despite the simplifications, using a new proof technique, we prove that the convergence rate of MCGD is strictly faster than FW-T.

Notations: For any $m \in \mathbb{N}$, $[m] := \{1, \dots, m\}$. The operator $\mathcal{P}_\Omega(\cdot) : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times p}$ is the projection operator on the set of entries in $\Omega \subset [n] \times [p]$, and $(\cdot)_+ : \mathbb{R} \rightarrow \mathbb{R}_+$ is the projection operator on the non-negative orthant $(x)_+ := \max\{0, x\}$. For matrices, we denote by $\|\cdot\|_F$ the Frobenius norm, $\|\cdot\|_*$ the nuclear norm, $\|\cdot\|$ the operator norm, and $\|\cdot\|_\infty$ the entry-wise infinity norm. For vectors, we denote by $\|\cdot\|_1$ is the ℓ_1 -norm, $\|\cdot\|_2$ the Euclidean norm, $\|\cdot\|_\infty$ the infinity norm, and $\|\cdot\|_0$ the number of non zero coefficients. The binary operator $\langle \mathbf{X}, \mathbf{Y} \rangle$ denotes the Frobenius

inner product. A function $f : \mathbb{R}^q \rightarrow \mathbb{R}$ is said to be σ -smooth if f is continuously differentiable and $\|\nabla f(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta}')\|_2 \leq \sigma \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2$ for all $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^q$.

2 Problem Formulation

Heterogenous Data Model Let (Y, X) be a probability space equipped with a σ -finite measure μ . The canonical exponential family distribution $\{\text{Exp}_{h,g}(m), m \in X\}$ with base measure $h : Y \rightarrow \mathbb{R}^+$, link function $g : X \rightarrow \mathbb{R}$, and scalar parameter, $m \in X$, has a density given by

$$f_m(y) = h(y) \exp(y m - g(m)). \quad (1)$$

The exponential family is a flexible framework to model different types of data. For example, $(Y = \mathbb{R}, g(m) = m^2 \sigma^2 / 2, h(y) = (2\pi \sigma^2)^{-1/2} \exp(-y^2 / 2\sigma^2))$ yields a Gaussian distribution with mean m and variance σ^2 for numerical data; $(Y = \{0, 1\}, g(m) = \log(1 + \exp(m)), h(y) = 1)$ yields a Bernoulli distribution with success probability $1/(1 + \exp(-m))$ for binary data; $(Y = \mathbb{N}, g(m) = \exp(am), h(y) = 1/y!)$ where $a \in \mathbb{R}$ yields a Poisson distribution with intensity $\exp(am)$ for count data. In these cases, the parameter space is $X = \mathbb{R}$.

Let $\{(Y_j, g_j, h_j), j \in [p]\}$ be a collection of observation spaces, base and link functions corresponding to the column types of a data frame $\mathbf{Y} = [\mathbf{Y}_{ij}]_{(i,j) \in [n] \times [p]} \in \mathbb{Y}_1^n \times \dots \times \mathbb{Y}_p^n$. For each $i \in [n]$ and $j \in [p]$, we denote by \mathbf{M}_{ij}^0 the target parameter minimizing the Kullback-Leibler divergence between the distribution of \mathbf{Y}_{ij} and the exponential family $\text{Exp}_{h_j, g_j}, j \in [p]$, given by

$$\mathbf{M}_{ij}^0 = \arg \max_m \mathbb{E}_{\mathbf{Y}_{ij}} [\log(h_j(\mathbf{Y}_{ij})) + \mathbf{Y}_{ij} m - g_j(m)]. \quad (2)$$

We propose the following model to estimate $\mathbf{M}^0 = [\mathbf{M}_{ij}^0]_{(i,j) \in [n] \times [p]}$ in the presence of additive effects and interactions.

Low-rank Interaction with Sparse additive effects (LORIS) model For every entry \mathbf{Y}_{ij} , assume a vector of covariates $\mathbf{x}_{ij} \in \mathbb{R}^q$ is also available, e.g., user information and item characteristics. Denote $\mathbf{x}_{ij}(k), k \in [q]$ the k -th component of \mathbf{x}_{ij} and define the matrix $\mathbf{X}(k) = [\mathbf{x}_{ij}(k)]_{(i,j) \in [n] \times [p]}$. We introduce the following decomposition of the parameter matrix \mathbf{M}^0 :

$$\mathbf{M}^0 = \sum_{k=1}^q \alpha_k^0 \mathbf{X}(k) + \boldsymbol{\Theta}^0. \quad (3)$$

We call (3) the LORIS model, where $\boldsymbol{\alpha} \in \mathbb{R}^q$ is a sparse vector with unknown support modeling additive effects and $\boldsymbol{\Theta}^0 \in \mathbb{R}^{n \times p}$ a low-rank matrix modeling the interactions.

In fact, LORIS is a generalization of *robust* matrix completion [5], where the parameter matrix can be decomposed as the sum of two matrices, one is low-rank and the other has some complementary low-dimensional structure such as entry-wise or column-wise sparsity. Statistical recoverability results in robust matrix estimation under a noiseless setting can be found in [32, 5, 6, 16]; the additive noise setting can be found in a recent work [21].

Estimation Problem Denote $\Omega = \{(i, j) \in [n] \times [p] : \mathbf{Y}_{ij} \text{ is observed}\}$ as the observation set. For $\mathbf{M} \in \mathbb{R}^{n \times p}$, $\mathcal{L}(\mathbf{M})$ is the negative log-likelihood of the observed data (\mathbf{Y}, Ω) parameterized by \mathbf{M} . Up to an additive constant,

$$\mathcal{L}(\mathbf{M}) = \sum_{(i,j) \in \Omega} \{-\mathbf{Y}_{ij} \mathbf{M}_{ij} + g_j(\mathbf{M}_{ij})\}. \quad (4)$$

For $a > 0$, we consider the following estimation problem:

$$(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}) \in \underset{\substack{\|\boldsymbol{\alpha}\|_\infty \leq a \\ \|\boldsymbol{\Theta}\|_\infty \leq a}}{\text{argmin}} \mathcal{L} \left(\sum_{k=1}^q \alpha_k \mathbf{X}(k) + \boldsymbol{\Theta} \right) + \lambda_S \|\boldsymbol{\alpha}\|_1 + \lambda_L \|\boldsymbol{\Theta}\|_\star. \quad (5)$$

We denote by $\hat{\mathbf{M}} = \sum_{k=1}^q \hat{\alpha}_k \mathbf{X}(k) + \hat{\boldsymbol{\Theta}}$ the estimated parameter matrix. The ℓ_1 and nuclear norm penalties are convex relaxations of the sparsity and low-rank constraints, and the regularization parameters λ_S and λ_L serve as trade-offs between fitting the data and enforcing sparsity of $\boldsymbol{\alpha}$ and controlling the "effective rank" of $\boldsymbol{\Theta}$.

Statistical Guarantees Here we establish convergence rates for the joint estimation of α^0 and Θ^0 ; the proofs can be found in the supplementary material. Consider the following assumptions.

H1 $\|\Theta^0\|_\infty \leq a$, $\|\alpha^0\|_\infty \leq a$ and for all $k \in [q]$ such that $\alpha_k^0 \neq 0$, $\langle \Theta^0, \mathbf{X}(k) \rangle = 0$.

In particular, **H1** guarantees the uniqueness of the decomposition in the LORIS model (3).

H2 For $\nu > 0$, all $k \in [q]$ and $(i, j) \in [n] \times [p]$, $\mathbf{X}(k)_{ij} \in [-1, 1]$. Furthermore for all $(i, j) \in [n] \times [p]$, $\sum_{k=1}^q |\mathbf{X}(k)_{ij}| \leq \nu$.

In particular, **H2** guarantees that for all (Θ, α) satisfying **H1**, the matrix $\mathbf{M} = \sum_{k=1}^q \alpha_k \mathbf{X}(k) + \Theta$ satisfies $\|\mathbf{M}\|_\infty \leq (1 + \nu)a$. Let \mathbf{G} be the $q \times q$ Gram matrix of the dictionary $(\mathbf{X}(1), \dots, \mathbf{X}(q))$ defined by $\mathbf{G} = [\langle \mathbf{X}(k), \mathbf{X}(l) \rangle]_{(k,l) \in [q] \times [q]}$.

H3 For $\kappa > 0$ and all $\alpha \in \mathbb{R}^q$, $\alpha^\top \mathbf{G} \alpha \geq \kappa^2 \|\alpha\|_2^2$.

Note we do not consider the case where the Gram matrix is singular, e.g., $q > np$. For $0 < \sigma_- \leq \sigma_+ < +\infty$ and $0 < \gamma < \infty$ consider the following assumption on the link functions g_j :

H4 The functions g_j are twice differentiable, and for all $x \in [-(1 + \nu)a - \gamma, (1 + \nu)a + \gamma]$,

$$\sigma_-^2 \leq g_j''(x) \leq \sigma_+^2, \quad j \in [p].$$

H4 implies the data fitting term $\mathcal{L}(\mathbf{M})$ is smooth and satisfies a restricted strong convexity property.

H5 For all $(i, j) \in [n] \times [p]$, Y_{ij} is a sub-exponential random variable with scale and variance parameters $1/\gamma$ and σ_+^2 .

If the random variables Y_{ij} are actually distributed according to an exponential family distribution of the form (1), then **H4** implies **H5**.

H6 For $(i, j) \in [n] \times [p]$, the events $\omega_{ij} = \{(i, j) \in \Omega\}$ are independent with occurrence probability π_{ij} . Furthermore, there exists $0 < \pi \leq 1$ such that for all $(i, j) \in [n] \times [p]$, $\pi_{ij} \geq \pi$.

H6 implies a data missing-at-random scenario where \mathbf{Y}_{ij} is observed with probability at least π .

Theorem 1 Assume **H1-6**. Set

$$\lambda_L = 2C\sigma_+ \sqrt{\pi \max(n, p) \log(n + p)}, \quad \text{and} \quad \lambda_S = 24 \max_k \|\mathbf{X}(k)\|_1 \log(n + p) / \gamma, \quad (6)$$

where C is a positive constant. Assume that $\max(n, p) \geq 4\sigma_+^2 / \gamma^6 \log^2(\sqrt{\min(n, p)} / (\pi\gamma\sigma_-)) + 2 \exp(\sigma_+^2 / \gamma^2 + 2\sigma_+^2 \gamma a)$. Then, with probability at least $1 - 9(n + p)^{-1}$,

$$\begin{aligned} \|\hat{\alpha} - \alpha^0\|_2^2 &\leq C_1 \frac{s \max_k \|\mathbf{X}(k)\|_1 \log(n + p)}{\kappa^2 \pi} + D_\alpha, \\ \|\hat{\Theta} - \Theta^0\|_F^2 &\leq C_2 \left(\frac{r \max(n, p)}{\pi} + \frac{s \max_k \|\mathbf{X}(k)\|_1}{\pi} \right) \log(n + p) + D_\Theta. \end{aligned} \quad (7)$$

In (23), $s := \|\alpha^0\|_0$, $r := \text{rank}(\Theta^0)$. C_1 and C_2 are positive constants and D_α and D_Θ are residuals of lower order whose exact values are given in Appendix A.

The proof can be found in Appendix A. In **Theorem 1**, the rate obtained for α^0 is the same as the bound obtained in [21] in the special case of robust matrix completion. Examples satisfying $\max_k \|\mathbf{X}(k)\|_1 / \kappa^2 = \mathcal{O}(1)$ include the case where the elements of the dictionary are matrices are all zeros except a row or a column of one, (to model row and column effects) and the number of rows n and columns p are of the same order; or when the covariates \mathbf{x}_{ij} are categorical and the categories are balanced, i.e., the number of samples per category is of the same order.

The rate obtained for Θ^0 is the sum of the standard low-rank matrix completion rate of order $r \max(n, p) / \pi$, e.g., [19], and of a term which boils down to sparse vector estimation rate as long as $\max_k \|\mathbf{X}(k)\|_1 = \mathcal{O}(1)$. Again, the latter can be satisfied by the special case of robust matrix completion, for which our rates match the results of [21].

3 A Mixed Coordinate Gradient Descent Method for LORIS

This section introduces a mixed coordinate gradient descent (MCGD) method to solve the LORIS estimation problem (5). We assume that a is sufficiently large such that the constraints $\|\alpha\|_\infty \leq a$, $\|\Theta\|_\infty \leq a$ are always inactive. To simplify notation, we denote the log-likelihood function as $\mathcal{L}(\alpha, \Theta) := \mathcal{L}(\sum_{k=1}^q \alpha_k \mathbf{X}(k) + \Theta)$. We assume

H7 (a) $\mathcal{L}(\alpha, \Theta)$ is σ_Θ -smooth w.r.t. Θ_{ij} for $(i, j) \in \Omega$ and (b) σ_α -smooth w.r.t. α ; (c) the gradient $\nabla_\alpha \mathcal{L}(\alpha, \Theta)$ is $\hat{\sigma}_\Theta$ -Lipschitz w.r.t. Θ . Moreover, the gradient $\nabla_\Theta \mathcal{L}(\alpha, \Theta)$ is bounded as long as α, Θ are bounded.

The above is implied by H4 for bounded (α, Θ) . We consider the augmented objective function:

$$F(\alpha, \Theta, R) := \mathcal{L}(\alpha, \Theta) + \lambda_S \|\alpha\|_1 + \lambda_L R. \quad (8)$$

For some $R_{\text{UB}} \geq 0$, if an optimal solution $(\hat{\alpha}, \hat{\Theta})$ to (5) satisfies $\|\hat{\Theta}\|_* \leq R_{\text{UB}}$, then any optimal solution to the following problem

$$P(R_{\text{UB}}) : \min_{\alpha \in \mathbb{R}^q, \Theta \in \mathbb{R}^{n \times p}, R \in \mathbb{R}_+} F(\alpha, \Theta, R) \text{ s.t. } R_{\text{UB}} \geq R \geq \|\Theta\|_*, \quad (9)$$

will also be optimal to (5). For example, $(\hat{\alpha}, \hat{\Theta}, \hat{R})$ with $\hat{R} = \|\hat{\Theta}\|_*$ is an optimal solution to (9). We have defined the problem as $P(R_{\text{UB}})$ to emphasize its dependence on the upper bound R_{UB} . Later we shall describe a simple strategy to estimate R_{UB} . We fix the set $\Xi \subseteq [n] \times [p]$ where $\Omega \subseteq \Xi$ is the target coordinate set for the low rank matrix $\hat{\Theta}$ that we are interested in.

Proposed Method A natural way to exploit structure in $P(R_{\text{UB}})$ is to apply coordinate gradient descent to update α and (Θ, R) separately. While the trace-norm constraint on (Θ, R) can be handled by the conditional gradient (CG) method [17], the ℓ_1 norm penalization on α is more efficiently tackled by the proximal gradient method in practice. In addition, we tighten the upper bound R_{UB} on-the-fly as the algorithm proceeds. The MCGD method goes as follows. At the t th iteration, we are given the previous iterate $(\alpha^{(t-1)}, \Theta^{(t-1)}, R^{(t-1)})$ and the upper bound $R_{\text{UB}}^{(t)}$ is computed. The first block α is updated with a proximal gradient step:

$$\begin{aligned} \alpha^{(t)} &= \text{prox}_{\gamma \lambda_S \|\cdot\|_1}(\alpha^{(t-1)} - \gamma \nabla_\alpha \mathcal{L}(\alpha^{(t-1)}, \Theta^{(t-1)})) \\ &= \mathbb{T}_{\gamma \lambda_S}(\alpha^{(t-1)} - \gamma \nabla_\alpha \mathcal{L}(\alpha^{(t-1)}, \Theta^{(t-1)})). \end{aligned} \quad (10)$$

In (10), $\nabla_\alpha \mathcal{L}(\cdot)$ is the gradient of the log-likelihood function taken w.r.t. α , $\gamma > 0$ is a pre-defined step size parameter and $\mathbb{T}_\lambda(\mathbf{x}) := \text{sign}(\mathbf{x}) \odot (\mathbf{x} - \lambda \mathbf{1})_+$ is the component-wise soft thresholding operator. Alternatively, we can exactly solve the problem

$$\alpha^{(t)} \in \arg \min_{\alpha \in \mathbb{R}^q} F(\alpha, \Theta^{(t-1)}, R^{(t-1)}), \quad (11)$$

for which closed-form solution can be obtained in certain special cases (see below).

The second block (Θ, R) is updated with a CG step

$$(\Theta^{(t)}, R^{(t)}) = (\Theta^{(t-1)}, R^{(t-1)}) + \beta_t (\hat{\Theta}^{(t)} - \Theta^{(t-1)}, \hat{R}^{(t)} - R^{(t-1)}), \quad (12)$$

where $\beta_t \in [0, 1]$ is a step size to be defined later. $(\hat{\Theta}^{(t)}, \hat{R}^{(t)})$ is a direction evaluated as

$$(\hat{\Theta}^{(t)}, \hat{R}^{(t)}) \in \arg \min_{Z, R} \langle Z, \nabla_\Theta \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)}) \rangle + \lambda_1 R \text{ s.t. } \|Z\|_* \leq R \leq R_{\text{UB}}^{(t)}, \quad (13)$$

and $\nabla_\Theta \mathcal{L}(\cdot)$ is the gradient of $\mathcal{L}(\cdot)$ taken w.r.t. Θ . If $(\Theta^{(t-1)}, R^{(t-1)})$ is feasible to $P(R_{\text{UB}}^{(t)})$, then $(\Theta^{(t)}, R^{(t)})$ must also be feasible to $P(R_{\text{UB}}^{(t)})$. Furthermore, if we let $\mathbf{u}_1, \mathbf{v}_1$ be the top left and right singular vectors of the gradient matrix $\nabla_\Theta \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)})$ and $\sigma_1(\nabla_\Theta \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)}))$ be the top singular value, then $(\hat{\Theta}^{(t)}, \hat{R}^{(t)})$ admits a simple closed form solution:

$$(\hat{\Theta}^{(t)}, \hat{R}^{(t)}) = \begin{cases} (\mathbf{0}, 0), & \text{if } \lambda_L \geq \sigma_1(\nabla_\Theta \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)})), \\ (-R_{\text{UB}}^{(t)} \mathbf{u}_1 \mathbf{v}_1^\top, R_{\text{UB}}^{(t)}), & \text{if } \lambda_L < \sigma_1(\nabla_\Theta \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)})). \end{cases} \quad (14)$$

Lastly, the step size β_t is determined by:

$$\beta_t = \min \left\{ 1, \frac{\langle \Theta^{(t-1)} - \hat{\Theta}^{(t)}, \nabla_{\Theta} \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)}) \rangle + \lambda_L (R^{(t-1)} - \hat{R}^{(t)})}{\sigma_{\Theta} \|\mathcal{P}_{\Omega}(\hat{\Theta}^{(t)} - \Theta^{(t-1)})\|_{\mathbb{F}}^2} \right\}. \quad (15)$$

The step size strategy ensures decrease in the objective value between successive iterations. This is essential for establishing convergence of the proposed method [cf. Theorem 2]. We remark that the arithmetics in the MCGD method are not affected when we restrict the update of $\Theta^{(t)}$ in (12) to the entries in Ξ only. This is due to $\mathcal{L}(\mathbf{X}) = \mathcal{L}(\mathcal{P}_{\Omega}(\mathbf{X}))$ and the CG update direction (13) only involves the gradient of $\nabla_{\Theta} \mathcal{L}(\alpha^{(t)}, \Theta^{(t-1)})$ w.r.t. entries of Θ in Ω , where $\Omega \subseteq \Xi$.

Computing the Upper Bound $R_{\text{UB}}^{(t)}$ We describe a strategy for computing a valid upper bound $R_{\text{UB}}^{(t)}$ for \hat{R} and $\|\hat{\Theta}\|_{\star}$ during the updates in the MCGD method. Let us assume that:

H8 For all Θ and α , we have $\mathcal{L}(\alpha, \Theta) \geq 0$.

The above can be enforced as the log-likelihood function is lower bounded [cf. H4]. From (5) and using the above assumption, it is obvious that

$$F_0(\mathbf{0}, \mathbf{0}) = \mathcal{L}(\mathbf{0}, \mathbf{0}) \geq \mathcal{L}(\hat{\alpha}, \hat{\Theta}) + \lambda_S \|\hat{\alpha}\|_1 + \lambda_L \|\hat{\Theta}\|_{\star} \geq \lambda_L \|\hat{\Theta}\|_{\star}, \quad (16)$$

and thus $R_{\text{UB}}^0 := \lambda_L^{-1} \mathcal{L}(\mathbf{0} + f_U(\mathbf{0}))$ is a valid upper bound to $\|\hat{\Theta}\|_{\star}$; furthermore it can be tightened as we progress in the MCGD method. In particular, observe that $(\hat{\alpha}, \hat{\Theta}, \hat{R})$ with $\hat{R} = \|\hat{\Theta}\|_{\star}$ is an optimal solution to $\text{P}(R_{\text{UB}}^0)$, we have

$$F(\alpha, \Theta, R) \geq F(\hat{\alpha}, \hat{\Theta}, \hat{R}) = \mathcal{L}(\hat{\alpha}, \hat{\Theta}) + \lambda_S \|\hat{\alpha}\|_1 + \lambda_L \hat{R} \geq \lambda_L \hat{R}. \quad (17)$$

In other words, for all feasible (α, Θ, R) to $\text{P}(R_{\text{UB}}^0)$, $\lambda_L^{-1} F(\alpha, \Theta, R)$ is an upper bound to \hat{R} and $\|\hat{\Theta}\|_{\star}$. The above motivates us to select $R_{\text{UB}}^{(t)} := \lambda_L^{-1} F(\alpha^{(t)}, \Theta^{(t-1)}, R^{(t-1)})$ at iteration t , where we observe that $R_{\text{UB}}^{(t)} \geq R^{(t-1)}$. That is, $(\alpha^{(t)}, \Theta^{(t-1)}, R^{(t-1)})$ is feasible to both $\text{P}(R_{\text{UB}}^{(t)})$ and $\text{P}(R_{\text{UB}}^{(t-1)})$. Lastly, we summarize the MCGD method in Algorithm 1.

Computation Complexity Consider the MCGD method in Algorithm 1. Observe that line 3 requires computing the gradient w.r.t. α which involves $|\Omega|q$ Floating Points Operations (FLOPS) and the soft thresholding operator involves $\mathcal{O}(q)$ FLOPS. As the log-likelihood function $\mathcal{L}(\cdot)$ is evaluated element-wisely on Θ , evaluating the objective value and the derivative w.r.t. Θ requires $\mathcal{O}(|\Omega|)$ FLOPS. As such, line 4 can be evaluated in $\mathcal{O}(|\Omega|)$ FLOPS and line 5 requires $\mathcal{O}(|\Omega| \max\{n, p\} \log(1/\delta))$ FLOPS where the additional complexity is due to the top SVD computation and δ is a preset accuracy level of SVD computation. Lastly, line 6 requires $\mathcal{O}(|\Xi|)$ FLOPS since we only need to update the entries of Θ in Ξ [cf. see the remark after (15)]. The overall per-iteration complexity is $\mathcal{O}(|\Xi| + |\Omega|(\max\{n, p\} \log(1/\delta) + q))$.

Algorithm 1 MCGD Method for (9).

-
- 1: **Initialize:** — $\Theta^{(0)}, \alpha^{(0)}, R^{(0)}$. E.g., $\Theta^{(0)}, \alpha^{(0)}, R^{(0)} = (\mathbf{0}, \mathbf{0}, 0)$.
 - 2: **for** $t = 1, 2, \dots, T$ **do**
 - 3: // *Update for α* //
Compute the proximal update using (10) [or exact update via (11)] to obtain $\alpha^{(t)}$.
 - 4: // *Update for (Θ, R)* //
Compute the upper bound as $R_{\text{UB}}^{(t)} := \lambda_L^{-1} F(\alpha^{(t)}, \Theta^{(t-1)}, R^{(t-1)})$.
 - 5: Compute the update direction, $(\hat{\Theta}^{(t)}, \hat{R}^{(t)})$, using Eq. (14).
 - 6: Compute the CG update using (12), where the step size β_t is set as Eq. (15).
 - 7: **end for**
 - 8: **Return:** $\Theta^{(T)}, \alpha^{(T)}, R^{(T)}$.
-

From the above, the per-iteration computation complexity of the MCGD method scales linearly with the problem dimension $\max\{n, p\}$ and $|\Omega|$. This is comparable to [27, 11], where the former focuses only on the least square loss case. The following theorem, whose proof can be found in Appendix C, shows that the MCGD method converges at a sublinear rate.

Theorem 2 Assume H7 and H8. Define the quantity

$$C(t) := \max \left\{ \frac{24(Q^{(t)})^2}{\gamma}, \frac{24\hat{\sigma}_{\Theta}^2(Q^{(t)})^2}{\sigma_{\Theta}} + \max\{6R_{\text{UB}}^{(t)}(\lambda_L + M^{(t)}), 24\sigma_{\Theta}(R_{\text{UB}}^{(t)})^2\} \right\}, \quad (18)$$

where $Q^{(t)} := \lambda_S^{-1} F(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\Theta}^{(t)}, R^{(t)})$, $M^{(t)} := \|\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\Theta}^{(t-1)})\|_2$ and $R_{\text{UB}}^{(t)} := \lambda_L^{-1} F(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\Theta}^{(t-1)}, R^{(t-1)})$. If we choose the step sizes as $\gamma \leq 1/\sigma_{\boldsymbol{\alpha}}$ and β_t as in (15), then (i) the above quantity is upper bounded as $C(t) \leq \bar{C}$ for all $t \geq 1$, where

$$\bar{C} := \max \left\{ \frac{24(Q^{(0)})^2}{\gamma}, \frac{24\hat{\sigma}_{\boldsymbol{\Theta}}^2(Q^{(0)})^2}{\sigma_{\boldsymbol{\Theta}}} + \max\{6R_{\text{UB}}^{(0)}(\lambda_L + \bar{M}), 24\sigma_{\boldsymbol{\Theta}}(R_{\text{UB}}^{(0)})^2\} \right\}, \quad (19)$$

such that \bar{M} is an upper bound to $M^{(t)}$, and (ii) the MCGD method converges to an ϵ -optimal solution to (5) in T iterations, i.e., $F_0(\boldsymbol{\alpha}^{(T)}, \boldsymbol{\Theta}^{(T)}) - F_0(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}) \leq \epsilon$, where

$$T \geq \bar{C}(T) \left(\frac{1}{\epsilon} - \frac{1}{F_0(\boldsymbol{\alpha}^{(0)}, \boldsymbol{\Theta}^{(0)}) - F_0(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}})} \right)_+ \quad \text{with} \quad \bar{C}(T) := \left(\frac{1}{T} \sum_{t=1}^T \frac{1}{C(t)} \right)^{-1}. \quad (20)$$

In particular, as $\bar{C}(T) \leq \bar{C}$, at most $\bar{C}(\epsilon^{-1} - (F_0(\boldsymbol{\alpha}^{(0)}, \boldsymbol{\Theta}^{(0)}) - F_0(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}))^{-1})_+$ iterations are required for the MCGD method to reach an ϵ -optimal solution to (5).

Detailed Comparison to Prior Algorithms Previous contributions have focused on the special case of (5) where $q = np$, the dictionary $(\mathbf{X}(1), \dots, \mathbf{X}(q))$ is the canonical basis of $\mathbb{R}^{n \times p}$, and the link functions are quadratic. In this particular case, (5) becomes the estimation problem solved in sparse plus low-rank matrix decomposition. Popular examples are the alternating direction method of multiplier [25, 29] or the projected gradient method on a reformulated problem [7]. These methods either require computing a complete SVD or knowing the optimal rank number of $\boldsymbol{\Theta}$ a priori. When $n, p \gg 1$, it is computationally prohibitive to evaluate the complete SVD since each iteration would require $\mathcal{O}(\max\{n^2 p, p^2 n\})$ FLOPS. Other related work rely on factorizing the low-rank component, yielding nonconvex problems [14]; see also [33] and references therein.

Similar to the development of MCGD, a natural alternative is to apply algorithms based on the CG (a.k.a. Frank-Wolfe) method [17], whose iterations only require the computation of a top SVD. The present work is closely related to the efforts in [27, 11] which focused on the quadratic setting. Mu et al. [27] combines the CG method with proximal update as a two-steps procedure; Garber et al. [11] combines a CD method with CG updates on both the sparse and low-rank components. The work in [11] is also related to [23, 4] which combine CD with CG updates for solving constrained problems, instead of penalized problems like (5). Sublinear convergence rates are proven for the above methods. Finally, Fithian and Mazumder [10] also suggested to apply CD on (5), yet the convergence properties were not discussed.

In fact, when the MCGD's result is specialized to the same setting as [27], our worst-case bound on iteration number computed with \bar{C} match the bound in [27]. As shown in the supplementary material, we have $C(t) \rightarrow C^*$, where C^* depends on the optimal objective value of (9) and is smaller than \bar{C} . Since the quantity $\bar{C}(T)$ in (20) is an average of $\{C(t)\}_{t=1}^T$, this implies that the MCGD method requires less number of iterations for convergence than that is required by [27]. Such reduction is possible due to the on-the-fly update for $R_{\text{UB}}^{(t)}$. Moreover, our analysis in Theorem 2 holds when the MCGD method is implemented with a few practical modifications.

Exact Partial Minimization for $\boldsymbol{\alpha}$ Consider the special case of (5) where the link functions are either quadratic or exponential and the dictionary matrices satisfy:

$$\text{supp}(\mathbf{X}(k)) \cap \text{supp}(\mathbf{X}(k')) = \emptyset, \quad k \neq k' \quad \text{and} \quad [\mathbf{X}(k)]_{i,j} = c_k, \quad \forall (i, j) \in \text{supp}(\mathbf{X}(k)). \quad (21)$$

In this case, the partial minimization (11) can be decoupled into q scalar optimizations involving one coordinate of $\boldsymbol{\alpha}$, which can be solved in closed form. Note that this modification to the MCGD method is supported by Theorem 2 and the sublinear convergence rate holds. On the contrary, closed form update of $\boldsymbol{\alpha}$ is not supported by prior works such as [27, 11, 23, 4].

Distributed MCGD Optimization Consider the case where the observed data entries are stored across K workers, each of them communicating with a central server. It is natural to distribute the MCGD optimization over these workers to offload computation burden, or for privacy protection. Formally, we divide Ω into K disjoint partitions such that $\Omega = \Omega_1 \cup \dots \cup \Omega_K$ and worker k holds Ω_k . In this way, $\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\Theta}) = \sum_{k=1}^K \mathcal{L}_k(\boldsymbol{\alpha}, \boldsymbol{\Theta})$, where $\mathcal{L}_k(\boldsymbol{\alpha}, \boldsymbol{\Theta})$ is defined by replacing the summation over Ω with Ω_k in (4). Clearly, when $\boldsymbol{\alpha}$ and $\mathcal{P}_{\Omega_k}(\boldsymbol{\Theta})$ are given to the k th worker, the worker will be

able to evaluate the *local* loss function and its gradient.

As shown in Appendix D, the MCGD method can be easily extended to utilize distributed computation. The proximal update in line 3 is replaced by the following procedure. First, the *local* gradients computed by the workers are aggregated, then the soft thresholding operation is performed at the central server. Meanwhile, as the CG update in line 5 essentially requires computing the top singular vectors of the gradient matrix $\nabla_{\Theta} L(\alpha, \Theta) = \sum_{k=1}^K \nabla_{\Theta} \mathcal{L}_k(\alpha, \mathcal{P}_{\Omega_k}(\Theta))$, the latter can be implemented through a distributed version of the power method exploiting the decomposable structure of the gradient, such as described in [34]. It only requires $\mathcal{O}(\log(1/\delta))$ power iterations to compute a top SVD solution of accuracy δ . Thus, for a sufficiently small $\delta > 0$, the overall per-iteration complexity of the distributed method at the t th iteration is reduced to $\mathcal{O}(|\Xi| + \max\{n, p\} \log(1/\delta))$ at the central server, and $\mathcal{O}(|\Omega_k|(\max\{n, p\} \log(1/\delta) + q))$ at the k th worker.

4 Numerical Experiments

Experimental Setup We first generate the target parameter \mathbf{M}^0 according to the LORIS model in (3). For the sparse additive effects component, we consider $q = pn/5$ where we set $(\mathbf{X}(k))_{ij} = 1$ if $j(n-1) + i \in \{5(k-1) + 1, \dots, 5k\}$. This models a categorical variable containing $n/5$ categories. Furthermore, the target sparse component α^0 has a sparsity level of 10%. For the low-rank component, the target parameter Θ^0 is generated as a rank-4 matrix formed by the outer product of random orthogonal vectors. Notice that due to the structure of sparse additive effects, the surveyed prior methods [25, 14, 7] cannot be applied directly.

Gaussian Design To compare our framework to a reasonable benchmark, we focus on a homogeneous setting with numerical data modeled with the quadratic link function $g(m) = m^2$. We set the regularization parameters λ_S and λ_L to the theoretical values given in Theorem 1. We compare our result with a common two-step procedure where the components α_{kj} are first estimated in a preprocessing step as the means of the variables taken by group; then Θ is estimated using the softImpute method proposed in [15]. The regularization parameter for [15] is set to the same value λ_L . We compare the results in terms of estimation error and computing time in Table 1, after letting the two methods converge to the same precision of 10^{-5} . We observe the two methods perform equally well in terms of estimating Θ . LORIS yields constant estimation errors of α^0 as the dimension increases and the support of α^0 is kept constant, contrary to the two-step procedure for which the estimation error of α^0 increases with the dimension. As expected, the two-step method is faster for small data sets, whereas for large data sizes LORIS is superior in computational time. The above results are consistent with our theoretical findings.

problem size ($n \times p$)	time (secs)		$\ \Theta^0 - \hat{\Theta}\ _F^2$		$\ \alpha^0 - \hat{\alpha}\ _2^2$	
	LORIS	two-step	LORIS	two-step	LORIS	two-step
150 × 30	0.17	0.02	52	52	1.8	3.0
1,500 × 300	13.8	10.7	175.5	234	0.95	17.1
15,000 × 300	130.2	136.6	675	720	0.95	16.2
15,000 × 3,000	348	528	2.7×10^3	2.6×10^3	2.34	180

Table 1: Comparison of proposed method with a two-step method in terms of computation time and estimation error for increasing dimensions (averaged over 10 experiments).

Survey data To test the efficacy of our framework with heterogeneous data, we examine a survey conducted by the French National Institute of Statistics (Insee: <http://www.insee.fr/>) concerning the hobbies of French people. The data set contains $n = 8,403$ individuals and $p = 19$ binary and quantitative variables, indicating whether or not the person has been involved in different activities (reading, fishing, etc.), the number of hours spent watching TV and the overall number of hobbies of the individuals. Individuals are grouped by age category (15 – 25, 25 – 35, etc.): this categorical variable is used as a predictor of the survey responses in the subsequent experiment. We introduce 30% of missing values in the data set, and compare the imputation error of LORIS with a mixed data model (using a quadratic loss for numeric columns, a logistic loss for binary columns and a Poisson loss for counts) and LORIS with a Gaussian data model, with the imputation error of softImpute. The

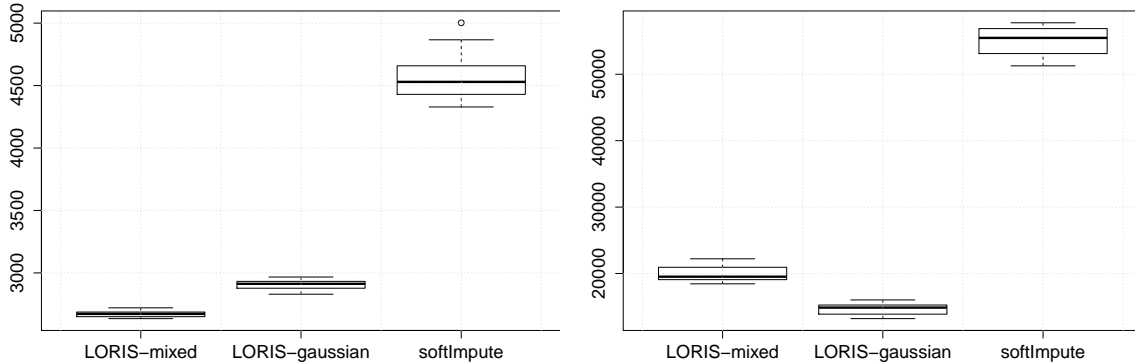


Figure 2: Imputation error of LORIS with mixed data model and Gaussian data model, and softImpute (10 replications) for categorical variables (left) and quantitative variables (right).

results are given in Figure 2 across 10 replications of the experiment, and show that, for this example, both LORIS models improve on the baseline softImpute by a factor 2. We also observe that modeling explicitly the binary variables leads to better imputation.

Finally, we apply LORIS with a mixed data model to the original data set. A subset of the resulting α vector is given in Table 2. There is a coefficient in α_{kj} for every age category k and every variable j . The coefficients in Table 2 indicate that young individuals engage in activities such as music and sport more than older people, and the opposite trend for collecting, knitting and fishing. Some coefficients are set to zero, indicating the absence of effect of the age category on the variable. We also observe that younger people engage overall in more activities than older people.

Age category	Music	Sport	Collecting	Mechanic	Knitting	Fishing	Nb activities
25-35	2.2	0.4	-2.1	0	-1.7	-1.9	10.0
35-45	2.0	0.3	-2.7	0	-2.3	-2.3	13.0
45-55	1.1	-0.8	-2.1	0	-2.7	-2.7	13.8
55-65	0	-2.2	-1.9	0	-1.0	-1.6	8.8
65-75	0	-2.1	-1.4	-1.1	-0.7	-1.3	5.5
75-85	-0.1	-0.9	-0.6	-0.5	-0.1	-0.6	2.2

Table 2: Estimated age category effects (α).

Conclusion In this paper, we proposed a new framework for handling large data frames with heterogeneous data and missing values which incorporates additive effects. It consists of a doubly penalized quasi-maximum likelihood estimator and a new optimization algorithm to implement the estimator. We examined both the statistical and computational efficiency of the framework and derived worst case bounds of its performance. Future work includes the incorporation of qualitative features with more than two categories and of missing values in the dictionary matrices.

5 Acknowledgement

The authors would like to thank for the useful comments from three anonymous reviewers. HTW’s work was supported by the grant NSF CCF-BSF 1714672.

References

- [1] A. Agresti. *Categorical Data Analysis, 3rd Edition*. Wiley, 2013.
- [2] J.-P. Aubin and I. Ekeland. *Applied nonlinear analysis*. Pure and applied mathematics. John Wiley, New-York, 1984. A Wiley-Interscience publication.

- [3] A. Beck and L. Tretuashvili. On the convergence of block coordinate descent type methods. *SIAM journal on Optimization*, 23(4):2037–2060, 2013.
- [4] A. Beck, E. Pauwels, and S. Sabach. The cyclic block conditional gradient method for convex optimization problems. *SIAM Journal on Optimization*, 25(4):2024–2049, 2015.
- [5] E. J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *J. ACM*, 58(3):11:1–11:37, June 2011. ISSN 0004-5411. doi: 10.1145/1970392.1970395. URL <http://doi.acm.org/10.1145/1970392.1970395>.
- [6] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky. Rank-sparsity incoherence for matrix decomposition. *SIAM Journal on Optimization*, 21(2):572–596, 2011. doi: 10.1137/090761793. URL <https://doi.org/10.1137/090761793>.
- [7] Y. Chen and M. J. Wainwright. Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. *CoRR*, abs/1509.03025, 2015.
- [8] J. de Leeuw. Principal component analysis of binary data by iterated singular value decomposition. *Comput. Stat. Data Anal.*, 50(1):21–39, Jan. 2006. ISSN 0167-9473. doi: 10.1016/j.csda.2004.07.010. URL <http://dx.doi.org/10.1016/j.csda.2004.07.010>.
- [9] A. Feuerwerker, Y. He, and S. Khatri. Statistical significance of the netflix challenge. *Statist. Sci.*, 27(2): 202–231, 05 2012. doi: 10.1214/11-STS368. URL <http://dx.doi.org/10.1214/11-STS368>.
- [10] W. Fithian and R. Mazumder. Flexible Low-Rank Statistical Modeling with Missing Data and Side Information. *Statistical Science*, 33(2):238–260, 2018.
- [11] D. Garber, S. Sabach, and A. Kaplan. Fast generalized conditional gradient method with applications to matrix recovery problems. *arXiv preprint arXiv:1802.05581*, 2018.
- [12] G. Gidel, F. Pedregosa, and S. Lacoste-Julien. Frank-wolfe splitting via augmented lagrangian method. In *OPTML 2017: 10th NIPS Workshop on Optimization for Machine Learning (NIPS 2017)*, page 21, 2017. URL http://opt-ml.org/papers/OPT2017_paper_21.pdf.
- [13] G. H. Golub and C. F. Van Loan. *Matrix computations*, volume 3. JHU Press, 2012.
- [14] Q. Gu, Z. W. Wang, and H. Liu. Low-rank and sparse structure pursuit via alternating minimization. In A. Gretton and C. C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 600–609, Cadiz, Spain, 09–11 May 2016. PMLR. URL <http://proceedings.mlr.press/v51/gu16.html>.
- [15] T. Hastie, R. Mazumder, J. Lee, and R. Zadeh. Matrix Completion and Low-Rank SVD via Fast Alternating Least Squares. *The Journal of Machine Learning Research*, 16:3367–3402, jan 2015.
- [16] D. Hsu, S. M. Kakade, and T. Zhang. Robust matrix decomposition with sparse corruptions. *EEE Transactions on Information Theory*, 57(11):7221–7234, 2011.
- [17] M. Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *ICML (1)*, pages 427–435, 2013.
- [18] H. A. L. Kiers. Simple structure in component analysis techniques for mixtures of qualitative and quantitative variables. *Psychometrika*, 56(2):197–212, Jun 1991. ISSN 1860-0980. doi: 10.1007/BF02294458. URL <https://doi.org/10.1007/BF02294458>.
- [19] O. Klopp. Noisy low-rank matrix completion with general sampling distribution. *Bernoulli*, 20(1):282–303, 2014.
- [20] O. Klopp. Matrix completion by singular value thresholding: sharp bounds. *Electronic journal of statistics*, 9(2):2348–2369, 2015. URL <https://hal.archives-ouvertes.fr/hal-01111757>.
- [21] O. Klopp, K. Lounici, and A. B. Tsybakov. Robust matrix completion. *Probability Theory and Related Fields*, 169(1):523–564, Oct 2017. doi: 10.1007/s00440-016-0736-y. URL <https://doi.org/10.1007/s00440-016-0736-y>.
- [22] N. K. Kumar and J. Schneider. Literature survey on low rank approximation of matrices. *Linear and Multilinear Algebra*, 65(11):2212–2244, 2017. doi: 10.1080/03081087.2016.1267104. URL <https://doi.org/10.1080/03081087.2016.1267104>.

- [23] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate frank-wolfe optimization for structural svms. In *Proceedings of the 30th International Conference on International Conference on Machine Learning-Volume 28*, pages 1–53. JMLR. org, 2013.
- [24] A. J. Landgraf and Y. Lee. Generalized principal component analysis: Projection of saturated model parameters. Technical report, The Ohio State University, Department of Statistics, 06 2015.
- [25] Z. Lin, R. Liu, and Z. Su. Linearized alternating direction method with adaptive penalty for low-rank representation. In *Advances in neural information processing systems*, pages 612–620, 2011.
- [26] L. T. Liu, E. Dobriban, and A. Singer. *e* pca: High dimensional exponential family pca. *Annals of Applied Statistics*, to appear, 2018.
- [27] C. Mu, Y. Zhang, J. Wright, and D. Goldfarb. Scalable robust matrix recovery: Frank–wolfe meets proximal methods. *SIAM Journal on Scientific Computing*, 38(5):A3291–A3317, 2016. doi: 10.1137/15M101628X. URL <https://doi.org/10.1137/15M101628X>.
- [28] J. Pagès. *Multiple factor analysis by example using R*. Chapman and Hall/CRC, 2014.
- [29] M. Tao and X. Yuan. Recovering low-rank and sparse components of matrices from incomplete and noisy observations. *SIAM Journal on Optimization*, 21(1):57–81, 2011.
- [30] M. Udell, C. Horn, R. Zadeh, and S. Boyd. Generalized low rank models. *Foundations and Trends in Machine Learning*, 9(1), 2016. doi: 10.1561/22000000055. URL <http://dx.doi.org/10.1561/22000000055>.
- [31] G. A. Watson. Characterization of the subdifferential of some matrix norms. *Linear algebra and its applications*, 170:33–45, 1992.
- [32] H. Xu, C. Caramanis, and S. Sanghavi. Robust pca via outlier pursuit. In *Proceedings of the 23rd International Conference on Neural Information Processing Systems, NIPS’10*, pages 2496–2504, USA, 2010. Curran Associates Inc. URL <http://dl.acm.org/citation.cfm?id=2997046.2997174>.
- [33] X. Zhang, L. Wang, and Q. Gu. A unified framework for nonconvex low-rank plus sparse matrix recovery. In *International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*, pages 1097–1107, 2018. URL <http://proceedings.mlr.press/v84/zhang18c.html>.
- [34] W. Zheng, A. Bellet, and P. Gallinari. A distributed frank-wolfe framework for learning low-rank matrices with the trace norm. *arXiv preprint arXiv:1712.07495*, 2017.

A Statistical guarantees

A.1 Main result

We recall the convergence rates for the Frobenius norm of the errors $\Delta\Theta = \hat{\Theta} - \Theta^0$ and $\Delta\alpha = \hat{\alpha} - \alpha^0$ given in Section 2. Define $d_{\mathbf{X}} = \max_k \|\mathbf{X}(k)\|_1$ and the following quantities:

$$D_{\alpha} = \frac{\|\alpha^0\|_1 \log(n+p)}{\pi \sigma_-^2 \gamma} + \left(\frac{a}{\pi}\right)^2 \log(n+p),$$

$$D_{\Theta} = D_{\alpha} + d_{\mathbf{X}} \|\alpha^0\|_1 \left\{ \frac{12\pi \sqrt{\log(n+p)}}{\gamma(1+\nu)a\sigma_+ \sqrt{\beta}} + \frac{1}{\pi\sigma_-^2} \left(\frac{\log(n+p)}{\gamma} \right) + 1 \right\}.$$

We assume that $M = (n \vee p)$ is large enough, that is

$$M \geq \left\{ \frac{4\sigma_+^2}{\gamma^6} \log^2 \left(\frac{\sqrt{n \wedge p}}{p\gamma\sigma_-} \right) \vee 2 \exp(\sigma_+^2/\gamma^2 \vee \sigma_+^2\gamma(1+\nu a)) \right\}.$$

Theorem 3 Assume H1-6. Set

$$\lambda_L = 2C\sigma_+ \sqrt{\pi \max(n,p) \log(n+p)}, \text{ and } \lambda_S \geq 24 \max_k \|\mathbf{X}(k)\|_1 \log(n+p)/\gamma, \quad (22)$$

where C is a positive constant. Assume that $\max(n, p) \geq 4\sigma_+^2/\gamma^6 \log^2(\sqrt{\min(n, p)/(\pi\gamma\sigma_-)}) + 2\exp(\sigma_+^2/\gamma^2 + 2\sigma_+^2\gamma a)$. Then, with probability at least $1 - 9(n + p)^{-1}$,

$$\begin{aligned}\|\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^0\|_2^2 &\leq C_1 \frac{sd_{\mathbf{X}} \log(n + p)}{\kappa^2 \pi} + D_{\boldsymbol{\alpha}}, \\ \|\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^0\|_F^2 &\leq C_2 \left(\frac{r \max(n, p)}{\pi} + \frac{sd_{\mathbf{X}}}{\pi} \right) \log(n + p) + D_{\boldsymbol{\Theta}}.\end{aligned}\tag{23}$$

In (23), $s := \|\boldsymbol{\alpha}^0\|_0$, $r := \text{rank}(\boldsymbol{\Theta}^0)$. C_1 and C_2 are positive constants and $D_{\boldsymbol{\alpha}}$ and $D_{\boldsymbol{\Theta}}$ are residuals of lower order whose exact values are given in Appendix A.

Denoting by \lesssim the inequality up to constant and logarithmic factors, the order of magnitude of the bounds are therefore:

$$\begin{aligned}\|\Delta_{\boldsymbol{\alpha}}\|_2^2 &\lesssim \frac{sd_{\mathbf{X}}}{p\kappa^2}, \\ \|\Delta_{\boldsymbol{\Theta}}\|_F^2 &\lesssim \frac{r\beta}{p^2} + \frac{sd_{\mathbf{X}}}{p},\end{aligned}$$

where $s = \|\boldsymbol{\alpha}^0\|_0$ and $r = \text{rank}(\boldsymbol{\Theta}^0)$. In the case of almost uniform sampling, i.e. $c_1\pi \leq \pi_{ij} \leq c_2\pi$ for all $(i, j) \in [n] \times [p]$ and two positive constants c_1 and c_2 , we obtain that $\beta \leq c_2(n \vee p)\pi$, which yields the following simplified bound:

$$\|\Delta_{\boldsymbol{\Theta}}\|_F^2 \lesssim \frac{rM}{\pi} + \frac{sd_{\mathbf{X}}}{\pi}.\tag{24}$$

The rate given in (24) is the sum of the usual low-rank convergence rate rM/p and, when $d_{\mathbf{X}}$ is a constant, of the usual sparse vector convergence rate.

A.2 Sketch of the proof

Let $\{\epsilon_{ij}\}$ be an i.i.d. Rademacher sequence independent of Y and Ω . We define

$$\Sigma_R = \sum_{i=1}^n \sum_{j=1}^p \omega_{ij} \epsilon_{ij} E_{ij}.$$

In Theorem 4 we give a general result under some assumptions on the regularization parameters λ_L and λ_S , which depend on the random matrices $\nabla \mathcal{L}(\mathbf{M}^0)$ and Σ_R . Then, Lemma 4 and 5 allow us to compute values of λ_L and λ_S that satisfy the assumptions of Theorem 4 with high probability. Finally we combining these results yield Theorem 3. Define

$$\Psi_{\boldsymbol{\alpha}} = \frac{\|\boldsymbol{\alpha}^0\|_1}{\pi} \left\{ \frac{\lambda_S}{\sigma_-^2} + a^2 d_{\mathbf{X}} \mathbb{E} \|\Sigma_R\|_{\infty} \right\} + \left(\frac{a}{\pi} \right)^2 \log(n + p),\tag{25}$$

$$\Psi_{\boldsymbol{\Theta}} = \frac{r}{\pi^2} \mathbb{E} \|\Sigma_R\|^2 + \frac{\|\boldsymbol{\alpha}\|_1}{\pi} \left\{ \frac{\lambda_S}{(1 + \nu)a\lambda_L} + d_{\mathbf{X}} \mathbb{E} \|\Sigma_R\|_{\infty} \right\} + \Psi_{\boldsymbol{\alpha}}.\tag{26}$$

Theorem 4 *Let*

$$\lambda_L \geq 2 \|\nabla \mathcal{L}(\mathbf{M}^0)\|, \quad \lambda_S \geq 2d_{\mathbf{X}} \left(\|\nabla \mathcal{L}(\mathbf{M}^0)\|_{\infty} + 2\sigma_+^2(1 + \nu)a \right),$$

and assumptions H 2-6 hold. Then, with probability at least $1 - 8(n + p)^{-1}$

$$\begin{aligned}(i) \quad &\|\Delta_{\boldsymbol{\alpha}}\|_2^2 \leq \frac{C}{\kappa^2} \Psi_{\boldsymbol{\alpha}}, \text{ and} \\ (ii) \quad &\|\Delta_{\boldsymbol{\Theta}}\|_F^2 \leq C \left\{ \frac{r\lambda_L^2}{\pi^2\sigma_-^4} + (1 + \nu)a\Psi_{\boldsymbol{\Theta}} \right\}.\end{aligned}\tag{27}$$

Denote $\Delta \mathbf{M} = \hat{\mathbf{M}} - \mathbf{M}^0$. We first derive an upper bound on the Frobenius error restricted to the observed entries $\|\mathcal{P}_{\Omega}(\Delta \mathbf{M})\|_F^2$. Then we show some restricted strong convexity property, meaning that $\mathbb{E} \|\mathcal{P}_{\Omega}(\Delta \mathbf{M})\|_F^2$ is upper bounded by $\|\mathcal{P}_{\Omega}(\Delta \mathbf{M})\|_F^2$ up to a residual term defined later.

Upper bound on $\|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2$. By definition of $\hat{\Theta}$ and $\hat{\alpha}$:

$$\mathcal{L}(\hat{\mathbf{M}}) - \mathcal{L}(\mathbf{M}^0) \leq \lambda_L \left(\|\Theta^0\|_* - \|\hat{\Theta}\|_* \right) + \lambda_S \left(\|\alpha^0\|_1 - \|\hat{\alpha}\|_1 \right).$$

Recall that, for $\alpha \in \mathbb{R}^q$, we use the notation $f_U(\alpha) = \sum_{k=1}^q \alpha_k \mathbf{X}(k)$. Adding $\langle \nabla \mathcal{L}(\mathbf{M}^0), \Delta\mathbf{M} \rangle$ on both sides of the last inequality, we get

$$\begin{aligned} \mathcal{L}(\hat{\mathbf{M}}) - \mathcal{L}(\mathbf{M}^0) + \langle \nabla \mathcal{L}(\mathbf{M}^0), \Delta\mathbf{M} \rangle &\leq \lambda_L \left(\|\Theta^0\|_* - \|\hat{\Theta}\|_* \right) - \langle \nabla \mathcal{L}(\mathbf{M}^0), \Delta\Theta \rangle \\ &\quad + \lambda_S \left(\|\alpha^0\|_1 - \|\hat{\alpha}\|_1 \right) - \langle \nabla \mathcal{L}(\mathbf{M}^0), f_U(\Delta\alpha) \rangle. \end{aligned}$$

The strong convexity of the link functions g_j , $j \in [p]$, allows us to lower bound the left hand side term and obtain

$$\begin{aligned} \frac{\sigma^2}{2} \|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 &\leq \lambda_L \left(\|\Theta^0\|_* - \|\hat{\Theta}\|_* \right) - \langle \nabla \mathcal{L}(\mathbf{M}^0), \Delta\Theta \rangle \\ &\quad + \lambda_S \left(\|\alpha^0\|_1 - \|\hat{\alpha}\|_1 \right) - \langle \nabla \mathcal{L}(\mathbf{M}^0), f_U(\Delta\alpha) \rangle. \end{aligned}$$

We now upper bound the right hand side using the following three arguments: the duality of the norms $\|\cdot\|_*$ and $\|\cdot\|$ on the one hand and of the norms $\|\cdot\|_1$ and $\|\cdot\|_\infty$ on the other hand, the triangular inequality and the following assumptions:

$$\lambda_L \geq 2 \|\nabla \mathcal{L}(\mathbf{M}^0)\|, \quad \lambda_S \geq 2 \|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty d_{\mathbf{X}}.$$

We obtain

$$\|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 \leq \frac{3\lambda_L}{\sigma^2} \sqrt{2 \text{rank}(\mathbf{M}^0)} \|\Delta\Theta\|_F + \frac{3\lambda_S}{\sigma^2} \|\alpha^0\|_1. \quad (28)$$

Restricted strong convexity We now show that when the errors $\Delta\Theta$ and $\Delta\alpha$ belong to a subspace \mathcal{C} and for a residual D - both defined later on - the following holds with high probability:

$$\|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 \geq \mathbb{E} \|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 - D. \quad (29)$$

We start by defining the set \mathcal{C} and prove that it contains the errors $\Delta\Theta$ and $\Delta\alpha$ with high probability (Lemma 1-2); then we show that restricted strong convexity holds on this subspace (Lemma 3).

For non-negative constants d_1 , d_Π , $\rho < m$ and ε that will be specified later on, define the two following sets:

$$\mathcal{A}(d_1, d_\Pi) = \left\{ \alpha \in \mathbb{R}^q : \|\alpha\|_1 \leq d_1, \|\mathcal{P}_\Omega(f_U(\alpha))\|_F^2 \leq d_\Pi \right\}. \quad (30)$$

The constants d_1 and d_Π define the constraints on the ℓ_1 norm of α and weighted Frobenius norm of $f_U(\alpha)$.

$$\begin{aligned} \mathcal{L}(\rho, \varepsilon) &= \left\{ \Theta \in \mathbb{R}^{n \times p}, \alpha \in \mathbb{R}^q : \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 \geq \frac{72 \log(n+p)}{\pi \log(6/5)}, \right. \\ &\quad \left. \|\Theta + f_U(\alpha)\|_\infty \leq 1, \|\Theta\|_* \leq \sqrt{\rho} \|\Theta\|_F + \varepsilon \right\} \end{aligned} \quad (31)$$

Condition $\|\Theta\|_* \leq \sqrt{\rho} \|\Theta\|_F + \varepsilon$ is a relaxed form of the condition $\|\Theta\|_* \leq \sqrt{\rho} \|\Theta\|_F$ satisfied for matrices of rank ρ . Finally, we define the constrained set of interest:

$$\mathcal{C}(d_1, d_\Pi, \rho, \varepsilon) = \mathcal{L}(\rho, \varepsilon) \cap \{\mathbb{R}^{n \times p} \times \mathcal{A}(d_1, d_\Pi)\}.$$

Let

$$\begin{aligned} d_1 &= 4 \|\alpha\|_1, \\ d_\Pi &= \frac{3\lambda_S}{\sigma^2} \|\alpha^0\|_1 + 64a^2 d_{\mathbf{X}} \mathbb{E} \|\Sigma_R\|_\infty \|\alpha\|_1 + 3072a^2 \pi^{-1} + \frac{72a^2 \log(n+p)}{\log(6/5)}. \end{aligned}$$

The following Lemma, proved in Appendix B.1 states that with high probability, $\Delta\alpha \in \mathcal{A}(d_1, d_\Pi)$.

Lemma 1 Let $\lambda_S \geq 2d_{\mathbf{X}} (\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + d_{\mathbf{X}})a)$ and assume **H 2-6** hold. Then, with probability at least $1 - 8(n + p)^{-1}$,

$$\Delta \boldsymbol{\alpha} \in \mathcal{A}(d_1, d_\Pi);$$

Lemma 1 (proved in Appendix B.2) implies (i) of Theorem 4. Thus, we only need to prove (ii).

Lemma 2 Let

$$\lambda_L \geq 2 \|\nabla \mathcal{L}(\mathbf{M}^0)\|, \quad \lambda_S \geq 2d_{\mathbf{X}} (\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + d_{\mathbf{X}})a),$$

and assumption **H 4** hold. Then, for $\rho = 32r$ and $\varepsilon = 3\lambda_S/\lambda_L \|\boldsymbol{\alpha}^0\|_1$,

$$\|\Delta \boldsymbol{\Theta}\|_* \leq \sqrt{\rho} \|\Delta \boldsymbol{\Theta}\|_F + \varepsilon.$$

A proof of Lemma 2 can be found in Appendix B.2. As a consequence, under the conditions on the regularization parameters λ_L and λ_S given in Lemma 2 and whenever

$$\mathbb{E} \|\mathcal{P}_\Omega(\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha}))\|_F^2 \geq \frac{72 \log(n + p)}{\pi \log(6/5)},$$

the error terms $(\Delta \boldsymbol{\Theta}, \Delta \boldsymbol{\alpha})$ belong to the constrained set $\mathcal{C}(d_1, d_\Pi, \rho, \varepsilon)$ with high probability. We therefore consider the two possible cases: $\mathbb{E} \|\mathcal{P}_\Omega(\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha}))\|_F^2 < \frac{72 \log(n + p)}{\pi \log(6/5)}$ and $\mathbb{E} \|\mathcal{P}_\Omega(\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha}))\|_F^2 \geq \frac{72 \log(n + p)}{\pi \log(6/5)}$.

Case 1: Suppose $\mathbb{E} \|\mathcal{P}_\Omega(\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha}))\|_F^2 < \frac{72 \log(n + p)}{\pi \log(6/5)}$. Then, Lemma 1 combined with the fact that $\|\mathbf{M}\|_F^2 \leq \pi^{-1} \|\mathcal{P}_\Omega(\mathbf{M})\|_F^2$ for all \mathbf{M} , and the identity $(a + b)^2 \geq a^2/4 - 4b^2$ ensures that

$$\|\Delta \boldsymbol{\Theta}\|_F^2 \leq 4 \|\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha})\|_F^2 + 16 \|f_U(\Delta \boldsymbol{\alpha})\|_F^2,$$

therefore

$$\|\Delta \boldsymbol{\Theta}\|_F^2 \leq \frac{288a^2 \log(n + p)}{\log(6/5)} + 16\Phi_\alpha,$$

which implies (ii) of Theorem 4.

Case 2: Suppose $\mathbb{E} \|\mathcal{P}_\Omega(\Delta \boldsymbol{\Theta} + f_U(\Delta \boldsymbol{\alpha}))\|_F^2 \geq \frac{72 \log(n + p)}{\pi \log(6/5)}$. Then, Lemma 1 and 2 yield that with probability at least $1 - 8(n + p)^{-1}$,

$$\left(\frac{\Delta \boldsymbol{\Theta}}{2(1 + \nu)a}, \frac{\Delta \boldsymbol{\alpha}}{2(1 + \nu)a} \right) \in \mathcal{C}(d'_1, d'_\Pi, \rho', \varepsilon'), \text{ with}$$

$$\begin{aligned} d'_1 &= \frac{d_1}{2(1 + \nu)a}, & d'_\Pi &= \frac{d_\Pi}{4(1 + \nu)^2 a^2}, \\ \rho' &= \rho, & \varepsilon' &= \frac{\varepsilon}{2(1 + \nu)a}, \end{aligned}$$

where d_1, d_Π, ρ and ε are defined in Lemma 1 and 2. We use the following result, proved in Appendix B.3. Define the set $\tilde{\mathcal{A}}(d_1)$ as follows:

$$\tilde{\mathcal{A}}(d_1) = \left\{ \boldsymbol{\alpha} \in \mathbb{R}^q : \|\boldsymbol{\alpha}\|_\infty \leq 1; \quad \|\boldsymbol{\alpha}\|_1 \leq d_1; \quad \|\mathcal{P}_\Omega(f_U \boldsymbol{\alpha})\|_F^2 \geq \frac{18 \log(n + p)}{\pi \log(6/5)} \right\}.$$

Let d_1, d_Π, ρ and ε be positive constants, and

$$\begin{aligned} D_\alpha &= 8\nu d_1 d_{\mathbf{X}} \mathbb{E} \|\Sigma_R\|_\infty + 768\pi^{-1}, \\ D_X &= \frac{112\rho}{\pi} \mathbb{E} \|\Sigma_R\|^2 + 8\nu\varepsilon \mathbb{E} \|\Sigma_R\| + 8\nu d_1 d_{\mathbf{X}} \mathbb{E} \|\Sigma_R\|_\infty + d_\Pi + 768\pi^{-1}. \end{aligned} \tag{32}$$

Lemma 3 Assume **H 6**. Then, the following properties hold:

(i) For any $\alpha \in \tilde{\mathcal{A}}(d_1)$, with probability at least $1 - 8(n+p)^{-1}$,

$$\|\mathcal{P}_\Omega(f_U(\alpha))\|_F^2 \geq \frac{1}{2} \mathbb{E} \|\mathcal{P}_\Omega(f_U(\alpha))\|_F^2 - D_\alpha.$$

(ii) For any pair $(\Theta, \alpha) \in \mathcal{C}(d_1, d_\Pi, \rho, \varepsilon)$, with probability at least $1 - 8(n+p)^{-1}$

$$\|\mathcal{P}_\Omega(\Delta\Theta + f_U(\Delta\alpha))\|_F^2 \geq \frac{1}{2} \mathbb{E} \|\mathcal{P}_\Omega(\Delta\Theta + f_U(\Delta\alpha))\|_F^2 - D_X. \quad (33)$$

Lemma 3 is proved in Appendix B.3. We apply Lemma 3 (ii) to $\left(\frac{\Delta\Theta}{2(1+\nu)a}, \frac{\Delta\alpha}{2(1+\nu)a}\right)$ which implies that with probability at least $1 - 8(n+p)^{-1}$, $\mathbb{E} \|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 \leq 2 \|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2 + 2(1+\nu)a\Psi_\Theta$. Combined with (28) and $\|\Delta\mathbf{M}\|_F^2 \leq \pi^{-1} \mathbb{E} \|\mathcal{P}_\Omega(\Delta\mathbf{M})\|_F^2$, it implies that

$$\|\Delta\mathbf{M}\|_F^2 \leq \frac{6\sqrt{2r}\lambda_L}{p\sigma_-^2} \|\Delta\Theta\|_F + \frac{6\lambda_S}{\pi\sigma_-^2} \|\alpha^0\|_1 + 2(1+\nu)a\Psi_\Theta.$$

Now using $\|\Delta\mathbf{M}\|_F^2 \geq \frac{\|\Delta\Theta\|_F^2}{2} - \|f_U(\Delta\alpha)\|_F^2$ and $\frac{6\sqrt{2r}\lambda_L}{\pi\sigma_-^2} \|\Delta\Theta\|_F \leq \frac{\|\Delta\Theta\|_F^2}{4} + \frac{288r\lambda_L^2}{p^2\sigma_-^4}$, we obtain

$$\|\Delta\Theta\|_F^2 \leq \frac{1152r\lambda_L^2}{p^2\sigma_-^4} + \frac{24\lambda_S}{\pi\sigma_-^2} \|\alpha^0\|_1 + 2(1+\nu)a\Psi_\Theta + 4\Psi_\alpha,$$

which gives the result of Theorem 4 (ii).

We now give deterministic upper bounds on $\mathbb{E}\|\Sigma_R\|$ and $\mathbb{E}\|\Sigma_R\|_\infty$, and probabilistic upper bounds on $\|\nabla\mathcal{L}(\mathbf{M}^0)\|$ and $\|\nabla\mathcal{L}(\mathbf{M}^0)\|_\infty$. We will use them to select values of λ_L and λ_S which satisfy the assumptions of Theorem 4 and compute the corresponding upper bounds.

Lemma 4 [21, Lemma 10] *Let assumption H 6 hold. Then, there exists an absolute constant C^* such that the two following inequalities hold*

$$\begin{aligned} \mathbb{E}\|\Sigma_R\|_\infty &\leq 1, \text{ and} \\ \mathbb{E}\|\Sigma_R\| &\leq C^* \left\{ \sqrt{\beta} + \sqrt{\log(\min(n,p))} \right\}. \end{aligned}$$

Lemma 5 [21, Lemma 10] *Let assumptions H 1-6 hold. Then, there exists an absolute constant c^* such that the following two inequalities hold with probability at least $1 - (n+p)^{-1}$.*

$$\|\nabla\mathcal{L}(\mathbf{M}^0)\|_\infty \leq 6 \max \left\{ \sigma_+ \sqrt{\log(n+p)}, \frac{\log(n+p)}{\gamma} \right\}, \quad (34)$$

$$\|\nabla\mathcal{L}(\mathbf{M}^0)\| \leq c^* \max \left\{ \sigma_+ \sqrt{\beta \log(n+p)}, \frac{\log(n+p)}{\gamma} \log \left(\frac{1}{\sigma_-} \sqrt{\frac{np}{\beta}} \right) \right\}. \quad (35)$$

From Theorem 4, Lemma 4 and 5 combined with a union bound argument, we deduce result given in Section 2.

B Technical results

B.1 Proof of Lemma 1

We start by proving $\|\Delta\alpha\|_1 \leq 4 \|\alpha^0\|_1$. By the optimality conditions over a convex set [2, Chapter 4, Section 2, Proposition 4], there exist two subgradients \hat{f}_Θ in the subdifferential of $\|\cdot\|_*$ taken at $\hat{\Theta}$ and \hat{f}_α in the subdifferential of $\|\cdot\|_1$ taken at $\hat{\alpha}$, such that for all feasible pairs (Θ, α) we have

$$\langle \nabla\mathcal{L}(\hat{\mathbf{M}}), \Theta - \hat{\Theta} + \sum_{k=1}^q (\alpha_k - \hat{\alpha}_k) \mathbf{X}(k) \rangle + \lambda_L \langle \hat{f}_\Theta, \Theta - \hat{\Theta} \rangle + \lambda_S \langle \hat{f}_\alpha, \alpha - \hat{\alpha} \rangle \geq 0. \quad (36)$$

Applying inequality (36) to the pair $(\hat{\Theta}, \alpha^0)$ we obtain

$$\langle \nabla \mathcal{L}(\hat{\mathbf{M}}), \sum_{k=1}^q \Delta \alpha_k \mathbf{X}(k) \rangle + \lambda_S \langle \hat{f}_\alpha, \Delta \alpha \rangle \geq 0.$$

Denote $\tilde{\mathbf{M}} = \hat{\Theta} + \sum_{k=1}^q \alpha_k^0 \mathbf{X}(k)$. The last inequality is equivalent to

$$\underbrace{\langle \nabla \mathcal{L}(\mathbf{M}^0), f_U(\Delta \alpha) \rangle}_{B_1} + \underbrace{\langle \nabla \mathcal{L}(\tilde{\mathbf{M}}) - \nabla \mathcal{L}(\mathbf{M}^0), f_U(\Delta \alpha) \rangle}_{B_2} + \underbrace{\langle \nabla \mathcal{L}(\hat{\mathbf{M}}) - \nabla \mathcal{L}(\tilde{\mathbf{M}}), f_U(\Delta \alpha) \rangle}_{B_3} + \lambda_S \langle \hat{f}_\alpha, \Delta \alpha \rangle \geq 0.$$

We now derive upper bounds on the three terms B_1 , B_2 and B_3 separately. Recall that we denote $d_{\mathbf{X}} = \max_k \|\mathbf{X}(k)\|_1$ and bound B_1 as follows:

$$B_1 \leq \|\Delta \alpha\|_1 \|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty d_{\mathbf{X}}. \quad (37)$$

Similarly, the duality between $\|\cdot\|_\infty$ and $\|\cdot\|_1$ gives

$$B_2 \leq \|\Delta \alpha\|_1 \|\nabla \mathcal{L}(\tilde{\mathbf{M}}) - \nabla \mathcal{L}(\mathbf{M}^0)\|_\infty d_{\mathbf{X}}.$$

Moreover, $\nabla \mathcal{L}(\tilde{\mathbf{M}}) - \nabla \mathcal{L}(\mathbf{M}^0)$ is a matrix with entries $g'_j(\tilde{\mathbf{M}}_{ij}) - g'_j(\mathbf{M}^0_{ij})$, therefore assumption **H 4** ensures

$$\|\nabla \mathcal{L}(\tilde{\mathbf{M}}) - \nabla \mathcal{L}(\mathbf{M}^0)\|_\infty \leq 2\sigma_+^2(1 + \nu)a,$$

and finally we obtain

$$B_2 \leq \|\Delta \alpha\|_1 2\sigma_+^2(1 + \nu)ad_{\mathbf{X}}. \quad (38)$$

We finally bound B_3 as follows. We have that

$$B_3 = \sum_{i=1}^n \sum_{j=1}^p \omega_{ij} \left(g'_j(\hat{\mathbf{M}}_{ij}) - g'_j(\tilde{\mathbf{M}}_{ij}) \right) \left(\tilde{\mathbf{M}}_{ij} - \hat{\mathbf{M}}_{ij} \right).$$

Now, for all $j \in [p]$, g'_j is increasing therefore

$$\left(g'_j(\hat{\mathbf{M}}_{ij}) - g'_j(\tilde{\mathbf{M}}_{ij}) \right) \left(\tilde{\mathbf{M}}_{ij} - \hat{\mathbf{M}}_{ij} \right) \leq 0,$$

which implies $B_3 \leq 0$. Combined with (37) and (38) this yields

$$\lambda_S \langle \hat{f}_\alpha, \hat{\alpha} - \alpha \rangle \leq \|\Delta \alpha\|_1 d_{\mathbf{X}} \left(\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + \nu)a \right).$$

Besides, the convexity of $\|\cdot\|_1$ gives $\langle \hat{f}_\alpha, \hat{\alpha} - \alpha \rangle \geq \|\hat{\alpha}\|_1 - \|\alpha\|_1$, therefore

$$\left\{ \lambda_S - d_{\mathbf{X}} \left(\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + \nu)a \right) \right\} \|\hat{\alpha}\|_1 \leq \left\{ \lambda_S + d_{\mathbf{X}} \left(\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + \nu)a \right) \right\} \|\alpha\|_1,$$

and the condition $\lambda_S \geq 2 \left\{ d_{\mathbf{X}} \left(\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1 + \nu)a \right) \right\}$ gives $\|\hat{\alpha}\|_1 \leq 3 \|\alpha\|_1$ and finally

$$\|\Delta \alpha\|_1 \leq 4 \|\alpha\|_1. \quad (39)$$

We consider the two following cases.

Case 1: $\mathbb{E} \|\mathcal{P}_\Omega(f_U(\Delta \alpha))\|_F^2 < \frac{72a^2 \log(n+p)}{\pi \log(6/5)}$. Then the result holds trivially.

Case 2: $\mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 \geq \frac{72a^2 \log(n+p)}{\pi \log(6/5)}$. For $d_1 > 0$ recall the definition of the set

$$\tilde{\mathcal{A}}(d_1) = \left\{ \boldsymbol{\alpha} \in \mathbb{R}^q : \|\boldsymbol{\alpha}\|_\infty \leq 1; \quad \|\boldsymbol{\alpha}\|_1 \leq d_1; \quad \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 \geq \frac{18 \log(n+p)}{\pi \log(6/5)} \right\}.$$

Inequality (39) and $\|\Delta\boldsymbol{\alpha}\|_\infty \leq 2a$ imply that

$$\frac{\Delta\boldsymbol{\alpha}}{2a} \in \tilde{\mathcal{A}}\left(\frac{2\|\boldsymbol{\alpha}\|_1}{a}\right).$$

Therefore we can apply Lemma 3(i) and obtain that with probability at least $1 - 8(n+p)^{-1}$,

$$\mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|^2 \leq 2 \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 + 64\nu a \|\boldsymbol{\alpha}\|_1 d_{\mathbf{X}} \mathbb{E} [\|\Sigma_R\|_\infty] + 3072a^2 p^{-1}. \quad (40)$$

We now must upper bound the quantity $\|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2$. Recall that $\tilde{\mathbf{M}} = \sum_{k=1}^q \boldsymbol{\alpha}_k \mathbf{X}(k) + \hat{\mathbf{M}}$. By definition,

$$\mathcal{L}(\hat{\mathbf{X}}) + \lambda_L \|\hat{\boldsymbol{\Theta}}\|_* + \lambda_S \|\hat{\boldsymbol{\alpha}}\|_1 \leq \mathcal{L}(\tilde{\mathbf{M}}) + \lambda_L \|\hat{\boldsymbol{\Theta}}\|_* + \lambda_S \|\boldsymbol{\alpha}\|_1,$$

i.e.

$$\mathcal{L}(\hat{\mathbf{M}}) - \mathcal{L}(\tilde{\mathbf{M}}) \leq \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\hat{\boldsymbol{\alpha}}\|_1).$$

Subtracting $\langle \nabla \mathcal{L}(\tilde{\mathbf{M}}), \hat{\mathbf{M}} - \tilde{\mathbf{M}} \rangle$ on both sides and by strong convexity of \mathcal{L} we obtain

$$\begin{aligned} \frac{\sigma^2}{2} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|^2 &\leq \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\hat{\boldsymbol{\alpha}}\|_1) + \langle \nabla \mathcal{L}(\tilde{\mathbf{M}}), \mathbf{f}_U(\Delta\boldsymbol{\alpha}) \rangle \\ &\leq \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\hat{\boldsymbol{\alpha}}\|_1) + \underbrace{\langle \nabla \mathcal{L}(\mathbf{M}^0), \mathbf{f}_U(\Delta\boldsymbol{\alpha}) \rangle}_{C_1} \\ &\quad + \underbrace{\langle \nabla \mathcal{L}(\mathbf{M}^0) - \nabla \mathcal{L}(\tilde{\mathbf{M}}), \mathbf{f}_U(\Delta\boldsymbol{\alpha}) \rangle}_{C_2}. \end{aligned} \quad (41)$$

The duality of $\|\cdot\|_1$ and $\|\cdot\|_\infty$ yields $C_1 \leq \|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty d_{\mathbf{X}} \|\Delta\boldsymbol{\alpha}\|_1$, and

$$C_2 \leq \|\nabla \mathcal{L}(\mathbf{M}^0) - \nabla \mathcal{L}(\tilde{\mathbf{M}})\|_\infty d_{\mathbf{X}} \|\Delta\boldsymbol{\alpha}\|_1.$$

Furthermore,

$$\|\nabla \mathcal{L}(\mathbf{M}^0) - \nabla \mathcal{L}(\tilde{\mathbf{M}})\|_\infty \leq 2\sigma_+^2 a,$$

since for all $(i, j) \in [n] \times [p]$ $|\tilde{\mathbf{M}}_{ij} - \mathbf{M}_{ij}^0| \leq 2a$ and $g_j''(\tilde{\mathbf{M}}_{ij}) \leq \sigma_+^2$. The last three inequalities plugged in (41) give

$$\frac{\sigma^2}{2} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 \leq \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\hat{\boldsymbol{\alpha}}\|_1) + d_{\mathbf{X}} \|\Delta\boldsymbol{\alpha}\|_1 \{ \|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2 a \}.$$

The triangular inequality gives

$$\begin{aligned} \frac{\sigma^2}{2} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 &\leq \{ d_{\mathbf{X}} (\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2 a) + \lambda_S \} \|\boldsymbol{\alpha}\|_1 \\ &\quad + \{ d_{\mathbf{X}} (\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2 a) - \lambda_S \} \|\hat{\boldsymbol{\alpha}}\|_1. \end{aligned}$$

Then, the assumption $\lambda_S \geq 2d_{\mathbf{X}} (\|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1+\nu)a)$ gives

$$\|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 \leq \frac{3\lambda_S}{\sigma_-^2} \|\boldsymbol{\alpha}\|_1.$$

Plugged into (40), this last inequality implies that with probability at least $1 - 8(n+p)^{-1}$

$$\mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\Delta\boldsymbol{\alpha}))\|_F^2 \leq \frac{3\lambda_S}{\sigma_-^2} \|\boldsymbol{\alpha}\|_1 + 64\nu a \|\boldsymbol{\alpha}\|_1 d_{\mathbf{X}} \mathbb{E} [\|\Sigma_R\|_\infty] + 3072a^2 p^{-1}. \quad (42)$$

Combining (39) and (42) gives the result.

B.2 Proof of Lemma 2

Using (36) for $L = \Theta^0$ and $\alpha = \alpha$ we obtain

$$\langle \nabla \mathcal{L}(\hat{\mathbf{M}}), \Delta \Theta \rangle + \sum_{k=1}^q (\Delta \alpha_k) \mathbf{X}(k) + \lambda_L \langle \hat{f}_L, \Delta \Theta \rangle + \lambda_S \langle \hat{f}_\alpha, \Delta \alpha \rangle \geq 0.$$

Then, the convexity of $\|\cdot\|_*$ and $\|\cdot\|_1$ imply that

$$\begin{aligned} \|\Theta^0\|_* &\geq \|\hat{\Theta}\|_* + \langle \partial \|\hat{\Theta}\|_*, \Delta \Theta \rangle, \\ \|\alpha\|_1 &\geq \|\hat{\alpha}\|_1 + \langle \partial \|\hat{\alpha}\|_1, \Delta \alpha \rangle. \end{aligned}$$

The last three inequalities yield

$$\begin{aligned} \lambda_L \left(\|\hat{\Theta}\|_* - \|\Theta^0\|_* \right) + \lambda_S (\|\hat{\alpha}\|_1 - \|\alpha\|_1) &\leq \langle \nabla \mathcal{L}(\hat{\mathbf{M}}), \Delta \Theta \rangle \\ &+ \langle \nabla \mathcal{L}(\hat{\mathbf{M}}), \sum_{k=1}^q (\Delta \alpha_k) \mathbf{X}(k) \rangle \\ &\leq \|\nabla \mathcal{L}(\hat{\mathbf{M}})\| \|\Delta \Theta\|_* + d_{\mathbf{X}} \|\nabla \mathcal{L}(\hat{\mathbf{M}})\|_\infty \|\Delta \alpha\|_1. \end{aligned}$$

Using the conditions

$$\lambda_L \geq 2 \|\nabla \mathcal{L}(\mathbf{M}^0)\|, \quad \lambda_S \geq 2d_{\mathbf{X}} \{ \|\nabla \mathcal{L}(\mathbf{M}^0)\|_\infty + 2\sigma_+^2(1+\nu)a \},$$

we get

$$\begin{aligned} \lambda_L (\|P_{\Theta^0}^\perp(\Delta \Theta)\|_* - \|P_{\Theta^0}(\Delta \Theta)\|_*) + \lambda_S (\|\hat{\alpha}\|_1 - \|\alpha\|_1) &\leq \\ \frac{\lambda_L}{2} (\|P_{\Theta^0}^\perp(\Delta \Theta)\|_* + \|P_{\Theta^0}(\Delta \Theta)\|_*) + \frac{\lambda_S}{2} \|\Delta \alpha\|_1, \end{aligned}$$

which implies

$$\|P_{\Theta^0}^\perp(\Delta \Theta)\|_* \leq 3 \|P_{\Theta^0}(\Delta \Theta)\|_* + 3\lambda_S/\lambda_L \|\alpha\|_1.$$

Now, using

$$\|\Delta \Theta\|_* \leq \|P_{\Theta^0}^\perp(\Delta \Theta)\|_* + \|P_{\Theta^0}(\Delta \Theta)\|_*, \quad \|P_{\Theta^0}(\Delta \Theta)\|_F \leq \|\Delta \Theta\|_F$$

and $\text{rank}(P_{\Theta^0}(\Delta \Theta)) \leq 2r$, we get

$$\|\Delta \Theta\|_* \leq \sqrt{32r} \|\Delta \Theta\|_F + 3\lambda_S/\lambda_L \|\alpha\|_1.$$

This completes the proof of Lemma 2.

B.3 Proof of Lemma 3

Proof of (i): Recall

$$D_\alpha = 8\nu d_1 d_{\mathbf{X}} \mathbb{E} [\|\Sigma_R\|_\infty] + 768p^{-1}$$

and

$$\tilde{\mathcal{A}}(d_1) = \left\{ \alpha \in \mathbb{R}^q : \|\alpha\|_\infty \leq 1; \quad \|\alpha\|_1 \leq d_1; \quad \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \geq \frac{18 \log(n+p)}{\pi \log(6/5)} \right\}.$$

We will show that the probability of the following event is small:

$$\mathcal{B} = \left\{ \exists \alpha \in \tilde{\mathcal{A}}(d_1) \text{ such that } \left| \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \right| > \frac{1}{2} \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 + D_\alpha \right\}.$$

Indeed, \mathcal{B} contains the complement of the event we are interested in. We use a peeling argument to upper bound the probability of event \mathcal{B} . Let $\nu = \frac{18 \log(n+p)}{\pi \log(6/5)}$ and $\eta = 6/5$. For $l \in \mathbb{N}$ set

$$\mathcal{S}_l = \left\{ \alpha \in \tilde{\mathcal{A}}(d_1) : \eta^{l-1}\nu \leq \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \leq \eta^l \nu \right\}.$$

Under the event \mathcal{B} , there exists $l \geq 1$ and $\alpha \in \tilde{\mathcal{A}}(d_1) \cap S_l$ such that

$$\begin{aligned} \left| \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \right| &> \frac{1}{2} \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 + D_\alpha \\ &> \frac{1}{2} \eta^{l-1} \nu + D_\alpha \\ &= \frac{5}{12} \eta^l \nu + D_\alpha. \end{aligned} \quad (43)$$

For $T > \nu$, consider the set of vectors

$$\tilde{\mathcal{A}}(d_1, T) = \left\{ \alpha \in \tilde{\mathcal{A}}(d_1) : \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \leq T \right\}$$

and the event

$$\mathcal{B}_l = \left\{ \exists \alpha \in \tilde{\mathcal{A}}(d_1, \eta^l \nu) : \left| \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \right| > \frac{5}{12} \eta^l \nu + D_\alpha \right\}.$$

If \mathcal{B} holds, then (43) implies that \mathcal{B}_l holds for some $l \leq 1$. Therefore, $\mathcal{B} \subset \cup_{l=1}^{+\infty} \mathcal{B}_l$, and it is enough to estimate the probability of the events \mathcal{B}_l and then apply the union bound. Such an estimation is given in the following Lemma, adapted from Lemma 10 in [20].

Lemma 6 Define $Z_T = \sup_{\alpha \in \tilde{\mathcal{A}}(d_1, T)} \left| \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\mathbf{f}_U(\alpha))\|_F^2 \right|$. Then,

$$\mathbb{P} \left(Z_T \geq D_\alpha + \frac{5}{12} T \right) \leq 4e^{-\pi T/18}.$$

Lemma 6 gives that $\mathbb{P}(\mathcal{B}_l) \leq 4 \exp(-\pi \eta^l \nu / 18)$. Applying the union bound we obtain

$$\begin{aligned} \mathbb{P}(\mathcal{B}) &\leq \sum_{l=1}^{\infty} \mathbb{P}(\mathcal{B}_l) \\ &\leq 4 \sum_{l=1}^{\infty} \exp(-\pi \eta^l \nu / 18) \\ &\leq 4 \sum_{l=1}^{\infty} \exp(-\pi \log(\eta) l \nu / 18), \end{aligned}$$

where we used $e^x \geq x$. Finally, for $\nu = \frac{18 \log(n+p)}{\pi \log(6/5)}$ we obtain

$$\mathbb{P}(\mathcal{B}) \leq \frac{4 \exp(-\pi \nu \log(\eta) / 18)}{1 - \exp(-\pi \nu \log(\eta) / 18)} \leq \frac{4 \exp(-\log(n+p))}{1 - \exp(-\log(n+p))} \leq \frac{8}{n+p},$$

since $d-1 \geq (n+p)/2$, which concludes the proof of (i).

Proof of (ii): The proof is very similar to that of (i); we recycle some of the notations for simplicity. Recall

$$D_X = \frac{112\rho}{\pi} \mathbb{E} [\|\Sigma_R\|]^2 + 8\nu \varepsilon \mathbb{E} [\|\Sigma_R\|] + 8\nu d_1 d_{\mathbf{X}} \mathbb{E} [\|\Sigma_R\|_\infty] + d_\Pi + 768p^{-1}.$$

Let

$$\mathcal{B} = \left\{ \exists (\Theta, \alpha) \in \mathcal{C}(d_1, d_\Pi, \rho, \varepsilon); \right. \\ \left. \left| \|\mathcal{P}_\Omega(\Theta + \mathbf{f}_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\Theta + \mathbf{f}_U(\alpha))\|_F^2 \right| > \frac{1}{2} \mathbb{E} \|\mathcal{P}_\Omega(\Theta + \mathbf{f}_U(\alpha))\|_F^2 + D_X \right\},$$

$\nu = \frac{72 \log(n+p)}{\pi \log(6/5)}$, $\eta = \frac{6}{5}$ and for $l \in \mathbb{N}$

$$S_l = \left\{ (\Theta, \alpha) \in \mathcal{C}(d_1, d_\Pi, \rho, \varepsilon) : \eta^{l-1} \nu \leq \mathbb{E} \|\mathcal{P}_\Omega(\Theta + \mathbf{f}_U(\alpha))\|_F^2 \leq \eta^l \nu \right\}.$$

As before, if \mathcal{B} holds, then there exist $l \geq 2$ and $(\Theta, \alpha) \in \mathcal{C}(d_1, d_\Pi, \rho, \varepsilon) \cap S_l$ such that

$$\left| \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 \right| > \frac{5}{12} \eta^l \nu + D_X. \quad (44)$$

For $T > \nu$, consider the set $\tilde{\mathcal{C}}(T) = \left\{ (\Theta, \alpha) \in \mathcal{C}(d_1, d_\Pi, \rho, \varepsilon) : \mathbb{E} \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 \leq T \right\}$, and the event

$$\mathcal{B}_l = \left\{ \exists (\Theta, \alpha) \in \tilde{\mathcal{C}}(\eta^l \nu) : \left| \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 \right| > \frac{5}{12} \eta^l \nu + D_X \right\}.$$

Then, (44) implies that \mathcal{B}_l holds and $\mathcal{B} \subset \cup_{l=1}^{\infty} \mathcal{B}_l$. Thus, we estimate in Lemma 7 the probability of the events \mathcal{B}_l , and then apply the union bound.

Lemma 7 Let $W_T = \sup_{(\Theta, \alpha) \in \tilde{\mathcal{C}}(T)} \left| \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 - \mathbb{E} \|\mathcal{P}_\Omega(\Theta + f_U(\alpha))\|_F^2 \right|$.

$$\mathbb{P} \left(W_T \geq D_X + \frac{5}{12} T \right) \leq 4e^{-\pi T/72}.$$

Lemma 7 gives that $\mathbb{P}(\mathcal{B}_l) \leq 4 \exp(-\pi \eta^l \nu / 72)$. Applying the union bound we obtain

$$\begin{aligned} \mathbb{P}(\mathcal{B}) &\leq \sum_{l=1}^{\infty} \mathbb{P}(\mathcal{B}_l) \\ &\leq 4 \sum_{l=1}^{\infty} \exp(-\pi \eta^l \nu / 72) \\ &\leq 4 \sum_{l=1}^{\infty} \exp(-\pi \log(\eta) l \nu / 72), \end{aligned}$$

where we used $e^x \geq x$. Finally, for $\nu = \frac{72 \log(n+p)}{\pi \log(6/5)}$ we obtain

$$\mathbb{P}(\mathcal{B}) \leq \frac{4 \exp(-\pi \nu \log(\eta) / 72)}{1 - \exp(-\pi \nu \log(\eta) / 72)} \leq \frac{4 \exp(-\log(n+p))}{1 - \exp(-\log(n+p))} \leq 8(n+p)^{-1},$$

since $n+p-1 \geq (n+p)/2$, which concludes the proof of (ii).

C Proof of Theorem 2

To prove the theorem, we first lower bound on the progress made by the algorithm at the two blocks between the iterations. With a slight abuse of notations, in the following we shall denote the iterates without the bracket in the superscripts, e.g., we denote $\alpha^{(t)}, \Theta^{(t)}, R^{(t)}$ by α^t, Θ^t, R^t , respectively, to simplify our discussions.

For the first block on α , in Section C.1 we show that

$$F(\alpha^t, \Theta^{t-1}, R^{t-1}) \leq F(\alpha^{t-1}, \Theta^{t-1}, R^{t-1}) - \frac{\gamma (g_\alpha(\alpha^{t-1}, \Theta^{t-1}; Q^{t-1}))^2}{2(2Q^{t-1})^2}, \quad (45)$$

where $Q^{t-1} := \lambda_S^{-1} F(\alpha^{t-1}, \Theta^{t-1}, R^{t-1})$ as defined in the main paper and

$$g_\alpha(\alpha^{t-1}, \Theta^{t-1}; Q^{t-1}) := \langle \nabla_\alpha \mathcal{L}(\alpha^{t-1}, \Theta^{t-1}), \alpha^{t-1} - \hat{\alpha}^{t-1} \rangle + \lambda_S (\|\alpha^{t-1}\|_1 - \|\hat{\alpha}^{t-1}\|_1), \quad (46)$$

such that

$$\hat{\alpha}^{t-1} := \arg \min_{\alpha} (\langle \nabla_\alpha \mathcal{L}(\alpha^{t-1}, \Theta^{t-1}), \alpha \rangle + \lambda_S \|\alpha\|_1) \text{ s.t. } \|\alpha\|_1 \leq Q^{t-1}. \quad (47)$$

For the second block on (Θ, R) , Section C.2 shows that

$$F(\alpha^t, \Theta^t, R^t) \leq F(\alpha^t, \Theta^{t-1}, R^{t-1}) - \frac{(g_\Theta(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{\max\{2R_{\text{UB}}^t(\lambda_L + M^t), 8\sigma_\Theta(R_{\text{UB}}^t)^2\}}, \quad (48)$$

where $M^t := \|\nabla_{\Theta}(\alpha^t, \Theta^{t-1})\|_2$ and we recall that $R_{\text{UB}}^t := \lambda_L^{-1}F(\alpha^t, \Theta^{t-1}, R^{t-1})$ and we have defined

$$g_{\Theta}(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t) := \langle \Theta^{t-1} - \hat{\Theta}^t, \nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1}) \rangle + \lambda_L(R^{t-1} - \hat{R}^t). \quad (49)$$

Moreover, Section C.2 shows that

$$F(\alpha^t, \Theta^t, R^t) - F(\alpha^t, \Theta^{t-1}, R^{t-1}) \leq -\frac{\sigma_{\Theta}}{2} \|\mathcal{P}_{\Omega}(\Theta^t - \Theta^{t-1})\|_F^2. \quad (50)$$

Statement (i). The above results show that the objective values for the iterates produced by the MCGD method are non-increasing, *i.e.*,

$$F(\alpha^t, \Theta^t, R^t) \leq F(\alpha^t, \Theta^{t-1}, R^{t-1}) \leq F(\alpha^{t-1}, \Theta^{t-1}, R^{t-1}) \quad (51)$$

Now, consider the time varying part in the quantity $C(t)$ [cf. (18)] — $Q^t, R_{\text{UB}}^t, M^t$. The first two quantities are defined from the objective values and are thus bounded by $\lambda_S^{-1}F(\alpha^0, \Theta^0, R^0)$, $\lambda_L^{-1}F(\alpha^0, \Theta^0, R^0)$, respectively. Moreover, from the monotonicity of $F(\alpha^t, \Theta^{t-1}, R^{t-1})$, we have $\lambda_L \|\Theta^{t-1}\|_* + \lambda_S \|\alpha^t\|_1 \leq F(\alpha^t, \Theta^{t-1}, R^{t-1}) \leq F(\alpha^0, \Theta^0, R^0)$ for all $t \geq 1$. As the gradient $\nabla_{\Theta} \mathcal{L}(\alpha, \Theta)$ is bounded whenever α, Θ are bounded, we conclude that M^t is bounded, e.g., $M^t \leq \bar{M}$. Finally, this shows for all $t \geq 1$ that

$$C(t) \leq \bar{C} := \max \left\{ \frac{24(Q^0)^2}{\gamma}, \frac{24\hat{\sigma}_{\Theta}^2(Q^0)^2}{\sigma_{\Theta}} + \max\{6R_{\text{UB}}^0(\lambda_L + \bar{M}), 24\sigma_{\Theta}(R_{\text{UB}}^0)^2\} \right\}. \quad (52)$$

Statement (ii). To characterize the convergence rate of the MCGD method, let us consider the Lyapunov function, $g^t(Q^t, R_{\text{UB}}^t)$, defined as:

$$g^t(Q^t, R_{\text{UB}}^t) := g_{\alpha}(\alpha^t, \Theta^{t-1}; Q^t) + g_{\Theta}(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t). \quad (53)$$

Note that as the loss function $\mathcal{L}(\alpha, \Theta)$ is convex and $\|\hat{\alpha}\|_1 \leq Q^t$, $\|\hat{\Theta}\|_* \leq R_{\text{UB}}^t$, it is possible to lower bound $g^t(Q^t, R_{\text{UB}}^t)$ by:

$$g^t(Q^t, R_{\text{UB}}^t) \geq F(\alpha^t, \Theta^{t-1}, R^{t-1}) - F(\hat{\alpha}, \hat{\Theta}, \hat{R}). \quad (54)$$

Furthermore, we can obtain an upper bound to $g^t(Q^t, R_{\text{UB}}^t)$ in terms of the objective values:

$$\begin{aligned} g_{\alpha}(\alpha^t, \Theta^{t-1}; Q^t) &= \max_{\|\alpha\| \leq Q^t} \langle \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^{t-1}), \alpha^t - \alpha \rangle + \lambda_S(\|\alpha^t\|_1 - \|\alpha\|_1) \\ &= \max_{\|\alpha\|_1 \leq Q^t} \langle \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^t), \alpha^t - \alpha \rangle + \langle \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^{t-1}) - \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^t), \alpha^t - \alpha \rangle \\ &\quad + \lambda_S(\|\alpha^t\|_1 - \|\alpha\|_1) \\ &\leq \max_{\|\alpha\|_1 \leq Q^t} \langle \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^t), \alpha^t - \alpha \rangle + \lambda_S(\|\alpha^t\|_1 - \|\alpha\|_1) \\ &\quad + \|\nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^{t-1}) - \nabla_{\alpha} \mathcal{L}(\alpha^t, \Theta^t)\|_2 \|\alpha^t - \alpha\|_2 \\ &\leq g_{\alpha}(\alpha^t, \Theta^t; Q^t) + 2\hat{\sigma}_{\Theta} Q^t \|\mathcal{P}_{\Omega}(\Theta^{t-1} - \Theta^t)\|_F. \end{aligned} \quad (55)$$

Consequently, we have

$$\begin{aligned} &(g^t(Q^t, R_{\text{UB}}^t))^2 \\ &\leq 3((g_{\Theta}(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2 + (g_{\alpha}(\alpha^t, \Theta^t; Q^t))^2 + 4\hat{\sigma}_{\Theta}^2(Q^t)^2 \|\mathcal{P}_{\Omega}(\Theta^{t-1} - \Theta^t)\|_F^2) \\ &\leq 3(C_1^t(F(\alpha^t, \Theta^t, R^t) - F(\alpha^{t+1}, \Theta^t, R^t)) + C_2^t(F(\alpha^t, \Theta^{t-1}, R^{t-1}) - F(\alpha^t, \Theta^t, R^t))) \end{aligned}$$

where

$$C_1^t := \frac{8(Q^t)^2}{\gamma}, \quad C_2^t = \frac{8\hat{\sigma}_{\Theta}^2(Q^t)^2}{\sigma_{\Theta}} + \max\{2R_{\text{UB}}^t(\lambda_L + M), 8\sigma_{\Theta}(R_{\text{UB}}^t)^2\} \quad (56)$$

Observe that $C(t)$ is defined by $C(t) = 3 \max\{C_1^t, C_2^t\}$ as the upper bound of the above constants, we get

$$(g^t(Q^t, R_{\text{UB}}^t))^2 \leq C(t)(F(\alpha^t, \Theta^{t-1}, R^{t-1}) - F(\alpha^{t+1}, \Theta^t, R^t)). \quad (57)$$

Using the shorthand notation $\Delta^t := F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) - F(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}, \hat{R})$ and notice that $(g^t(Q^t, R_{\text{UB}}^t))^2 \geq (\Delta^t)^2$, we arrive at the following inequality:

$$\Delta^{t+1} \leq \Delta^t - \frac{1}{C(t)}(\Delta^t)^2 \quad (58)$$

Applying Lemma 8 in Section C.3, we can show that

$$\Delta^{t+1} \leq \frac{1}{(\Delta^1)^{-1} + \sum_{i=1}^t \frac{1}{C(i)}}, \quad (59)$$

Note that $\Delta^1 \leq \tilde{\Delta}^0 := F(\boldsymbol{\alpha}^0, \boldsymbol{\Theta}^0, R^0) - F(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}, \hat{R})$, we have

$$\Delta^{t+1} \leq \frac{1}{(\tilde{\Delta}^0)^{-1} + \sum_{i=1}^t \frac{1}{C(i)}} \leq \frac{1}{(\tilde{\Delta}^0)^{-1} + t\bar{C}(t)}, \quad \forall t \geq 0. \quad (60)$$

The proof is concluded by the straightforward inequality $F_0(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\Theta}^{t+1}) - F_0(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}) \leq F(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\Theta}^{t+1}, R^{t+1}) - F(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}, \hat{R}) \leq \Delta^{t+1}$.

Comment on $\lim_{t \rightarrow \infty} C(t)$. Since both $F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^t, R^t)$ and $F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1})$ converge to $F^* := F(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}, \hat{R})$, *i.e.*, the optimal objective value. It is clear that $Q^t \rightarrow \hat{Q} := \lambda_S^{-1} F^*$ and $R_{\text{UB}}^t \rightarrow \hat{R}_{\text{UB}} := \lambda_L^{-1} F^*$ as well. Furthermore, by continuity of the gradient, we have $M^t \rightarrow \|\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}})\|_2$. This shows that the limit $C^* = \lim_{t \rightarrow \infty} C(t)$ exists.

To obtain a computable bound for C^* , note that $(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}})$ is also an optimal solution to (5) and the optimality condition shows that

$$\mathbf{0} \in \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}}) + \lambda_L \partial \|\hat{\boldsymbol{\Theta}}\|_* \quad (61)$$

By [31, P. 41], we know that $\partial \|\hat{\boldsymbol{\Theta}}\|_* = \{\mathbf{U}_1 \mathbf{V}_1^\top + \mathbf{W} : \|\mathbf{W}\|_2 \leq 1, \mathbf{U}_1^\top \mathbf{W} = \mathbf{0}, \mathbf{W} \mathbf{V}_1 = \mathbf{0}\}$ such that $\mathbf{U}_1 \in \mathbb{R}^{m_1 \times r}$, $\mathbf{V}_1 \in \mathbb{R}^{m_2 \times r}$ are the left/right singular vectors of $\hat{\boldsymbol{\Theta}}$ corresponding the $r := \text{rank}(\hat{\boldsymbol{\Theta}})$ non-zero singular values of $\hat{\boldsymbol{\Theta}}$. Importantly, this implies that $\|\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\Theta}})\|_2 \leq 2\lambda_L$ and

$$C^* \leq \bar{C}^* := \max \left\{ \frac{24(\hat{Q})^2}{\gamma}, \frac{24\hat{\sigma}_{\boldsymbol{\Theta}}^2(\hat{Q})^2}{\sigma_{\boldsymbol{\Theta}}} + \max\{18\hat{R}_{\text{UB}}\lambda_L, 24\sigma_{\boldsymbol{\Theta}}(\hat{R}_{\text{UB}})^2\} \right\}. \quad (62)$$

C.1 Proof of Eq. (45)

Suppose $\boldsymbol{\alpha}^t$ is obtained by the proximal update in (10), we observe that

$$\begin{aligned} F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) &\leq F(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}) + \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}), \boldsymbol{\alpha}^t - \boldsymbol{\alpha}^{t-1} \rangle \\ &\quad + \frac{\sigma_{\boldsymbol{\alpha}}}{2} \|\boldsymbol{\alpha}^t - \boldsymbol{\alpha}^{t-1}\|_2^2 + \lambda_S (\|\boldsymbol{\alpha}^t\|_1 - \|\boldsymbol{\alpha}^{t-1}\|_1). \end{aligned} \quad (63)$$

On the other hand, when $\boldsymbol{\alpha}^t$ is obtained by the exact minimization in (11), denoted by $\boldsymbol{\alpha}_{\text{exact}}^t$ to avoid confusion, we have $F(\boldsymbol{\alpha}_{\text{exact}}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) \leq F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1})$ since the latter is an exact minimizer. Thus, $F(\boldsymbol{\alpha}_{\text{exact}}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1})$ is upper bounded by the right hand side in the above inequality.

Using the property of the proximal operator, it can be shown that

$$\boldsymbol{\alpha}^t \in \arg \min_{\boldsymbol{\alpha}} \left(\langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}), \boldsymbol{\alpha} - \boldsymbol{\alpha}^{t-1} \rangle + \frac{1}{2\gamma} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}^{t-1}\|_2^2 + \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\boldsymbol{\alpha}^{t-1}\|_1) \right) \quad (64)$$

Due to our choice of step size, we have $\sigma_{\boldsymbol{\alpha}} \leq 1/\gamma$. Combining this with the above inequality implies that

$$\begin{aligned} F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) &\leq F(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}) + \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}), \boldsymbol{\alpha} - \boldsymbol{\alpha}^{t-1} \rangle \\ &\quad + \frac{1}{2\gamma} \|\boldsymbol{\alpha} - \boldsymbol{\alpha}^{t-1}\|_2^2 + \lambda_S (\|\boldsymbol{\alpha}\|_1 - \|\boldsymbol{\alpha}^{t-1}\|_1), \quad \forall \boldsymbol{\alpha} \in \mathbb{R}^K. \end{aligned} \quad (65)$$

Furthermore, for all $b \in \mathbb{R}$ it holds that

$$\begin{aligned}
F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) &\leq F(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}) + b \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}), \hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1} \rangle \\
&\quad + \frac{b^2}{2\gamma} \|\hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1}\|_2^2 + \lambda_S (\|b\hat{\boldsymbol{\alpha}}^{t-1} + (1-b)\boldsymbol{\alpha}^{t-1}\|_1 - \|\boldsymbol{\alpha}^{t-1}\|_1) \\
&\leq F(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}) + b \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}), \hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1} \rangle \\
&\quad + \frac{b^2}{2\gamma} \|\hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1}\|_2^2 + b\lambda_S (\|\hat{\boldsymbol{\alpha}}^{t-1}\|_1 - \|\boldsymbol{\alpha}^{t-1}\|_1),
\end{aligned} \tag{66}$$

where we have limited our search space from $\boldsymbol{\alpha} \in \mathbb{R}^K$ to $\boldsymbol{\alpha} = b\hat{\boldsymbol{\alpha}} + (1-b)\boldsymbol{\alpha}^{t-1}$ for $b \in \mathbb{R}$. Minimizing the right hand side of the above with respect to b yields

$$\begin{aligned}
&F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) - F(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}) \\
&\leq -\frac{\gamma}{2} \frac{(g_{\boldsymbol{\alpha}}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}; Q^{t-1}))^2}{\|\hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1}\|_2^2} \leq -\frac{\gamma}{2} \frac{(g_{\boldsymbol{\alpha}}(\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}; Q^{t-1}))^2}{(2Q^{t-1})^2},
\end{aligned} \tag{67}$$

where we have used $\|\hat{\boldsymbol{\alpha}}^{t-1} - \boldsymbol{\alpha}^{t-1}\|_2^2 \leq (2Q^{t-1})^2$ in the last inequality.

C.2 Proof of Eq. (48) and (50)

Let us observe that

$$\begin{aligned}
F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^t, R^t) &= F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) - \beta_t g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) \\
&\quad + \frac{\beta_t^2}{2} \begin{pmatrix} \text{vec}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1}) \\ \hat{R}^t - R^{t-1} \end{pmatrix}^\top \nabla_{\boldsymbol{\Theta}, R}^2(\boldsymbol{\xi}) \begin{pmatrix} \text{vec}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1}) \\ \hat{R}^t - R^{t-1} \end{pmatrix},
\end{aligned} \tag{68}$$

where $\boldsymbol{\xi}$ is any point that lies on the line $[(\text{vec}(\boldsymbol{\Theta}^{t-1}); R^{t-1}), (\text{vec}(\boldsymbol{\Theta}^t); R^t)]$. From the property of F , we observe that

$$\nabla_{\boldsymbol{\Theta}, R}^2(\boldsymbol{\xi}) \preceq \begin{pmatrix} \sigma_{\boldsymbol{\Theta}} \text{Diag}(\mathcal{P}_{\Omega}(\mathbf{J})) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}, \tag{69}$$

where \mathbf{J} is the $m_1 \times m_2$ all-ones matrix. The above implies that

$$\begin{aligned}
F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^t, R^t) &\leq F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) - \beta_t g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) \\
&\quad + \frac{\beta_t^2 \sigma_{\boldsymbol{\Theta}}}{2} \|\mathcal{P}_{\Omega}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1})\|_F^2.
\end{aligned} \tag{70}$$

Recall that $\beta_t = \min\{1, g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) / (\sigma_{\boldsymbol{\Theta}} \|\mathcal{P}_{\Omega}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1})\|_F^2)\}$. If $g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) \geq \sigma_{\boldsymbol{\Theta}} \|\mathcal{P}_{\Omega}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1})\|_F^2$, then we choose $\beta_t = 1$ and observe:

$$\begin{aligned}
&F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^t, R^t) - F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) \\
&\leq -\frac{1}{2} g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) = -\frac{1}{2} \frac{(g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t)} \\
&\leq -\frac{1}{2} \frac{(g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{R_{\text{UB}}^t (\lambda_L + 2M^t)},
\end{aligned} \tag{71}$$

where we have used the upper bound to $g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t)$ as follows:

$$\begin{aligned}
g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) &\leq \lambda_L R_{\text{UB}}^t + \langle \boldsymbol{\Theta}^{t-1} - \hat{\boldsymbol{\Theta}}^t, \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}) \rangle \\
&\leq R_{\text{UB}}^t (\lambda_L + 2M^t),
\end{aligned} \tag{72}$$

with $M^t := \|\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1})\|_2$ being the spectral norm of the gradient.

Otherwise, we choose $\beta_t = g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t) / (\sigma_{\boldsymbol{\Theta}} \|\mathcal{P}_{\Omega}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1})\|_F^2)$ and observe:

$$\begin{aligned}
&F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^t, R^t) - F(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}) \\
&\leq -\frac{1}{2} \frac{(g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{\sigma_{\boldsymbol{\Theta}} \|\mathcal{P}_{\Omega}(\hat{\boldsymbol{\Theta}}^t - \boldsymbol{\Theta}^{t-1})\|_F^2} \leq -\frac{1}{2} \frac{(g_{\boldsymbol{\Theta}}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{\sigma_{\boldsymbol{\Theta}} (2R_{\text{UB}}^t)^2},
\end{aligned} \tag{73}$$

where we have used $\|\mathcal{P}_\Omega(\hat{\Theta}^t - \Theta^{t-1})\|_F^2 \leq \|\hat{\Theta}^t - \Theta^{t-1}\|_F^2 \leq \|\hat{\Theta}^t - \Theta^{t-1}\|_*^2 \leq (2R_{\text{UB}}^t)^2$.

To prove (50), we observe that

$$\|\mathcal{P}_\Omega(\Theta^t - \Theta^{t-1})\|_F^2 = \beta_t^2 \|\mathcal{P}_\Omega(\hat{\Theta}^t - \Theta^{t-1})\|_F^2. \quad (74)$$

If $\beta_t = 1$, then we have $g_\Theta(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t) \geq \sigma_\Theta \|\mathcal{P}_\Omega(\hat{\Theta}^t - \Theta^{t-1})\|_F^2$ and therefore we can upper bound $\|\mathcal{P}_\Omega(\Theta^t - \Theta^{t-1})\|_F^2$ by:

$$\frac{1}{\sigma_\Theta} g_\Theta(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t) \leq \frac{2}{\sigma_\Theta} \left(F(\alpha^t, \Theta^{t-1}, R^{t-1}) - F(\alpha^t, \Theta^t, R^t) \right) \quad (75)$$

where the last inequality follows from (72). Otherwise, we choose $\beta_t = g_\Theta(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t) / \sigma_\Theta \|\mathcal{P}_\Omega(\hat{\Theta}^t - \Theta^{t-1})\|_F^2$ and therefore,

$$\begin{aligned} \|\mathcal{P}_\Omega(\Theta^t - \Theta^{t-1})\|_F^2 &= \frac{1}{\sigma_\Theta} \frac{(g_\Theta(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\text{UB}}^t))^2}{\sigma_\Theta \|\mathcal{P}_\Omega(\hat{\Theta}^t - \Theta^{t-1})\|_F^2} \\ &\leq \frac{2}{\sigma_\Theta} \left(F(\alpha^t, \Theta^{t-1}, R^{t-1}) - F(\alpha^t, \Theta^t, R^t) \right), \end{aligned} \quad (76)$$

where the last inequality follows from (73).

C.3 Additional Lemma

The following lemma is modified from [3, Lemma 3.5].

Lemma 8 *Let $\{A_k\}_{k \geq 1}$ be a non-negative sequence satisfying:*

$$A_{k+1} \leq A_k - \gamma_k A_k^2, \quad k \geq 1, \quad (77)$$

where γ_k is some positive number for all $k \geq 1$. Then,

$$A_{k+1} \leq \frac{1}{\frac{1}{A_1} + \sum_{i=1}^k \gamma_i}, \quad k \geq 1. \quad (78)$$

Proof: Consider the following chain of inequality:

$$\frac{1}{A_{k+1}} - \frac{1}{A_k} = \frac{A_k - A_{k+1}}{A_k A_{k+1}} \geq \gamma_k \frac{A_k}{A_{k+1}} \geq \gamma_k, \quad (79)$$

where the last inequality is due to the fact that $A_{k+1} \leq A_k$. Consequently, we have

$$\frac{1}{A_{k+1}} - \frac{1}{A_1} = \sum_{i=1}^k \left(\frac{1}{A_{i+1}} - \frac{1}{A_i} \right) \geq \sum_{i=1}^k \gamma_i. \quad (80)$$

Reshuffling terms shows the desired result in (78). **Q.E.D.**

D Distributed MCGD Optimization

Similar to the previous section, in the following we shall denote the iterates without the bracket in the superscripts, e.g., we denote $\alpha^{(t)}, \Theta^{(t)}, R^{(t)}$ by α^t, Θ^t, R^t , respectively, to simplify our discussions.

Let us describe a distributed version of the MCGD method under a *master-slave* architecture setting where there exists K workers and each of them is connected to a central server. Our goal is to offload the computation required by MCGD method to the workers, while protecting the privacy sensitive data owned by the workers. To describe our setting, the set of observed data \mathbf{Y}_{ij} , $(i, j) \in \Omega$ are stored in K different workers, where the k th worker holds \mathbf{Y}_{ij} with $(i, j) \in \Omega_k \subset \Omega$. Particularly, we have $\Omega = \Omega_1 \cup \dots \cup \Omega_K$ with $\Omega_k \cap \Omega_{k'} = \emptyset$ for all $k \neq k'$. In this way, we can write

$$\mathcal{L}(\alpha, \Theta) = \sum_{(i,j) \in \Omega} \{-\mathbf{Y}_{ij} \mathbf{M}_{ij} + g_j(\mathbf{M}_{ij})\} = \underbrace{\sum_{k=1}^K \sum_{(i,j) \in \Omega_k} \{-\mathbf{Y}_{ij} \mathbf{M}_{ij} + g_j(\mathbf{M}_{ij})\}}_{:= \mathcal{L}_k(\alpha, \Theta)} \quad (81)$$

such that the log-likelihood function can be decomposed as $\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\Theta}) := \sum_{k=1}^K \mathcal{L}_k(\boldsymbol{\alpha}, \boldsymbol{\Theta})$. Moreover, notice that $\mathcal{L}_k(\boldsymbol{\alpha}, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta})) = \mathcal{L}_k(\boldsymbol{\alpha}, \boldsymbol{\Theta})$ since the k th local function is evaluated only on the entries in Ω_k . For simplicity, we assume that computation can be done synchronously among the workers.

We can implement the MCGD method in a distributed setting as follows. We focus on the t th iteration where $\boldsymbol{\alpha}^{t-1}, \boldsymbol{\Theta}^{t-1}, R^{t-1}$ have been previously computed and worker k now holds $\boldsymbol{\alpha}^{t-1}, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1}), R^{t-1}$.

Firstly, the proximal update step of line 3 is replaced by a natural distributed implementation where the workers compute and transmit the local gradients of the log-likelihood function, $\nabla_{\boldsymbol{\alpha}} \mathcal{L}_k(\boldsymbol{\alpha}^{t-1}, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1}))$, to the master node; the master node can then *aggregate* the received local gradients to form the update in (10), yielding $\boldsymbol{\alpha}^t$ which is then transmitted back to the workers.

Secondly, the CG update of line 5 requires the top SVD of $\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1})$ whose complexity is $\mathcal{O}(|\Omega| \max\{n, p\} \log(1/\delta))$ using a centralized implementation, where $\delta > 0$ is the desired accuracy of SVD. In a distributed setting, we can replace the step by a *distributed power method* for offloading the complexity. Importantly, we observe that the top singular vectors of $\nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1})$ can be *approximated* by the following power method recursions:

Algorithm 2 Distributed Power Method for MCGD.

- 1: **Initialize:** initialization — $\mathbf{u}(0) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \in \mathbb{R}^n$, and the parameter $P \in \mathbb{Z}$.
- 2: **for** $p = 1, 2, \dots, P$ **do**
- 3: The central server sends the vector $\mathbf{u}(p-1)$ to workers.
- 4: For all k , worker k computes the vector:

$$\mathbf{v}_k(p) = \nabla_{\boldsymbol{\Theta}} \mathcal{L}_k(\boldsymbol{\alpha}^t, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1})) \mathbf{u}(p-1) \quad (82)$$

and transmit it to the central server.

- 5: The central server forms the next iterate by $\mathbf{v}(p) = \sum_{k=1}^K \mathbf{v}_k(p)$ and sends the vector $\mathbf{v}(p)$ to workers.
- 6: For all k , worker k computes the vector:

$$\mathbf{u}_k(p) = \nabla_{\boldsymbol{\Theta}} \mathcal{L}_k(\boldsymbol{\alpha}^t, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1}))^\top \mathbf{v}(p) \quad (83)$$

and transmit it to the central server.

- 7: The central server forms the next iterate by $\mathbf{u}(p) = \sum_{k=1}^K \mathbf{u}_k(p)$.
 - 8: **end for**
 - 9: At the central server, compute the top left and right singular vector as $\mathbf{u}_{(1)}^t = \mathbf{u}(P) / \|\mathbf{u}(P)\|$ and $\mathbf{v}_{(1)}^t = \mathbf{v}(P) / \|\mathbf{v}(P)\|$.
 - 10: **Return:** the top singular vectors $\mathbf{u}_{(1)}^t, \mathbf{v}_{(1)}^t$.
-

Line 4 and 5 in the above pseudo code implement the following power iterations:

$$\mathbf{v}(p) = \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}) \mathbf{u}(p-1) = \sum_{k=1}^K \nabla_{\boldsymbol{\Theta}} \mathcal{L}_k(\boldsymbol{\alpha}^t, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1})) \mathbf{u}(p-1) \quad (84)$$

$$\mathbf{u}(p) = \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1})^\top \mathbf{v}(p) = \sum_{k=1}^K \nabla_{\boldsymbol{\Theta}} \mathcal{L}_k(\boldsymbol{\alpha}^t, \mathcal{P}_{\Omega_k}(\boldsymbol{\Theta}^{t-1}))^\top \mathbf{v}(p), \quad (85)$$

where we have exploited the decomposable structure of the log-likelihood function in the distributed setting. Upon computing $\mathbf{u}_{(1)}^t, \mathbf{v}_{(1)}^t$, we can estimate the top singular value by $(\mathbf{v}_{(1)}^t)^\top \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}) \mathbf{u}_{(1)}^t$ which can also be computed distributively using similar scheme as in the above. Consequently, the update direction $(\hat{\boldsymbol{\Theta}}^t, \hat{R}^t)$ can be computed at the central server using

$$(\hat{\boldsymbol{\Theta}}^t, \hat{R}^t) = \begin{cases} (\mathbf{0}, 0), & \text{if } \lambda_L \geq (\mathbf{v}_{(1)}^t)^\top \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}) \mathbf{u}_{(1)}^t, \\ (-R_{\text{UB}}^t \mathbf{u}_{(1)}^t (\mathbf{v}_{(1)}^t)^\top, R_{\text{UB}}^t), & \text{if } \lambda_L < (\mathbf{v}_{(1)}^t)^\top \nabla_{\boldsymbol{\Theta}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\Theta}^{t-1}) \mathbf{u}_{(1)}^t. \end{cases} \quad (86)$$

Lastly, to compute the step size β_t required in line 6, an efficient way is to observe the following decomposition of the inner product:

$$\langle \Theta^{t-1} - \hat{\Theta}^t, \nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1}) \rangle = \sum_{k=1}^K \langle \mathcal{P}_{\Omega_k}(\Theta^{t-1} - \hat{\Theta}^t), \nabla_{\Theta} \mathcal{L}_k(\alpha^t, \mathcal{P}_{\Omega_k}(\Theta^{t-1})) \rangle. \quad (87)$$

This implies that the inner product on the left hand side can be computed by aggregating the K terms on the right hand side, where each of the K terms can be computed at the k th worker once $\mathcal{P}_{\Omega_k}(\hat{\Theta}^t)$ is available. As such, the central server also sends $\mathcal{P}_{\Omega_k}(\hat{\Theta}^t)$ to the workers after (86). Consequently, the step size is given by:

$$\beta_t = \min \left\{ 1, \frac{(\hat{g}_{\Theta}(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\cup B}^t))_+}{\sigma_{\Theta} \|\mathcal{P}_{\Omega}(\hat{\Theta}^t - \Theta^{t-1})\|_F^2} \right\}, \quad (88)$$

where

$$\hat{g}_{\Theta}(\alpha^t, \Theta^{t-1}, R^{t-1}; R_{\cup B}^t) := \langle \Theta^{t-1} - \hat{\Theta}^t, \nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1}) \rangle + \lambda_L (R^{t-1} - \hat{R}^t). \quad (89)$$

Note that unlike the function $g_{\Theta}(\cdot)$ defined in (49), the function $\hat{g}_{\Theta}(\cdot)$ can be negative since the matrix $\hat{\Theta}^t$ herein is computed from an inexact pair of top singular vectors.

Several remarks are in order. Throughout the optimization, the central server is unaware of the local gradient matrix *w.r.t.* Θ , instead only its corresponding matrix-vector products are transmitted from the workers to the server. In this way, the privacy-sensitive data from the workers will not be revealed to the server.

For any $\delta > 0$, it is well known that in high probability (with respect to the random initialization), the power method in Algorithm 2 converges [13] to an δ -accurate top SVD solution in $P = \mathcal{O}(\log(1/\delta))$ steps¹. Therefore, for the distributed MCGD method, the overall complexity required per iteration is $\mathcal{O}(|\Xi| + \max\{n, p\} \log(1/\delta))$ at the central server, and it is $\mathcal{O}(|\Omega_k| \max\{n, p\} \log(1/\delta))$ for the k th worker. The overall complexity is lower than a centralized implementation especially when $|\Omega_k| \ll |\Omega|$, e.g., when the number of workers increases.

¹For example, an δ -accurate top SVD solution satisfies

$$\left| (\mathbf{u}_{(1)}^t)^\top \nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1}) \mathbf{v}_{(1)}^t - \sigma_1(\nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1})) \right| \leq \delta. \quad (90)$$

In the complexity measure, we have hidden the dependency on the spectral gap $\Delta := \sigma_2(\nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1})) / \sigma_1(\nabla_{\Theta} \mathcal{L}(\alpha^t, \Theta^{t-1})) \leq 1$ in the big-O notation.