

Statistica Sinica Preprint No: SS-2019-0196

Title	Matrix Completion under Low-Rank Missing Mechanism
Manuscript ID	SS-2019-0196
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202019.0196
Complete List of Authors	Xiaojun Mao, Raymond K. W. Wong and Song Xi Chen
Corresponding Author	Song Xi Chen
E-mail	csx@gsm.pku.edu.cn
Notice: Accepted version subject to English editing.	

Matrix Completion under Low-Rank Missing Mechanism

Xiaojun Mao, Raymond K. W. Wong and Song Xi Chen

Fudan University, Texas A&M University, Peking University

Abstract: Matrix completion is a modern missing data problem where both the missing structure and the underlying parameter are high dimensional. Although missing structure is a key component to any missing data problems, existing matrix completion methods often assume a simple uniform missing mechanism. In this work, we study matrix completion from corrupted data under a novel low-rank missing mechanism. The probability matrix of observation is estimated via a high dimensional low-rank matrix estimation procedure, and further used to complete the target matrix via inverse probabilities weighting. Due to both high dimensional and extreme (i.e., very small) nature of the true probability matrix, the effect of inverse probability weighting requires careful study. We derive optimal asymptotic convergence rates of the proposed estimators for both the observation probabilities and the target matrix.

Key words and phrases: Low-rank; Missing; Nuclear-norm; Regularization.

1. Introduction

The problem of recovering a high-dimensional matrix $\mathbf{A}_* \in \mathbb{R}^{n_1 \times n_2}$ from very few (noisy) observations of its entries is commonly known as matrix completion, whose applications include, collaborative filtering, computer visions and

Song Xi Chen is the corresponding author.

positioning. From a statistical viewpoint, it is a high-dimensional missing data problem where a high percentage of matrix entries are missing. As in many missing data problems, the underlying missing mechanism plays an important role. Most existing work (e.g., Candès and Recht, 2009; Keshavan et al., 2009; Recht, 2011; Rohde and Tsybakov, 2011; Koltchinskii et al., 2011) adopt a uniform observation mechanism, where each entry has the same marginal probability of being observed. This leads to significant simplifications, and enables the domain to move forward rapidly with various theoretical breakthroughs in the last decade. However, the uniform mechanism is often unrealistic. Recent works (Foygel et al., 2011; Negahban and Wainwright, 2012; Klopp, 2014; Cai and Zhou, 2016; Cai et al., 2016; Bi et al., 2017; Mao et al., 2019) have been devoted to relaxing such an restrictive assumption by adopting other missing structures. The usage of these settings hinges on strong prior knowledge of the underlying problems. At a high level, many of them utilize some special forms of low-rank structures for missing mechanism. For instance, Foygel et al. (2011) and Negahban and Wainwright (2012) both adopt a rank-1 structure based on the estimated marginal probabilities. In this paper, we aim at recovering the target matrix A_* under a flexible high-dimensional low-rank sampling structure. This is achieved by a weighted empirical risk minimization, with application of inverse probability weighting (e.g., Schnabel et al., 2016; Mao et al., 2019) to adjust for the effect of non-uniform missingness.

Data arising in many applications of matrix completion, such as recom-

mender systems, usually possesses complex “sampling” structure which is largely unknown. For a movie recommender system, users tend to rate movies that they prefer or dislike most, while often remain “silent” to other movies. Another example of the complex sampling regime is in the online merchandising, where some users may purchase certain items regularly without often rating them, but often evaluate products that they rarely buy. Similar to the widely adopted model that ratings are generated from a small number of hidden factors, it is reasonable to believe that the missingness is also governed by a small and possibly different set of hidden factors, which leads to a low-rank model the missing structure.

Inspired by generalized linear models, we model the probabilities of observation $\Theta_{\star} = (\theta_{\star,ij})_{i,j=1}^{n_1,n_2} \in (0,1)^{n_1 \times n_2}$ by a high-dimensional low-rank matrix $M_{\star} = (m_{\star,ij})_{i,j=1}^{n_1,n_2} \in \mathbb{R}^{n_1 \times n_2}$ through a known function f . That means, on the entry-wise level, we have $\theta_{\star,ij} = f(m_{\star,ij})$. In generalized linear models, the linear predictor $m_{\star,ij}$ is further modeled as a linear function of observed covariates. However, to reflect difficulties to attain (appropriate and adequate) covariate information and the complexity in the modeling of Θ_{\star} in some situations of the matrix completion, the predictor matrix M_{\star} is assumed completely hidden in this study. Despite M_{\star} being hidden, as demonstrated in this work, the low-rankness of M_{\star} together with the high dimensionality of M_{\star} allows both identification and consistent estimation of Θ_{\star} , which facilitates inverse probability weighting based matrix completion. Motivated by the nature of matrix completion, we propose a novel parametrization $M_{\star} = \mu_{\star} \mathbf{1}_{n_1} \mathbf{1}_{n_2}^T + Z_{\star}$ where

\mathbf{Z}_* satisfies $\mathbf{1}_{n_1}^\top \mathbf{Z}_* \mathbf{1}_{n_2} = 0$. Our proposal extends the work of Davenport et al. (2014), which aims to solve a binary matrix completion problem and pursues a different goal. Compared with Davenport et al. (2014), the proposed method does not regularize the estimation of μ_* , but only regularize the nuclear norm of the estimation of \mathbf{Z}_* , which require different algorithmic treatment to avoid bias caused by the nuclear-norm penalty.

There are three fundamental aspects that set our work aside from the existing works of matrix completion as we consider: (i) the high-dimensional probability matrix Θ , whose dimensions n_1, n_2 go to infinity in our asymptotic setting; (ii) the diminishing lower bound of the observation probabilities (as n_1, n_2 go to infinity), and added issue to the instability of inverse probability weighting; (iii) the effects of estimation error in inverse probability weighting to the matrix completion procedure. Aspects (i) and (ii) are unique to our problem, and not found in the literature of missing data. The work related to Aspect (iii) is sparse in the literature of matrix completion. It is noted that Mao et al. (2019) focused on a low-dimensional parametric modeling of inverse probability weighting with observable covariates.

We develop non-asymptotic upper bounds of the mean squared errors for the proposed estimators of the observation probabilities and the target matrix. To sustain the convergence rate of the target matrix under the high-dimensionality of \mathbf{M}_* and low levels of observation probabilities, we propose to re-estimate \mathbf{Z}_* by constraining the magnitude of its entries to a smaller threshold. Our anal-

ysis shows that the proposed constrained inverse probability weighting estimator achieves the optimal rate (up to a logarithmic factor in estimation of target matrix). We also compare the inverse probability weighting based completion based on the proposed constrained estimation, with the completion based on direct weight trimming (or winsorization), a known practice in the conventional missing value literature (e.g., Rubin, 2001; Kang and Schafer, 2007; Schafer and Kang, 2008) and show that the constrained estimation has both theoretical and empirical advantages.

2. Model and Method

2.1 General Setup

Let $\mathbf{A}_\star = (a_{\star,ij})_{i,j=1}^{n_1,n_2} \in \mathbb{R}^{n_1 \times n_2}$ be an unknown high-dimensional matrix of interest, and $\mathbf{Y} = (y_{ij})_{i,j=1}^{n_1,n_2}$ be a contaminated version of \mathbf{A}_\star according to the following additive noise model:

$$y_{ij} = a_{\star,ij} + \epsilon_{ij}, \quad \text{for } i = 1, \dots, n_1; j = 1, \dots, n_2, \quad (2.1)$$

where $\{\epsilon_{ij}\}$ are independently distributed random errors with zero mean and finite variance. In the setting of matrix completion, only a portion of $\{y_{ij}\}$ is observed. For the (i, j) th entry, define the sampling indicator $w_{ij} = 1$ if y_{ij} is observed, and 0 otherwise, and assume $\{\epsilon_{ij}\}$ are independent of $\{w_{ij}\}$.

As for the sampling mechanism, we adopt a Bernoulli model where $\{w_{ij}\}$ are independent Bernoulli random variables with observation probabilities $\{\theta_{\star,ij}\}$,

collectively denoted by a matrix $\Theta_\star = (\theta_{\star,ij})_{i,j=1}^{n_1,n_2} \in (0, 1)^{n_1 \times n_2}$. Similar to generalized linear models, the observation probabilities can be expressed in terms of an unknown matrix $M_\star = (m_{\star,ij})_{i,j=1}^{n_1,n_2} \in \mathbb{R}^{n_1 \times n_2}$ and a pre-specified monotone and differentiable function $f : \mathbb{R} \rightarrow [0, 1]$, i.e., $\theta_{\star,ij} = f(m_{\star,ij})$ for all i, j . The matrix M_\star plays the same role as a linear predictor in the generalized linear model, while the function f is an inverse link function. Two popular choices of f are inverse logit function $g(m) = e^m / (1 + e^m)$ (logistic model) and the standard normal cumulative distribution function (probit model).

2.2 Low-rank Modeling of A_\star and M_\star

The above setup is general. Without additional assumption, it is virtually impossible to recover the hidden feature matrix M_\star and also the target matrix A_\star . A common and powerful assumption is that A_\star is a low-rank matrix, i.e., $\text{rank}(A_\star) \ll \min\{n_1, n_2\}$. Take the Yahoo! Webscope data set (to be analyzed in Section 7) as an example. This data set contains a partially observed matrix of ratings from 15,400 users to 1000 songs, and the goal is to complete the rating matrix. The low-rank assumption reflects the belief that users' ratings are generated by a small number of factors, representing several standard preference profiles for songs. This viewpoint has been proven useful in the modeling of recommender systems (e.g., Candès and Plan, 2010; Cai et al., 2010).

The same idea could be adapted to the missing pattern, despite that the factors that induce the missingness may be different from those that generate the

ratings. To this end, we assume M_* is also low-rank. Next, we consider decomposing M_* as

$$M_* = \mu_* \mathbf{J} + \mathbf{Z}_* \quad \text{where} \quad \mathbf{1}_{n_1}^T \mathbf{Z}_* \mathbf{1}_{n_2} = 0 \quad (2.2)$$

with $\mathbf{1}_n$ being an n -vector of ones, and $\mathbf{J} = \mathbf{1}_{n_1} \mathbf{1}_{n_2}^T$. Here μ_* is the mean of M_* , i.e., $\mu_* = \mathbf{1}_{n_1}^T M_* \mathbf{1}_{n_2} / (n_1 n_2)$. Note that this decomposition holds for any matrix M by setting $\mu = (n_1 n_2)^{-1} \mathbf{1}_{n_1}^T M \mathbf{1}_{n_2}$ and $\mathbf{Z} = M - \mu \mathbf{J}$. Moreover, the decomposition is unique due to the constraint that $\mathbf{1}_{n_1}^T \mathbf{Z}_* \mathbf{1}_{n_2} = 0$. The key here is to reparametrize M_* in terms of μ_* and \mathbf{Z}_* , which require different treatments in their estimation. See Section 3 for details. Further, the low-rankness of M_* can be translated to the low-rankness of \mathbf{Z}_* .

We note that the rank of M_* is not the same as that of Θ_* due to the nonlinear transformation f . In general, the low-rank structure of M_* implies a specific low-dimensional nonlinear structure of Θ_* . For a common high missingness scenario, most entries of M_* are significantly negative, where many common choices of the inverse link function can be well-approximated by a linear function. So our modeling can be regarded as a low-rank modeling of Θ_* .

There are some related but more specialized models. Srebro and Salakhutdinov (2010) and Negahban and Wainwright (2012) utilize an independent row and column sampling mechanism, leading to a rank-1 structure for Θ_* . Cai et al. (2016) consider a matrix block structure for Θ_* and hence M_* , which can be regarded as a special case of the low-rank modeling. Mao et al. (2019) considered the case when the missingness depends on observable covariates, and adopted a

low-rank modeling with a known row space of M_* . The proposal in this paper is for the situation when the missingness is dependent on some hidden factors, which reflects situations when obvious covariates are unknown or not available.

2.3 Inverse Probability Weighting Based Matrix Completion: Motivations and Challenges

Write the Hadamard product as \circ and the Frobenius norm as $\|\cdot\|_F$. To recover the target matrix A_* , many existing matrix completion techniques assume uniform missing structure and hence utilize an unweighted/uniform empirical risk function $\widehat{R}_{\text{UNI}}(\mathbf{A}) = (n_1 n_2)^{-1} \|\mathbf{W} \circ (\mathbf{A} - \mathbf{Y})\|_F^2$ (e.g., Candès and Plan, 2010; Koltchinskii et al., 2011; Mazumder et al., 2010), which is an unbiased estimator of the risk $R(\mathbf{A}) = \mathbf{E}(\|\mathbf{A} - \mathbf{Y}\|_F^2)/(n_1 n_2)$ (up to a multiplicative constant) under uniform missingness. The work of Klopp (2014) is a notable exception that considers the use of \widehat{R}_{UNI} under non-uniform missingness.

For any matrix $\mathbf{B} = (b_{ij})_{i,j=1}^{n_1, n_2}$, we denote $\mathbf{B}^\dagger = (b_{ij}^{-1})_{i,j=1}^{n_1, n_2}$ and $\mathbf{B}^\ddagger = (b_{ij}^{-1/2})_{i,j=1}^{n_1, n_2}$. Under general missingness (uniform or non-uniform), one can show that, for any $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$,

$$R(\mathbf{A}) = \frac{1}{n_1 n_2} \mathbf{E}(\|\mathbf{A} - \mathbf{Y}\|_F^2) = \frac{1}{n_1 n_2} \mathbf{E}\left(\|\mathbf{W} \circ \Theta_*^\ddagger \circ (\mathbf{A} - \mathbf{Y})\|_F^2\right).$$

Clearly, A_* uniquely minimizes R . If Θ_* were known, an unbiased estimator of R would be

$$\widehat{R}(\mathbf{A}) = \frac{1}{n_1 n_2} \|\mathbf{W} \circ \Theta_*^\ddagger \circ (\mathbf{A} - \mathbf{Y})\|_F^2, \quad (2.3)$$

which motivates the use of inverse probability weighting in matrix completion as in Mao et al. (2019). In addition, our theoretical analysis shows that the nuclear-norm-regularized empirical risk estimator (to be defined in details later) based on \widehat{R} (assuming the use of true observation probabilities) improves upon existing error upper bound of corresponding estimator based on \widehat{R}_{UNI} achieved by Klopp (2014) as shown in Section 5.3. However, the inverse probability weights Θ_{\star}^{\dagger} are often unknown and have to be estimated, which has to be conducted carefully in the context of matrix completion.

Despite the popularity of inverse probability weighting in missing data literature, it is known to produce unstable estimation due to occurrences of small probabilities (e.g., Rubin, 2001; Kang and Schafer, 2007; Schafer and Kang, 2008). This problematic scenario is common for matrix completion problems for recovering a target matrix from very few observations. Theoretically, a reasonable setup should allow some $\theta_{\star,ij}$ to go to zero as $n_1, n_2 \rightarrow \infty$, leading to diverging weights and a non-standard setup of inverse probability weighting. Due to these observations, a careful construction of the estimation procedure is required.

For uniform sampling ($\theta_{\star,ij} \equiv \theta_0$ for some probability θ_0), one only has to worry about a small common probability θ_0 (or that θ_0 diminishes in an asymptotic sense.) Although small θ_0 increases the difficulty of estimation, $\widehat{R}(\mathbf{A})$ changes only up to a multiplicative constant. However, for non-uniform setting, it is not as straightforward due to the heterogeneity among $\{\theta_{\star,ij}\}$. To demon-

strate the issue, we now briefly look at the Yahoo! Webscope dataset described in Section 7. A sign of the strong heterogeneity in $\{\theta_{*,ij}\}$ is a large θ_U/θ_L , where $\theta_L = \min_{i,j} \theta_{*,ij}$ and $\theta_U = \max_{i,j} \theta_{*,ij}$. We found that the corresponding ratio of estimated probabilities $\widehat{\theta}_U/\widehat{\theta}_L$ based on the rank-1 structure of Negahban and Wainwright (2012) was 25656.2, and that based on our proposed method (without re-estimation, to be described below) was 23988.0. This strong heterogeneity can jeopardize the convergence rate of our estimator, which will be properly addressed in our framework.

In the following section, we propose an estimation approach for Θ_* in Section 3.1 and an appropriate modification in Section 3.3 which, when substituted into the empirical risk \widehat{R} , allows us to construct a stable estimator for \mathbf{A}_* .

3. Estimation of Θ_*

3.1 Regularized Maximum Likelihood Estimation

We develop the estimation of Θ_* based upon the framework of regularized maximum likelihood. Given the inverse of link function f , the log-likelihood function with respect to the indicator matrix $\mathbf{W} = (w_{ij}) \in \mathbb{R}^{n_1 \times n_2}$ is

$$\ell_{\mathbf{W}}(\mathbf{M}) = \sum_{i,j} [\mathbb{I}_{[w_{ij}=1]} \log \{f(m_{ij})\} + \mathbb{I}_{[w_{ij}=0]} \log \{1 - f(m_{ij})\}],$$

for any $\mathbf{M} = (m_{ij})_{i,j=1}^{n_1, n_2} \in \mathbb{R}^{n_1 \times n_2}$, where $\mathbb{I}_{\mathcal{A}}$ is the indicator of an event \mathcal{A} . Due to the low-rank assumption of \mathbf{M}_* , one natural candidate of estimators to \mathbf{M}_* is the maximizer of the regularized log-likelihood $\ell_{\mathbf{W}}(\mathbf{M}) - \lambda \|\mathbf{M}\|_*$, where $\|\cdot\|_*$

represents the nuclear norm and $\lambda > 0$ is a tuning parameter. It is also common to enforce an additional max-norm constraint $\|\mathbf{M}\|_\infty \leq \alpha$ for some $\alpha > 0$ in the maximization (e.g., Davenport et al., 2014). Note that the nuclear norm penalty favors $\mathbf{M} = 0$, corresponding to that $\Pr(w_{ij} = 1) = 0.5$ for all i, j . Nevertheless, this would not align well with common settings of matrix completion under which the average probability of observations is quit small, and hence would result in a large bias. In view of this, we instead adopt a parametrization $\mathbf{M}_\star = \mu_\star \mathbf{J} + \mathbf{Z}_\star$ and consider the following estimator of $(\mu_\star, \mathbf{Z}_\star)$:

$$\left(\hat{\mu}, \hat{\mathbf{Z}}\right) = \arg \max_{(\mu, \mathbf{Z}) \in \mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)} \ell_{\mathbf{W}}(\mu \mathbf{J} + \mathbf{Z}) - \lambda \|\mathbf{Z}\|_*, \text{ where} \quad (3.1)$$

$$\mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2) = \{(\mu, \mathbf{Z}) \in \mathbb{R} \times \mathbb{R}^{n_1 \times n_2} : |\mu| \leq \alpha_1, \|\mathbf{Z}\|_\infty \leq \alpha_2, \mathbf{1}_{n_1}^\top \mathbf{Z} \mathbf{1}_{n_2} = 0\}.$$

Note that the mean μ of the linear predictor $\mu \mathbf{J} + \mathbf{Z}$ is not penalized. The constraint $\mathbf{1}_{n_1}^\top \mathbf{Z} \mathbf{1}_{n_2} = 0$ ensures the identifiability of μ and \mathbf{Z} . Apparently, the constraints in $\mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)$ are analogous to $\|\mathbf{M}\|_\infty \leq \alpha_0$, where $\alpha_0 = \alpha_1 + \alpha_2$, but on the parameters μ and \mathbf{Z} respectively. With $(\hat{\mu}, \hat{\mathbf{Z}})$, the corresponding estimator of \mathbf{M}_\star is $\hat{\mathbf{M}} = \hat{\mu} \mathbf{J} + \hat{\mathbf{Z}}$.

Davenport et al. (2014) considered a regularized maximum likelihood approach for a binary matrix completion problem. Their goal was different, as they aimed at recovering a binary rating matrix in lieu of the missing structure, and considered a regularization on \mathbf{M} (instead of \mathbf{Z}) via $\|\mathbf{M}\|_* \leq \alpha' \{\text{rank}(\mathbf{M}_\star) n_1 n_2\}^{1/2}$. As for the scaling parameter α' , Davenport et al. (2014) considered an α' independent of the dimensions n_1 and n_2 to restrict the “spikiness” of \mathbf{M} . As

explained earlier, in our framework, θ_L should be allowed to go to zero as $n_1, n_2 \rightarrow \infty$. To this end, we allow α_1 and α_2 to depend on the dimensions n_1 and n_2 . See more details in Section 5.

3.2 Computational algorithm and tuning parameter selection

To solve the optimization (3.1), we begin with the observation that $\ell_{\mathbf{W}}$ is a smooth concave function, which allows the usage of an iterative algorithm called accelerated proximal gradient algorithm (Beck and Teboulle, 2009). Given a pair $(\mu_{\text{old}}, \mathbf{Z}_{\text{old}})$ from a previous iteration, a quadratic approximation of the objective function $-\ell_{\mathbf{W}}(\mu\mathbf{J} + \mathbf{Z}) + \lambda\|\mathbf{Z}\|_*$ is formed:

$$\begin{aligned} P_L \{(\mu, \mathbf{Z}), (\mu_{\text{old}}, \mathbf{Z}_{\text{old}})\} = & -\ell_{\mathbf{W}}(\mu_{\text{old}}\mathbf{J} + \mathbf{Z}_{\text{old}}) \\ & + (\mu - \mu_{\text{old}}) \mathbf{1}_{n_1}^T \{-\nabla_{\mu} \ell_{\mathbf{W}}(\mu_{\text{old}}\mathbf{J} + \mathbf{Z}_{\text{old}})\} \mathbf{1}_{n_2} + \frac{Ln_1n_2}{2} (\mu - \mu_{\text{old}})^2 \\ & + \langle \mathbf{Z} - \mathbf{Z}_{\text{old}}, -\nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}}\mathbf{J} + \mathbf{Z}_{\text{old}}) \rangle + \frac{L}{2} \|\mathbf{Z} - \mathbf{Z}_{\text{old}}\|_F^2 + \lambda \|\mathbf{Z}\|_* , \end{aligned}$$

where $L > 0$ is an algorithmic parameter determining the step size of the proximal gradient algorithm, and is chosen by a backtracking method (Beck and Teboulle, 2009). Here $\langle \mathbf{B}, \mathbf{C} \rangle = \sum_{i,j} b_{ij}c_{ij}$ for any matrices $\mathbf{B} = (b_{ij})$ and $\mathbf{C} = (c_{ij})$ of same dimensions.

In this iterative algorithm, a successive update of (μ, \mathbf{Z}) can be obtained by

$$\arg \min_{(\mu, \mathbf{Z}) \in \mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)} P_L \{(\mu, \mathbf{Z}), (\mu_{\text{old}}, \mathbf{Z}_{\text{old}})\},$$

where the optimization with respect to μ and \mathbf{Z} can be performed separately.

For μ , one can derive a closed-form update

$$\min \left[\alpha_1, \max \left[-\alpha_1, \mu_{\text{old}} + (Ln_1n_2)^{-1} \mathbf{1}_{n_1}^T \{-\nabla_{\mu} \ell_{\mathbf{W}}(\mu_{\text{old}}\mathbf{J} + \mathbf{Z}_{\text{old}})\} \mathbf{1}_{n_2} \right] \right].$$

As for \mathbf{Z} , we need to perform the minimization

$$\arg \min_{\|\mathbf{Z}\|_\infty \leq \alpha_2, \mathbf{1}_{n_1}^\top \mathbf{Z} \mathbf{1}_{n_2} = 0} \langle \mathbf{Z} - \mathbf{Z}_{\text{old}}, -\nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \rangle + \frac{L}{2} \|\mathbf{Z} - \mathbf{Z}_{\text{old}}\|_F^2 + \lambda \|\mathbf{Z}\|_*,$$

which is equivalent to

$$\arg \min_{\|\mathbf{Z}\|_\infty \leq \alpha_2, \mathbf{1}_{n_1}^\top \mathbf{Z} \mathbf{1}_{n_2} = 0} \frac{1}{2} \left\| \mathbf{Z} - \mathbf{Z}_{\text{old}} - \frac{1}{L} \nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \right\|_F^2 + \frac{\lambda}{L} \|\mathbf{Z}\|_* . \quad (3.2)$$

We apply a three-block extension of the alternative direction method of multipliers (Chen et al., 2016) to an equivalent form of (3.2):

$$\arg \min_{\mathbf{Z}=\mathbf{G}_1=\mathbf{G}_2, \mathbf{1}_{n_1}^\top \mathbf{G}_1 \mathbf{1}_{n_2}=0, \|\mathbf{G}_2\|_\infty \leq \alpha_2} \frac{\lambda}{L} \|\mathbf{Z}\|_* + \frac{1}{2} \left\| \mathbf{G}_2 - \mathbf{Z}_{\text{old}} - \frac{1}{L} \nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \right\|_F^2 . \quad (3.3)$$

Write $\mathbf{H} = (\mathbf{H}_1, \mathbf{H}_2)$. The augmented Lagrangian for (3.3) is

$$\begin{aligned} \mathcal{L}_u(\mathbf{Z}, \mathbf{G}_1, \mathbf{G}_2; \mathbf{H}) &= \frac{\lambda}{L} \|\mathbf{Z}\|_* + \frac{1}{2} \left\| \mathbf{G}_2 - \mathbf{Z}_{\text{old}} - \frac{1}{L} \nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \right\|_F^2 \\ &\quad - \langle \mathbf{H}_1, \mathbf{Z} - \mathbf{G}_1 \rangle - \langle \mathbf{H}_2, \mathbf{Z} - \mathbf{G}_2 \rangle + \frac{u}{2} \|\mathbf{Z} - \mathbf{G}_1\|_F^2 + \frac{u}{2} \|\mathbf{Z} - \mathbf{G}_2\|_F^2 \\ &\quad + \mathbb{I}_{[\mathbf{1}_{n_1}^\top \mathbf{G}_1 \mathbf{1}_{n_2} = 0]} + \mathbb{I}_{[\|\mathbf{G}_2\|_\infty \leq \alpha_2]}, \end{aligned}$$

where $u > 0$ is an algorithmic parameter and, $\mathbb{I}_{\mathcal{A}} = 0$ if the constraint \mathcal{A} holds and ∞ otherwise. The detailed algorithm to solve (3.3) is summarized in Algorithm 1 in the supplementary material. It is noted that, in general, the multi-block alternative direction method of multipliers may fail to converge for some $u > 0$ (Chen et al., 2016). In those cases, an appropriate selection of u is crucial. However, we are able to show that the form of our algorithm belongs to a special class (Chen et al., 2016) in which convergence is guaranteed for any $u > 0$.

Therefore, we simply set $u = 1$. We summarize the corresponding convergence result in the following theorem whose proof is provided in the supplementary material.

Theorem 1. *The sequence $\{\mathbf{Z}^{(k)}, \mathbf{G}_1^{(k)}, \mathbf{G}_2^{(k)}\}$, generated by Algorithm 1 in the supplementary material, converges to the global optimum of (3.3).*

Notice that the alternative direction method of multipliers algorithm is nested within the proximal gradient algorithm. But, from our practical experiences, both the number of inner iterations (alternative direction method of multipliers) and outer iterations (proximal gradient) are small, usually less than twenty in our numerical experiments. Similarly, we summarize the corresponding convergence result of the overall proximal gradient algorithm in the following theorem whose proof is provided in the supplementary material.

Theorem 2. *The estimator $(\hat{\mu}, \hat{\mathbf{Z}})$ generated by the proximal gradient algorithm, converges to the global optimum of (3.1).*

The tuning parameters α_1 and α_2 can be chosen according to prior knowledge of the problem setup, if available. When a prior knowledge is not available, one can choose large values for these parameters. Once these parameters are large enough, our method is not sensitive to their specific values. A more principled way to tune α_1 and α_2 is a challenging problem and beyond the scope of this work. As for λ , we adopt Akaike information criterion (AIC) where the degree of freedom is approximated by $r_{\widehat{M}}(n_1 + n_2 - r_{\widehat{M}})$.

3.3 Constrained estimation

To use \widehat{R} of (2.3), a naive idea is to obtain $\widehat{\Theta} = (\widehat{\theta}_{ij})_{i,j=1}^{n_1, n_2} = \mathcal{F}(\widehat{\mathbf{M}})$, where $\mathcal{F}(\mathbf{M}) = (f(m_{ij}))_{i,j=1}^{n_1, n_2} \in \mathbb{R}^{n_1 \times n_2}$ for any $\mathbf{M} = (m_{ij})_{i,j=1}^{n_1, n_2} \in \mathbb{R}^{n_1 \times n_2}$, and then replace Θ_{\star}^{\dagger} by $\widehat{\Theta}^{\dagger} = (\widehat{\theta}_{ij}^{-1/2})_{i,j=1}^{n_1, n_2}$. However, this direct implementation is not robust to extremely small probabilities of observation, and our theoretical analysis shows that this could lead to a slower convergence rate in the estimation of \mathbf{A}_{\star} . In the literature of missing data, a simple solution to robustify is winsorizing the small probabilities (Potter, 1990; Scharfstein et al., 1999).

In the estimation of $\widehat{\Theta}$ defined in (3.1) that assumes $\|\mathbf{Z}_{\star}\|_{\infty} \leq \alpha_2$, a large α_2 has an adverse effect on the estimation. In the setting of diverging α_2 (due to diminishing θ_L), the convergence rate of $\widehat{\mathbf{Z}}$ becomes slower and the estimator obtained after direct winsorization will also be affected. That is, even though the extreme probabilities could be controlled by winsorizing, the unchanged entries of $\widehat{\mathbf{Z}}$ (in the procedure of winsorizing) may already suffer from a slower rate of convergence. This results in a larger estimation error under certain settings of missingness, as revealed in Section 5.

A seemingly better strategy is to impose a tighter constraint directly in the minimization problem (3.1). That is to adopt the constraint $\|\mathbf{Z}\|_{\infty} \leq \beta$ where $0 \leq \beta \leq \alpha_2$. Theoretically, one can better control the errors on those entries of magnitude smaller than β . However, the mean-zero constraint of \mathbf{Z} no longer makes sense as the constraint $\|\mathbf{Z}\|_{\infty} \leq \beta$ may have shifted the mean.

We propose a re-estimation of \mathbf{Z}_* with a different constraint level β :

$$\widehat{\mathbf{Z}}_\beta = \arg \max_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \ell_{\mathbf{W}}(\widehat{\mu} \mathbf{J} + \mathbf{Z}) - \lambda' \|\mathbf{Z}\|_* \quad \text{subject to} \quad \|\mathbf{Z}\|_\infty \leq \beta. \quad (3.4)$$

Note that we only re-compute \mathbf{Z} but not μ , which allows us to drop the mean-zero constraint. Thus, $\widehat{\mathbf{M}}_\beta = \widehat{\mu} \mathbf{J} + \widehat{\mathbf{Z}}_\beta$. The corresponding algorithm for optimization (3.4) can be derived similarly as in Davenport et al. (2014), and is provided in the supplementary material. In what follows, we write $\widehat{\Theta} = \mathcal{F}(\widehat{\mathbf{M}})$ and $\widehat{\Theta}_\beta = \mathcal{F}(\widehat{\mathbf{M}}_\beta)$.

4. Estimation of \mathbf{A}_*

Now, we come back to (2.3) and replace Θ_*^\dagger by $\widehat{\Theta}_\beta^\dagger$ to obtain a modified empirical risk:

$$\widetilde{R}(\mathbf{A}) = \frac{1}{n_1 n_2} \left\| \mathbf{W} \circ \widehat{\Theta}_\beta^\dagger \circ (\mathbf{A} - \mathbf{Y}) \right\|_F^2, \quad (4.1)$$

where $\widehat{\Theta}_\beta^\dagger = (\widehat{\theta}_{ij,\beta}^{-1/2}) \in \mathbb{R}^{n_1 \times n_2}$. Since \mathbf{A} is a high-dimensional parameter, a direct minimization of \widetilde{R}^* often results in over-fitting. To circumvent it, we consider a regularized version:

$$\widetilde{R}(\mathbf{A}) + \tau \|\mathbf{A}\|_*, \quad (4.2)$$

where $\tau > 0$ is a regularization parameter. Again, the nuclear norm regularization encourages low-rank solution. Based on (4.2), our estimator of \mathbf{A}_* is

$$\widehat{\mathbf{A}}_\beta = \arg \min_{\|\mathbf{A}\|_\infty \leq a} \left\{ \frac{1}{n_1 n_2} \left\| \mathbf{W} \circ \widehat{\Theta}_\beta^\dagger \circ (\mathbf{A} - \mathbf{Y}) \right\|_F^2 + \tau \|\mathbf{A}\|_* \right\}, \quad (4.3)$$

where a is an upper bound on $\|\mathbf{A}_\star\|_\infty$. The above $\widehat{\mathbf{A}}_\beta$ contains as special cases (i) the matrix completion $\widehat{\mathbf{A}}_{\alpha_2}$, with unconstrained probability estimator $\widehat{\Theta}$, by setting $\beta = \alpha_2$ and (ii) the estimator $\widehat{\mathbf{A}}_\beta$, with constrained probability estimator $\widehat{\Theta}_\beta$, when $\beta < \alpha_2$.

We use an accelerated proximal gradient algorithm (Beck and Teboulle, 2009) to solve (4.3). For the choice of tuning parameter τ in (4.3), we adopt a 5-fold cross-validation to select the remaining tuning parameters. Due to the non-uniform missing mechanism, we use a weighted version of the validation errors. The specific details are given in Algorithm 3 in the supplementary material.

5. Theoretical Properties

5.1 Probabilities of observation

Let $\|\mathbf{B}\| = \sigma_{\max}(\mathbf{B})$, $\|\mathbf{B}\|_\infty = \max_{i,j} |b_{ij}|$ and $\|\mathbf{B}\|_{\infty,2} = (\max_i \sum_j b_{ij}^2)^{1/2}$ be the spectral norm, the maximum norm and $l_{\infty,2}$ -norm of a matrix \mathbf{B} respectively.

We use the symbol \asymp to represent the asymptotic equivalence in order, i.e., $a_n \asymp b_n$ if $a_n = O(b_n)$ and $b_n = O(a_n)$. The average squared distance between two matrices $\mathbf{B}, \mathbf{C} \in \mathbb{R}^{n_1 \times n_2}$ is $d^2(\mathbf{B}, \mathbf{C}) = \|\mathbf{B} - \mathbf{C}\|_F^2 / (n_1 n_2)$. The average squared errors of $\widehat{\mathbf{M}}_\beta$ and $\widehat{\Theta}_\beta^\dagger$ are then $d^2(\widehat{\mathbf{M}}_\beta, \mathbf{M}_\star)$ and $d^2(\widehat{\Theta}_\beta^\dagger, \Theta_\star^\dagger)$ respectively. We adopt the Hellinger distance for any two matrices $S, T \in [0, 1]^{n_1 \times n_2}$, $d_H^2(S, T) = (n_1 n_2)^{-1} \sum_{i,j=1}^{n_1, n_2} d_H^2(s_{ij}, t_{ij})$ where $d_H^2(s, t) = (s^{1/2} - t^{1/2})^2 + \{(1 -$

$s)^{1/2} - (1 - t)^{1/2}\}^2$ for $s, t \in [0, 1]$. In the literature of matrix completion, most discussions related to optimal convergence rate are only up to certain polynomial orders of $\log n$. For convenience, we use $\text{polylog}(n)$ for polynomials of $\log n$.

To investigate the asymptotic properties of $\widehat{\mathbf{M}}_\beta$ and $\widehat{\Theta}_\beta^\dagger$ defined in Section 3, we introduce the following conditions on the missing structure.

C1. The indicators $\{w_{ij}\}_{i,j=1}^{n_1, n_2}$ are mutually independent, and independent of $\{\epsilon_{ij}\}_{i,j=1}^{n_1, n_2}$. For $i = 1, \dots, n_1$ and $j = 1, \dots, n_2$, w_{ij} follows a Bernoulli distribution with probability of success $\theta_{*,ij} = f(m_{*,ij}) \in (0, 1)$. Furthermore, f is monotonic increasing and differentiable.

C2. The hidden feature matrix $\mathbf{M}_* = \mu_* \mathbf{J} + \mathbf{Z}_*$ where $\mathbf{1}_{n_1}^T \mathbf{Z}_* \mathbf{1}_{n_2} = 0$, $|\mu_*| \leq \alpha_1 < \infty$ and $\|\mathbf{Z}_*\|_\infty \leq \alpha_2 < \infty$. Here α_1 and α_2 are allowed to depend on the dimensions n_1 and n_2 . This also implies that there exists a lower bound $\theta_L \in (0, 1)$ (allowed to depend on n_1, n_2) such that $\min_{i,j} \{\theta_{ij}\} \geq \theta_L \geq f(-\alpha_1 - \alpha_2) > 0$.

For the convenience of the theoretical analysis, we consider an equivalent estimator of (μ_*, \mathbf{Z}_*) defined by the constrained maximization problem (5.1) instead of the Lagrangian form (3.1). For $r_{\mathbf{Z}_*} \leq \min\{n_1, n_2\}$ and $\alpha_1, \alpha_2 \geq 0$,

$$\begin{aligned} (\widehat{\mu}, \widehat{\mathbf{Z}}) &= \arg \max_{(\mu, \mathbf{Z}) \in \widetilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2)} \ell_{\mathbf{W}}(\mu \mathbf{J} + \mathbf{Z}), \text{ where} & (5.1) \\ \widetilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2) &= \{(\mu, \mathbf{Z}) \in \mathbb{R} \times \mathbb{R}^{n_1 \times n_2} : |\mu| \leq \alpha_1, \|\mathbf{Z}\|_\infty \leq \alpha_2, \\ &\quad \|\mathbf{Z}\|_* \leq \alpha_2 \sqrt{r_{\mathbf{Z}_*} n_1 n_2}, \mathbf{1}_{n_1}^T \mathbf{Z} \mathbf{1}_{n_2} = 0\}. \end{aligned}$$

It is easy to see that we have $(\mu_*, \mathbf{Z}_*) \in \widetilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2)$ once $(\mu_*, \mathbf{Z}_*) \in \mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)$

holds. For the ease of presentation, we assume $n_1 = n_2 = n$ and choose the logit function as the inverse link function f in the rest of Section 5, while corresponding results under general settings of n_1, n_2 and f are delegated to Section S1.3 in the supplementary material. We first establish the convergence results for $\hat{\mu}$, $\hat{\mathbf{Z}}$ and $\hat{\mathbf{M}}$, respectively. To simplify notations, let $\alpha_0 = \alpha_1 + \alpha_2$, $h_{\alpha_1, \beta} = (1 + e^{\alpha_1 + \beta})^{-1}$ and $\Gamma_n = e^{\alpha_0}(\alpha_1 + \alpha_2 r_{\mathbf{Z}_*}^{1/2})n^{-1/2}$.

Lemma 1. *Suppose Conditions C1-C2 hold, and $(\mu_*, \mathbf{Z}_*) \in \mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)$. Consider $\hat{\mathbf{M}} = \hat{\mu}\mathbf{J} + \hat{\mathbf{Z}}$ where $(\hat{\mu}, \hat{\mathbf{Z}})$ is the solution to (5.1). There exist positive constants C_1, C_2 such that we have with probability at least $1 - C_1/n$,*

$$\begin{aligned} (\mu_* - \hat{\mu})^2 &\leq C_2 (\alpha_1^2 \wedge \Gamma_n), \quad \frac{1}{n^2} \left\| \hat{\mathbf{Z}} - \mathbf{Z}_* \right\|_F^2 \leq C_2 (\alpha_2^2 \wedge \Gamma_n) \\ \text{and } \frac{1}{n^2} \left\| \hat{\mathbf{M}} - \mathbf{M}_* \right\|_F^2 &\leq C_2 (\alpha_0^2 \wedge \Gamma_n). \end{aligned} \quad (5.2)$$

The upper bounds in (5.2) all consist of trivial bounds α_j^2 and a more dedicated bound Γ_n . The trivial upper bounds α_1^2, α_2^2 and α_0^2 can be easily derived from the constraint set $\mathcal{C}_{n_1, n_2}(\alpha_1, \alpha_2)$. For extreme settings of increasing α_0 , the more dedicated bound Γ_n is diverging and the trivial bounds may provide better control. The term Γ_n can be controlled by the rank of \mathbf{Z}_* . For a range of non-extreme scenarios, i.e., $\alpha_0 \leq 1/2 \log n$ or $\theta_L \geq n^{-1/2}$, the second term in Γ_n attains the convergence order once $r_{\mathbf{Z}_*} = O(1)$.

Similarly, we study the theoretical results of the re-estimation of \mathbf{Z}_* in terms

of the constrained optimization:

$$\begin{aligned} \widehat{\mathbf{Z}}_\beta &= \arg \max_{\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}} \ell_{\mathbf{W}}(\widehat{\boldsymbol{\mu}}\mathbf{J} + \mathbf{Z}) \\ \text{subject to } \quad &\|\mathbf{Z}\|_\infty \leq \beta, \quad \|\mathbf{Z}\|_* \leq \beta \sqrt{r_{\mathcal{T}_\beta(\mathbf{Z}_*)} n_1 n_2}. \end{aligned} \quad (5.3)$$

We now consider the constrained estimation for \mathbf{Z}_* , \mathbf{M}_* and $\boldsymbol{\Theta}_*^\dagger$. For any matrix $\mathbf{B} = (b_{ij})_{i,j=1}^{n_1, n_2}$, define the winsorizing operator \mathcal{T}_β by $\mathcal{T}_\beta(\mathbf{B}) = (T_\beta(b_{ij}))$ where

$$T_\beta(b_{ij}) = b_{ij} \mathbb{I}_{[-\beta \leq b_{ij} \leq \beta]} + \beta \mathbb{I}_{[b_{ij} > \beta]} - \beta \mathbb{I}_{[b_{ij} < -\beta]} \quad \text{for any } \beta \geq 0. \quad (5.4)$$

Write $\mathbf{M}_{*,\beta} = \mu_* \mathbf{J} + \mathcal{T}_\beta(\mathbf{Z}_*)$ and $\widehat{\mathbf{M}}_{*,\beta} = \widehat{\boldsymbol{\mu}}\mathbf{J} + \mathcal{T}_\beta(\mathbf{Z}_*)$, and $\boldsymbol{\Theta}_{*,\beta} = \mathcal{F}(\mathbf{M}_{*,\beta})$ and $\widehat{\boldsymbol{\Theta}}_{*,\beta} = \mathcal{F}(\widehat{\mathbf{M}}_{*,\beta})$ respectively. It is noted that $\widehat{\mathbf{M}}_{*,\beta}$ serves as a ‘‘bridge’’ between the underlying $\mathbf{M}_{*,\beta}$ and the empirical $\widehat{\mathbf{M}}_\beta$. Write $N_\beta = \sum_{i,j} (\mathbb{I}_{[z_{*,ij} > \beta]} + \mathbb{I}_{[z_{*,ij} < -\beta]})$ as the number of extreme values in \mathbf{Z}_* at level β . The convergence rates of $d^2(\widehat{\mathbf{M}}_\beta, \mathbf{M}_*)$ and $d^2(\widehat{\boldsymbol{\Theta}}_\beta^\dagger, \boldsymbol{\Theta}_*^\dagger)$ are investigated in the next theorem. Define $\Lambda_n = \min[\beta^2, \widetilde{\Gamma}_n + h_{\alpha_1, \beta}^{-1} n^{-2} \beta \{8N_\beta + (n^2 - N_\beta)|\mu_* - \widehat{\boldsymbol{\mu}}|\}]$ where $\widetilde{\Gamma}_n = h_{\alpha_1, \beta}^{-1} (\alpha_1 + \beta r_{\mathcal{T}_\beta(\mathbf{Z}_*)}^{1/2}) n^{-1/2}$.

Theorem 3. *Assume that Conditions C1-C2 hold, and $(\mu_*, \mathbf{Z}_*) \in \widetilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2)$. Consider $\widehat{\mathbf{M}}_\beta = \widehat{\boldsymbol{\mu}}\mathbf{J} + \widehat{\mathbf{Z}}_\beta$ where $\widehat{\mathbf{Z}}_\beta$ is the solution to (5.3) and $\beta \geq 0$, there exist some positive constants C_1 , C_2 and C_3 such that we have with probability at least $1 - 2C_1/n$,*

$$\begin{aligned} d^2 \left\{ \widehat{\mathbf{Z}}_\beta, \mathcal{T}_\beta(\mathbf{Z}_*) \right\} &\leq C_3 \Lambda_n, \quad d^2 \left(\widehat{\mathbf{M}}_\beta, \mathbf{M}_* \right) \leq C_2 (\alpha_1^2 \wedge \Gamma_n) + C_3 \Lambda_n + \frac{2(\alpha_2 - \beta)_+^2 N_\beta}{n^2} \\ \text{and } d^2 \left(\widehat{\boldsymbol{\Theta}}_\beta^\dagger, \boldsymbol{\Theta}_*^\dagger \right) &\leq \frac{C_2}{h_{\alpha_1, \beta}^2} (\alpha_1^2 \wedge \Gamma_n) + \frac{C_3 \Lambda_n}{h_{\alpha_1, \beta}^2} + \frac{8N_\beta}{n^2 \theta_L^2}. \end{aligned} \quad (5.5)$$

We can derive an upper bound $4\beta^2$ for $d^2(\widehat{\mathbf{Z}}_{\text{Win},\beta}, \mathcal{T}_\beta(\mathbf{Z}_*))$ from the second term in Theorem 1 where $\widehat{\mathbf{Z}}_{\text{Win},\beta} = \mathcal{T}_\beta(\widehat{\mathbf{Z}})$ is directly winsorized from $\widehat{\mathbf{Z}}$. Obviously, the order of this upper bound is larger than or equal to Λ_n . Moreover, there are scenarios where Λ_n is a smaller order of β^2 . To illustrate, assume that both $\alpha_1 \asymp 1$ and $\beta \asymp 1$, we have $h_{\alpha_1,\beta} \asymp 1$. Once we have $N_\beta = o(n)$, $r_{\mathcal{T}_\beta(\mathbf{Z}_*)} = o(n)$ and $|\widehat{\mu} - \mu_*| = o(1)$, then $\Lambda_n = o(\beta^2)$.

With a more dedicated investigation of (5.5), one can derive an upper bound for $d^2(\widehat{\Theta}_\beta^\dagger, \widehat{\Theta}_{*,\beta}^\dagger)$, which will be used in Section 5.2. Denote $k'_{\alpha_1,\alpha_2,n} = \min\{\alpha_1^2, e^{\alpha_0}(\alpha_1 + \alpha_2 r_{\mathbf{Z}_*}^{1/2})n^{-1/2}\}$, such an upper bound is of order $k_{\alpha_1,\alpha_2,\beta,n} h_{\alpha_1,\beta}^{-2}$ where

$$k_{\alpha_1,\alpha_2,\beta,n} \asymp \min \left[\beta^2, h_{\alpha_1,\beta}^{-1} \beta \left\{ 8N_\beta + (n^2 - N_\beta) k_{\alpha_1,\alpha_2,n}' \right\} n^{-2} + h_{\alpha_1,\beta}^{-1} n^{-1/2} (\alpha_1 + \beta r_{\mathcal{T}_\beta(\mathbf{Z}_*)}^{1/2}) \right].$$

5.2 Target matrix

To study the convergence of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$, we require the following conditions regarding the random errors ϵ and the target matrix \mathbf{A}_* . Recall that $\widehat{\mathbf{A}}_\beta$ includes both the estimations obtained with the unconstrained estimator $\widehat{\Theta}$ and the constrained estimator $\widehat{\Theta}_\beta$ as $\widehat{\mathbf{A}}(\widehat{\Theta}) = \widehat{\mathbf{A}}_{\alpha_2}$ with $\beta = \alpha_2$.

C3. (a) The random errors $\{\epsilon_{ij}\}$ in Model (2.1) are independently distributed random variables such that $\mathbf{E}(\epsilon_{ij}) = 0$ and $\mathbf{E}(\epsilon_{ij}^2) = \sigma_{ij}^2 < \infty$ for all i, j . (b) For some finite positive constants c_σ and η , $\max_{i,j} \mathbf{E}|\epsilon_{ij}|^l \leq \frac{1}{2} l! c_\sigma^2 \eta^{l-2}$ for any positive integer $l \geq 2$.

C4. There exists a positive constant a_0 such that $\|\mathbf{A}_*\|_\infty \leq a_0$.

Denote $h_{(1),\beta} = \max_{i,j}(\theta_{\star,i,j}^{-1}\theta_{\star,i,j,\beta})$ and

$$\Delta = \max \left\{ \frac{(c_\sigma \vee a) e^{-\mu_\star/2 + \alpha_2 - \beta + |\alpha_2/2 - \beta|} (n \log n)^{1/2}}{n^2}, \frac{\eta e^{\mu_\star/2 + \alpha_1 + |\alpha_2/2 - \beta|} k_{\alpha_1, \alpha_2, \beta, n}^{1/2} \log^{3/2} n}{h_{\alpha_1, \beta} n} \right\}. \quad (5.6)$$

The following theorem established a general upper bound for $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$.

Theorem 4. *Assume Conditions C1-C4 hold. For $\beta \geq 0$, there exist some positive constants C_4, C_5, C_6 and C_7 , both independent of β , such that for $h_{(1),\beta}\tau \geq C_4\Delta$, we have with probability at least $1 - 3/(2n)$,*

$$d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star) \leq \max \left\{ C_6 n^2 h_{(1),\beta}^2 r_{\mathbf{A}_\star} \tau^2 + C_7 h_{(1),\beta}^2 h_{(2),\beta}^2 r_{\mathbf{A}_\star} n^{-1} \log(n), C_5 h_{(1),\beta} h_{(3),\beta} n^{-1} \log^{1/2}(n) \right\}. \quad (5.7)$$

As for the estimator of the target matrix based on direct winsorization $\widehat{\Theta}_{\text{Win},\beta} = \mathcal{F}(\widehat{\mu}\mathbf{J} + \widehat{\mathbf{Z}}_{\text{Win},\beta})$ where $\widehat{\mathbf{Z}}_{\text{Win},\beta} = \mathcal{T}_\beta(\widehat{\mathbf{Z}})$, an upper bound can be derived using Theorem 3. As noted in a remark after Theorem 3, $d^2(\widehat{\mathbf{Z}}_{\text{Win},\beta}, \mathcal{T}_\beta(\mathbf{Z}_\star))$ converges at a slower rate β^2 which will cause a larger error bound for the target matrix.

Now, we discuss the rates of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$ under various missing structures. For simplicity, the following discussion focuses on the low-rank linear predictor (\mathbf{M}_\star) setting such that $r_{\mathbf{M}_\star} \asymp 1$.

Uniform missingness. Under the uniform missingness, i.e., $\theta_{ij} \equiv \theta_0$, it has been shown in Koltchinskii et al. (2011) that $\theta_0^{-1}n^{-1}\text{polylog}(n)$ is the optimal rate for $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$. Therefore it is reasonable to require $\alpha_1 + \alpha_2 = \alpha_0 = O(\text{polylog}(n))$ for the convergence of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$. Under the uniform missingness, we have $\alpha_2 = 0$, $\alpha_0 = \alpha_1$ and $e^{\mu_\star} \asymp \theta_0$. For $\beta = 0$, our estimator $\widehat{\mathbf{A}}_\beta$ degenerates to the estimator based on the unweighted empirical risk func-

tion. Theorem 4 shows that achieves the optimal rate $\theta_0^{-1}n^{-1}\text{polylog}(n)$. As for $\beta > 0$, by taking $\beta \rightarrow 0$ such that $k_{\alpha_1, \alpha_2, \beta, n} = O(e^{\mu_* - 2\alpha_1 - 2\beta}n^{-1}\log^{-2}n)$, the estimator can also reach the optimal rate. Of interest here is that β is allowed to be strictly positive to achieve the same rate.

Non-uniform missingness. Under the non-uniform missingness, suppose the lower and upper bounds of observation probability satisfy $\theta_L \asymp e^{\mu_* - \alpha_2}$ and $\theta_U \asymp e^{\mu_* + \alpha_2}$. For the non-constrained case of $\beta = \alpha_2$ and $h_{\alpha_1, \beta} \asymp e^{-\alpha_1 - \alpha_2}$, the second term of Δ in (5.6) dominates due to the fact that

$$e^{-\mu_*/2 + \alpha_2/2}n^{-3/2}\log^{1/2}n = o(e^{\mu_*/2 + 5\alpha_1/2 + 3\alpha_2/2}n^{-5/4}\log^{3/2}n).$$

Thus, the convergence rate of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$ is $e^{\mu_* + 5\alpha_1 + 3\alpha_2}n^{-1/2}\log^3n$. To guarantee convergence, as $e^{\mu_*/2 + 5\alpha_1/2 + 3\alpha_2/2} \leq e^{3\alpha_1 + 3\alpha_2/2}$, it requires that $\alpha_1 + \alpha_2/2 < (1/12)\log n$ which implies that $\theta_L^{-1} = O(n^{1/6})$.

However, the above range of $\theta_L^{-1} = O(n^{1/6})$ excludes $\theta_L \equiv (n^{-1}\text{polylog}(n))$, the case that results in the number of the observed matrix entries at the order of $n \text{polylog}(n)$ which represents the most sparse case of observation where the matrix can still be recovered (Candès and Recht, 2009; Candès and Plan, 2010; Koltchinskii et al., 2011; Negahban and Wainwright, 2012). We will show in the following that with an appropriately chosen β , the constrained estimator $\widehat{\Theta}_\beta$ can accommodate the case of $\theta_L^{-1} = O(n \log^{-1}n)$.

Case (I): $\beta = 0$. To demonstrate this, we start with the absolute constrained case, i.e., $\beta = 0$, which forces the estimated probabilities to be uniform and implies $e^{-\mu_*/2 + \alpha_2 - \beta + |\alpha_2/2 - \beta|} = e^{-\mu_*/2 + 3\alpha_2/2} \asymp \theta_U^{1/2}\theta_L^{-1}$. Then, accord-

ing to Theorem 4, $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$ attains the convergence rate $\theta_U \theta_L^{-2} n^{-1} \log(n)$, which converges to 0 provided $\theta_U \theta_L^{-2} = o(n \log^{-1} n)$. , the condition $\theta_U \theta_L^{-2} = o(n \log^{-1} n)$ includes the extreme case of $\theta_L^{-1} = O(n \log^{-1} n)$ and n polylog(n) observations.

Case (II): $\beta > 0$. For the more interesting setting $\beta > 0$, to simplify the discussion, we concentrate on the case when the first term in $k_{\alpha_1, \alpha_2, \beta, n}$ is of a smaller order, which can be achieved by choosing $\beta = O(e^{-\mu_\star - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$. Then, according to Theorem 4,

$$d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star) = O_p(e^{-\mu_\star + 2\alpha_2 - 2\beta + 2|\alpha_2/2 - \beta|} n^{-1} \log n) = O_p(e^{\alpha_1/2 + 3\alpha_2/2} n^{-1} \log n),$$

since $e^{-\mu_\star/2 + \alpha_2 - \beta + |\alpha_2/2 - \beta|} \leq e^{\alpha_1/2 + 3\alpha_2/2}$. In the following we consider two further cases: (i) $\alpha_2 = O((\log \log n)^{-1} \alpha_1)$ and (ii) $\alpha_1 = o(\alpha_2 \log \log n)$. Note that for either cases, $e^{-\mu_\star + 2\alpha_2 - 2\beta + 2|\alpha_2/2 - \beta|} \asymp \theta_U \theta_L^{-2}$ which leads to

$$d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star) = O_p(\theta_U \theta_L^{-2} n^{-1} \log n).$$

If $\alpha_2 = O\{(\log \log n)^{-1} \alpha_1\}$, we require $\alpha_1 < (1 + 3 \log \log n)^{-1} (\log n - \log \log n)$ to guarantee convergence, which implies that $\theta_L = O(n^{-1})$. Thus, we only lose a polylog(n) factor when compared with the most extreme but feasible setting of $\theta_L^{-1} = O[n\{\text{polylog}(n)\}^{-1}]$. Also $\beta = O(e^{-\mu_\star - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$ implies that $\beta = O(n^{-1/2} \log^{-1} n)$.

If $\alpha_1 = o\{(\log \log n) \alpha_2\}$, we require that $\alpha_2 < \{3 + (\log \log n)^{-1}\}^{-1} (\log n - \log \log n)$ which leads to $\theta_L^{-1} = O(n^{1/3})$. Also $\beta = O(e^{-\mu_\star - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$ implies that $\beta = O(n^{-1/6} \log^{-1} n)$. However, to make $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_\star)$ convergent,

the attained rate for θ_L^{-1} has to be $O(n^{1/3})$, which excludes the most extreme heterogeneity case of $\theta_L^{-1} = O\{n(\text{polylog}(n))^{-1}\}$. The reason for not being able to cover the most extreme case of $\theta_L^{-1} = O\{n(\text{polylog}(n))^{-1}\}$ is that the current Case (ii) allows more heterogeneity in \mathbf{Z}_* as reflected by having a larger α_2 than that prescribed under Case (i). As μ_* is jointly estimated with \mathbf{Z}_* in the unconstrained estimation (Section 3.1), stronger heterogeneity slows down the convergence rate in the estimation of μ_* , which becomes a bottleneck for further improvement. If μ_* was observable, the problem would not be as serious despite the adverse effect of stronger heterogeneity on the estimation of \mathbf{Z}_* .

To summarize, under the uniform missing and Case (I), (II)(i) in the non-uniform missing, we can achieve the optimal rate up to a $\text{polylog}(n)$ order. For Case (II)(ii), when the missingness is not extreme, with an appropriately chosen $\beta > 0$, the proposed estimator can also attain the optimal rate up to the $\text{polylog}(n)$ order.

5.3 Comparison with Uniform Objective Function

Recall that the unweighted empirical risk function $\widehat{R}_{\text{UNI}}(\mathbf{A}) = n^{-2} \|\mathbf{W} \circ (\mathbf{A} - \mathbf{Y})\|_F^2$ is adopted by many existing matrix completion techniques (Klopp, 2014). An interesting question is whether there is any benefit in adopting the proposed weighted empirical risk function for matrix completion. In this subsection, we aim to shed some light on this aspect by comparing the non-asymptotic error bounds of the corresponding estimators. Due to the additional complication from

the estimation error of the observation probability matrix, we only focus on the weighted empirical risk function with true inverse probability weighting in this section. We will demonstrate empirically in Sections 6 and 7 the benefits of the weighted objective function with estimated weights.

Most existing work with unweighted empirical risk function assume the true missingness is uniform (Candès and Plan, 2010; Koltchinskii et al., 2011). One notable exception is Klopp (2014), where unweighted empirical risk function is studied under possibly non-uniform missing structure. The estimator of Klopp (2014) is equivalent to our estimator when $\beta = 0$, which is denoted by $\widehat{\mathbf{A}}^{\text{UNI}}$. Thus, according to Theorem 4, we have with probability at least $1 - 3/(2n)$,

$$d^2(\widehat{\mathbf{A}}^{\text{UNI}}, \mathbf{A}_\star) \leq \min \left\{ (C_6 + C_7)r_{\mathbf{A}_\star} \theta_U \theta_L^{-2} n^{-1} \log n, C_5 \theta_U^{1/2} \theta_L^{-1} n^{-1/2} \log^{1/2} n \right\} = U^{\text{UNI}},$$

which is the same upper bound obtained in Klopp (2014). Define $\widehat{\mathbf{A}}^{\text{KNOWN}}$ as the estimator which minimizes the known weighted empirical risk function (2.3). Then,

$$d^2(\widehat{\mathbf{A}}^{\text{KNOWN}}, \mathbf{A}_\star) \leq \min \left\{ (C_6 + C_7)r_{\mathbf{A}_\star} \theta_L^{-1} n^{-1} \log n, C_5 \theta_L^{-1/2} n^{-1/2} \log^{1/2} n \right\} = U^{\text{KNOWN}}.$$

The improvement in the upper bounds of the weighted objective function \widehat{R} lies in that, under non-uniform missingness, $\theta_U \theta_L^{-1} > 1$ which implies that $U^{\text{KNOWN}} < U^{\text{UNI}}$ as summarized below.

Theorem 5. *Assume Conditions C1-C4 holds, and take $\tau_{\text{KNOWN}} = C_3 \theta_L^{-1/2} n^{-3/2} \log^{1/2} n$ and $\tau_{\text{UNI}} = C_3 \theta_U^{1/2} f^{-1}(\mu_\star) n^{-3/2} \log^{1/2} n$. The upper bound of $d^2(\widehat{\mathbf{A}}^{\text{UNI}}, \mathbf{A}_\star)$ is the same as U^{UNI} and the upper bound of $d^2(\widehat{\mathbf{A}}^{\text{KNOWN}}, \mathbf{A}_\star)$ is the same as U^{KNOWN} . In addition, $U^{\text{KNOWN}} \leq U^{\text{UNI}}$, and $U^{\text{KNOWN}} < U^{\text{UNI}}$ if $\theta_U > \theta_L$, i.e., the true missing mechanism is non-uniform.*

Our approach draws inspiration from the missing value literature, for instance in Chen et al. (2008), which showed that using the estimated parameters in the inverse probability weighting can actually reduce the variance of the parameter of interest; see Theorem 1 of the paper. Given the results of Chen et al. (2008), we would expect using the estimated parameters $\hat{\Theta}_\beta$ in the weighting probability would not be inferior to the version with the true parameter $\hat{\Theta}_*$. It was verified by numerical studies at least, although the exact theoretical proof cannot be achieved.

6. Simulation Study

6.1 Missingness

This section reports results from simulation experiments which were designed to evaluate the numerical performance of the proposed methodologies. We first evaluate the estimation performances of the observation probabilities in Section 6.1 and then those of the target matrix in Section 6.2. In the simulation, the true observation probabilities Θ_* and the target matrix \mathbf{A}_* were randomly generated once and kept fixed for each simulation setting to be described below. To generate Θ_* , we first generated $\mathbf{U}_{M_*} \in \mathbb{R}^{n_1 \times (r_{M_*} - 1)}$ and $\mathbf{V}_{M_*} \in \mathbb{R}^{(r_{M_*} - 1) \times n_2}$ as random Gaussian matrices with independent entries each following $\mathcal{N}(-0.4, 1)$. We then obtained $\mathbf{M}_* = \mathbf{U}_{M_*} \mathbf{V}_{M_*}^\top - \bar{m}_{n_1, n_2, r_{M_*}} \mathbf{J}$ where $\bar{m}_{n_1, n_2, r_{M_*}}$ is a scalar chosen to ensure the average observation rate is 0.2 in each simulation setting.

We finally set $\Theta_\star = \mathcal{F}(M_\star)$ where the inverse link function f is a logistic function.

In our study, we set $r_{M_\star} = 11$, (or $r_{Z_\star} = 10$) and chose $n_1 = n_2$ with four sizes: 600, 800, 1000 and 1200, and the number of simulation runs for each settings was 500. For the purpose of benchmarking, we compared various estimators of the missingness:

1. the non-constrained estimator $\hat{\Theta}_\alpha$ defined in (3.1);
2. the constrained estimator $\hat{\Theta}_\beta$ defined in (3.4);
3. the directly winsorized estimator $\hat{\Theta}_{\text{Win},\beta} = \mathcal{F}\{\hat{\mu}\mathbf{J} + \mathcal{T}_\beta(\hat{\mathbf{Z}})\}$;
4. the 1-bit estimator $\hat{\Theta}_{1\text{-bit},\alpha}$ proposed in Davenport et al. (2014) and its corresponding constrained and winsorized versions $\hat{\Theta}_{1\text{-bit},\beta}$ and $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$; (note that the 1-bit estimator $\hat{\Theta}_{1\text{-bit},\alpha}$ imposes the nuclear-norm regularization on the whole M instead of Z , when compared to $\hat{\Theta}_\alpha$)
5. the rank-1 probability estimator $\hat{\Theta}_{\text{NW}}$ used in Negahban and Wainwright (2012) where $g_{i\cdot} = n_2^{-1} \sum_{j=1}^{n_2} w_{ij}$, $g_{\cdot j} = n_1^{-1} \sum_{i=1}^{n_1} w_{ij}$ and $\theta_{ij,\text{NW}} = g_{i\cdot} g_{\cdot j}$;
6. the uniform estimator $\hat{\Theta}_{\text{UNI}} = N/(n_1 n_2) \mathbf{J}$.

For the non-constrained estimator $\hat{\Theta}_\alpha$ and the 1-bit estimator $\hat{\Theta}_{1\text{-bit},\alpha}$, the parameter α is set according to the knowledge of the true M_\star . For the constrained estimators $\hat{\Theta}_\beta$ and $\hat{\Theta}_{\text{Win},\beta}$, the constraint level β was chosen so that either 5% or 10% of the elements in $\hat{\mathbf{Z}}_\alpha$ were winsorized. Similarly for $\hat{\Theta}_{1\text{-bit},\beta}$

and $\widehat{\Theta}_{1\text{-bit}, \text{Win}, \beta}$.

To quantify the estimation performance of linear predictor M_* and observation probabilities Θ_* , we considered the empirical root mean squared errors $\text{RMSE}(\mathbf{B}, \mathbf{C})$ with respect to any two matrices \mathbf{B} and \mathbf{C} of dimension $n_1 \times n_2$, and the Hellinger distance $d_H^2(\widehat{\Theta}, \Theta_*)$ between $\widehat{\Theta}$ and Θ_* defined as follows:

$$\text{RMSE}(\mathbf{B}, \mathbf{C}) = \frac{\|\mathbf{B} - \mathbf{C}\|_F}{(n_1 n_2)^{1/2}} \quad \text{and} \quad d_H^2(\widehat{\Theta}, \Theta_*) = \frac{\sum_{i,j=1}^{n_1, n_2} d_H^2(\widehat{\theta}_{ij}, \theta_{*,ij})}{(n_1 n_2)^{1/2}}.$$

As the estimators $\mathcal{F}^{-1}(\widehat{\Theta}_\alpha)$ and $\mathcal{F}^{-1}(\widehat{\Theta}_{1\text{-bit}, \alpha})$ are both low-rank, we also report their corresponding ranks.

Table 1 summarizes the simulation results for the missingness. The most visible aspect of the results is that the proposed estimators $\widehat{\Theta}_\alpha$ and $\widehat{\Theta}_{1\text{-bit}, \alpha}$ both have superior performance than the two existing estimators $\widehat{\Theta}_{\text{NW}}$ and $\widehat{\Theta}_{\text{UNI}}$ by having smaller root mean square errors with respect to \widehat{M} , Hellinger distances $d_H^2(\widehat{\Theta}, \Theta_*)$ and more accuracy estimated rank of M_* . Without the separation of μ_* from M_* , $\widehat{\Theta}_{1\text{-bit}, \alpha}$ has larger error and Hellinger distance than the proposed estimators. The performance of $\widehat{\Theta}_{\text{NW}}$ is roughly between the proposed estimators and the uniform estimator $\widehat{\Theta}_{\text{UNI}}$. Estimator $\widehat{\Theta}_{\text{UNI}}$ is a benchmark which captures no variation of the observation probabilities.

6.2 Target matrix

To generate a target matrix A_* , we first generated $U_{A_*} \in \mathbb{R}^{n_1 \times (r_{A_*} - 1)}$ and $V_{A_*} \in \mathbb{R}^{(r_{A_*} - 1) \times n_2}$ as random matrices with independent Gaussian entries dis-

Table 1: Root mean squared errors $\text{RMSE}(\widehat{M}, M_*)$, Hellinger distance $d_H^2(\widehat{\Theta}, \Theta_*)$, rank of linear predictor \widehat{M} and estimated $\widehat{\Theta}$ and their standard errors (in parentheses) under the low rank missing observation mechanism, with $(n_1, n_2) = (600, 600), (800, 800), (1000, 1000), (1200, 1200)$ and $r_{M_*} = 11$, for the proposed estimators $\widehat{\Theta}_\alpha, \widehat{\Theta}_{1\text{-bit},\alpha}$ and the two existing estimators ($\widehat{\Theta}_{\text{NW}}$ and $\widehat{\Theta}_{\text{UNI}}$).

	600	$\widehat{\Theta}_\alpha$	$\widehat{\Theta}_{1\text{-bit},\alpha}$	$\widehat{\Theta}_{\text{NW}}$	$\widehat{\Theta}_{\text{UNI}}$	
$\text{RMSE}(\widehat{M}, M_*)$	2.6923	(0.0342)	2.9155	(0.0295)	-	-
$d_H^2(\widehat{\Theta}, \Theta_*)$	0.0369	(0.0015)	0.0450	(0.0016)	0.1233	(1e-04)
$r_{\widehat{M}}$	12.45	(0.50)	12.69	(0.46)	-	-
$r_{\widehat{\Theta}}$	600.00	(0.00)	600.00	(0.00)	-	-
	800	$\widehat{\Theta}_\alpha$	$\widehat{\Theta}_{1\text{-bit},\alpha}$	$\widehat{\Theta}_{\text{NW}}$	$\widehat{\Theta}_{\text{UNI}}$	
$\text{RMSE}(\widehat{M}, M_*)$	2.5739	(0.0116)	2.7796	(0.0033)	-	-
$d_H^2(\widehat{\Theta}, \Theta_*)$	0.0317	(5e-04)	0.0379	(1e-04)	0.1219	(1e-04)
$r_{\widehat{M}}$	12.04	(0.20)	12.03	(0.17)	-	-
$r_{\widehat{\Theta}}$	800.00	(0.00)	800.00	(0.00)	-	-
	1000	$\widehat{\Theta}_\alpha$	$\widehat{\Theta}_{1\text{-bit},\alpha}$	$\widehat{\Theta}_{\text{NW}}$	$\widehat{\Theta}_{\text{UNI}}$	
$\text{RMSE}(\widehat{M}, M_*)$	2.4870	(0.0212)	2.7731	(0.0015)	-	-
$d_H^2(\widehat{\Theta}, \Theta_*)$	0.0266	(8e-04)	0.0351	(1e-04)	0.1246	(1e-04)
$r_{\widehat{M}}$	12.68	(0.53)	12.00	(0.00)	-	-
$r_{\widehat{\Theta}}$	1000.00	(0.00)	1000.00	(0.00)	-	-
	1200	$\widehat{\Theta}_\alpha$	$\widehat{\Theta}_{1\text{-bit},\alpha}$	$\widehat{\Theta}_{\text{NW}}$	$\widehat{\Theta}_{\text{UNI}}$	
$\text{RMSE}(\widehat{M}, M_*)$	2.3809	(0.0018)	2.6470	(0.0012)	-	-
$d_H^2(\widehat{\Theta}, \Theta_*)$	0.0242	(1e-04)	0.0314	(1e-04)	0.1211	(1e-04)
$r_{\widehat{M}}$	12.00	(0.00)	12.00	(0.00)	-	-
$r_{\widehat{\Theta}}$	1200.00	(0.00)	1200.00	(0.00)	-	-

tributed as $\mathcal{N}(0, \sigma_{A_*}^2)$ and obtained $A_* = 2.5J + U_{A_*} V_{A_*}^T$. Here we set the standard deviation of the entries in the matrix product $U_{A_*} V_{A_*}^T$ to be 2.5 to mimic the Yahoo! Webscope data set described in Section 7. To achieve this, $\sigma_{A_*} = (2.5^2 / (r_{A_*} - 1))^{1/4}$. The contaminated version of A_* was then generated as $Y = A_* + \epsilon$, where $\epsilon \in \mathbb{R}^{n_1 \times n_2}$ has i.i.d. mean zero Gaussian entries

$\epsilon_{ij} \sim \mathcal{N}(0, \sigma_\epsilon^2)$. The σ_ϵ^2 is chosen such that $\text{SNR} = (\mathbf{E}\|\mathbf{A}_*\|_F^2 / \mathbf{E}\|\epsilon\|_F^2)^{1/2} = 1$, where $\mathbf{E}\|\mathbf{A}_*\|_F^2 = n_1 n_2 (r_{\mathbf{A}_*} - 1 + 2.5^2)$ implies $\sigma_\epsilon = 0.5(r_{\mathbf{A}_*} - 1 + 2.5^2)^{1/2}$.

For the estimation of the target matrix, we evaluated ten versions of the proposed estimators $\text{Prop}_{\hat{\Theta}_{\beta-t}}$, $\text{Prop}_{\hat{\Theta}_{\text{Win},\beta-t}}$, $\text{Prop}_{\hat{\Theta}_\alpha}$, $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\beta-t}}$, $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\text{Win},\beta-t}}$ and $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\alpha}}$. Here Prop indicates the estimators are obtained by solving problem (4.3), while $\hat{\Theta}_\beta$, $\hat{\Theta}_{\text{Win},\beta}$, $\hat{\Theta}_\alpha$, $\hat{\Theta}_{1\text{-bit},\beta}$, $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ and $\hat{\Theta}_{1\text{-bit},\alpha}$ represents the probability estimators used in (4.3), as described in Section 6.1, and $t = 0.05$ or 0.1 denote the winsorized proportion for which β is chosen. In addition, same as Mao et al. (2019), we also compared them with three existing matrix completion techniques: the methods proposed in Negahban and Wainwright (2012) (NW), Koltchinskii et al. (2011) (KLT) and Mazumder et al. (2010) (MHT). Among these three methods, NW is the only one that adjusts for non-uniform missingness. All three methods require tuning parameter selection, for which cross-validation is adopted. See Mao et al. (2019) for more details.

To quantify the performance of the matrix completion, in addition to the empirical root mean squared errors with respect to $\hat{\mathbf{A}}_\beta$ and \mathbf{A}_* , we used one more measure: Test Error = $\|\mathbf{W}^* \circ (\hat{\mathbf{A}}_\beta - \mathbf{A}_*)\|_F^2 / \|\mathbf{W}^* \circ \mathbf{A}_*\|_F^2$, where \mathbf{W}^* is the matrix of missing indicator with the (i, j) th entry being $(1 - w_{ij})$. The test error measures the relative estimation error of the unobserved entries to their signal strength. The estimated ranks of $\hat{\mathbf{A}}_\beta$ are also reported.

Tables 2 summarize the simulation results for different dimensions $n_1=n_2$ ranges from 600 to 800 and two different settings of $r_{\mathbf{A}_*} = 11$. The results of

$r_{A_*} = 11$ for different dimensions $n_1=n_2$ ranges from 1000 to 1200 are delegated to Table S1 and the results of $r_{A_*} = 31$ are delegated to Tables S2-S3 of Section S1.5 in the supplementary material. From the tables, we notice that the ten versions of the proposed methods possess superior performance than the three existing methods by having smaller root mean squared errors and Test Errors. Among the first five proposed methods in the tables, $\text{Prop}_{\hat{\Theta}_\beta}$ is better than $\text{Prop}_{\hat{\Theta}_\alpha}$ for most of the time. It is because that the constrained estimator $\hat{\Theta}_\beta$ has much smaller ratio $\hat{\theta}_U/\hat{\theta}_L$ than $\hat{\Theta}_\alpha$ which improve the stability of prediction and the accuracy. Another observation is that $\text{Prop}_{\hat{\Theta}_\beta.0.1}$ performs better than $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\alpha}}$ at most times.

7. Real data application

In this section we demonstrate the proposed methodology by analyzing the Yahoo! Webscope dataset (ydata-ymusic-user-artist-ratings-v1.0) available at http://research.yahoo.com/Academic_Relations. It contains (incomplete) ratings from 15,400 users on 1000 songs. The dataset consists of two subsets, a training set and a test set. The training set records approximately 300,000 ratings given by the aforementioned 15,400 users. Each song has at least 10 ratings. The test set was constructed by surveying 5,400 out of these 15,400 users, each rates exactly 10 songs that are not rated in the training set. The missing rates are 0.9763 overall, 0.3520 to 0.9900 across users, and 0.6372 to 0.9957 across songs. The non-uniformity of the missingness is shown in Figure S1 of Section S1.6 in

the supplementary material. In this experiment, we applied those methods as described in Section 6 to the training set and evaluated the test errors based on the corresponding test set. As there is no prior knowledge about true parameters α_1 and α_2 , we suggest to choose α_1 and α_2 large enough, say $\alpha_1 = 100$ and $\alpha_2 = 100$, to ensure that the range covers all the missing probabilities. It was noted that $\hat{\Theta}_\alpha$ is not sensitive to larger α .

Table 3 reports the root mean squared prediction errors, where $\text{RMSPE} = \|\mathbf{W}^{test} \circ (\hat{\mathbf{A}}_\beta - \mathbf{Y})\|_F / (\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} w_{ij}^{test})^{1/2}$ and \mathbf{W}^{test} is the indicator matrix of test set with the (i, j) th entry being w_{ij}^{test} . Note that $\text{Prop}_{\hat{\Theta}_\beta-0.05}$ performs the best among all ten versions of proposed methods. Besides, $\text{Prop}_{\hat{\Theta}_\alpha}$ also has much smaller root mean squared prediction error than the other eight versions of proposed methods. This may indicate that only slight constraint is required for the probabilities estimator for this dataset. Note that we cannot guarantee the optimal convergence rate or even asymptotic convergence in certain setting of missingness for $\text{Prop}_{\hat{\Theta}_\alpha}$, see Section 5.2 for details.

With the separation of μ , $\text{Prop}_{\hat{\Theta}_\alpha}$ is better than $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\alpha}}$; analogously, $\text{Prop}_{\hat{\Theta}_\beta-t}$ is better than $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\beta-t}}$ with different constraint level t , same to $\text{Prop}_{\hat{\Theta}_{\text{Win},\beta-s}}$ and $\text{Prop}_{\hat{\Theta}_{1\text{-bit},\text{Win},\beta-s}}$ with different winsorization level s .

As compared with the existing methods NW, KLT and MHT, our proposed methods perform significantly better in terms of root mean squared prediction errors, and achieve as much as 25% improvement when compared with Mazumder, Hastie and Tibshirani's method (the best among the three existing methods).

This suggests that a more flexible modeling of missing structure improves the prediction power.

8. Concluding Remarks

When the matrix entries are heterogeneously observed due to selection bias, this heterogeneity should be taken into account. This paper focuses on the problem of matrix completion under low-rank missing structure. In the recovery of probabilities of observation, we adopt a generalized linear model with a low-rank linear predictor matrix. To avoid unnecessary bias, we introduce a separation of the mean effect μ . As the extreme values of probabilities may lead to unstable estimation of target matrix, we propose an inverse probability weighting based method with constrained probability estimates and demonstrate the improvements in empirical perspectives. Our theoretical result shows that the estimator of the high dimensional probability matrix can be embedded into the inverse probability weighting framework without compromising the rate of convergence of the target matrix (for an appropriately tuned $\beta > 0$), and reveals a possible regime change in the tuning of the constraint parameter ($\beta > 0$ vs. $\beta = 0$). In addition, corresponding computational algorithms are developed, and a related algorithmic convergence result is established. Empirical studies show the attractive performance of the proposed methods as compared with existing matrix completion methods.

Supplementary Materials

The online supplementary materials contain some useful lemmas, the proofs of the main theorems and some additional numerical studies.

Acknowledgment

The authors thank the Editors and two reviewers for constructive comments which led to improvement in the presentation of the paper. Xiaojun Mao's research is partially supported by Shanghai Sailing Program 19YF1402800. Raymond K.W. Wong's research is partially supported by the National Science Foundation under Grants DMS-1806063, DMS-1711952 (subcontract) and CCF-1934904. Song Xi Chen acknowledge supports from National Key Research Special Program of China grant 2016YFC0207701, National Natural Science Foundation of China grant 71532001 and LMEQF at Peking University.

References

- Beck, A. and M. Teboulle (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* 2(1), 183–202.
- Bi, X., A. Qu, J. Wang, and X. Shen (2017). A group-specific recommender system. *Journal of the American Statistical Association* 112(519), 1344–1353.
- Cai, J.-F., E. J. Candès, and Z. Shen (2010). A singular value thresholding algorithm for matrix completion. *SIAM Journal on Optimization* 20(4), 1956–1982.
- Cai, T., T. T. Cai, and A. Zhang (2016). Structured matrix completion with applications to genomic data

REFERENCES

- integration. *Journal of the American Statistical Association* 111(514), 621–633.
- Cai, T. T. and W.-X. Zhou (2016). Matrix completion via max-norm constrained optimization. *Electronic Journal of Statistics* 10(1), 1493–1525.
- Candès, E. J. and Y. Plan (2010). Matrix completion with noise. *Proceedings of the IEEE* 98(6), 925–936.
- Candès, E. J. and B. Recht (2009). Exact matrix completion via convex optimization. *Foundations of Computational Mathematics* 9(6), 717–772.
- Chen, C., B. He, Y. Ye, and X. Yuan (2016). The direct extension of admm for multi-block convex minimization problems is not necessarily convergent. *Mathematical Programming* 155(1-2), 57–79.
- Chen, S. X., D. H. Leung, and J. Qin (2008). Improving semiparametric estimation by using surrogate data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 70(4), 803–823.
- Davenport, M. A., Y. Plan, E. van den Berg, and M. Wootters (2014). 1-bit matrix completion. *Information and Inference* 3(3), 189–223.
- Foygel, R., O. Shamir, N. Srebro, and R. R. Salakhutdinov (2011). Learning with the weighted trace-norm under arbitrary sampling distributions. In *Advances in Neural Information Processing Systems*, pp. 2133–2141.
- Kang, J. D. and J. L. Schafer (2007). Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical science* 22(4), 523–539.
- Keshavan, R. H., A. Montanari, and S. Oh (2009). Matrix completion from noisy entries. In *Advances in Neural Information Processing Systems*, pp. 952–960.
- Klopp, O. (2014). Noisy low-rank matrix completion with general sampling distribution. *Bernoulli* 20(1), 282–303.

REFERENCES

- Koltchinskii, V., K. Lounici, and A. B. Tsybakov (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics* 39(5), 2302–2329.
- Mao, X., S. X. Chen, and R. K. Wong (2019). Matrix completion with covariate information. *Journal of the American Statistical Association* 114(525), 198–210.
- Mazumder, R., T. Hastie, and R. Tibshirani (2010). Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research* 11, 2287–2322.
- Negahban, S. and M. J. Wainwright (2012). Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *Journal of Machine Learning Research* 13(1), 1665–1697.
- Potter, F. J. (1990). A study of procedures to identify and trim extreme sampling weights. In *Proceedings of the American Statistical Association, Section on Survey Research Methods*, Volume 225230.
- Recht, B. (2011). A simpler approach to matrix completion. *Journal of Machine Learning Research* 12, 3413–3430.
- Rohde, A. and A. B. Tsybakov (2011). Estimation of high-dimensional low-rank matrices. *The Annals of Statistics* 39(2), 887–930.
- Rubin, D. B. (2001). Using propensity scores to help design observational studies: Application to the tobacco litigation. *Health Services and Outcomes Research Methodology* 2(3-4), 169–188.
- Schafer, J. L. and J. Kang (2008). Average causal effects from nonrandomized studies: a practical guide and simulated example. *Psychological methods* 13(4), 279.
- Scharfstein, D. O., A. Rotnitzky, and J. M. Robins (1999). Adjusting for nonignorable drop-out using semiparametric nonresponse models. *Journal of the American Statistical Association* 94(448), 1096–1120.

REFERENCES

Schnabel, T., A. Swaminathan, A. Singh, N. Chandak, and T. Joachims (2016). Recommendations as treatments: Debiasing learning and evaluation. *arXiv preprint arXiv:1602.05352*.

Srebro, N. and R. R. Salakhutdinov (2010). Collaborative filtering in a non-uniform world: Learning with the weighted trace norm. In *Advances in Neural Information Processing Systems*, pp. 2056–2064.

School of Data Science, Fudan University, Shanghai 200433, China.

E-mail: maoxj@fudan.edu.cn

Department of Statistics, Texas A&M University, College Station, Texas 77843, U.S.A.

E-mail: raywong@stat.tamu.edu

Department of Business Statistics and Econometrics, Guanghua School of Management and Center for Statistical Science, Peking University, Beijing 100651, China.

E-mail: csx@gsm.pku.edu.cn

REFERENCES

Table 2: Root mean squared errors, test errors, estimated ranks $r_{\hat{A}_\beta}$ and their standard deviations (in parentheses) under the low rank missing observation mechanism, for three existing methods and ten versions of the proposed methods where Prop indicates the estimators are obtained by solving problem (4.3), while $\hat{\Theta}_\beta$, $\hat{\Theta}_{\text{Win},\beta}$, $\hat{\Theta}_\alpha$, $\hat{\Theta}_{1\text{-bit},\beta}$, $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ and $\hat{\Theta}_{1\text{-bit},\alpha}$ represents the probability estimators used in (4.3), as described in Section 6.1, and $t = 0.05$ or 0.1 denote the winsorized proportion for which β is chosen.

$(n_1, n_2) = (600, 600)$	RMSE($\hat{A}_\beta, \mathbf{A}_\star$)	Test Error	$r_{\hat{A}_\beta}$
Prop- $\hat{\Theta}_{\text{Win},\beta}$ -0.05	1.5615 (0.0147)	0.3005 (0.0062)	65.28 (5.72)
Prop- $\hat{\Theta}_\beta$ -0.05	1.5548 (0.0085)	0.2996 (0.0034)	54.98 (3.01)
Prop- $\hat{\Theta}_{\text{Win},\beta}$ -0.1	1.5621 (0.0111)	0.3013 (0.0046)	63.68 (5.36)
Prop- $\hat{\Theta}_\beta$ -0.1	1.5509 (0.0085)	0.2983 (0.0034)	53.13 (2.72)
Prop- $\hat{\Theta}_\alpha$	1.5637 (0.0147)	0.3010 (0.0061)	65.63 (5.89)
Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ -0.05	1.5664 (0.0093)	0.3028 (0.0037)	62.76 (5.96)
Prop- $\hat{\Theta}_{1\text{-bit},\beta}$ -0.05	1.5573 (0.0089)	0.2996 (0.0036)	61.80 (5.34)
Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ -0.1	1.5669 (0.0092)	0.3032 (0.0037)	62.78 (2.68)
Prop- $\hat{\Theta}_{1\text{-bit},\beta}$ -0.1	1.5540 (0.0089)	0.2987 (0.0036)	60.79 (3.01)
Prop- $\hat{\Theta}_{1\text{-bit},\alpha}$	1.5612 (0.0097)	0.3005 (0.0040)	62.12 (4.76)
NW	1.9896 (0.2814)	0.4676 (0.1341)	167.67 (54.78)
KLT	2.2867 (0.0073)	0.5951 (0.0026)	1.00 (0.00)
MHT	1.6543 (0.0097)	0.3432 (0.0041)	51.20 (2.61)
$(n_1, n_2) = (800, 800)$	RMSE($\hat{A}_\beta, \mathbf{A}_\star$)	Test Error	$r_{\hat{A}_\beta}$
Prop- $\hat{\Theta}_{\text{Win},\beta}$ -0.05	1.4754 (0.0107)	0.2669 (0.0041)	88.58 (10.81)
Prop- $\hat{\Theta}_\beta$ -0.05	1.4797 (0.0080)	0.2714 (0.0030)	71.79 (4.12)
Prop- $\hat{\Theta}_{\text{Win},\beta}$ -0.1	1.4724 (0.0108)	0.2664 (0.0042)	86.25 (10.34)
Prop- $\hat{\Theta}_\beta$ -0.1	1.4763 (0.0082)	0.2704 (0.0031)	67.08 (4.22)
Prop- $\hat{\Theta}_\alpha$	1.4783 (0.0115)	0.2676 (0.0041)	88.92 (11.70)
Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ -0.05	1.4917 (0.0078)	0.2743 (0.0030)	83.51 (1.45)
Prop- $\hat{\Theta}_{1\text{-bit},\beta}$ -0.05	1.4804 (0.0080)	0.2705 (0.0031)	82.60 (3.47)
Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta}$ -0.1	1.4972 (0.0080)	0.2765 (0.0031)	81.64 (7.23)
Prop- $\hat{\Theta}_{1\text{-bit},\beta}$ -0.1	1.4800 (0.0078)	0.2708 (0.0030)	74.89 (3.54)
Prop- $\hat{\Theta}_{1\text{-bit},\alpha}$	1.4790 (0.0099)	0.2685 (0.0039)	88.57 (9.56)
NW	1.9515 (0.3625)	0.4585 (0.1593)	215.61 (82.24)
KLT	2.3447 (0.0064)	0.6081 (0.0020)	1.00 (0.00)
MHT	1.6067 (0.0086)	0.3245 (0.0036)	63.68 (3.02)

¹ With $r_{M_\star} = 11$, $r_{\mathbf{A}_\star} = 11$, $(n_1, n_2) = (600, 600)$, $(800, 800)$ and $\text{SNR} = 1$.
 The three existing methods are proposed respectively in Negahban and Wainwright (2012)(NW), Koltchinskii et al. (2011)(KLT) and Mazumder et al. (2010)(MHT)

REFERENCES

Table 3: Root mean squared prediction errors based on Yahoo! Webscope dataset for the ten versions of the proposed method and the three existing methods proposed respectively in Negahban and Wainwright (2012)(NW), Koltchinskii et al. (2011)(KLT) and Mazumder et al. (2010)(MHT).

	Prop- $\hat{\Theta}_{\text{Win},\beta-0.05}$	Prop- $\hat{\Theta}_{\beta-0.05}$	Prop- $\hat{\Theta}_{\text{Win},\beta-0.1}$
RMSPE	1.0396	1.0381	1.0476
	Prop- $\hat{\Theta}_{\beta-0.1}$	Prop- $\hat{\Theta}_{\alpha}$	Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta-0.05}$
RMSPE	1.0490	1.0383	1.0831
	Prop- $\hat{\Theta}_{1\text{-bit},\beta-0.05}$	Prop- $\hat{\Theta}_{1\text{-bit},\text{Win},\beta-0.1}$	Prop- $\hat{\Theta}_{1\text{-bit},\beta-0.1}$
RMSPE	1.1091	1.0760	1.0523
	Prop- $\hat{\Theta}_{1\text{-bit},\alpha}$	NW	KLT
RMSPE	1.1065	1.7068	3.6334
	MHT		
RMSPE	1.3821		