

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

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 in which both the missing structure and the underlying parameter are high dimensional. Despite the missing structure being a key component of any missing-data problem, 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 observation probability matrix is estimated using a high-dimensional, low-rank matrix-estimation procedure, and then used to complete the target matrix via inverse probability weighting. Owing to the high-dimensional and extreme (i.e., very small) nature of the true probability matrix, the effect of inverse probability weighting requires careful study. Lastly, 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, with applications including collaborative filtering, computer vision and

Song Xi Chen is the corresponding author.

positioning. From a statistical viewpoint, it is a high-dimensional, missing-data problem in which a high percentage of matrix entries are missing. As in many missing-data problems, the underlying missing mechanism plays an important role. Most existing works (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 has enabled the domain to move forward rapidly, with various theoretical breakthroughs occurring in the last decade. However, a 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 tried to relaxing this restrictive assumption by adopting other missing structures. The use of these settings hinges on having strong prior knowledge of the underlying problem. At a high level, many of these works use a special form of low-rank structure for the missing mechanism. For instance, Foygel et al. (2011) and Negahban and Wainwright (2012) both adopt a rank-one structure based on the estimated marginal probabilities. In this study, we aim to recover the target matrix \mathbf{A}_* under a flexible, high-dimensional, low-rank sampling structure. This is achieved by using a weighted empirical risk minimization and by applying 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 possess a complex “sampling” structure that is largely unknown. In a movie recommender system, users tend to rate movies that they prefer or dislike the most, while often remaining “silent” about other movies. Another example of a complex sampling regime is that of online merchandising. Here some users may purchase certain items regularly without often rating them, but often evaluate products that they rarely buy. Similarly 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, leading to a low-rank model of the missing structure.

Inspired by generalized linear models, we model the probability of observation $\Theta_* = (\theta_{*,ij})_{i,j=1}^{n_1,n_2} \in (0, 1)^{n_1 \times n_2}$ using a high-dimensional, low-rank matrix $\mathbf{M}_* = (m_{*,ij})_{i,j=1}^{n_1,n_2} \in \mathbb{R}^{n_1 \times n_2}$ with a known function f . Therefore, at the entry level, we have $\theta_{*,ij} = f(m_{*,ij})$. In generalized linear models, the linear predictor $m_{*,ij}$ is further modeled as a linear function of the observed covariates. However, to reflect the difficulty of attaining (appropriate and adequate) covariate information and the complexity of modeling Θ_* in some situations of the matrix completion, we assume the predictor matrix \mathbf{M}_* is completely hidden. Despite \mathbf{M}_* being hidden, the low rank and high dimensionality of \mathbf{M}_* allow both identification and consistent estimation of Θ_* , which facilitates matrix completion based on inverse probability weighting. Motivated by the nature of matrix completion, we propose a novel parametrization $\mathbf{M}_* = \mu_* \mathbf{1}_{n_1} \mathbf{1}_{n_2}^T + \mathbf{Z}_*$, where \mathbf{Z}_*

satisfies $\mathbf{1}_{n_1}^T \mathbf{Z}_* \mathbf{1}_{n_2} = 0$. Our proposal extends the work of Davenport et al. (2014), who solve a binary matrix-completion problem and pursue a different goal. In contrast to that of Davenport et al. (2014), the proposed method does not regularize the estimation of μ_* . Instead, it regularizes the nuclear norm of the estimation of \mathbf{Z}_* , which requires a different algorithmic treatment to avoid bias caused by the nuclear-norm penalty.

Three fundamental aspects that set our work apart from existing works on matrix completion: (i) the high-dimensional probability matrix Θ , the dimensions of which, 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 the estimation error in the inverse probability weighting of the matrix completion procedure. Aspects (i) and (ii) are unique to our problem, and not found in the literature on missing data. Works related to aspect (iii) are sparse in the literature on matrix completion. Noted that Mao et al. (2019) focused on the low-dimensional parametric modeling of inverse probability weighting with observable covariates.

We develop non-asymptotic upper bounds for the mean squared errors of 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 the low levels of the observation probabilities, we propose re-estimating \mathbf{Z}_* by constraining the magnitude of its entries to a smaller threshold.

Our analysis shows that the proposed constrained inverse probability weighting estimator achieves the optimal rate (up to a logarithmic factor in the estimation of the target matrix). We also compare the completion based on inverse probability weighting and the proposed constrained estimation with that 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). As such, we 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 context 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 zero otherwise, and assume $\{\epsilon_{ij}\}$ are independent of $\{w_{ij}\}$.

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 $\Theta_\star = (\theta_{\star,ij})_{i,j=1}^{n_1,n_2} \in (0, 1)^{n_1 \times n_2}$. Similarly 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 prespecified monotone and differentiable function $f : \mathbb{R} \rightarrow [0, 1]$; that is, $\theta_{\star,ij} = f(m_{\star,ij})$, for all i, j . The matrix M_\star plays the same role as a linear predictor in the generalized linear model. The function f is an inverse link function. Two popular choices of f are the 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 assumptions, it is virtually impossible to recover the hidden feature matrix M_\star or the target matrix A_\star . A common and powerful assumption is that A_\star is a low-rank matrix; that is, $\text{rank}(A_\star) \ll \min\{n_1, n_2\}$. For example, consider the Yahoo! Webscope data set (see Section 7). This data set contains a partially observed matrix of ratings from 15,400 users for 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 proven useful in modeling recommender systems (e.g., Candès and Plan, 2010; Cai et al., 2010).

The same idea can be adapted to the missing pattern, despite the factors that induce the missingness possibly differing from those that generate the ratings.

To this end, we assume M_* is of low rank. Next, we decompose M_* as

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

where $\mathbf{1}_n$ is an n -vector of ones, and $J = \mathbf{1}_{n_1} \mathbf{1}_{n_2}^T$. Here, μ_* is the mean of M_* ; that is, $\mu_* = \mathbf{1}_{n_1}^T M_* \mathbf{1}_{n_2} / (n_1 n_2)$. 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 $Z = M - \mu J$. Moreover, the decomposition is unique, owing to the constraint that $\mathbf{1}_{n_1}^T Z_* \mathbf{1}_{n_2} = 0$. The key here is to reparametrize M_* in terms of μ_* and Z_* , which require different treatments in their estimations. See Section 3 for details. Furthermore, the low-rankness of M_* can be translated to the low-rankness of Z_* .

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

There are several related, but more specialized models. Srebro and Salakhutdinov (2010) and Negahban and Wainwright (2012) use an independent row and column sampling mechanism, leading to a rank-one structure for Θ_* . Cai et al. (2016) consider a matrix block structure for Θ_* and hence M_* , that can be regarded as a special case of low-rank modeling. Mao et al. (2019) examine the case when the missingness depends on observable covariates, and adopt a low-rank modeling with a known row space of M_* . In this study, we focus on the

situation in which the missingness is dependent on some hidden factors, reflecting 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 \mathbf{A}_* , many existing matrix completion techniques assume a uniform missing structure. Hence, they use 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 $R(\mathbf{A}) = \mathbb{E}(\|\mathbf{A} - \mathbf{Y}\|_F^2)/(n_1 n_2)$ (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 $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2}$,

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

Clearly, \mathbf{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 (defined later) based on \hat{R} (assuming the use of true observation probabilities) improves upon the existing error upper bound of the corresponding estimator based on \hat{R}_{UNI} achieved by Klopp (2014), as shown in Section 5.3. However, the inverse probability weights $\Theta_{\star}^{\ddagger}$ 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 the missing-data literature, it is known to produce unstable estimations, owing to the occurrence of small probabilities (e.g., Rubin, 2001; Kang and Schafer, 2007; Schafer and Kang, 2008). This problematic scenario is common in matrix-completion problems in which we attempt to recover 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 nonstandard setup of inverse probability weighting. Therefore, 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 a small θ_0 increases the difficulty of the estimation, $\hat{R}(\mathbf{A})$ changes only up to a multiplicative constant. However, in a non-uniform setting, this is not as straightforward, owing to the heterogeneity among $\{\theta_{\star,ij}\}$. To demonstrate the issue, we examine the Yahoo! Webscope data set 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 $\theta_U = \max_{i,j} \theta_{*,ij}$. The corresponding ratio of estimated probabilities $\widehat{\theta}_U/\widehat{\theta}_L$ based on the rank-one structure of Negahban and Wainwright (2012) is 25656.2, and that based on our proposed method (without re-estimation, to be described below) is 23988.0. This strong heterogeneity can jeopardize the convergence rate of our estimator, which was address properly 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 A_* .

3. Estimation of Θ_*

3.1 Regularized maximum likelihood estimation

We develop the estimation of Θ_* based on the framework of a regularized maximum likelihood. Given the inverse of the 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} .

Owing to the low-rank assumption of \mathbf{M}_* , a natural candidate as an estimator for \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 $\Pr(w_{ij} = 1) = 0.5$, for all i, j . Nevertheless, this does not align well with common settings of matrix completion, under which the average probability of observations is quite small, resulting 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 $(\mu_\star, \mathbf{Z}_\star)$:

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

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

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} . 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} . 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, because 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 $\|\mathbf{M}\|_* \leq \alpha' \{\text{rank}(\mathbf{M}_\star) n_1 n_2\}^{1/2}$. 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 Section 5.

3.2 Computational algorithm and tuning parameter selection

To solve the optimization given in (3.1), we begin with the observation that $\ell_{\mathbf{W}}$ is a smooth concave function, which allows us to use an iterative algorithm called the 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, as follows:

$$\begin{aligned} P_L \{(\mu, \mathbf{Z}), (\mu_{\text{old}}, \mathbf{Z}_{\text{old}})\} = & -\ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \\ & + (\mu - \mu_{\text{old}}) \mathbf{1}_{n_1}^T \{-\nabla_\mu \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}})\} \mathbf{1}_{n_2} + \frac{Ln_1 n_2}{2} (\mu - \mu_{\text{old}})^2 \\ & + \langle \mathbf{Z} - \mathbf{Z}_{\text{old}}, -\nabla_{\mathbf{Z}} \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \rangle + \frac{L}{2} \|\mathbf{Z} - \mathbf{Z}_{\text{old}}\|_F^2 + \lambda \|\mathbf{Z}\|_*, \end{aligned}$$

where $L > 0$ is an algorithmic parameter determining the step size of the proximal gradient algorithm, and is chosen using 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 the 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

$$\min \left[\alpha_1, \max \left[-\alpha_1, \mu_{\text{old}} + (Ln_1 n_2)^{-1} \mathbf{1}_{n_1}^T \{-\nabla_\mu \ell_{\mathbf{W}}(\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}})\} \mathbf{1}_{n_2} \right] \right].$$

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

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

$$\arg \min_{\mathbf{Z} = \mathbf{G}_1 = \mathbf{G}_2, \mathbf{1}_{n_1}^\top \mathbf{G}_1 \mathbf{1}_{n_2} = 0, \|\mathbf{G}_2\|_\infty \leq \alpha_2} \frac{\lambda}{L} \|\mathbf{Z}\|_* + \frac{1}{2} \left\| \mathbf{G}_2 - \mathbf{Z}_{\text{old}} - \frac{1}{L} \nabla_{\mathbf{Z}} \ell_{\mathbf{W}} (\mu_{\text{old}} \mathbf{J} + \mathbf{Z}_{\text{old}}) \right\|_F^2.$$
(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 $\mathbb{I}_{\mathcal{A}}$ is equal to zero if the constraint \mathcal{A} holds, and ∞ otherwise. The detailed algorithm to solve (3.3) is provided as Algorithm 1 in the Supplementary Material. 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 such 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, the proof of which 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).*

Note that the alternative direction method of multipliers algorithm is nested within the proximal gradient algorithm. However, our numerical experiments show that the numbers of inner iterations (alternative direction method of multipliers) and outer iterations (proximal gradient) are both small, usually less than 20. We summarize the corresponding convergence result of the overall proximal gradient algorithm in the following theorem, the proof of which 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 based on 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 sufficiently large, 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. For λ , we adopt the Akaike information criterion (AIC), where we approximate the degrees of freedom 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 to replace Θ_*^\dagger with $\widehat{\Theta}^\dagger = (\widehat{\theta}_{ij}^{-1/2})_{i,j=1}^{n_1, n_2}$. However, this direct implementation is not robust to extremely small probabilities of observation, and our theoretical analysis shows that it could lead to a slower convergence rate in the estimation of \mathbf{A}_* . In the literature on missing data, a simple solution is to winsorize the small probabilities (Potter, 1990; Scharfstein et al., 1999).

In the estimation of $\widehat{\Theta}$ defined in (3.1) that assumes $\|\mathbf{Z}_*\|_\infty \leq \alpha_2$, 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 is affected. That is, even though the extreme probabilities can be controlled by winsorizing, the unchanged entries of $\widehat{\mathbf{Z}}$ (in the winsorizing procedure) 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 given in (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 entries of magnitude smaller than β . However, the mean-zero constraint of \mathbf{Z} no longer makes sense, because 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 :

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

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 to that 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_*^\ddagger with $\widehat{\Theta}_\beta^\ddagger$ to obtain a modified empirical risk:

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

where $\widehat{\Theta}_\beta^\ddagger = (\widehat{\theta}_{ij,\beta}^{-1/2}) \in \mathbb{R}^{n_1 \times n_2}$. Because \mathbf{A} is a high-dimensional parameter, a direct minimization of \widehat{R}^* often results in over-fitting. To circumvent this, 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 a low-rank solution. Based on (4.2), our estimator of \mathbf{A}_* is

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

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

We use an accelerated proximal gradient algorithm (Beck and Teboulle, 2009) to solve (4.3). For the choice of the tuning parameter τ in (4.3), we adopt a five-fold cross-validation to select the remaining tuning parameters. Owing to the non-uniform missing mechanism, we use a weighted version of the validation errors; see 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, maximum norm, and $l_{\infty,2}$ -norm, respectively, of a matrix \mathbf{B} . We use the symbol \asymp to represent the asymptotic equivalence in order; that is, $a_n \asymp b_n$ 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 on matrix completion, most discussions related to the optimal convergence rate are only up to certain poly-

nomial 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 convenience in 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 $\alpha_1, \alpha_2 \geq 0$,

$$\left(\widehat{\mu}, \widehat{\mathbf{Z}}\right) = \arg \max_{(\mu, \mathbf{Z}) \in \widetilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2)} \ell_{\mathbf{W}}(\mu \mathbf{J} + \mathbf{Z}), \text{ where} \quad (5.1)$$

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

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 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; the 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 $\widehat{\mu}$, $\widehat{\mathbf{Z}}$, and $\widehat{\mathbf{M}}$. To simplify the notation, 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 $\widehat{\mathbf{M}} = \widehat{\mu}\mathbf{J} + \widehat{\mathbf{Z}}$, where $(\widehat{\mu}, \widehat{\mathbf{Z}})$ is the solution to (5.1). There exist positive constants C_1, C_2 , such that we have with probability at least $1 - C_1/n$,

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

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 derived easily 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, that is, $\alpha_0 \leq 1/2 \log n$ 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:

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

We now consider the constrained estimation for \mathbf{Z}_* , \mathbf{M}_* , and Θ_*^\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 $\widehat{\mathbf{M}}_{*,\beta} = \widehat{\mu} \mathbf{J} + \mathcal{T}_\beta(\mathbf{Z}_*)$, and $\Theta_{*,\beta} = \mathcal{F}(\mathbf{M}_{*,\beta})$ and $\widehat{\Theta}_{*,\beta} = \mathcal{F}(\widehat{\mathbf{M}}_{*,\beta})$. 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 $d^2(\widehat{\Theta}_{*,\beta}^\dagger, \Theta_*^\dagger)$ are investigated in the next theorem. Define $\Lambda_n = \min[\beta^2, \widetilde{\Gamma}_n + h_{\alpha_1, \beta}^{-1} n^{-2} \beta \{8N_\beta + (n^2 - N_\beta) |\mu_* - \widehat{\mu}|]\}$, 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 $(\mu_*, \mathbf{Z}_*) \in \tilde{\mathcal{C}}_{n_1, n_2}(\alpha_1, \alpha_2)$. Consider $\widehat{\mathbf{M}}_\beta = \widehat{\mu} \mathbf{J} + \widehat{\mathbf{Z}}_\beta$, where $\widehat{\mathbf{Z}}_\beta$ is the solution to (5.3) and $\beta \geq 0$. Then there exist some positive constants C_1 , C_2 , and C_3 , such that we have with probability at least $1 - 2C_1/n$,

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

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 winsorized directly from $\widehat{\mathbf{Z}}$. Obviously, the order of this upper bound is larger than or equal to Λ_n . Moreover, there are scenarios in which Λ_n is a smaller order of β^2 . To illustrate, assume that $\alpha_1 \asymp 1$ and $\beta \asymp 1$, and 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 is 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

$$k_{\alpha_1,\alpha_2,\beta,n} \asymp \min \left[\beta^2, h_{\alpha_1,\beta}^{-1} \beta \left\{ 8N_\beta + (n^2 - N_\beta) k_{\alpha_1,\alpha_2,n}^{1/2} \right\} n^{-2} + h_{\alpha_1,\beta}^{-1} n^{-1/2} (\alpha_1 + \beta r_{\mathcal{T}_\beta(\mathbf{Z}_*)}^{1/2}) \right].$$

5.2 Target matrix

To study the convergence of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$, we require the following conditions on the random errors ϵ and the target matrix \mathbf{A}_* . Recall that $\widehat{\mathbf{A}}_\beta$ includes the estimations obtained with the unconstrained estimator $\widehat{\Theta}$ and those with the constrained estimator $\widehat{\Theta}_\beta$, because $\widehat{\mathbf{A}}(\widehat{\Theta}) = \widehat{\mathbf{A}}_{\alpha_2}$, 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 $h_{(1),\beta} = \max_{i,j}(\theta_{*,ij}^{-1}\theta_{*,ij,\beta})$ and

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

The following theorem establishes a general upper bound for $d^2(\hat{\mathbf{A}}_\beta, \mathbf{A}_*)$.

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 , all independent of β , such that for $h_{(1),\beta}\tau \geq C_4\Delta$, we have with probability at least $1 - 3/(2n)$,*

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

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

Now, we discuss the rates of $d^2(\hat{\mathbf{A}}_\beta, \mathbf{A}_*)$ under various missing structures.

For simplicity, the following discussion focuses on the low-rank linear predictor (\mathbf{M}_*) setting, such that $r_{\mathbf{M}_*} \asymp 1$.

Uniform missingness. Under uniform missingness (i.e., $\theta_{ij} \equiv \theta_0$), Koltchinskii et al. (2011) show that $\theta_0^{-1} n^{-1} \text{polylog}(n)$ is the optimal rate for $d^2(\hat{\mathbf{A}}_\beta, \mathbf{A}_*)$. Therefore, it is reasonable to require $\alpha_1 + \alpha_2 = \alpha_0 = O(\text{polylog}(n))$ for the convergence of $d^2(\hat{\mathbf{A}}_\beta, \mathbf{A}_*)$. Under uniform missingness, we have $\alpha_2 = 0$, $\alpha_0 = \alpha_1$, and $e^{\mu_*} \asymp \theta_0$. For $\beta = 0$, our estimator $\hat{\mathbf{A}}_\beta$ degenerates to that based on the unweighted empirical risk function. Theorem 4 shows that this

achieves the optimal rate $\theta_0^{-1}n^{-1}\text{polylog}(n)$. 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 non-uniform missingness, suppose the lower and upper bounds of the observation probability satisfy $\theta_L \asymp e^{\mu_* - \alpha_2}$ and $\theta_U \asymp e^{\mu_* + \alpha_2}$, respectively. 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 because

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

Thus, the convergence rate of $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$ is $e^{\mu_* + 5\alpha_1 + 3\alpha_2} n^{-1/2} \log^3 n$. Because $e^{\mu_*/2 + 5\alpha_1/2 + 3\alpha_2/2} \leq e^{3\alpha_1 + 3\alpha_2/2}$, guaranteeing convergence 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 $\theta_L \equiv (n^{-1}\text{polylog}(n))$. This is the case that results in the number of observed matrix entries being of 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 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, according to

Theorem 4, $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$ attains the convergence rate $\theta_U \theta_L^{-2} n^{-1} \log(n)$, which converges to zero provided that $\theta_U \theta_L^{-2} = o(n \log^{-1} n)$. 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 $n \text{polylog}(n)$ 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 $\beta = O(e^{-\mu_* - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$.

Then, according to Theorem 4,

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

because $e^{-\mu_*/2 + \alpha_2 - \beta + |\alpha_2/2 - \beta|} \leq e^{\alpha_1/2 + 3\alpha_2/2}$. 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 both cases, $e^{-\mu_* + 2\alpha_2 - 2\beta + 2|\alpha_2/2 - \beta|} \asymp \theta_U \theta_L^{-2}$, which leads to

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

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 $\text{polylog}(n)$ factor compared with the most extreme, but feasible setting of $\theta_L^{-1} = O[n\{\text{polylog}(n)\}^{-1}]$. In addition, $\beta = O(e^{-\mu_* - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$ 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 $\theta_L^{-1} = O(n^{1/3})$. Furthermore, $\beta = O(e^{-\mu_* - 2\alpha_1 + \alpha_2} n^{-1/2} \log^{-1} n)$ implies that $\beta = O(n^{-1/6} \log^{-1} n)$. However, to make $d^2(\widehat{\mathbf{A}}_\beta, \mathbf{A}_*)$ 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 Z_* , as reflected by having a larger α_2 than that prescribed under Case (i). Because μ_* is jointly estimated with 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_* were observable, the problem would not be as serious, despite the adverse effect of the stronger heterogeneity on the estimation of Z_* .

To summarize, under uniform missingness and Case (I) and (II)(i) under non-uniform missingness, 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 the 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 to adopting the proposed weighted empirical risk function for matrix completion. In this subsection, we examine this aspect by comparing the non-asymptotic error bounds of the corresponding estimators. Owing to the additional complication from the estimation

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

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

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

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

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

The improvement in the upper bound 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 hold, and take $\tau_{\text{KNOWN}} = C_3 \theta_L^{-1/2} n^{-3/2} \log^{1/2} n$ and $\tau_{\text{UNI}} = C_3 \theta_U^{1/2} f^{-1}(\mu_*) n^{-3/2} \log^{1/2} n$. The upper bound of $d^2(\widehat{\mathbf{A}}^{\text{UNI}}, \mathbf{A}_*)$ is the same as U^{UNI} , and the upper bound of $d^2(\widehat{\mathbf{A}}^{\text{KNOWN}}, \mathbf{A}_*)$ is the same as*

U^{KNOWN} . In addition, $U^{KNOWN} \leq U^{UNI}$ and $U^{KNOWN} < U^{UNI}$ if $\theta_U > \theta_L$; that is, the true missing mechanism is non-uniform.

Our approach draws inspiration from the missing-value literature. For instance, Chen et al. (2008) show that using the estimated parameters in the inverse probability weighting can actually reduce the variance of the parameter of interest; see Theorem 1 of their paper. Given the results of Chen et al. (2008), we expect that using the estimated parameters $\widehat{\Theta}_\beta$ in the weighting probability will not be inferior to the version with the true parameter $\widehat{\Theta}_*$. However, although this is verified using numerical studies, a theoretical proof cannot be achieved.

6. Simulation Study

6.1 Missingness

This section reports the results of our simulation experiments, which were designed to evaluate the numerical performance of the proposed methodologies.

We first evaluate the estimation performance of the observation probabilities in Section 6.1, and then do so for the target matrix in Section 6.2. In the simulation, the true observation probabilities Θ_* and the target matrix A_* were randomly generated once and kept fixed for each simulation setting. To generate Θ_* , we first generated $U_{M_*} \in \mathbb{R}^{n_1 \times (r_{M_*}-1)}$ and $V_{M_*} \in \mathbb{R}^{(r_{M_*}-1) \times n_2}$ as random Gaussian matrices, with independent entries each following $\mathcal{N}(-0.4, 1)$. We then obtained $M_* = U_{M_*} V_{M_*}^T - \bar{m}_{n_1, n_2, r_{M_*}} J$, 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 $\Theta_* = \mathcal{F}(M_*)$, where the inverse link function f is a logistic function.

In our study, we set $r_{M_*} = 11$ (or $r_{Z_*} = 10$) and choose $n_1 = n_2$, with four sizes: 600, 800, 1000, and 1200. The number of simulation runs for each setting is 500. For the purpose of benchmarking, we compare various missingness estimators:

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 $\hat{\Theta}_{Win,\beta} = \mathcal{F}\{\hat{\mu}\mathbf{J} + \mathcal{T}_\beta(\hat{\mathbf{Z}})\}$;
4. the one-bit estimator $\hat{\Theta}_{1-bit,\alpha}$ proposed in Davenport et al. (2014), and its corresponding constrained and winsorized versions $\hat{\Theta}_{1-bit,\beta}$, and $\hat{\Theta}_{1-bit,Win,\beta}$, respectively (note that the one-bit estimator $\hat{\Theta}_{1-bit,\alpha}$ imposes the nuclear-norm regularization on the whole M instead of Z , in contrast to $\hat{\Theta}_\alpha$);
5. the rank-one probability estimator $\hat{\Theta}_{NW}$ 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 $\theta_{ij,NW} = g_{i\cdot}g_{\cdot j}$;
6. the uniform estimator, $\hat{\Theta}_{UNI} = N/(n_1 n_2) \mathbf{J}$.

For the non-constrained estimator $\hat{\Theta}_\alpha$ and the one-bit estimator $\hat{\Theta}_{1-bit,\alpha}$, the parameter α is set based on our knowledge of the true M_* . For the constrained estimators $\hat{\Theta}_\beta$ and $\hat{\Theta}_{Win,\beta}$, the constraint level β is chosen so that either 5%

or 10% of the elements in $\widehat{\mathbf{Z}}_\alpha$ are winsorized. The parameters for $\widehat{\Theta}_{1\text{-bit},\beta}$ and $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}$ are set in a similar manner.

To quantify the estimation performance of the linear predictor M_* and the observation probabilities Θ_* , we consider 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:

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

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

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

Table 1: The root mean squared errors $\text{RMSE}(\widehat{\boldsymbol{M}}, \boldsymbol{M}_*)$, Hellinger distance $d_H^2(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}_*)$, rank of linear predictor $\widehat{\boldsymbol{M}}$, and estimated $\widehat{\boldsymbol{\Theta}}$ 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 $r_{\boldsymbol{M}_*} = 11$, for the proposed estimators $\widehat{\boldsymbol{\Theta}}_\alpha$, $\widehat{\boldsymbol{\Theta}}_{1\text{-bit},\alpha}$ and the two existing estimators $\widehat{\boldsymbol{\Theta}}_{\text{NW}}$ and $\widehat{\boldsymbol{\Theta}}_{\text{UNI}}$.

	600	$\widehat{\boldsymbol{\Theta}}_\alpha$	$\widehat{\boldsymbol{\Theta}}_{1\text{-bit},\alpha}$	$\widehat{\boldsymbol{\Theta}}_{\text{NW}}$	$\widehat{\boldsymbol{\Theta}}_{\text{UNI}}$
$\text{RMSE}(\widehat{\boldsymbol{M}}, \boldsymbol{M}_*)$	2.6923 (0.0342)	2.9155 (0.0295)	-	-	-
$d_H^2(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}_*)$	0.0369 (0.0015)	0.0450 (0.0016)	0.1233 (1e-04)	0.1729 (1e-04)	
$r_{\widehat{\boldsymbol{M}}}$	12.45 (0.50)	12.69 (0.46)	-	-	
$r_{\widehat{\boldsymbol{\Theta}}}$	600.00 (0.00)	600.00 (0.00)	-	-	
	800	$\widehat{\boldsymbol{\Theta}}_\alpha$	$\widehat{\boldsymbol{\Theta}}_{1\text{-bit},\alpha}$	$\widehat{\boldsymbol{\Theta}}_{\text{NW}}$	$\widehat{\boldsymbol{\Theta}}_{\text{UNI}}$
$\text{RMSE}(\widehat{\boldsymbol{M}}, \boldsymbol{M}_*)$	2.5739 (0.0116)	2.7796 (0.0033)	-	-	-
$d_H^2(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}_*)$	0.0317 (5e-04)	0.0379 (1e-04)	0.1219 (1e-04)	0.1767 (1e-04)	
$r_{\widehat{\boldsymbol{M}}}$	12.04 (0.20)	12.03 (0.17)	-	-	
$r_{\widehat{\boldsymbol{\Theta}}}$	800.00 (0.00)	800.00 (0.00)	-	-	
	1000	$\widehat{\boldsymbol{\Theta}}_\alpha$	$\widehat{\boldsymbol{\Theta}}_{1\text{-bit},\alpha}$	$\widehat{\boldsymbol{\Theta}}_{\text{NW}}$	$\widehat{\boldsymbol{\Theta}}_{\text{UNI}}$
$\text{RMSE}(\widehat{\boldsymbol{M}}, \boldsymbol{M}_*)$	2.4870 (0.0212)	2.7731 (0.0015)	-	-	-
$d_H^2(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}_*)$	0.0266 (8e-04)	0.0351 (1e-04)	0.1246 (1e-04)	0.1767 (1e-04)	
$r_{\widehat{\boldsymbol{M}}}$	12.68 (0.53)	12.00 (0.00)	-	-	
$r_{\widehat{\boldsymbol{\Theta}}}$	1000.00 (0.00)	1000.00 (0.00)	-	-	
	1200	$\widehat{\boldsymbol{\Theta}}_\alpha$	$\widehat{\boldsymbol{\Theta}}_{1\text{-bit},\alpha}$	$\widehat{\boldsymbol{\Theta}}_{\text{NW}}$	$\widehat{\boldsymbol{\Theta}}_{\text{UNI}}$
$\text{RMSE}(\widehat{\boldsymbol{M}}, \boldsymbol{M}_*)$	2.3809 (0.0018)	2.6470 (0.0012)	-	-	-
$d_H^2(\widehat{\boldsymbol{\Theta}}, \boldsymbol{\Theta}_*)$	0.0242 (1e-04)	0.0314 (1e-04)	0.1211 (1e-04)	0.1761 (1e-04)	
$r_{\widehat{\boldsymbol{M}}}$	12.00 (0.00)	12.00 (0.00)	-	-	
$r_{\widehat{\boldsymbol{\Theta}}}$	1200.00 (0.00)	1200.00 (0.00)	-	-	

6.2 Target matrix

To generate a target matrix \boldsymbol{A}_* , we first generated $\boldsymbol{U}_{\boldsymbol{A}_*} \in \mathbb{R}^{n_1 \times (r_{\boldsymbol{A}_*}-1)}$ and $\boldsymbol{V}_{\boldsymbol{A}_*} \in \mathbb{R}^{(r_{\boldsymbol{A}_*}-1) \times n_2}$ as random matrices with independent Gaussian entries distributed as $\mathcal{N}(0, \sigma_{\boldsymbol{A}_*}^2)$, obtaining $\boldsymbol{A}_* = 2.5\boldsymbol{J} + \boldsymbol{U}_{\boldsymbol{A}_*}\boldsymbol{V}_{\boldsymbol{A}_*}^\top$. Here, we set the standard deviation of the entries in the matrix product $\boldsymbol{U}_{\boldsymbol{A}_*}\boldsymbol{V}_{\boldsymbol{A}_*}^\top$ to 2.5 to mimic the Yahoo! Web-

scope 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 independent and identically distributed mean-zero Gaussian entries $\epsilon_{ij} \sim \mathcal{N}(0, \sigma_\epsilon^2)$. Here, σ_ϵ^2 is chosen such that $\text{SNR} = (\mathbf{E}\|A_*\|_F^2/\mathbf{E}\|\epsilon\|_F^2)^{1/2} = 1$, where $\mathbf{E}\|A_*\|_F^2 = n_1 n_2 (r_{A_*} - 1 + 2.5^2)$ implies $\sigma_\epsilon = 0.5(r_{A_*} - 1 + 2.5^2)^{1/2}$.

For the estimation of the target matrix, we evaluated 10 versions of the proposed estimators $\text{Prop-}\widehat{\Theta}_{\beta-t}$, $\text{Prop-}\widehat{\Theta}_{\text{Win},\beta-t}$, $\text{Prop-}\widehat{\Theta}_\alpha$, $\text{Prop-}\widehat{\Theta}_{1\text{-bit},\beta-t}$, $\text{Prop-}\widehat{\Theta}_{1\text{-bit},\text{Win},\beta-t}$, and $\text{Prop-}\widehat{\Theta}_{1\text{-bit},\alpha}$. Here, Prop indicates the estimators are obtained by solving problem (4.3), $\widehat{\Theta}_\beta$, $\widehat{\Theta}_{\text{Win},\beta}$, $\widehat{\Theta}_\alpha$, $\widehat{\Theta}_{1\text{-bit},\beta}$, $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}$, and $\widehat{\Theta}_{1\text{-bit},\alpha}$ represent the probability estimators used in (4.3), as described in Section 6.1, and $t = 0.05$ or 0.1 denotes the winsorized proportion for which β is chosen. In addition, as in Mao et al. (2019), we compare these 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**). Of 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 \widehat{A}_β and A_* , we use one more measure: Test Error = $\|\mathbf{W}^* \circ (\widehat{A}_\beta - A_*)\|_F^2 / \|\mathbf{W}^* \circ A_*\|_F^2$, where \mathbf{W}^* is the matrix of the missing indicator, with the (i, j) th entry being $(1 - w_{ij})$. The test error measures the estimation error of the unobserved entries relative to their

signal strength. The estimated ranks of $\widehat{\mathbf{A}}_\beta$ are also reported.

Table 2 summarizes the simulation results for dimensions $n_1=n_2$ ranging from 600 to 800 and two different settings of $r_{\mathbf{A}_*} = 11$. The results for $r_{\mathbf{A}_*} = 11$ for dimensions $n_1=n_2$ ranging from 1000 to 1200 are relegated to Table S1, and the results for $r_{\mathbf{A}_*} = 31$ are relegated to Tables S2–S3 of Section S1.5 in the Supplementary Material. The tables show that the 10 versions of the proposed methods outperform the three existing methods by having smaller root mean squared errors and Test Error values. Among the first five proposed methods, $\text{Prop-}\widehat{\Theta}_\beta$ is better than $\text{Prop-}\widehat{\Theta}_\alpha$ in most cases. This is because the constrained estimator $\widehat{\Theta}_\beta$ has a much smaller ratio $\widehat{\theta}_U/\widehat{\theta}_L$ than $\widehat{\Theta}_\alpha$, which improves the stability of the prediction and the accuracy. Furthermore, $\text{Prop-}\widehat{\Theta}_\beta\text{-}0.1$ performs better than $\text{Prop-}\widehat{\Theta}_{1\text{-bit},\alpha}$ in most cases.

7. Real-data application

In this section, we demonstrate the proposed methodology by analyzing the Yahoo! Webscope data set (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 data set 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 of whom rated exactly 10 songs that were not rated in the training set. The missing

rates are 0.9763 overall, ranging from 0.3520 to 0.9900 across users, and from 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 the methods described in Section 6 to the training set, and evaluated the test errors based on the corresponding test set. Because there is no prior knowledge about the true parameters α_1 and α_2 , we suggest choosing α_1 and α_2 sufficiently large, say $\alpha_1 = 100$ and $\alpha_2 = 100$, to ensure that the range covers all missing probabilities. Noted that $\widehat{\Theta}_\alpha$ is not sensitive to a larger α .

Table 3 reports the root mean squared prediction errors, where RMSPE = $\|\mathbf{W}^{test} \circ (\widehat{\mathbf{A}}_\beta - \mathbf{Y})\|_F / (\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} w_{ij}^{test})^{1/2}$, and \mathbf{W}^{test} is the indicator matrix of the test set with an (i, j) th entry of w_{ij}^{test} . Note that Prop $_{\widehat{\Theta}_\beta}$ -0.05 performs the best of the 10 versions of the proposed methods. In addition, Prop $_{\widehat{\Theta}_\alpha}$ has a much smaller root mean squared prediction error than those of the other eight versions of the proposed methods. This may indicate that only a slight constraint is required for the probability estimator for this data set. Note that we cannot guarantee the optimal convergence rate or even asymptotic convergence in certain settings of missingness for Prop $_{\widehat{\Theta}_\alpha}$; see Section 5.2 for details.

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

Compared with the existing methods NW, KLT, and MHT, our proposed

methods perform significantly better in terms of the root mean squared prediction errors, achieving as much as a 25% improvement over the method of Mazumder, Hastie, and Tibshirani (the best of the three existing methods). This suggests that a more flexible model of a missing structure improves the model's prediction power.

8. Conclusion

When matrix entries are heterogeneously observed owing to a selection bias, this heterogeneity needs to be taken into account. This study focuses on the problem of matrix completion under a low-rank missing structure. To recover the observation probabilities, we adopted a generalized linear model with a low-rank linear predictor matrix. To avoid unnecessary bias, we introduced a separation of the mean effect μ . Because extreme values of the probabilities may lead to an unstable estimation of the target matrix, we propose an inverse probability weighting-based method with constrained probability estimates, and demonstrate its improvements empirically. 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 compare

MATRIX COMPLETION UNDER LOW-RANK MISSING

the proposed methods with existing matrix-completion methods, demonstrating their appealing performance.

Supplementary Material

The online Supplementary Material contains useful lemmas, the proofs of the main theorems and some additional numerical studies.

Acknowledgments

The authors thank the editors and two reviewers for their constructive comments and suggestions. Xiaojun Mao's research was partially supported by Shanghai Sailing Program 19YF1402800, and the Science and Technology Commission of Shanghai Municipality grant 20dz1200600. Raymond K.W. Wong's research was partially supported by the National Science Foundation under Grants DMS-1806063, DMS-1711952 (subcontract), and CCF-1934904. Song Xi Chen acknowledges support from the National Key Research Special Program of China grant 2016YFC0207701, National Natural Science Foundation of China grant 71532001, 71973005, 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.

REFERENCES

- 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 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*

REFERENCES

- Neural Information Processing Systems*, pp. 952–960.
- Klopp, O. (2014). Noisy low-rank matrix completion with general sampling distribution. *Bernoulli* **20**(1), 282–303.
- 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**(80), 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**(53), 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**. American Statistical Association Washington, DC.
- 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

REFERENCES

- 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.
- Schnabel, T., A. Swaminathan, A. Singh, N. Chandak, and T. Joachims (2016). Recommendations as treatments: Debiasing learning and evaluation. Volume **48** of *Proceedings of Machine Learning Research*, New York, New York, USA, pp. 1670–1679. PMLR.
- 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 $r_{\hat{\mathbf{A}}_\beta}$, and their standard deviations (in parentheses) under the low-rank, missing-observation mechanism for three existing methods and 10 versions of the proposed methods, where Prop indicates the estimators are obtained by solving problem (4.3), while $\widehat{\Theta}_\beta$, $\widehat{\Theta}_{\text{Win},\beta}$, $\widehat{\Theta}_\alpha$, $\widehat{\Theta}_{1\text{-bit},\beta}$, $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}$, and $\widehat{\Theta}_{1\text{-bit},\alpha}$ represent the probability estimators used in (4.3), as described in Section 6.1, and $t = 0.05$ or 0.1 denotes the winsorized proportion for which β is chosen.

$(n_1, n_2) = (600, 600)$	RMSE($\hat{\mathbf{A}}_\beta, \mathbf{A}_*$)	Test Error	$r_{\hat{\mathbf{A}}_\beta}$
Prop- $\widehat{\Theta}_{\text{Win},\beta}-0.05$	1.5615 (0.0147)	0.3005 (0.0062)	65.28 (5.72)
Prop- $\widehat{\Theta}_\beta-0.05$	1.5548 (0.0085)	0.2996 (0.0034)	54.98 (3.01)
Prop- $\widehat{\Theta}_{\text{Win},\beta}-0.1$	1.5621 (0.0111)	0.3013 (0.0046)	63.68 (5.36)
Prop- $\widehat{\Theta}_\beta-0.1$	1.5509 (0.0085)	0.2983 (0.0034)	53.13 (2.72)
Prop- $\widehat{\Theta}_\alpha$	1.5637 (0.0147)	0.3010 (0.0061)	65.63 (5.89)
Prop- $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}-0.05$	1.5664 (0.0093)	0.3028 (0.0037)	62.76 (5.96)
Prop- $\widehat{\Theta}_{1\text{-bit},\beta}-0.05$	1.5573 (0.0089)	0.2996 (0.0036)	61.80 (5.34)
Prop- $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}-0.1$	1.5669 (0.0092)	0.3032 (0.0037)	62.78 (2.68)
Prop- $\widehat{\Theta}_{1\text{-bit},\beta}-0.1$	1.5540 (0.0089)	0.2987 (0.0036)	60.79 (3.01)
Prop- $\widehat{\Theta}_{1\text{-bit},\alpha}$	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($\hat{\mathbf{A}}_\beta, \mathbf{A}_*$)	Test Error	$r_{\hat{\mathbf{A}}_\beta}$
Prop- $\widehat{\Theta}_{\text{Win},\beta}-0.05$	1.4754 (0.0107)	0.2669 (0.0041)	88.58 (10.81)
Prop- $\widehat{\Theta}_\beta-0.05$	1.4797 (0.0080)	0.2714 (0.0030)	71.79 (4.12)
Prop- $\widehat{\Theta}_{\text{Win},\beta}-0.1$	1.4724 (0.0108)	0.2664 (0.0042)	86.25 (10.34)
Prop- $\widehat{\Theta}_\beta-0.1$	1.4763 (0.0082)	0.2704 (0.0031)	67.08 (4.22)
Prop- $\widehat{\Theta}_\alpha$	1.4783 (0.0115)	0.2676 (0.0041)	88.92 (11.70)
Prop- $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}-0.05$	1.4917 (0.0078)	0.2743 (0.0030)	83.51 (1.45)
Prop- $\widehat{\Theta}_{1\text{-bit},\beta}-0.05$	1.4804 (0.0080)	0.2705 (0.0031)	82.60 (3.47)
Prop- $\widehat{\Theta}_{1\text{-bit},\text{Win},\beta}-0.1$	1.4972 (0.0080)	0.2765 (0.0031)	81.64 (7.23)
Prop- $\widehat{\Theta}_{1\text{-bit},\beta}-0.1$	1.4800 (0.0078)	0.2708 (0.0030)	74.89 (3.54)
Prop- $\widehat{\Theta}_{1\text{-bit},\alpha}$	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 $r_{M_*} = 11$, $r_{\mathbf{A}_*} = 11$, $(n_1, n_2) = (600, 600)$, $(800, 800)$, 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 data set for the 10 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- $\widehat{\Theta}_{Win,\beta-0.05}$	Prop- $\widehat{\Theta}_{\beta-0.05}$	Prop- $\widehat{\Theta}_{Win,\beta-0.1}$
RMSPE	1.0396	1.0381	1.0476
	Prop- $\widehat{\Theta}_{\beta-0.1}$	Prop- $\widehat{\Theta}_{\alpha}$	Prop- $\widehat{\Theta}_{1-bit,Win,\beta-0.05}$
RMSPE	1.0490	1.0383	1.0831
	Prop- $\widehat{\Theta}_{1-bit,\beta-0.05}$	Prop- $\widehat{\Theta}_{1-bit,Win,\beta-0.1}$	Prop- $\widehat{\Theta}_{1-bit,\beta-0.1}$
RMSPE	1.1091	1.0760	1.0523
	Prop- $\widehat{\Theta}_{1-bit,\alpha}$	NW	KLT
RMSPE	1.1065	1.7068	3.6334
	MHT		
RMSPE	1.3821		