

Tripartite Collaborative Filtering with Observability and Selection for Debiasing Rating Estimation on Missing-Not-at-Random Data

Qi Zhang,^{1,2} Longbing Cao,² Chongyang Shi,^{1,*} Liang Hu^{3,2}

¹ Beijing Institute of Technology, ² University of Technology Sydney, ³ DeepBlue Academy of Science Shanghai
zhangqi_cs@bit.edu.cn, lonbing.cao@uts.edu.au, cy_shi@bit.edu.cn, rainmilk@gmail.com

Abstract

Most collaborative filtering (CF) models estimate missing ratings with an implicit assumption that the ratings are *missing-at-random*, which may cause the biased rating estimation and degraded performance since recent deep exploration shows that ratings may likely be *missing-not-at-random* (MNAR). To debias MNAR rating estimation, we introduce item observability and user selection to depict the generation of MNAR ratings and propose a tripartite CF (TCF) framework to jointly model the triple aspects of rating generation: item observability, user selection, and ratings, and to estimate the MNAR ratings. An item observability variable is introduced to a *complete observability model* to infer whether an item is observable to a user. TCF also conducts a *complete rating model* for rating generation and utilizes a *user selection model* dependent on the item observability and rating values to model user selection of the observable items. We further elaborately instantiate TCF as a Tripartite Probabilistic Matrix Factorization model (TPMF) by leveraging the probabilistic matrix factorization. Besides, TPMF introduces multifaceted dependency between user selection and ratings to model the influence of user selection on ratings. Extensive experiments on synthetic and real-world datasets show that modeling item observability and user selection effectively debias MNAR rating estimation, and TPMF outperforms the state-of-the-art methods in estimating the MNAR ratings.

Introduction

The research on recommender systems continues with major challenges on more precisely estimating ratings where a large proportion of ratings were missing (Liu et al. 2017; Wang et al. 2019; Zhang et al. 2019). A general and evolving approach is to build collaborative filtering (CF) models (Lin et al. 2014; Zhang et al. 2016). Those models typically assume that rating data is *missing-at-random* (MAR), i.e., the process that generates the available ratings is independent of the values of missing ratings (Hernández-Lobato et al. 2014). In reality, such an assumption may not hold. Taking movie recommendation as an example, users tend to rate preferred movies but rarely rate movies they dislike, rendering usually lower-valued ratings missed and obtaining biased results when estimating missing ratings for

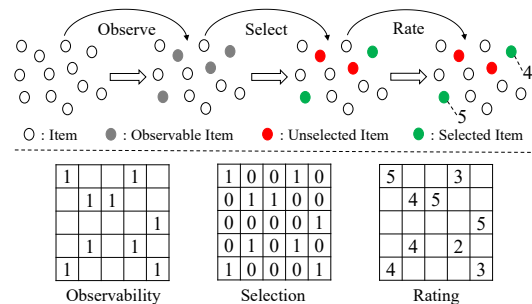


Figure 1: The influence of item observability and user selection on the rating generation.

specific users/items by averaging over the ratings available. This analysis of the potential rating generation process and the dependencies between the missing ratings and the generation process indicates that ratings may be *missing-not-at-random* (MNAR) instead of following the MAR assumption. The MAR-based rating estimation may cause biased rating estimation and degraded performance.

Some recent studies further explore the MNAR rating issue (Yang et al. 2015; Schnabel et al. 2016; Wang et al. 2019) to debias the rating estimation. For example, a classic debiasing approach (Little and Rubin 1986) to the MNAR data is the probabilistic theory of missing data. Such methods (Marlin and Zemel 2009; Ling et al. 2012; Hernández-Lobato et al. 2014; Chen et al. 2018) treat the problem as missing data imputation based on the joint likelihood of the missing rating model and the complete rating model, where the missed ratings (i.e., non-selections) are dependent on the rating values. The intuition behind these methods is that all ratings are firstly generated by the complete rating model and the missing rating model then estimates which entries to be selected (or missed) according to their rating values.

Beyond the dependency on rating values, we argue that the generation process of ratings may be actually more complicated with the MNAR ratings. Revisiting movie recommendation, movie recommenders often suggest those movies that they believe interesting to users, e.g., popular movies, but rarely suggest movies potentially less interesting. Meanwhile, users cannot select and rate those movies unobservable to them. The observability of movies to users

*Corresponding author

influences the user selection of movies. The MNAR perspective indicates the *item observability* to users and the *user selection* of items may jointly influence the rating generation (Melucci 2016; Beauxis-Aussalet and Hardman 2017; Yang et al. 2020), as shown in Figure 1. The figure reflects that the missing ratings contain both preferred yet unobserved entries caused by poor item observability and non-preferred (also called negative) entries. The above debiasing methods neglect the impact of item observability but unreasonably treat all missing entries as non-preferred, which may not conform to the generation process of ratings and lead to a biased modeling of actual user selection and missing ratings.

We argue to simultaneously model item observability and user selection to debias rating estimation, which however is challenging especially when only on the rating data as these aspects are often coupled and co-influence each other (Ohsawa, Obara, and Osogami 2016; Cao 2016). Such modeling needs to properly infer the relationships between the triple aspects while excessively complex modeling may render overfitting. To tackle these challenges, we take two new perspectives: (1) ratings are influenced by factors describing user selection; and (2) item observability depicts the scope of user selection and could correct the probability of the missing entries being negative. New CF models built on the two aspects have potential to address the MNAR nature of rating data and avoid modeling to be biased and skewed to the available ratings.

The above aspects motivate us to develop a tripartite collaborative filtering (TCF) framework by incorporating both item observability and user selection into rating estimation to cater for the MNAR rating data and to tackle the rating estimation bias. We further instantiate the framework by Probabilistic Matrix Factorization (PMF) (Salakhutdinov and Mnih 2007) and propose a Tripartite Probabilistic Matrix Factorization model (TPMF) to infer the three corresponding variables in three sub-models: (1) a *complete rating model* to factorize the ratings with multifaceted factors and model the dependency of ratings on user selection by factorizing the two aspects into shared subspaces simultaneously; (2) a *complete observability model* to introduce a Bernoulli distribution to model the item observability, which determines whether an item is observable to a user and assigns each missing entry a confidence of being truly negative; and (3) a *user selection model* to treat user selection by following a Gaussian distribution whose mean is a function of the corresponding rating value and determining which observable items will be selected by the user.

To the best of our knowledge, this work represents the first attempt to address the MNAR ratings by exploring the complex dependencies between item observability, user selection, and ratings. Extensive empirical results show that modeling item observability and user selection is essential and can effectively debias rating estimation in the MNAR data, and our model outperforms the existing state-of-the-art methods for the MNAR data.

Related Work

As this work explores the impact of item observability and user selection on the rating formation and the bias in esti-

ating missing ratings of recommendation (Schnabel et al. 2016), below we review the related work on modeling item observability and user selection.

Recently, some researchers believe missing ratings reflect both non-preferred (negative) missing ratings and unobserved missing ratings (Liang et al. 2016). They introduce a *user exposure* variable indicating whether an item was exposed to a user to joint probabilistic models and infer the exposure from user selection by the iterative estimation of user selection and the exposure (Liang et al. 2016; Wang et al. 2018a,b; Liu et al. 2020). These methods distribute a confidence of being truly negative to each missing entry and then down-weight the unobserved items to avoid simply treating them as negative that are accordingly not recommended. These methods outperform the state-of-the-art CF methods for the MNAR data, but they only model the dependency between user selection and item observability and are tailored for recommendation with implicit feedback.

Existing models dealing with the MNAR data follow the theory of missing data in (Little and Rubin 1986), which introduces a parametric joint probability distribution on the ratings and selection indicator. For example, CPT-v and Logitvd (Marlin et al. 2007; Marlin and Zemel 2009) use a Mixture of Multinomials (MM) to generate user rating values and model user selection based on these values. More recently, RAPMF (Ling et al. 2012), MF-MNAR (Hernández-Lobato et al. 2014) and SPMF-MNAR (Chen et al. 2018) leverage the powerful probabilistic matrix factorization (PMF) to model user ratings and selection, and SPMF-MNAR further applies social influence rather than just the rating to generate user selection. However, these models neglect the influence of item observability on user selection and treat that all missing entries equally as unselected. This treatment may introduce bias as the missing values actually contain both non-preferred and unobserved entries. Furthermore, these models only consider the dependency of user selection on rating values but fail to reveal the intrinsic multifaceted correlation embodied between user ratings and selection.

In addition, some methods address the MNAR problem by computing an estimated error of the prediction error of imputed values on missing entries (Steck 2011, 2013; Wang et al. 2019). These methods often have a large bias due to imputation inaccuracy, which is then propagated into training and degrade the performance. Some other recent methods (Swaminathan and Joachims 2015; Schnabel et al. 2016; Yang et al. 2018; Joachims, Swaminathan, and Schnabel 2017; Saito 2020) leverage causal inference to handle the MNAR problem. These methods leverage the inverse propensity score (IPS) for each observed entry to propose an unbiased estimator for model training and evaluation. They are suitable for recommendation of either explicit or implicit feedback and have been theoretically and are empirically demonstrated effective and robust. However, IPS-based methods, different from our method and the aforementioned missing theory-based methods, often suffer from the high variance of the propensities (Thomas and Brunskill 2016) and extra metadata may be necessary for estimating the propensity. To the best of our knowledge, no deep learn-

extra metadata (e.g., user demographic or item features) is available, it can be used to infer item observability and differentiate item observability for different users.

User Selection Model. We adopt matrix factorization to factorize variable \mathbf{S} and model the variable as a function of \mathbf{R} and \mathbf{O} , see Figure 2. Specifically, we treat $s_{ij}|o_{ij} = 0$ following constant distribution (denoted by ρ_0) since we have $p(s_{ij} = 0|o_{ij} = 0) = 1$, and we further model $s_{ij}|o_{ij} = 1$ with a Gaussian distribution (note that the Bernoulli distribution is also suitable but brings difficulty in inference). Then, we have:

$$p(\mathbf{S}|\mathbf{R}, \mathbf{O}, \Omega_s) = \prod_{i=1}^n \prod_{j=1}^m \mathcal{N}(s_{ij}|\hat{s}_{ij}, \sigma_s^2)^{o_{ij}} \rho_0^{1-o_{ij}}, \quad (3)$$

$$\hat{s}_{ij} = \mathbf{G}_i^T \mathbf{H}_j + \sum_{l=1}^L (\zeta_{il}^{row} + \zeta_{jl}^{col}) \mathbb{I}[r_{ij} = l] + b_s, \quad (4)$$

where $\mathbb{I}(\cdot)$ denotes the indicator function, b_s is a bias term and Ω_s denotes \mathbf{G} , \mathbf{H} , ζ^{row} , ζ^{col} and b_s . Two matrices $\mathbf{G} \in (0, 1)^{d \times n}$ and $\mathbf{H} \in (0, 1)^{d \times m}$ with $d < \min(n, m)$ are used to factorize \mathbf{S} and follow truncated Gaussian distributions. $\zeta^r \in \mathbb{R}^{n \times L}$ and $\zeta^c \in \mathbb{R}^{m \times L}$ follow zero-mean spherical Gaussian. Note that we put the prior distribution formulas of the parameters into Supplementary. ζ_{il}^{row} and ζ_{jl}^{col} reflect the influence of rating value r_{ij} on s_{ij} . Intuitively, a larger value of $(\zeta_{il}^{row} + \zeta_{jl}^{col})$ when $r_{ij} = l$ implies a higher probability that $s_{ij} = 1$.

Complete Rating Model (CRM). We further factorize R by the inner product of two low-rank latent matrices $\mathbf{U} \in \mathbb{R}^{d \times n}$ and $\mathbf{V} \in \mathbb{R}^{d \times m}$, representing latent user preferences and item attraction respectively. Specifically, we assume Gaussian noise on the ratings below:

$$p(\mathbf{R}|\mathbf{U}, \mathbf{V}, \sigma) = \prod_{i=1}^n \prod_{j=1}^m \mathcal{N}(r_{ij}|\mathbf{U}_i^T \mathbf{V}_j, \sigma^2), \quad (5)$$

where \mathbf{U} and \mathbf{V} follow a zero-mean spherical Gaussian distribution. However, in addition to the influence of rating values on user selection, it is worthy noting how a user selection affects the user rating. We expect that user ratings are also influenced by the factors describing user selection. To model the factor correlation, we regularize the factorization of \mathbf{R} :

$$p(\mathbf{R}|\mathbf{U}, \mathbf{V}, \mathbf{G}, \mathbf{H}, \sigma) = \prod_{i=1}^n \prod_{j=1}^m \mathcal{N}(r_{ij}|\hat{r}_{ij}, \sigma^2), \quad (6)$$

where the estimated rating $\hat{r}_{ij} = \mathbf{U}_i^T \mathbf{\Gamma}_{ij} \mathbf{V}_j$ and we have:

$$\mathbf{\Gamma}_{ij} = \frac{d[\text{diag}(\mathbf{G}_i \circ \mathbf{H}_j) + \varepsilon \mathbf{I}]}{\mathbf{G}_i^T \mathbf{H}_j + d\varepsilon}, \quad (7)$$

where d is the latent dimension, operator \circ calculates the element-wise product and $\text{diag}(\cdot)$ denotes a function constructing a diagonal matrix with a vector. $0 \leq \varepsilon \leq 1$ is an adjustment factor that large ε reduces the influence of \mathbf{G}_i and \mathbf{H}_j on \hat{r}_{ij} and avoids the denominator being zero.

Let $\mathbf{\Gamma}_{ij}^{kk}$ be the k -th ($k \in [1, d]$) diagonal element, we have $\frac{d\varepsilon}{d-1+\varepsilon} < \mathbf{\Gamma}_{ij}^{kk} < \frac{d+d\varepsilon}{1+d\varepsilon}$ and $\mathbb{E}_k(\mathbf{\Gamma}_{ij}^{kk}) = \mathbb{E}_k\left(\frac{dg_{ik}h_{jk}+d\varepsilon}{\mathbf{G}_i^T \mathbf{H}_j + d\varepsilon}\right) =$

1. Hence, we can treat $\mathbf{\Gamma}_{ij}$ as a mask over the d multiplicative factors in calculating $\mathbf{U}_i^T \mathbf{V}_j$, and Equation (6) is equivalent to PMF when $\mathbf{\Gamma}_{ij}$ equals an identity matrix. A larger value of $g_{ik}h_{jk}$ contributes more to user selection, and it also upweights $u_{ik}v_{jk}$ in the estimation of user ratings. The above settings constrain that user preference and item feature show consistency to some extent on the estimation of user selection and rating.

Joint Model. Based on the three sub-models, we obtain the following log joint probability according to Equation (1):

$$\begin{aligned} & \log(\mathbf{R}, \mathbf{O}, \mathbf{S}|\Omega_o, \Omega_s, \Omega_r) \\ &= \sum_{i=1}^n \sum_{j=1}^m o_{ij} \log \mathcal{N}(s_{ij}|\hat{s}_{ij}, \sigma_s^2) + \log \mathcal{N}(r_{ij}|\hat{r}_{ij}, \sigma^2) \\ & \quad + \log \mathcal{B}(o_{ij}|\mu_{ij}) + (1 - o_{ij}) \log \mathbb{I}(s_{ij} = 0) + \mathcal{C} \end{aligned} \quad (8)$$

where \mathcal{C} denotes a constant independent of parameters.

Optimization

We use Expectation-Maximization (EM) (Dempster, Laird, and Rubin 1977), for convenience, to find the maximum a posterior estimates of the parameters of TPMF.

E-step. Both the rating matrix \mathbf{R} and the item observability matrix \mathbf{O} are partly available, we thus calculate the expectation of the ratings and item observability for missing entries, i.e., the entries with $s_{ij} = 0$. Note that we put rating expectation in the M-step via marginalizing $\mathbf{R}_{\mathcal{A}}$ for conveniently updating the latent factors.

Since the estimated rating values (i.e., \hat{r}_{ij}) for missing entries are continuous, we adopt a step function to scatter the values to $\{1, 2, \dots, L\}$ for the calculation of Equation (4). For simplicity, we partition \mathbb{R} into L contiguous intervals with boundaries b_0, b_1, \dots, b_L where $b_0 = -\infty, b_1 = 1, \dots, b_{L-1} = L-1, b_L = \infty$. r_{ij} is obtained according to the interval which the estimated rating belongs to: for example $r_{ij} = l$, if $b_{l-1} < \hat{r}_{ij} \leq b_l$. Since r_{ij} follows $\mathcal{N}(\hat{r}_{ij}, \sigma)$, we define:

$$p(r_{ij} = l|\hat{r}_{ij}) = \Phi\left(\frac{b_l - \hat{r}_{ij}}{\sigma}\right) - \Phi\left(\frac{b_{l-1} - \hat{r}_{ij}}{\sigma}\right), \quad (9)$$

where we denote $\phi(i, j, l) = p(r_{ij} = l|\hat{r}_{ij})$, and Φ is the cumulative distribution function for the standard Gaussian distribution:

$$\Phi(z) = \text{Pr}(\mathcal{N}(0, 1) \leq z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt. \quad (10)$$

Then, we obtain the expectation of o_{ij} below:

$$\begin{aligned} & \mathbb{E}(o_{ij}|\hat{r}_{ij}, \hat{s}_{ij}, \mu_{ij}, s_{ij} = 0) \\ &= \frac{\mu_{ij} \sum_{l=1}^L \phi(i, j, l) \mathcal{N}(0|\hat{s}_{ij}, \sigma_s^2)}{\mu_{ij} \sum_{l=1}^L \phi(i, j, l) \mathcal{N}(0|\hat{s}_{ij}, \sigma_s^2) + (1 - \mu_{ij})}. \end{aligned} \quad (11)$$

M-step. With respect to μ_j following the Beta distribution, we update μ_{ij} by finding the mode of the complete conditional $Beta(\alpha + \sum_i o_{ij}, \beta + n - \sum_i o_{ij})$ as below:

$$\mu_{ij} \leftarrow \frac{\alpha + \sum_i o_{ij} - 1}{\alpha + \beta + n - 2}. \quad (12)$$

Algorithm 1 Generalized EM for TPMF

- 1: **Input:** Rating dataset \mathcal{D}
 - 2: Obtain triple aspects of \mathcal{D} : rating matrix \mathbf{R} , item observability matrix \mathbf{O} , and user selection matrix \mathbf{S}
 - 3: Initialize $\Omega = \{\mathbf{U}, \mathbf{V}, \mathbf{G}, \mathbf{H}, \zeta^{row}, \zeta^{col}, b_s\}, \mu$
 - 4: **while** stopping criteria is not satisfied **do**
 - 5: Compute the expected item observability for missing entries, i.e., $\mathbf{O}_{\bar{\mathcal{A}}}$, by Equation (11)
 - 6: Update Ω by batch gradient ascent along the gradient $\nabla_{\Omega} \mathcal{L}(\Omega, \Theta)$
 - 7: Update μ by Equation (12)
 - 8: **end while**
-

To update the latent factors, we calculate the posterior probability given the estimated \mathbf{O} , i.e., $p(\mathbf{R}_{\mathcal{A}}, \mathbf{S}, \Omega | \mathbf{O}, \Theta)$, and separate the data into the available and missing parts to marginalize $\mathbf{R}_{\bar{\mathcal{A}}}$. Then, we calculate the log-likelihood of the probability and obtain the objective function below (see Supplementary for more details):

$$\begin{aligned} \mathcal{L}(\Omega, \Theta) = & \sum_{(i,j) \in \mathcal{A}} -\frac{1}{2\sigma^2} (\hat{r}_{ij} - r_{ij})^2 - \frac{1}{2\sigma_s^2} (\hat{s}_{ij} - 1)^2 \\ & + \sum_{(i,j) \in \bar{\mathcal{A}}} o_{ij} \sum_{l=1}^L \phi(i, j, l) \left(\frac{\hat{s}_{ij}^2}{2\sigma_s^2} + \rho \right) - \frac{\|\mathbf{U}\|_F}{2\sigma_u^2} - \frac{\|\mathbf{V}\|_F}{2\sigma_v^2} \\ & - \frac{\|\mathbf{G}\|_F}{2\sigma_g^2} - \frac{\|\mathbf{H}\|_F}{2\sigma_h^2} - \frac{\|\zeta^{row}\|_F}{2\sigma_r^2} - \frac{\|\zeta^{col}\|_F}{2\sigma_c^2} + \mathcal{C}. \end{aligned} \quad (13)$$

where $\rho = \log \frac{1}{\sqrt{2\pi\sigma_s^2}}$ is independent of Ω and \mathcal{C} is a constant. Our objective is to maximize $\mathcal{L}(\Omega, \Theta)$ to learn an optimal of $\Omega = \{\mathbf{U}, \mathbf{V}, \mathbf{G}, \mathbf{H}, \zeta^{row}, \zeta^{col}, b_s\}$ under the hyperparameter Θ . Since $\mathcal{L}(\Omega, \Theta)$ has no analytical solution, we take batch gradient ascent to update Ω following Yang et al. (2015).

The resulting optimization algorithm shown in Algorithm 1 belongs to the class of *generalized EM algorithms* guaranteed to converge to a (local) optimum of the log-likelihood (Wu 1983; Greff, van Steenkiste, and Schmidhuber 2017). Due to space limitation, we move the gradients of the parameters to Supplementary.

Discussion. Let us calculate the likelihood probability of an entry being missing (unselected), i.e., $s_{ij} = 0$ by marginalizing r_{ij} and o_{ij} :

$$\begin{aligned} & p(s_{ij} = 0 | \mu_{ij}, \Omega, \Theta) \\ &= \int_{o_{ij}} \int_{r_{ij}} p(s_{ij} = 0, o_{ij}, r_{ij} | \mu_{ij}, \Omega, \Theta) dr_{ij} do_{ij} \\ &= (1 - \mu_{ij}) + \mu_{ij} \sum_{l=1}^L \phi(i, j, l) \mathcal{N}(0 | \hat{s}_{ij}, \sigma_s^2) \end{aligned} \quad (14)$$

with regarding to maximizing the likelihood, we find that μ_{ij} downweights the probability of the missing entries being negative (i.e., $\hat{s}_{ij} = 0$), and the smaller μ_{ij} corresponds to the higher probability of \hat{s} not being 0.

Since the missing entries are partly attributed to the other entries being unobservable (i.e., $o_{ij} = 0$), when we set $\mu_{ij} = 1$ for all entries, the TPMF model degrades to the classic MNAR models (e.g., Logitvd and MF-MNAR) which treat all missing entries as negative ones, which is intuitively not the real case.

Experiments

Since it is difficult to obtain unavailable ratings for testing, we first generate synthetic data to mimic different types of MNAR data and conduct experiments to investigate the effectiveness and robustness of TPMF in handling MNAR ratings. We then compare TPMF against several state-of-the-art methods on four real-world datasets.

Datasets

Synthetic Datasets. The synthetic datasets are generated by a matrix factorization model. First, we set $n = m = 1,000$, $d = 10$ and $L = 5$ and generate the matrices $\mathbf{U}, \mathbf{V}, \zeta^r$ and ζ^c from the standard Gaussian and \mathbf{G} and \mathbf{H} from a uniform distribution within $[0, 1]$. Then, we generate the integer ratings by $r_{ij} = \lceil L \times \psi(\mathbf{U}_i \mathbf{G}_{ij} \mathbf{V}_j^T) \rceil$ and draw o_{ij} from $Bernoulli(\mu_j)$ where μ_j is drawn from $Beta(2, \beta_0)$. Accordingly, we assign the selection variable $s_{ij} = 1$ with a probability of $\rho_s \delta \left(\mathbf{G}_i \mathbf{H}_j^T + \sum_{l=1}^L z_l \mathbf{I}[r_{ij} = l] - 2 \right) / Z$ when $o_{ij} = 1$, and $s_{ij} = 0$ otherwise. Here, δ is a logistic function and $(z_1, \dots, z_5) = (-2, -2, -2, 2, 2)$ reflects items with high ratings are more likely to be selected and Z is used to normalize the probability. The ratings with $s_{ij} = 1$ are selected to construct the dataset. Here, β_0 and ρ_s are used to control the global observability (i.e., #observable items per nm , denoted $p_o \in (0, 1]$) and rating density (i.e., #ratings per nm , denoted $d_r \in (0, 1)$). Roughly, we have $\beta_0 = 2/p_o - 2$ and $\rho_s = d_r/p_o$. We denote this synthetic data as **DDC** which indicates the combination of *item-dependent* observability scheme, *rating-dependent* selection scheme, and *factor-correlated* rating scheme.

To investigate how different observability, selection and rating schemes affect the prediction performance of TPMF, we change scheme combination based on DDC and generate another three datasets: 1) **RDC** - using *random* observability scheme, i.e., $o_{ij} \sim Bernoulli(p_o)$; 2) **DDU** - changed to *factor-uncorrelated* rating scheme, i.e., $r_{ij} = \lceil L \times \psi(\mathbf{U}_i \mathbf{V}_j^T) \rceil$; 3) **DRU** - using *random* selection scheme, i.e., $s_{ij} | o_{ij} = 1 \sim Bernoulli(d_r/p_o)$, and the *factor-uncorrelated* rating scheme; and 4) **RRU** - using *random* observability and selection schemes and *factor-uncorrelated* rating scheme. During the generation, we tune β_0 and p_s to keep the global observability p_o and rating density d_r nearly the same. For all synthetic datasets, we randomly sample two test sets: a *standard set* sampled from the available ratings r_{ij} with $s_{ij} = 1$ and a *special set* sampled from the missing rating r_{ij} with $s_{ij} = 0$, and treat the rest of the available ratings as the training set.

Real-world Datasets. The evaluation of debiasing rating estimation should be verified on MAR ratings. Two real-world rating datasets with MAR ratings are considered: 1)

Dataset	Metric	Special Test Set				Standard Test Set			
		PMF	T-FO	T-NF	TPMF	PMF	T-FO	T-NF	TPMF
RRU	MAE	0.2779	0.2677	0.2651	0.2696	0.2792	0.2685	0.2651	0.2707
	RMSE	0.3357	0.32	0.3158	0.3229	0.3359	0.3207	0.3157	0.3239
DRU	MAE	0.2758	0.2667	0.2628	0.2649	0.2856	0.2723	0.2664	0.2703
	RMSE	0.33	0.3169	0.3114	0.3144	0.349	0.3271	0.3184	0.3243
DDU	MAE	0.2765	0.2614	0.2598	0.2573	0.2897	0.2773	0.2769	0.2758
	RMSE	0.3325	0.311	0.3102	0.3088	0.3504	0.3326	0.3324	0.3299
RDC	MAE	0.2873	0.2718	0.273	0.2699	0.289	0.2737	0.2745	0.2723
	RMSE	0.3461	0.3251	0.3254	0.3241	0.3527	0.3307	0.3319	0.3281
DDC	MAE	0.2901	0.2751	0.2742	0.2719	0.2949	0.2828	0.2822	0.2821
	RMSE	0.3498	0.3294	0.3288	0.3274	0.3591	0.3421	0.3414	0.3404

Table 1: Performance of TPMF compared against PMF and its variants on the five synthetic datasets ($p_o = 0.5$ and $d_r = 0.1$).

Yahoo R3 (denoted Yahoo) collects 311, 704 MNAR ratings and 45, 000 MAR ratings from 15,400 users on 1, 000 songs. 2) The Coat (Coat) has 6, 960 MNAR ratings and 4, 640 ratings of 290 users to 300 coats. And we collect another two real-world datasets that only have MNAR ratings: 3) MovieLens-1M (ML1M) contains about 1M MNAR ratings from 6, 040 users and 3, 706 movies. 4) The Movie Tweets (MTweet) collects 106, 337 MNAR ratings by 3, 972 users on 2, 043 movies from Twitter, where we rescale the original ratings from $[0; 10]$ to the interval $[1; 5]$. Refer to the Supplementary for the links of the four datasets. We use MNAR ratings for training and MAR ratings for testing on Yahoo and Coat, while we randomly split the dataset into training/test sets with 80/20 proportions on ML1M and MTweet. Since there are no MAR ratings in ML1M and MTweet, we set aside 5% of the MNAR ratings and use Naive Bayes to learn propensities.

Experimental Settings

Baseline Methods. We compare TPMF with one basic approach and four state-of-the-art debiasing approaches, including: 1) **PMF** (Salakhutdinov and Mnih 2007) which is based on MAR assumption; 2) **MF-MNAR** (Hernández-Lobato et al. 2014) which deals with the MNAR nature of rating data based on jointly learning the missing data model and the complete rating model; 3) **MF-IPS** (Schnabel et al. 2016) which develops an unbiased estimator for the MNAR rating data based on the Inverse-Propensity-Scoring (IPS); 4) **MF-JL**; and 5) **MF-DR-JL** (Wang et al. 2019) which propose a more robust unbiased estimator by integrating IPS and estimated imputed errors for the MNAR rating data. Besides, we introduce two variants of the proposed model: **T-FO** treating all items being fully observed, i.e., $o_{ij} = 1$, and **T-NF** neglecting the factor correlation between ratings and selection, i.e., prediction ratings by $\hat{r}_{ij} = \mathbf{U}_i^T \mathbf{V}_j$.

Parameter Settings. We utilize the mean absolute error (MAE) and root mean squared error (RMSE) to evaluate the experimental results. For a fair comparison, we tune the hyperparameters on validation sets by grid search and obtain the best for testing. Specifically, we

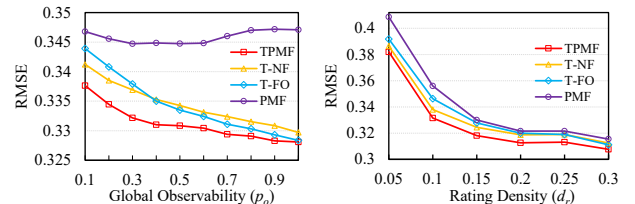


Figure 3: Evaluation on dataset DDC with varying global observability and rating density.

choose the latent dimension d in $\{10, 20, 30, 40\}$, learning rate in $\{0.01, 0.05, 0.1, 1\}$, and L_2 regularization rate in $\{0.01, 0.1, 1\}$ (if required) and keep other hyperparameters recommended from the source codes of the baselines. Regarding TPMF, we fix $\alpha = \beta = 1$ and $\varepsilon = 0.5$ for simplicity. Besides, we tune $\lambda_s = \sigma^2/\sigma_s^2$, $\lambda_v = \sigma^2/\sigma_v^2$, $\lambda_u = \sigma^2/\sigma_u^2$, $\lambda_g = \sigma^2/\sigma_g^2$, $\lambda_h = \sigma^2/\sigma_h^2$, $\lambda_r = \sigma^2/\sigma_r^2$ and $\lambda_c = \sigma^2/\sigma_c^2$ over $\{0.01, 0.1, 1, 10\}$, the learning rate over $\{0.005, 0.01, 0.05, 0.1\}$, and L_2 regularization rate over $\{0.1, 0.5, 1\}$. To guarantee a fast convergence and avoid overfitting, we further initialize \mathbf{U} and \mathbf{V} from a pretrained PMF model and initialize $\boldsymbol{\mu}$ with item frequency.

Experimental Results

Synthetic Experiments. To analyze the effectiveness of TPMF, we evaluate TPMF and its two variants on the five synthetic datasets. Results reporting MAE and RMSE on the special test data and standard test data are provided in Table 1. The results show that the proposed TPMF and its variants outperform the biased method PMF. Considering the characteristics of the datasets, we notice that TPMF performs the best under both metrics except on RRU and DRU. This is reasonable because: 1) TPMF models item observability and factor correlation to handle both simple (i.e., RDC) and complex (i.e., DDC and DDU) item observability schemes and are suitable for the cases with existence (i.e., DDC and RDC) and nonexistence (i.e., DDU) of factor correlation. 2) Relative to the other datasets, both DRU and RRU are sim-

Dataset	Metric	PMF	MF-MNAR	MF-IPS	MF-JL	MF-DR-JL	T-FO	T-NF	TPMF	<i>Imp.(%)</i>
Coat	MAE	0.736	0.704	0.735	0.69	0.701	0.697	0.679	0.67	2.99
	RMSE	0.934	0.899	0.927	0.883	0.897	0.893	0.869	0.857	3.03
Yahoo	MAE	0.973	0.956	0.918	0.903	0.804	0.907	0.821	0.771	4.28
	RMSE	1.223	1.196	1.215	1.182	1.177	1.186	1.172	1.165	3.23
ML1M	MAE	0.701	0.691	0.702	0.671	0.68	0.684	0.671	0.662	1.36
	RMSE	0.886	0.878	0.89	0.857	0.865	0.869	0.856	0.845	1.42
MTWeet	MAE	0.556	0.519	0.53	0.511	0.502	0.521	0.492	0.493	3.65
	RMSE	0.741	0.692	0.701	0.685	0.661	0.695	0.651	0.652	5.22

Table 2: Performance of TPMF compared against its variants and the state-of-the-art baselines on four real-world datasets. The best performance on each dataset is marked in bold. *Imp* reports the performance improvement of TPMF over the best baseline.

ple and randomly select ratings without adding factor correlation. In this case, TPMF may overfit these two datasets and degrade the prediction performance. Comparing the two tables, we see that all models perform better on standard test data than on special test data except on RRU which is an MAR dataset, confirming that the MNAR issues degrade the generalization of the model trained on the biased data to random data. Overall, the results indicate that TPMF can effectively model item observability, user selection and ratings, and infer the relationships between the triple aspects.

Robustness Study. We further investigate the performance of the proposed methods on DDC with varying global observability rates (i.e., $p_o \in \{0.1, 0.2, \dots, 1.0\}$) and rating density levels (i.e., $den \in \{0.05, 0.1, \dots, 0.25, 0.03\}$). Results reporting RMSE on special test sets are provided in Figure 3, where we observe that the proposed methods achieve higher prediction accuracy than PMF. In terms of item observability, higher p_o improves higher prediction accuracy for the proposed methods, which is attributed to the fact that higher item observability simplifies the dataset (note that PMF is not sensitive to the simplicity) and benefits the inference of the two sub-models USM and CRM (see the discussion in Inference). And T-FO performs worse than T-NF when the global observability p_o is small and catches up and exceeds T-NF when $p_o > 0.4$, which confirms that capturing factor correlation plays more important roles with item observability increasing. In addition, all methods obtain obvious improvement with the increase of rating density, which is attributable since more ratings intuitively facilitate the inference of rating generation.

Performance Comparison. To further investigate the effectiveness of TPMF, we report the performance of TPMF and baseline methods on real-world datasets in Table 2. Our TPMF outperforms the state-of-art methods under both metrics on all datasets. Note that MF-IPS performs worse than other debiasing methods and even PMF on the MovieLens dataset while MF-MNAR, MF-JL and MF-DR-JL achieve desirable performance on all the datasets. The results are well explainable. IPS-based methods debias rating estimates by inducing the knowledge of the selection bias and guarantee no bias (if the propensities are correct) but high variance. Meanwhile, the imputation-base methods, i.e., MF-MNAR, rely on modeling the entire generation process of rating to

counterfactually estimate ratings, which gives non-zero bias but very low/zero variance. MF-JL and MF-DR-JL get the best of both the worlds i.e. no bias when either of the models is unbiased and lower variance than IPS. Hence, one might expect that a method like MF-DR-JL using TPMF instead of the MF-MNAR would lead to better results.

In addition, TPMF shows clear advantages over the comparative methods on Coat and Yahoo (two MAR test sets) relative to its performance on ML1M and MTWeet. Deep insight behind the superior results lays that jointly considering item observability, user selection and ratings data facilitates debiasing the rating estimation on MNAR data, and, more importantly, TPMF effectively models the triple aspects. Table 2 also reports T-NF performs better than T-FO, indicating that item observability plays a more important role than factor correlation in debiasing rating estimation. This may be caused by the fact that a large number of items are unobservable to users in practical recommendation data.

Conclusions

We propose a new framework TCF to model the missing-not-at-random rating generation and estimate the MNAR ratings by deeply exploring the relations between rating missingness, item observability, and user selection. The proposed framework includes three sub-models for jointly inferring triple aspects: item observability, user selection and ratings. The newly-added latent variable observability distributes a confidence of being truly negative to each missing entry. We also instantiate the framework to a probabilistic model TPMF, which further introduces the factor dependency between user selection and ratings to model their multifaceted factor correlation. Extensive experiments on the synthetic datasets show that TPMF effectively model the triple aspects simultaneously and infer their relationships. Results on real-world datasets show that both item observability and factor dependency are critical to MNAR rating estimation and TPMF outperforms the state-of-the-art debiasing methods in rating prediction with respect to RMSE and MAE. Further work includes introducing extra metadata into modeling item observability, which may improve the estimate accuracy of item observability and alleviate overfitting issues and even cold-start issues.

Acknowledgments

This work is supported in part by Australian Research Council Discovery Grant (DP190101079), ARC Future Fellowship Grant (FT190100734), the National Key R&D Program of China (2019YFB1406302, 2018YFB1003903), National Natural Science Foundation of China (No. 61502033, 61472034, 61772071, 61272361 and 61672098), and the Fundamental Research Funds for the Central Universities.

References

- Beauxis-Aussalet, E.; and Hardman, L. 2017. Extended Methods to Handle Classification Biases. In *DSAA'2017*, 765–774.
- Cao, L. 2016. Non-IID Recommender Systems: A Review and Framework of Recommendation Paradigm Shifting. *Engineering* 2(2): 212–224.
- Chen, J.; Wang, C.; Ester, M.; Shi, Q.; Feng, Y.; and Chen, C. 2018. Social Recommendation with Missing Not at Random Data. In *IEEE ICDM*, 29–38.
- Dempster, A. P.; Laird, N. M.; and Rubin, D. B. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society* 39(1): 1–38.
- Greff, K.; van Steenkiste, S.; and Schmidhuber, J. 2017. Neural Expectation Maximization. In *NIPS*, 6691–6701.
- Hernández-Lobato, J. M.; Houlisby, N.; Ghahramani, Z.; and . 2014. Probabilistic Matrix Factorization with Non-random Missing Data. In *ICML*, 1512–1520.
- Joachims, T.; Swaminathan, A.; and Schnabel, T. 2017. Unbiased Learning-to-Rank with Biased Feedback. In *WSDM*, 781–789. ISBN 978-1-4503-4675-7.
- Liang, D.; Charlin, L.; McInerney, J.; and Blei, D. M. 2016. Modeling User Exposure in Recommendation. In *WWW*, 951–961.
- Lin, C. H.; Kamar, E.; ; and Horvitz, E. 2014. Signals in the Silence: Models of Implicit Feedback in a Recommendation System for Crowdsourcing. In *AAAI*, 908–915. AAAI Press.
- Ling, G.; Yang, H.; Lyu, M. R.; and King, I. 2012. Response Aware Model-based Collaborative Filtering. In *UAI*, 501–510. ISBN 978-0-9749039-8-9.
- Little, R. J. A.; and Rubin, D. B. 1986. *Statistical Analysis with Missing Data*. New York, NY, USA: John Wiley & Sons, Inc. ISBN 0-471-80254-9.
- Liu, Y.; Xiao, Y.; Wu, Q.; Miao, C.; Zhang, J.; Zhao, B.; and Tang, H. 2020. Diversified Interactive Recommendation with Implicit Feedback. In *AAAI*.
- Liu, Y.; Zhao, P.; Liu, X.; Wu, M.; Duan, L.; and Li, X.-L. 2017. Learning user dependencies for recommendation. In *IJCAI*, 2379–2385.
- Marlin, B. M.; and Zemel, R. S. 2009. Collaborative prediction and ranking with non-random missing data. In *ACM RecSys*, 5–12.
- Marlin, B. M.; Zemel, R. S.; Roweis, S.; and Slaney, M. 2007. Collaborative Filtering and the Missing at Random Assumption. In *UAI*, 267–275. ISBN 0-9749039-3-0.
- Melucci, M. 2016. Impact of Query Sample Selection Bias on Information Retrieval System Ranking. In *DSAA'2016*, 341–350.
- Ohsawa, S.; Obara, Y.; and Osogami, T. 2016. Gated Probabilistic Matrix Factorization: Learning Users' Attention from Missing Values. In *IJCAI*, 1888–1894.
- Saito, Y. 2020. Asymmetric Tri-training for Debiasing Missing-Not-At-Random Explicit Feedback. In *SIGIR*, 309–318. ACM.
- Salakhutdinov, R.; and Mnih, A. 2007. Probabilistic Matrix Factorization. In *NIPS*, 1257–1264.
- Schnabel, T.; Swaminathan, A.; Singh, A.; Chandak, N.; and Joachims, T. 2016. Recommendations as Treatments: Debiasing Learning and Evaluation. In *ICML*, 1670–1679.
- Steck, H. 2011. Item popularity and recommendation accuracy. In *ACM RecSys*, 125–132.
- Steck, H. 2013. Evaluation of recommendations: rating-prediction and ranking. In *ACM RecSys*, 213–220.
- Swaminathan, A.; and Joachims, T. 2015. The Self-Normalized Estimator for Counterfactual Learning. In *NIPS*, 3231–3239.
- Thomas, P. S.; and Brunskill, E. 2016. Data-Efficient Off-Policy Policy Evaluation for Reinforcement Learning. In *ICML*, 2139–2148.
- Wang, M.; Gong, M.; Zheng, X.; and Zhang, K. 2018a. Modeling Dynamic Missingness of Implicit Feedback for Recommendation. In *NIPS*, 6670–6679.
- Wang, M.; Zheng, X.; Yang, Y.; and Zhang, K. 2018b. Collaborative Filtering With Social Exposure: A Modular Approach to Social Recommendation. In *AAAI*, 2516–2523.
- Wang, X.; Zhang, R.; Sun, Y.; and Qi, J. 2019. Doubly Robust Joint Learning for Recommendation on Data Missing Not at Random. In *ICML*, 6638–6647.
- Wu, C. F. J. 1983. On the Convergence Properties of the EM Algorithm. *The Annals of Statistics* 11(1): 95–103.
- Yang, H.; Ling, G.; Su, Y.; Lyu, M. R.; and King, I. 2015. Boosting Response Aware Model-Based Collaborative Filtering. *IEEE Trans. Knowl. Data Eng.* 27(8): 2064–2077.
- Yang, L.; Cui, Y.; Xuan, Y.; Wang, C.; Belongie, S. J.; and Estrin, D. 2018. Unbiased offline recommender evaluation for missing-not-at-random implicit feedback. In *ACM RecSys*, 279–287.
- Yang, T.; Yao, R.; Yin, Q.; and Wu, O. 2020. Mitigating sentimental bias via a polar attention mechanism. *Int J Data Sci Anal* .
- Zhang, H.; Shen, F.; Liu, W.; He, X.; Luan, H.; and Chua, T. 2016. Discrete Collaborative Filtering. In *SIGIR*, 325–334.
- Zhang, S.; Yao, L.; Sun, A.; and Tay, Y. 2019. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* 52(1): 5:1–5:38.