

Hierarchical Bayesian Matrix Factorization with Side Information

Sunho Park¹, Yong-Deok Kim¹, Seungjin Choi^{1,2}

¹ Department of Computer Science and Engineering

² Division of IT Convergence Engineering

Pohang University of Science and Technology

77 Cheongam-ro, Nam-gu, Pohang 790-784, Korea

{titan,karma13,seungjin}@postech.ac.kr

Abstract

Bayesian treatment of matrix factorization has been successfully applied to the problem of collaborative prediction, where unknown ratings are determined by the predictive distribution, inferring posterior distributions over user and item factor matrices that are used to approximate the user-item matrix as their product. In practice, however, Bayesian matrix factorization suffers from cold-start problems, where inferences are required for users or items about which a sufficient number of ratings are not gathered. In this paper we present a method for Bayesian matrix factorization with side information, to handle cold-start problems. To this end, we place Gaussian-Wishart priors on mean vectors and precision matrices of Gaussian user and item factor matrices, such that mean of each prior distribution is regressed on corresponding side information. We develop variational inference algorithms to approximately compute posterior distributions over user and item factor matrices. In addition, we provide Bayesian Cramér-Rao Bound for our model, showing that the hierarchical Bayesian matrix factorization with side information improves the reconstruction over the standard Bayesian matrix factorization where the side information is not used. Experiments on MovieLens data demonstrate the useful behavior of our model in the case of cold-start problems.

1 Introduction

Matrix factorization refers to a method for uncovering a low-rank latent structure of data, approximating the data matrix as a product of two factor matrices. Matrix factorization is popular for collaborative prediction, where unknown ratings are predicted by user and item factor matrices which are determined to approximate a user-item matrix as their product [Rennie and Srebro, 2005; Salakhutdinov and Mnih, 2008b; Lim and Teh, 2007; Raiko *et al.*, 2007; Takács *et al.*, 2009; Koren *et al.*, 2009; Singh and Gordon, 2010].

Suppose that $\mathbf{X} \in \mathbb{R}^{I \times J}$ is a user-item rating matrix, the (i, j) -entry of which, X_{ij} , represents the rating of user i on item j . Matrix factorization determines factor matrices

$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_I] \in \mathbb{R}^{K \times I}$ and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_J] \in \mathbb{R}^{K \times J}$ (where K denotes the dimension of latent space), to approximate the rating matrix \mathbf{X} by $\mathbf{U}^\top \mathbf{V}$:

$$\mathbf{X} \approx \mathbf{U}^\top \mathbf{V}. \quad (1)$$

A popular approach is to minimize the regularized squared error loss defined as

$$\sum_{(i,j) \in \Omega} [(X_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2 + \lambda(\|\mathbf{u}_i\|^2 + \|\mathbf{v}_j\|^2)], \quad (2)$$

where Ω is a set of indices of observed entries in \mathbf{X} , and λ is the regularization parameter. Although this works even on large-scale datasets [Takács *et al.*, 2009; Koren *et al.*, 2009], it is prone to overfitting on the training data and require careful tuning of regularization parameter and the number of optimization steps.

Bayesian treatment of matrix factorization successfully alleviates the overfitting problem by integrating out all model parameter, thus allowing for complex models to be learned without requiring much parameter tuning [Lim and Teh, 2007; Raiko *et al.*, 2007; Salakhutdinov and Mnih, 2008a; Yoo and Choi, 2011; 2012; Kim and Choi, 2013]. In practice, however, Bayesian matrix factorization still suffers from *cold-start problem*, where inferences are required for users or items about which a sufficient number of ratings are not gathered. The cold-start problem commonly occurs in practical recommendation system based on collaborative prediction because new users or items are continuously added into the system. Moreover users usually do not rate unpopular items, thus these items will be fixed in long-tail part for a long time but for effective recommendation. To resolve the cold-start problem, it is crucial to use side information, such as user demographic information and item content information.

In this paper we present a method for Bayesian matrix factorization with side information, to handle cold-start problems. To this end, we place Gaussian-Wishart priors on mean vectors and precision matrices of Gaussian user and item factor matrices, such that mean of each prior distribution is regressed on corresponding side information. We develop variational inference algorithms to approximately compute posterior distributions over user and item factor matrices. In addition, we provide Bayesian Cramér-Rao Bound for our model, showing that the hierarchical Bayesian matrix factorization with side information improves the reconstruction over the

standard Bayesian matrix factorization where the side information is not used. Experiments on MovieLens data demonstrate the useful behavior of our model in the case of cold-start problems.

2 Related Work

We briefly review several matrix factorization models which incorporate side information. Suppose that side information is given in the form of feature matrices $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_I] \in \mathbb{R}^{D_u \times I}$ and $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_J] \in \mathbb{R}^{D_v \times J}$, where \mathbf{f}_i and \mathbf{g}_j are feature vector related with user i and item j respectively.

Matrix co-factorization is one of efficient ways to exploit the side information. Matrix co-factorization jointly decomposes multiple data matrices, where each deposition is coupled by sharing some factor matrices [Singh and Gordon, 2008]. For example, rating matrix \mathbf{X} , user side information matrix \mathbf{F} and item side information matrix \mathbf{G} are jointly decomposed as

$$\mathbf{F} = \mathbf{A}^\top \mathbf{U} + \mathbf{E}_F, \quad \mathbf{X} = \mathbf{U}^\top \mathbf{V} + \mathbf{E}_X, \quad \mathbf{G} = \mathbf{B}^\top \mathbf{V} + \mathbf{E}_G,$$

where factor matrices \mathbf{U} and \mathbf{V} are shared in the first two and last two decompositions, respectively and $\mathbf{E}_F, \mathbf{E}_X, \mathbf{E}_G$ represents Gaussian noise which reflect uncertainties. Recently Bayesian treatment of matrix co-factorization was developed, where the inference was performed using a sampling method [Singh and Gordon, 2010] or the variational method [Yoo and Choi, 2011; 2012].

In contrast to matrix co-factorization model, *Regression-based Latent Factor Model* (RLFM) [Agarwal and Chen, 2009] assumes that the user and item factor matrices are generated from the side information via linear regression,

$$\mathbf{U} = \mathbf{A}^\top \mathbf{F} + \mathbf{E}_U, \quad \mathbf{V} = \mathbf{B}^\top \mathbf{G} + \mathbf{E}_V, \quad (3)$$

and then rating matrix \mathbf{X} is generated by (1). More precisely, RLFM place isotropic Gaussian priors on user and item factor matrices \mathbf{U} and \mathbf{V} where mean of each prior is regressed on corresponding side information (see Figure 1-(b)),

$$p(\mathbf{U} | \mathbf{A}, \lambda_u) = \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i | \mathbf{A}^\top \mathbf{f}_i, \lambda_u^{-1}), \quad (4)$$

$$p(\mathbf{V} | \mathbf{B}, \lambda_v) = \prod_{j=1}^J \mathcal{N}(\mathbf{v}_j | \mathbf{B}^\top \mathbf{g}_j, \lambda_v^{-1}). \quad (5)$$

Bayesian matrix factorization with side information (BMFSI), proposed in [Proteous *et al.*, 2010], also utilizes side information via regression but it is differentiated from RLFM. In BMFSI model, the regression coefficient and side information are augmented in factor matrices such that the rating matrix \mathbf{X} is estimated by the sum of the collaborative prediction part and regression part against side information:

$$\mathbf{X} = \left[\mathbf{U}^\top \mathbf{A}^\top \mathbf{F}^\top \right] \left[\mathbf{V}^\top \mathbf{G}^\top \mathbf{B}^\top \right]^\top + \mathbf{E}_X \quad (6)$$

$$= \mathbf{U}^\top \mathbf{V} + \mathbf{A}^\top \mathbf{G} + \mathbf{F}^\top \mathbf{B} + \mathbf{E}_X. \quad (7)$$

Moreover BMFSI utilizes user-item specific features such as N ratings for item j , given by N nearest-neighborhoods of user i , which is not considered in matrix co-factorization and RLFM.

3 Model and Inference

In this section we present our hierarchical Bayesian modeling for matrix factorization which utilizes side information and explain how to our model brides Bayesian matrix factorization (BMF) [Salakhutdinov and Mnih, 2008a] and RLFM. We develop a variational inference algorithm to approximately compute the posterior distributions over factor matrices and provide rating prediction procedures for hold-out and fold-in cases also in the variational inference framework.

3.1 Hierarchical Bayesian Modeling

The graphical model representing our hierarchical Bayesian matrix factorization with side information (HBMFSI) is shown in Figure 1-(c). At first, the observed data X_{ij} is assumed to be generated by a sum of inner product of latent factors, bias terms regressed on the feature vectors, and noise:

$$X_{ij} = \mathbf{u}_i^\top \mathbf{v}_j + \mathbf{a}_0^\top \mathbf{f}_i + \mathbf{b}_0^\top \mathbf{g}_j + \mathcal{E}_{ij}, \quad (8)$$

for $(i, j) \in \Omega$, where Ω is a set of indices of observe entries in \mathbf{X} . The uncertainty in the model is absorbed by the noise \mathcal{E}_{ij} which is assumed to be Gaussian, i.e., $\mathcal{E}_{ij} \sim \mathcal{N}(\mathcal{E}_{ij} | 0, \rho^{-1})$ where ρ is the precision (the inverse of variance). Thus, the likelihood is given by

$$\begin{aligned} p(\mathbf{X} | \mathbf{U}, \mathbf{V}, \mathbf{a}_0, \mathbf{v}_0, \rho) \\ = \prod_{(i,j) \in \Omega} \mathcal{N}(X_{ij} | \mathbf{u}_i^\top \mathbf{v}_j + \mathbf{a}_0^\top \mathbf{f}_i + \mathbf{b}_0^\top \mathbf{g}_j, \rho^{-1}). \end{aligned} \quad (9)$$

The prior distributions over factor matrices \mathbf{U} and \mathbf{V} are assumed to be Gaussian:

$$p(\mathbf{U} | \Theta_u) = \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i | \mathbf{m}_i, \mathbf{\Lambda}_u^{-1}), \quad (10)$$

$$p(\mathbf{V} | \Theta_v) = \prod_{j=1}^J \mathcal{N}(\mathbf{v}_j | \mathbf{n}_j, \mathbf{\Lambda}_v^{-1}), \quad (11)$$

where $\Theta_u = \{\{\mathbf{m}_i\}_{i=1}^I, \mathbf{\Lambda}_u\}$ and $\Theta_v = \{\{\mathbf{n}_j\}_{j=1}^J, \mathbf{\Lambda}_v\}$. We further place Gaussian-Wishart priors on the $\Theta_u = \{\{\mathbf{m}_i\}_{i=1}^I, \mathbf{\Lambda}_u\}$ and $\Theta_v = \{\{\mathbf{n}_j\}_{j=1}^J, \mathbf{\Lambda}_v\}$:

$$p(\Theta_u | \mathbf{A}, \mathbf{m}_0, \gamma_0, \nu_0, \mathbf{W}_0) \quad (12)$$

$$= \prod_{i=1}^I \mathcal{N}(\mathbf{m}_i | \mathbf{m}_0 + \mathbf{A}^\top \mathbf{f}_i, (\gamma_0 \mathbf{\Lambda}_u)^{-1}) \mathcal{W}(\mathbf{\Lambda}_u | \mathbf{W}_0, \mu_0),$$

$$p(\Theta_v | \mathbf{B}, \mathbf{n}_0, \gamma_0, \nu_0, \mathbf{W}_0) \quad (13)$$

$$= \prod_{j=1}^J \mathcal{N}(\mathbf{n}_j | \mathbf{n}_0 + \mathbf{B}^\top \mathbf{g}_j, (\gamma_0 \mathbf{\Lambda}_v)^{-1}) \mathcal{W}(\mathbf{\Lambda}_v | \mathbf{W}_0, \mu_0),$$

and finally place Gaussian priors on regression coefficient matrices $\mathbf{A} \in \mathbb{R}^{D_u \times K}$, $\mathbf{B} \in \mathbb{R}^{D_v \times K}$ and vectors $\mathbf{a}_0 \in \mathbb{R}^{D_u}$,

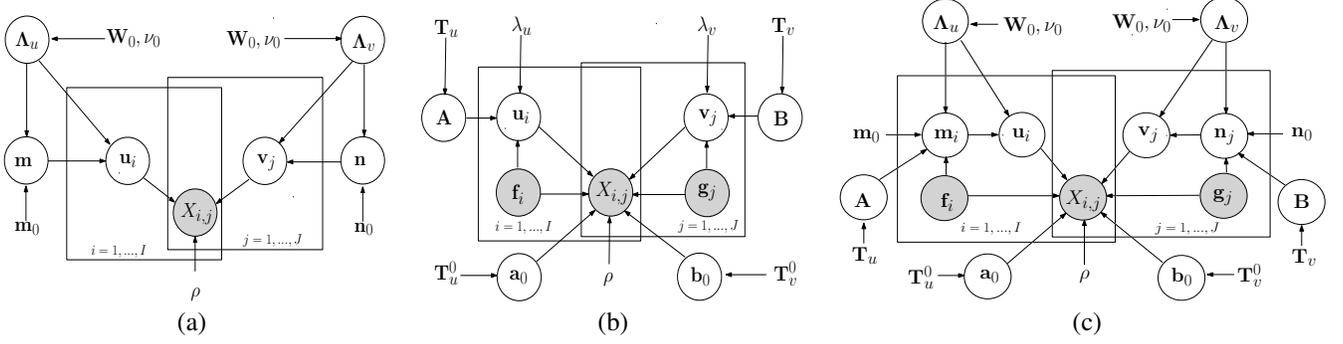


Figure 1: The probabilistic graphical models of (a) Standard Bayesian matrix factorization without side information, (b) Regression-based Latent Factor Model (RLFM) and (c) our hierarchical Bayesian matrix factorization with side information.

$\mathbf{b}_0 \in \mathbb{R}^{D_v}$:

$$p(\mathbf{A}) = \prod_{k=1}^K \mathcal{N}(\mathbf{a}_k | 0, \mathbf{T}_u^{-1}), \quad (14)$$

$$p(\mathbf{B}) = \prod_{k=1}^K \mathcal{N}(\mathbf{b}_k | 0, \mathbf{T}_v^{-1}), \quad (15)$$

$$p(\mathbf{a}) = \mathcal{N}(\mathbf{a}_0 | 0, (\mathbf{T}_u^0)^{-1}), \quad (16)$$

$$p(\mathbf{b}) = \mathcal{N}(\mathbf{b}_0 | 0, (\mathbf{T}_v^0)^{-1}), \quad (17)$$

where the precision matrices $\mathbf{T}_u, \mathbf{T}_u^0 \in \mathbb{R}^{D_u \times D_u}$ and $\mathbf{T}_v^0, \mathbf{T}_v \in \mathbb{R}^{D_v \times D_v}$ are assumed to be diagonal.

In standard BMF (Figure 1-(a)), means of the Gaussian-Wishart priors, \mathbf{m} and \mathbf{n} , are shared for all users and items and they generated from the non-informative constant vectors, \mathbf{m}_0 and \mathbf{n}_0 , which are trivially set to zero. However our HBMFSI places individual means for Gaussian-Wishart priors, $\{\mathbf{m}_i\}_{i=1}^I$ and $\{\mathbf{n}_j\}_{j=1}^J$, by associating with corresponding side information via regression:

$$\mathbf{m}_i = \mathbf{m}_0 + \mathbf{A}^\top \mathbf{f}_i + \mathbf{e}_i, \quad \text{for } i = 1, \dots, I, \quad (18)$$

$$\mathbf{n}_j = \mathbf{n}_0 + \mathbf{B}^\top \mathbf{g}_j + \mathbf{e}_j, \quad \text{for } j = 1, \dots, J. \quad (19)$$

Note that HBMFSI can be reduced to RLFM (Figure 1-(b)) if we marginalize out $\{\mathbf{m}_i\}_{i=1}^I$ and $\{\mathbf{n}_j\}_{j=1}^J$ by using the property of the linear Gaussian model. After marginalization, the conditional distributions over factor matrices \mathbf{U} and \mathbf{V} become

$$p(\mathbf{U} | \Lambda_u, \mathbf{A}, \mathbf{m}_0) = \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i | \mathbf{m}_0 + \mathbf{A}^\top \mathbf{f}_i, (\gamma \Lambda_u)^{-1}),$$

$$p(\mathbf{V} | \Lambda_v, \mathbf{B}, \mathbf{n}_0) = \prod_{j=1}^J \mathcal{N}(\mathbf{v}_j | \mathbf{n}_0 + \mathbf{B}^\top \mathbf{g}_j, (\gamma \Lambda_v)^{-1}),$$

where $\gamma = \gamma_0 / (\gamma_0 + 1)$. These conditional distributions are equivalent to priors over factor matrices in RLFM (4, 5), provided that $\mathbf{m}_0 = \mathbf{n}_0 = 0$ and precision matrices, Λ_u and Λ_v , are set to isotropic.

3.2 Variational Inference

We approximately computes posterior distributions over factor matrices by maximizing a lower-bound on the marginal

log-likelihood. Let \mathcal{Z} be a set of all latent variables:

$$\mathcal{Z} = \{\mathbf{U}, \mathbf{V}, \Lambda_u, \Lambda_v, \mathbf{A}, \mathbf{B}, \mathbf{a}_0, \mathbf{b}_0\}. \quad (20)$$

The marginal log-likelihood $\log p(\mathbf{X})$ is given by

$$\begin{aligned} \log p(\mathbf{X}) &= \log \int p(\mathbf{X}, \mathcal{Z}) d\mathcal{Z} \\ &\geq \int q(\mathcal{Z}) \log \frac{p(\mathbf{X}, \mathcal{Z})}{q(\mathcal{Z})} d\mathcal{Z} \\ &\equiv \mathcal{F}(q), \end{aligned} \quad (21)$$

where Jensen's inequality was used to obtain the *variational lower-bound* $\mathcal{F}(q)$ and $q(\mathcal{Z})$ is *variational distribution* that is assumed to be factorized:

$$q(\mathcal{Z}) = q(\mathbf{U})q(\mathbf{V})q(\Lambda_u)q(\Lambda_v)q(\mathbf{A})q(\mathbf{B})q(\mathbf{a}_0)q(\mathbf{b}_0). \quad (22)$$

Variational posterior distributions $q^*(\cdot)$, are determined by maximizing variational lower-bound $\mathcal{F}(q)$, leading to

$$q^*(\mathcal{Z}) \propto \exp \left\{ \langle \log p(\mathbf{X}, \mathcal{Z}) \rangle_{q(\mathcal{Z} \setminus \mathcal{Z})} \right\}, \quad \forall \mathcal{Z} \in \mathcal{Z}, \quad (23)$$

where the expectation is taken with respect to the variational distributions over all variables excluding \mathcal{Z} . Variational posterior distributions and corresponding updating equations for all latent variables $\mathcal{Z} \in \mathcal{Z}$ are summarized in Table 1.

We use the empirical Bayes estimation to update hyperparameters $\{\rho, \mathbf{m}_0, \mathbf{n}_0, \mathbf{T}_u^0, \mathbf{T}_u, \mathbf{T}_v^0, \mathbf{T}_v\}$. Setting derivatives of the variational lower-bound $\mathcal{F}(q)$ to zero yields the updating equations for hyperparameters:

$$\rho = |\Omega| / \sum_{(i,j) \in \Omega} \mathcal{E}_{ij}, \quad (24)$$

where

$$\mathcal{E}_{ij} = \left\langle \left(X_{ij} - \mathbf{u}_i^\top \mathbf{v}_j - \mathbf{a}_0^\top \mathbf{f}_i - \mathbf{b}_0^\top \mathbf{g}_j \right)^2 \right\rangle, \quad (25)$$

and

$$\mathbf{m}_0 = \frac{1}{I} \sum_{i=1}^I \left[\langle \mathbf{u}_i \rangle - \langle \mathbf{A} \rangle^\top \mathbf{f}_i \right], \quad (26)$$

$$\mathbf{n}_0 = \frac{1}{J} \sum_{j=1}^J \left[\langle \mathbf{v}_j \rangle - \langle \mathbf{B} \rangle^\top \mathbf{g}_j \right], \quad (27)$$

Table 1: Variational posteriors and corresponding updating equations are summarized. The symbol $\langle \cdot \rangle$ is the statistical expectation with respect to the corresponding variational distribution, Ω_i is a set of indices j for which $X_{i,j}$ is observed, $\mathbf{\Gamma}_u \in \mathbb{R}^{K \times K}$ is a diagonal matrix with K elements, $\boldsymbol{\lambda}_k \in \mathbb{R}^K$ is the k th column vector of $\langle \mathbf{\Lambda}_u \rangle$, $\mathbf{1}_u = [1, \dots, 1]^\top \in \mathbb{R}^I$, $\tilde{X}_{i,j} = X_{i,j} - \langle \mathbf{a}_0 \rangle^\top \mathbf{f}_i - \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_j$, and $\hat{\mathbf{C}}_u = \sum_{i=1}^I \mathbf{f}_i \mathbf{f}_i^\top$. Variational posterior for latent variables corresponding to items can be computed in similar way.

Variational distributions	Updating equations for variational parameters	Sufficient statistics
$q^*(\mathbf{U})$ $= \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i \boldsymbol{\psi}_i, \mathbf{L}_i^{-1})$	$\mathbf{L}_i = \gamma \langle \mathbf{\Lambda} \rangle + \rho \sum_{j \in \Omega_i} \langle \mathbf{v}_j \mathbf{v}_j^\top \rangle$ $\boldsymbol{\psi}_i = \mathbf{L}_i^{-1} \left\{ \gamma \langle \mathbf{\Lambda}^u \rangle (\mathbf{m}_0 + \langle \mathbf{A} \rangle^\top \mathbf{f}_i) + \rho \sum_{j \in \Omega_i} \tilde{X}_{i,j} \langle \mathbf{v}_j \rangle \right\}$	$\langle \mathbf{u}_i \rangle = \boldsymbol{\psi}_i$ $\langle \mathbf{u}_i \mathbf{u}_i^\top \rangle = \mathbf{L}_i^{-1} + \boldsymbol{\psi}_i \boldsymbol{\psi}_i^\top$
$q^*(\mathbf{\Lambda}_u)$ $= \mathcal{W}(\mathbf{\Lambda}_u \mathbf{W}_u, \nu_u)$	$\nu_u = \nu_0 + I$ $\mathbf{W}_u^{-1} = \mathbf{W}_0^{-1} + \gamma \left\{ \left(\langle \mathbf{U} \rangle - \langle \mathbf{A} \rangle^\top \mathbf{F} - \mathbf{m}_0 \mathbf{1}_u^\top \right) \cdot \left(\langle \mathbf{U} \rangle - \langle \mathbf{A} \rangle^\top \mathbf{F} - \mathbf{m}_0 \mathbf{1}_u^\top \right)^\top + \sum_{i=1}^I \mathbf{L}_i^{-1} + \mathbf{\Gamma}_u \right\}$	$\langle \mathbf{\Lambda}_u \rangle = \nu_u \mathbf{W}_u$ $\left([\mathbf{\Gamma}_u]_{k,k} = \text{Tr}(\hat{\mathbf{C}}_u \mathbf{H}_k^{-1}) \right)$
$q^*(\mathbf{A})$ $= \prod_{k=1}^K \mathcal{N}(\mathbf{a}_k \boldsymbol{\varphi}_k, \mathbf{H}_k^{-1})$	$\mathbf{H}_k = \mathbf{T}_u + \gamma [\langle \mathbf{\Lambda}_u \rangle]_{k,k} \hat{\mathbf{C}}_u$ $\boldsymbol{\varphi}_k = \gamma \mathbf{H}_k^{-1} \cdot \left\{ \sum_{i=1}^I \boldsymbol{\lambda}_k^\top (\langle \mathbf{u}_i \rangle - \mathbf{m}_0) \mathbf{f}_i - \hat{\mathbf{C}}_u \sum_{k' \neq k} [\boldsymbol{\lambda}_k]_{k'} \langle \mathbf{a}_{k'} \rangle \right\}$	$\langle \mathbf{a}_k \rangle = \boldsymbol{\varphi}_k$, $\langle \mathbf{a}_k \mathbf{a}_k^\top \rangle = \mathbf{H}_k^{-1} + \boldsymbol{\varphi}_k \boldsymbol{\varphi}_k^\top$
$q^*(\mathbf{a}_0)$ $= \mathcal{N}(\mathbf{a}_0 \boldsymbol{\varphi}_0, \mathbf{H}_0^{-1})$	$\mathbf{H}_0 = \mathbf{T}_u^0 + \rho \sum_{i=1}^I \left[\Omega_i \mathbf{f}_i \mathbf{f}_i^\top \right]$ $\boldsymbol{\varphi}_0 = \rho \mathbf{H}_0^{-1} \cdot \left\{ \sum_{i=1}^I \mathbf{f}_i \sum_{j \in \Omega_i} \left(X_{i,j} - \langle \mathbf{u}_i \rangle^\top \langle \mathbf{v}_j \rangle - \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_j \right) \right\}$	$\langle \mathbf{a}_0 \rangle = \boldsymbol{\varphi}_0$, $\langle \mathbf{a}_0 \mathbf{a}_0^\top \rangle = \mathbf{H}_0^{-1} + \boldsymbol{\varphi}_0 \boldsymbol{\varphi}_0^\top$

and

$$(\mathbf{T}_u^0)_{d,d} = 1 / \langle \mathbf{a}_0 \mathbf{a}_0^\top \rangle_{d,d}, \quad (28)$$

$$(\mathbf{T}_u)_{d,d} = K / \sum_{k=1}^K \langle \mathbf{a}_k \mathbf{a}_k^\top \rangle_{d,d}, \quad (29)$$

for $d = 1, \dots, D_u$ and

$$(\mathbf{T}_v^0)_{d,d} = 1 / \langle \mathbf{b}_0 \mathbf{b}_0^\top \rangle_{d,d}, \quad (30)$$

$$(\mathbf{T}_v)_{d,d} = K / \sum_{k=1}^K \langle \mathbf{b}_k \mathbf{b}_k^\top \rangle_{d,d}, \quad (31)$$

for $d = 1, \dots, D_v$.

3.3 Prediction

We consider two kinds of prediction tasks in the collaborative filtering: *hold-out* prediction and *fold-in* prediction. In the hold-out prediction, we aim to predict a missing entry X_{i^o, j^o} in the rating matrix \mathbf{X} , where $(i^o, j^o) \notin \Omega$. To do this, we simply make use of the posterior means of each variable:

$$\hat{X}_{i^o, j^o} = \langle \mathbf{u}_{i^o} \rangle^\top \langle \mathbf{v}_{j^o} \rangle + \langle \mathbf{a}_0 \rangle^\top \mathbf{f}_{i^o} + \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_{j^o} \quad (32)$$

In the fold-in prediction, we want to predict the rating value of a new user or new item. We assume that a user i^+ is newly added into the system and some rating entries \mathbf{x}_{i^+} made by the user i^+ are available. Then the prediction of an unknown entry of the user i^+ can be made similarly to (32):

$$\hat{X}_{i^+, j} = \langle \mathbf{u}_{i^+} \rangle^\top \langle \mathbf{v}_j \rangle + \langle \mathbf{a}_0 \rangle^\top \mathbf{f}_{i^+} + \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_j, \quad (33)$$

where $\langle \mathbf{u}_{i^+} \rangle$ is a mean of the predictive distribution over a new factor \mathbf{u}_{i^+} given \mathbf{x}_{i^+} , which can be defined by

$$p(\mathbf{u}_{i^+} | \mathbf{x}_{i^+}, \mathbf{X}) = \int p(\mathbf{u}_{i^+} | \mathbf{x}_{i^+}, \mathcal{Z}) p(\mathcal{Z} | \mathbf{X}) d\mathcal{Z}, \quad (34)$$

where \mathcal{Z} is the set of latent variables defined in (20). We approximate the predictive distribution (34) by maximizing a lower-bound on marginal log-posterior $p(\mathbf{x}_{i^+} | \mathbf{X})$:

$$\begin{aligned} \log p(\mathbf{x}_{i^+} | \mathbf{X}) &= \log \int \int p(\mathbf{x}_{i^+}, \mathbf{u}_{i^+}, \mathcal{Z} | \mathbf{X}) d\mathbf{u}_{i^+} d\mathcal{Z} \\ &\geq \int \int \tilde{q}(\mathbf{u}_{i^+}, \mathcal{Z}) \log \frac{p(\mathbf{x}_{i^+}, \mathbf{u}_{i^+}, \mathcal{Z} | \mathbf{X})}{\tilde{q}(\mathbf{u}_{i^+}, \mathcal{Z})} d\mathbf{u}_{i^+} d\mathcal{Z} \\ &\equiv \mathcal{F}(\tilde{q}). \end{aligned} \quad (35)$$

Variational distribution $\tilde{q}(\mathbf{u}_{i^+}, \mathcal{Z})$ is assumed to be factorized:

$$\tilde{q}(\mathbf{u}_{i^+}, \mathcal{Z}) = \tilde{q}(\mathbf{u}_{i^+}) \tilde{q}^*(\mathcal{Z}), \quad (36)$$

where $\tilde{q}^*(\mathcal{Z})$ is fixed to the pre-trained variational distribution in Table 1. For given pre-fixed $\tilde{q}^*(\mathcal{Z})$, the optimal variational posterior $\tilde{q}^*(\mathbf{u}_{i^+})$ maximizing (35) is computed with single update,

$$\tilde{q}^*(\mathbf{u}_{i^+}) \propto \exp \left\{ \langle \log p(\mathbf{x}_{i^+} | \mathbf{u}_{i^+}, \mathcal{Z}) p(\mathbf{u}_{i^+} | \mathcal{Z}) \rangle_{\tilde{q}^*(\mathcal{Z})} \right\}, \quad (37)$$

yielding to Gaussian distribution such that

$$\tilde{q}^*(\mathbf{u}_{i^+}) = \mathcal{N}(\mathbf{u}_{i^+} | \boldsymbol{\psi}_{i^+}, \mathbf{L}_{i^+}^{-1}), \quad (38)$$

where

$$\mathbf{L}_{i^+} = \gamma \langle \mathbf{\Lambda}_u \rangle + \langle \rho \rangle \sum_{j \in \Omega_{i^+}} \langle \mathbf{v}_j \mathbf{v}_j^\top \rangle, \quad (39)$$

$$\begin{aligned} \boldsymbol{\psi}_{i^+} &= \mathbf{L}_{i^+}^{-1} \left\{ \gamma \langle \mathbf{\Lambda}_u \rangle (\mathbf{m}_0 + \langle \mathbf{A} \rangle^\top \mathbf{f}_{i^+}) \right. \\ &\quad \left. + \langle \rho \rangle \sum_{j \in \Omega_{i^+}} \left(X_{i^+, j} - \langle \mathbf{a}_0 \rangle^\top \mathbf{f}_{i^+} - \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_j \right) \langle \mathbf{v}_j \rangle \right\}. \end{aligned}$$

Thus the prediction of unknown rating entry of the new user, $X_{i+,j}$ can be made based on (33) and (38):

$$\widehat{X}_{i+,j} = \boldsymbol{\psi}_{i+}^\top \mathbf{v}_j + \langle \mathbf{a}_0 \rangle^\top \mathbf{f}_{i+} + \langle \mathbf{b}_0 \rangle^\top \mathbf{g}_j. \quad (40)$$

The aforementioned prediction method can be similarly applied to predict missing ratings in the presence of a new item.

4 Bayesian Cramér-Rao Bound

Motivated by [Yoo and Choi, 2011], we provide Bayesian Cramér-Rao Bound for our model, showing that our HBMFSI improves the reconstruction over the standard BMF where the side information is not used.

The Bayesian Cramér-Rao bound or Posterior Cramér-Rao bound places a lower bound on the variance of any parametric estimator [Tichavsky *et al.*, 1998], as the inverse of the Fisher information matrix \mathcal{F} , which is written by,

$$\left\langle (\boldsymbol{\theta} - \boldsymbol{\theta}(\mathbf{x})) (\boldsymbol{\theta} - \boldsymbol{\theta}(\mathbf{x}))^\top \right\rangle \succeq \mathcal{F}^{-1}, \quad (41)$$

where $\boldsymbol{\theta}(\mathbf{x})$ is the estimate of the $\boldsymbol{\theta}$. Each element of the Fisher information matrix is defined by

$$\mathcal{F}_{i,j} = - \left\langle \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} \right\rangle_{p(\mathbf{x}, \boldsymbol{\theta})}, \quad (42)$$

where $p(\mathbf{x}, \boldsymbol{\theta})$ is the joint probability between the observations and the parameters. The benefit of using BCRB over classical Cramér-Rao bound is that the BCRB is known to provide a lower bound on the variance of any parametric estimators, even for the unbiased ones [Tichavsky *et al.*, 1998].

To compute the BCRB for HBMFSI, we rearrange the factor matrices and other parameters to be a vector form:

$$\boldsymbol{\theta} = [\mathbf{u}_1^\top \cdots \mathbf{u}_I^\top \mathbf{v}_1^\top \cdots \mathbf{v}_J^\top \mathbf{m}_1^\top \cdots \mathbf{m}_I^\top \mathbf{n}_1^\top \cdots \mathbf{n}_J^\top \text{vec}(\mathbf{A}_u)^\top \text{vec}(\mathbf{A}_v)^\top \text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B})^\top \mathbf{a}_0^\top \mathbf{b}_0^\top]^\top. \quad (43)$$

Each element of the Fisher information matrix is computed as (42). The second derivatives of the log joint probability with respect to elements of factor matrices are given by

$$\begin{aligned} \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial^2 U_{k,i}} &= -[\mathbf{A}_u]_{k,k} - \rho \sum_{j \in \Omega_i} V_{k,j}^2, \\ \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial U_{k,i} \partial U_{k^+,i}} &= -[\mathbf{A}_u]_{k,k^+} - \rho \sum_{j \in \Omega_i} V_{k,j} V_{k^+,j}, \\ \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial U_{k,i} \partial U_{k,i^+}} &= \frac{\partial^2 \log p(\mathbf{X}, \boldsymbol{\theta})}{\partial U_{k,i} \partial U_{k^+,i^+}} = 0, \\ \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial U_{k,i} \partial V_{k,j}} &= -\rho \{ (X_{i,j} - \mathbf{u}_i^\top \mathbf{v}_j) - U_{k,i} V_{k,j} \}, \\ \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial U_{k,i} \partial V_{k^+,j}} &= -\rho U_{k^+,i} V_{k,j}. \end{aligned}$$

In standard BMF where side information is not used, if we set \mathbf{m}_0 , \mathbf{n}_0 to zero vector and \mathbf{W}_0 to an identity matrix as in [Salakhutdinov and Mnih, 2008a], the most of the expectation above second derivatives are vanished and only nonzero values come from the differentiating with the same parameter,

$$- \left\langle \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial^2 U_{k,i}} \right\rangle_{p(\mathbf{x}, \boldsymbol{\theta})} = \nu_0 + \frac{\rho |\Omega_i| (\gamma_0 + 1) \nu_0}{\gamma_0 (\nu_0 - K - 1)}, \quad (44)$$

hence the Fisher information matrix becomes diagonal.

The Fisher information matrix of HBMFSI is also diagonal but larger than (44) because of extra term involved with side information,

$$\begin{aligned} & - \left\langle \frac{\partial^2 \log p(\mathbf{x}, \boldsymbol{\theta})}{\partial^2 U_{k,i}} \right\rangle_{p(\mathbf{x}, \boldsymbol{\theta})} \\ &= \frac{\nu_0 + \rho |\Omega_i| (\gamma_0 + 1) \nu_0}{\gamma_0 (\nu_0 - K - 1)} + \rho \sum_{j \in \Omega_i} \mathbf{g}_j^\top \mathbf{T}_v^{-1} \mathbf{g}_j, \end{aligned} \quad (45)$$

hence has lower BCRB (the inverse of the fisher information matrix) and also lower reconstruction error [Yoo and Choi, 2011].

5 Numerical Experiments

We evaluated the proposed HBMFSI for the collaborative prediction problem in both warm-start and cold-start situations, and compared the performance with that of BMF [Salakhutdinov and Mnih, 2008a], BMCF [Yoo and Choi, 2012], BMFSI [Proteous *et al.*, 2010] and RLFM [Agarwal and Chen, 2009] to show the benefit of the HBMFSI.

Datasets: We conducted experiments on two MovieLens datasets: MovieLens-100K, which consists of 100K ratings with 943 users and 1682 movies, and MovieLens-1M, which consists of 1M ratings with 6040 users and 3706 movies. MovieLens datasets provide side information for both users and items. The user information which consists of the user's age, gender and occupation was encoded into a binary valued vector with length 28. Similarly, the item information which consists of the 18 category of movie genre was encoded into a binary valued vector with length 18. Ratings were normalized to be zero-mean.

Evaluation Metrics: The performance of the collaborative prediction algorithm was evaluated in terms of the widely used root mean squared error (RMSE) defined by $\text{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \hat{x}_i)^2}$, and mean absolute error (MAE) defined by $\text{MAE} = \frac{1}{N} \sum_{i=1}^N |x_i - \hat{x}_i|$, where N is the total number of test data points, x_i and \hat{x}_i are the true rating and predicted rating of the i -th test data, respectively.

Model Training: Because each model was trained by different inference algorithm such as Markov Chain Monte Carlo (BMF, BMFSI), Monte Carlo Expectation Maximization (RLF), and the variational method (BMCF, HBMFSI), it is difficult to directly compare the effectiveness of the modeling for side information. For this reason, we decided to train all models by the variational method. We set the latent dimension K to 10 (compare to the Netflix data, relatively small K is sufficient for MovieLens datasets), and hyperparameters of Gaussian-Wishart priors to $\nu_0 = K$, $\gamma_0 = 1$ and $\mathbf{W}_0 = \mathbf{I}$.

Warm-start Performance: Firstly we investigate the predictive performance on the common warm-start situation with MovieLens-1M data. The data is randomly split around 50% training and 50% test. Experiments are repeated ten times, and we report their averages and standard deviations in Table 2. We confirmed that HBMFSI outperforms other methods in terms of both MAE and RMSE. Interestingly, BMFSI and BMCF showed slightly worse performance than BMF, even

Table 2: MAE and RMSE on MovieLens-1M data in warm-start situation. The number in parenthesis represents the standard deviation.

Method	MAE	RMSE
BMF	0.6761(0.0005)	0.8621(0.0005)
BMFSI	0.6785(0.0005)	0.8636(0.0005)
BMCF	0.6818(0.0009)	0.8666(0.0008)
RLMF	0.6718(0.0004)	0.8561(0.0004)
HBMFSI	0.6698(0.0005)	0.8536(0.0005)

though side information is utilized. In the case of BMFSI, similar result was reported in [Proteous *et al.*, 2010], where movie metadata which extracted from Wikipedia (e.g. movie director, actors, languages) do not improve performance measured by RMSE in Netflix Prize data.

Cold-start Performance: Next we investigate the predictive performance on the cold-start situation with MovieLens-100K data. We considered the cold-start problems which can occur in the hold-out prediction (prediction of missing ratings of the existing users) and also in the fold-in prediction (prediction of unknown ratings of the new users).

In the case of hold-out prediction, we randomly select 20% ratings as the test data and train each model on ten training sets, which consist of 10%, 20%, ..., 100% of the remained ratings. In the case of fold-in prediction, we randomly divide the users into two sets: a test set containing 20% of the users and training set containing the rest of the users. First each model is trained on all ratings by the training users. For each of the test users we train the model varying with the number of given ratings from 0 to 20 according to the procedure outlined in Section 3.3. Experiments are repeated ten times for both hold-out and fold-in predictions.

Figure 2 shows that HBMFSI outperforms BMF both in hold-out and fold-in prediction. Especially HBMFSI significantly improves the performance in fold-in prediction which reflects practical cold-start situation more than hold-out prediction. Experimental results also showed that side information modeling used in RLFM and HBMFSI, where side information is used to modeling the means of priors over factor matrices is more efficient to utilize side information than other methods including BMCF and BMFSI. Finally HBMFSI showed slightly better performance than RLFM, which is a special case of our model.

6 Conclusions

We have presented hierarchical Bayesian matrix factorization with side information, to handle cold-start problems. We used side information to properly regularize user and item factor matrices by placing Gaussian-Wishart priors on mean vectors and precision matrices of Gaussian user and item factor matrices, such that mean of each prior distribution is regressed on corresponding side information. We developed variational inference algorithms to approximately compute posterior distributions over user and item factor matrices. In addition, we provided Bayesian Cramér-Rao Bound for our model, showing that the hierarchical Bayesian matrix factorization with

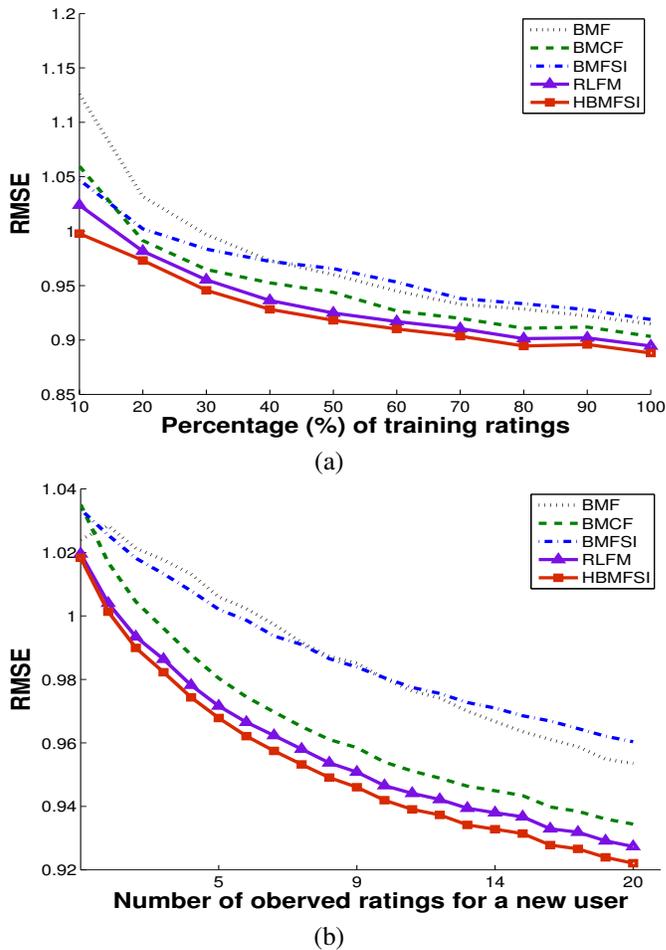


Figure 2: RMSE on MovieLens-100k data in cold-start situation: (a) hold-out prediction case, (b) fold-in prediction.

side information improves the reconstruction over the standard Bayesian matrix factorization where the side information is not used. Experiments on MovieLens data demonstrated the useful behavior of our model in the case of cold-start and also warm-start problems.

7 Acknowledgments

This work was supported by National Research Foundation (NRF) of Korea (2012-0005032), NIPA ITRC Support Program (NIPA-2013-H0301-13-3002), POSTECH Rising Star Program, and NRF World Class University Program (R31-10100).

References

- [Agarwal and Chen, 2009] D. Agarwal and B.-C. Chen. Regression-based latent factor models. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, Paris, France, 2009.
- [Kim and Choi, 2013] Y.-D. Kim and S. Choi. Variational Bayesian view of weighted trace norm regularization for

- matrix factorization. *IEEE Signal Processing Letters*, 20(3):261–264, 2013.
- [Koren *et al.*, 2009] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [Lim and Teh, 2007] Y. J. Lim and Y. W. Teh. Variational Bayesian approach to movie rating prediction. In *Proceedings of KDD Cup and Workshop*, San Jose, CA, 2007.
- [Proteous *et al.*, 2010] I. Proteous, A. Asuncion, and M. Welling. Bayesian matrix factorization with side information and Dirichlet process mixtures. In *Proceedings of the AAAI National Conference on Artificial Intelligence (AAAI)*, Atlanta, Georgia, USA, 2010.
- [Raiko *et al.*, 2007] T. Raiko, A. Ilin, and J. Karhunen. Principal component analysis for large scale problems with lots of missing values. In *Proceedings of the European Conference on Machine Learning (ECML)*, pages 691–698, Warsaw, Poland, 2007.
- [Rennie and Srebro, 2005] J. D. M. Rennie and N. Srebro. Fast maximum margin matrix factorization for collaborative prediction. In *Proceedings of the International Conference on Machine Learning (ICML)*, Bonn, Germany, 2005.
- [Salakhutdinov and Mnih, 2008a] R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using MCMC. In *Proceedings of the International Conference on Machine Learning (ICML)*, Helsinki, Finland, 2008.
- [Salakhutdinov and Mnih, 2008b] R. Salakhutdinov and A. Mnih. Probabilistic matrix factorization. In *Advances in Neural Information Processing Systems (NIPS)*, volume 20. MIT Press, 2008.
- [Singh and Gordon, 2008] A. P. Singh and G. J. Gordon. Relational learning via collective matrix factorization. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, Las Vegas, Nevada, 2008.
- [Singh and Gordon, 2010] A. P. Singh and G. J. Gordon. A Bayesian matrix factorization model for relational data. In *Proceedings of the Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, Catalina Island, CA, 2010.
- [Takács *et al.*, 2009] G. Takács, I. Pilászy, B. Németh, and D. Tikk. Scalable collaborative filtering approaches for large recommender systems. *Journal of Machine Learning Research*, 10:623–656, 2009.
- [Tichavsky *et al.*, 1998] P. Tichavsky, C. H. Muravchik, and A. Nehorai. Posterior Cramér-Rao bounds for discrete-time nonlinear filtering. *IEEE Transactions on Signal Processing*, 46(5):1386–1395, 1998.
- [Yoo and Choi, 2011] J. Yoo and S. Choi. Bayesian matrix co-factorization: Variational algorithm and Cramér-Rao bound. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD)*, Athens, Greece, 2011.
- [Yoo and Choi, 2012] J. Yoo and S. Choi. Hierarchical variational Bayesian matrix co-factorization. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Kyoto, Japan, 2012.