
Scalable Variational Bayesian Matrix Factorization with Side Information

Yong-Deok Kim and Seungjin Choi

Department of Computer Science and Engineering
Pohang University of Science and Technology
77 Cheongam-ro, Nam-gu, Pohang 790-784, Korea
{karma13, seungjin}@postech.ac.kr

Abstract

Bayesian matrix factorization (BMF) is a popular method for collaborative prediction, because of its robustness to overfitting as well as of being free from cross-validation for fine tuning of regularization parameters. In practice, however, due to its cubic time complexity with respect to the rank of factor matrices, existing variational inference algorithms for BMF are not well suited to web-scale datasets where billions of ratings provided by millions of users are available. The time complexity even increases when the side information, such as user binary implicit feedback or item content information, is incorporated into variational Bayesian matrix factorization (VBMF). For instance, a state of the arts in VBMF with side information, is to place Gaussian priors on user and item factor matrices, where mean of each prior is regressed on the corresponding side information. Since this approach introduces additional cubic time complexity with respect to the size of feature vectors, the use of rich side information in a form of high-dimensional feature vector is prohibited. In this paper, we present a scalable inference for VBMF with side information, the complexity of which is linear in the rank K of factor matrices. Moreover, the algorithm can be easily parallelized on multi-core systems. Experiments on large-scale datasets demonstrate the useful behavior of our algorithm such as scalability, fast learning, and prediction accuracy.

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 [1, 2, 3, 4, 5, 6]. 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 is the rank of factor matrices), 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. The problem (2) can be efficiently solved by alternating least squares or stochastic gradient descent methods [6, 7, 8], and parallelization of these methods are recently developed [9, 10, 11, 12]. However these approaches are prone to overfitting on the training data, hence require careful tuning of meta parameter such as regularization parameter, learning rate, and the number of iteration.

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 [2, 13, 3, 14, 15, 16]. In addition, Bayesian matrix factorization (BMF) easily incorporates side information such as user binary implicit feedback and item content information by placing Gaussian priors on user and item factor matrices, where mean of each prior is regressed on corresponding side information [17, 18].

The time and space complexity of VBMF depends on the mean-field assumption for variational distributions. When variational distributions are assumed to be matrix-wise independent, VBMF requires cubic time and quadratic space complexity with respect to K , existing variational inference algorithms for BMF are not well suited to web-scale datasets where billions of ratings provided by millions of users are available. When side information is incorporated into VBMF, it requires additional cubic time and quadratic space complexity with respect to the size of feature vectors formed by side information [18]. Thus, the use of rich side information in the form of high dimensional feature vector is prohibited. Assuming that variational distributions fully factorize to satisfy element-wise independence, a naive implementation of coordinate ascent algorithm requires quadratic time complexity with respect to K to update variational posterior means in element-wise fashion.

In this paper, we also assume that variational distributions fully factorize, but present a scalable algorithm for VBMF with side information, the time complexity of which is *linear* in K to update variational posterior means in element-wise fashion. Regression on side information also requires only linear time complexity with respect to the number of nonzero elements in feature vectors provided by side information. In addition, the proposed algorithm is easily parallelized on multi-core systems. Experiments on large-scale datasets (MovieLens10M, Douban, Netflix, and Yahoo! Music) with various side information (binary implicit feedback, neighborhood, social network, movie contents, music taxonomy) demonstrate the useful behavior of our algorithm such as scalability, fast learning, and prediction accuracy.

2 METHOD

In this section, we present our main contribution, a scalable algorithm for VBMF with side information, whose graphical representation is shown in Fig. 1.

2.1 Model

VBMF with side information assumes that the observed data X_{ij} is generated by

$$X_{ij} = \mathbf{u}_i^\top \mathbf{v}_j + e_{ij}, \quad (3)$$

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 e_{ij} which is assumed to be Gaussian, i.e., $e_{ij} \sim \mathcal{N}(e_{ij} | 0, \tau^{-1})$ where τ is the precision (the inverse of variance). Thus, the likelihood is given by

$$p(\mathbf{X} | \mathbf{U}, \mathbf{V}, \tau) = \prod_{(i,j) \in \Omega} \mathcal{N}(X_{ij} | \mathbf{u}_i^\top \mathbf{v}_j, \tau^{-1}). \quad (4)$$

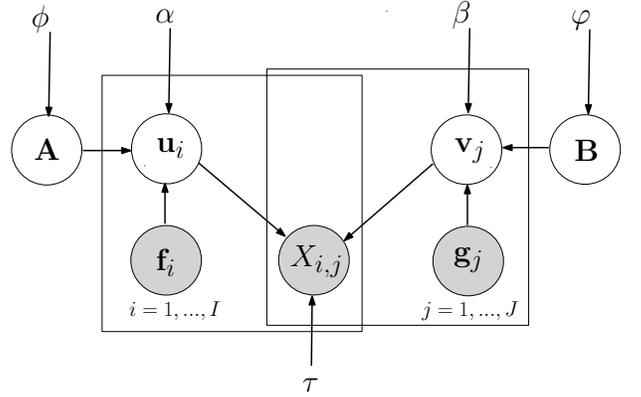


Figure 1: A graphical model for VBMF with side information.

Suppose that side information is given in the form of feature matrices $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_I] \in \mathbb{R}^{M \times I}$ and $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_J] \in \mathbb{R}^{N \times J}$, where \mathbf{f}_i and \mathbf{g}_j are feature vectors related with user i and item j respectively. VBMF with side information assumes that factor matrices \mathbf{U} and \mathbf{V} are generated from the side information via the linear regression,

$$\mathbf{U} = \mathbf{A}^\top \mathbf{F} + \mathbf{E}^u, \quad (5)$$

$$\mathbf{V} = \mathbf{B}^\top \mathbf{G} + \mathbf{E}^v, \quad (6)$$

where $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_K] \in \mathbb{R}^{M \times K}$ and $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_K] \in \mathbb{R}^{N \times K}$ are regression coefficient matrices. More precisely, Gaussian priors are placed on user and item factor matrices, where mean of each prior is regressed on corresponding side information:

$$p(\mathbf{U} | \mathbf{A}, \mathbf{F}, \alpha) = \prod_{k=1}^K \prod_{i=1}^I \mathcal{N}(u_{ki} | \mathbf{a}_k^\top \mathbf{f}_i, \alpha_k^{-1}), \quad (7)$$

$$p(\mathbf{V} | \mathbf{B}, \mathbf{G}, \beta) = \prod_{k=1}^K \prod_{j=1}^J \mathcal{N}(v_{kj} | \mathbf{b}_k^\top \mathbf{g}_j, \beta_k^{-1}). \quad (8)$$

Finally Gaussian priors are placed on regression coefficient matrices \mathbf{A} and \mathbf{B} :

$$p(\mathbf{A} | \phi) = \prod_{m=1}^M \prod_{k=1}^K \mathcal{N}(a_{mk} | 0, \phi_k^{-1}), \quad (9)$$

$$p(\mathbf{B} | \varphi) = \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(b_{nk} | 0, \varphi_k^{-1}), \quad (10)$$

with hyperparameters ϕ, φ .

2.2 Variational Inference

We approximately computes posterior distributions over factor matrices \mathbf{U}, \mathbf{V} and regression coefficient matrices

Table 1: Variational distributions and corresponding updating rules are summarized when matrix-wise independence (13) is assumed. In addition, the space and time complexity for updating rules are also summarized. Ω_i is a set of indices j for which X_{ij} is observed, $\text{diag}(\boldsymbol{\alpha})$ is a diagonal matrix with diagonal entries $\{\alpha_k\}$, \mathbf{I}_M is the identity matrix with dimension M , and $\mathbf{U}_k \in \mathbb{R}^{1 \times I}$ is the k^{th} row vector of \mathbf{U} .

Variational distributions	Updating rules	Space complexity	Time complexity
$q(\mathbf{U})$ $= \prod_{i=1}^I \mathcal{N}(\mathbf{u}_i \bar{\mathbf{u}}_i, \boldsymbol{\Lambda}_i^u)$	$\boldsymbol{\Lambda}_i^u = \left(\text{diag}(\boldsymbol{\alpha}) + \tau \sum_{j \in \Omega_i} (\bar{\mathbf{v}}_j \bar{\mathbf{v}}_j^\top + \boldsymbol{\Lambda}_j^v) \right)^{-1}$, $\bar{\mathbf{u}}_i = \boldsymbol{\Lambda}_i^u \left(\text{diag}(\boldsymbol{\alpha}) \mathbf{A}^\top \mathbf{f}_i + \tau \sum_{j \in \Omega_i} X_{ij} \bar{\mathbf{v}}_j \right)$	$I(K + K^2)$	$ \Omega K^2 + IK^3$
$q(\mathbf{V})$ $= \prod_{j=J}^I \mathcal{N}(\mathbf{v}_j \bar{\mathbf{v}}_j, \boldsymbol{\Lambda}_j^v)$	$\boldsymbol{\Lambda}_j^v = \left(\text{diag}(\boldsymbol{\beta}) + \tau \sum_{i \in \Omega_j} (\bar{\mathbf{u}}_i \bar{\mathbf{u}}_i^\top + \boldsymbol{\Lambda}_i^u) \right)^{-1}$, $\bar{\mathbf{v}}_j = \boldsymbol{\Lambda}_j^v \left(\text{diag}(\boldsymbol{\beta}) \mathbf{B}^\top \mathbf{g}_j + \tau \sum_{i \in \Omega_j} X_{ij} \bar{\mathbf{u}}_i \right)$	$J(K + K^2)$	$ \Omega K^2 + JK^3$
$q(\mathbf{A})$ $= \prod_{k=1}^K \mathcal{N}(\mathbf{a}_k \bar{\mathbf{a}}_k, \boldsymbol{\Lambda}_k^a)$	$\boldsymbol{\Lambda}_k^a = \left(\phi_k \mathbf{I}_M + \alpha_k \mathbf{F} \mathbf{F}^\top \right)^{-1}$, $\bar{\mathbf{a}}_k = \boldsymbol{\Lambda}_k^a \left(\alpha_k \mathbf{F} \bar{\mathbf{U}}_k^\top \right)$	$K(M + M^2)$	$K(IM + M^3)$
$q(\mathbf{B})$ $= \prod_{k=1}^K \mathcal{N}(\mathbf{b}_k \bar{\mathbf{b}}_k, \boldsymbol{\Lambda}_k^b)$	$\boldsymbol{\Lambda}_k^b = \left(\varphi_k \mathbf{I}_N + \beta_k \mathbf{G} \mathbf{G}^\top \right)^{-1}$, $\bar{\mathbf{b}}_k = \boldsymbol{\Lambda}_k^b \left(\beta_k \mathbf{G} \bar{\mathbf{V}}_k^\top \right)$	$K(N + N^2)$	$K(JN + N^3)$

\mathbf{A}, \mathbf{B} by maximizing a lower-bound on the marginal log-likelihood. Let \mathcal{Z} be a set of all latent variables:

$$\mathcal{Z} = \{\mathbf{U}, \mathbf{V}, \mathbf{A}, \mathbf{B}\}. \quad (11)$$

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

$$\begin{aligned} \log p(\mathbf{X}) &= \log \int p(\mathbf{X}, \mathcal{Z}) d\mathcal{Z} \\ &\geq \int q(\mathcal{Z}) \log \frac{p(\mathbf{X}, \mathcal{Z})}{q(\mathcal{Z})} d\mathcal{Z} \\ &\equiv \mathcal{F}(q), \end{aligned} \quad (12)$$

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

In our previous work on VBMF with side information [18], variational distributions are assumed to be matrix-wise independent as in VBMF [2]

$$q(\mathcal{Z}) = q(\mathbf{U}, \mathbf{V}, \mathbf{A}, \mathbf{B}) = q(\mathbf{U})q(\mathbf{V})q(\mathbf{A})q(\mathbf{B}). \quad (13)$$

The free-form optimization results in products of multivariate Gaussian distributions as summarized in Table 1. Numerical experiments with MovieLens datasets showed that prediction accuracy can be significantly improved by incorporating user demographic information and movie genre information.

In practice, however, due to the cubic time and quadratic space complexity with respect to K , the variational inference algorithms under matrix-wise independence assumption are not well suited to web-scale datasets billions of ratings provided by millions of users are available. To make matter worse, the use of rich side information in the form of high dimensional feature vector (e.g. $M, N \gg 10000$) is prohibited because of cubic time complexity with respect to M and N .

2.3 Scalable Algorithm

Instead of the matrix-wisely factorized variational distributions (13), we consider fully factorized variational distributions which satisfies an element-wise independence,

$$q(\mathbf{U}, \mathbf{V}, \mathbf{A}, \mathbf{B}) = \prod_{k,i} q(u_{ki}) \prod_{k,j} q(v_{kj}) \prod_{m,k} q(a_{mk}) \prod_{n,k} q(b_{nk}). \quad (14)$$

The free-form optimization results in products of univariate Gaussian distributions:

$$\begin{aligned} q(\mathbf{U}) &= \prod_{k=1}^K \prod_{i=1}^I \mathcal{N}(u_{ki} | \bar{u}_{ki}, s_{ki}^u), \\ q(\mathbf{V}) &= \prod_{k=1}^K \prod_{j=1}^J \mathcal{N}(v_{kj} | \bar{v}_{kj}, s_{kj}^v), \\ q(\mathbf{A}) &= \prod_{m=1}^M \prod_{k=1}^K \mathcal{N}(a_{mk} | \bar{a}_{mk}, s_{mk}^a), \\ q(\mathbf{B}) &= \prod_{n=1}^N \prod_{k=1}^K \mathcal{N}(b_{nk} | \bar{b}_{nk}, s_{nk}^b). \end{aligned} \quad (15)$$

We compute the variational lower-bound $\mathcal{F}(q)$:

$$\begin{aligned} \mathcal{F}(q) &= \mathbb{E}_q \{ \log p(\mathbf{X} | \mathbf{U}, \mathbf{V}) \\ &\quad + \log p(\mathbf{U} | \mathbf{A}, \mathbf{F}) + \log p(\mathbf{V} | \mathbf{B}, \mathbf{G}) \\ &\quad + \log p(\mathbf{A}) + \log p(\mathbf{B}) - \log q(\mathbf{U}, \mathbf{V}, \mathbf{A}, \mathbf{B}) \}, \end{aligned}$$

where $\mathbb{E}_q\{\cdot\}$ denotes the statistical expectation with respect to $q(\mathbf{U}, \mathbf{V}, \mathbf{A}, \mathbf{B})$ that is of the form (15). Invoking the model (4) with prior distributions (7), (8), (9) and (10) as well as the fully factorized variational distribution (15), we

calculate $\mathcal{F}(q)$ as

$$\begin{aligned} \mathcal{F}(q) = & \sum_{(i,j) \in \Omega} \mathcal{F}_{ij} + \sum_{k=1}^K \sum_{i=1}^I \mathcal{F}_{ki}^u + \sum_{k=1}^K \sum_{j=1}^J \mathcal{F}_{kj}^v \\ & + \sum_{m=1}^M \sum_{k=1}^K \mathcal{F}_{mk}^a + \sum_{n=1}^N \sum_{k=1}^K \mathcal{F}_{nk}^b, \end{aligned} \quad (16)$$

where

$$\begin{aligned} \mathcal{F}_{ij} &= \mathbb{E}_q \{ \log p(X_{ij} | \mathbf{U}, \mathbf{V}) \} \\ &= -\frac{\tau}{2} (R_{ij}^2 + W_{ij}) - \frac{1}{2} \log(2\pi\tau^{-1}), \\ R_{ij} &= X_{ij} - \sum_{k=1}^K \bar{u}_{ki} \bar{v}_{kj}, \\ W_{ij} &= \sum_{k=1}^K (\bar{u}_{ki}^2 s_{kj}^v + \bar{v}_{kj}^2 s_{ki}^u + s_{ki}^u s_{kj}^v), \\ \mathcal{F}_{ki}^u &= \mathbb{E}_q \{ \log p(u_{ki} | \mathbf{a}_k, \mathbf{f}_i) - \log q(u_{ki}) \} \\ &= -\frac{\alpha_k}{2} \left[(\bar{u}_{ki} - \bar{\mathbf{a}}_k^\top \mathbf{f}_i)^2 + s_{ki}^u + \sum_{m=1}^M s_{mk}^a f_{mi}^2 \right] \\ &\quad + \frac{1}{2} \log(s_{ki}^u \alpha_k) + \frac{1}{2}, \\ \mathcal{F}_{kj}^v &= \mathbb{E}_q \{ \log p(v_{kj} | \bar{\mathbf{b}}_k, \mathbf{g}_j) - \log q(v_{kj}) \} \\ &= -\frac{\beta_k}{2} \left[(\bar{v}_{kj} - \bar{\mathbf{b}}_k^\top \mathbf{g}_j)^2 + s_{kj}^v + \sum_{n=1}^N s_{nk}^b g_{nj}^2 \right] \\ &\quad + \frac{1}{2} \log(s_{kj}^v \beta_k) + \frac{1}{2}, \\ \mathcal{F}_{mk}^a &= \mathbb{E}_q \{ \log p(a_{mk}) - \log q(a_{mk}) \} \\ &= -\frac{\phi_k}{2} (\bar{a}_{mk}^2 + s_{mk}^a) + \frac{1}{2} \log(s_{mk}^a \phi_k) + \frac{1}{2}, \\ \mathcal{F}_{nk}^b &= \mathbb{E}_q \{ \log p(b_{nk}) - \log q(b_{nk}) \} \\ &= -\frac{\varphi_k}{2} (\bar{b}_{nk}^2 + s_{nk}^b) + \frac{1}{2} \log(s_{nk}^b \varphi_k) + \frac{1}{2}. \end{aligned}$$

The updating rules for variational parameters can be obtained by setting derivative the variational lower-bound (16) to zero. For example, in the case of $q(u_{ki})$, the updating rules for variance and mean are given by

$$s_{ki}^u = (\alpha_k + \tau \sum_{j \in \Omega_i} (\bar{v}_{kj}^2 + s_{kj}^v))^{-1}, \quad (17)$$

$$\bar{u}_{ki} = s_{ki}^u (\alpha_k \bar{\mathbf{a}}_k^\top \mathbf{f}_i + \tau \sum_{j \in \Omega_i} (X_{ij} - \sum_{k' \neq k} \bar{u}_{k'i} \bar{v}_{k'j}) \bar{v}_{kj}),$$

where Ω_i is a set of indices j for which X_{ij} is observed. Since updating of \bar{u}_{ki} requires $O(|\Omega_i|K)$, and similarly updating \bar{v}_{kj} requires $O(|\Omega_j|K)$, the time complexity for updating all variational posterior means of factor matrices is quadratic with respect to K :

$$O \left(\sum_{k=1}^K \sum_{i=1}^I |\Omega_i| K + \sum_{k=1}^K \sum_{j=1}^J |\Omega_j| K \right) = O(2|\Omega|K^2).$$

Although this naive implementation requires quadratic time complexity with respect to K , we can efficiently reduce the time complexity to linear with a simple trick. Let $R_{ij} = X_{ij} - \sum_{k=1}^K \bar{u}_{ki} \bar{v}_{kj}$ denote the residual on (i, j) observation. With the residual R_{ij} , the updating rule (17) can be rewritten as

$$\bar{u}_{ki} = s_{ki}^u \left(\alpha_k \bar{\mathbf{a}}_k^\top \mathbf{f}_i + \tau \sum_{j \in \Omega_i} (R_{ij} + \bar{u}_{ki} \bar{v}_{kj}) \bar{v}_{kj} \right)$$

and it requires constant time complexity $O(|\Omega_i|)$ with respect to K . When \bar{u}_{ki} is changed to \bar{u}'_{ki} , R_{ij} can be easily updated to R'_{ij} for all $j \in \Omega_i$ by:

$$R'_{ij} = R_{ij} - (\bar{u}'_{ki} - \bar{u}_{ki}) \bar{v}_{kj}. \quad (18)$$

Now time complexity per iteration (including updating of variational posterior variances and residuals) is reduced to linear with respect to K :

$$O \left(\sum_{k=1}^K \sum_{i=1}^I 3|\Omega_i| + \sum_{k=1}^K \sum_{j=1}^J 3|\Omega_j| \right) = O(6|\Omega|K). \quad (19)$$

Using the similar idea, variational posterior means of regression coefficient matrices can be efficiently updated. Let $\mathbf{R}^u = \bar{\mathbf{U}} - \bar{\mathbf{A}}^\top \mathbf{F}$ denote the residual on the variational posterior mean $\bar{\mathbf{U}}$. With the residual \mathbf{R}^u , the updating rules for s_{mk}^a and \bar{a}_{mk} can be written as

$$s_{mk}^a = (\phi_k + \alpha_k \|F_m\|^2)^{-1}, \quad (20)$$

$$\bar{a}_{mk} = s_{mk}^a \alpha_k (R_k^u + \bar{a}_{mk} F_m) F_m^\top, \quad (21)$$

where $R_k^u \in \mathbb{R}^{1 \times I}$ is the k^{th} row vector of \mathbf{R}^u and $F_m \in \mathbb{R}^{1 \times I}$ is the m^{th} row vector of \mathbf{F} . When \bar{a}_{mk} is updated to \bar{a}'_{mk} , R_k^u is updated to $R_k^{u'}$:

$$R_k^{u'} = R_k^u - (\bar{a}'_{mk} - \bar{a}_{mk}) F_m. \quad (22)$$

Note that updating s_{mk}^a and \bar{a}_{mk} takes $O(3\|F_m\|_0)$, where $\|\cdot\|$ denotes the L_0 norm which equals to number of non-zero elements. Similarly updating s_{nk}^b and \bar{b}_{nk} takes $O(3\|G_n\|_0)$, hence overall time complexity for updating variational posterior variances and means $\{\mathbf{S}^a, \bar{\mathbf{A}}, \mathbf{S}^b, \bar{\mathbf{B}}\}$ of regression coefficient matrices is linear in the number of nonzero elements in feature matrices:

$$\begin{aligned} O \left(\sum_{m=1}^M \sum_{k=1}^K 3\|F_m\|_0 + \sum_{n=1}^N \sum_{k=1}^K 3\|G_n\|_0 \right) \\ = O(3(\|\mathbf{F}\|_0 + \|\mathbf{G}\|_0)K). \end{aligned}$$

Algorithm 1 VBMFSI-CA: linear complexity learning algorithm for variational Bayesian matrix factorization with side information. VBMFSI-CA can be easily parallelized in column-by-column manner.

```

1: Initialize  $\bar{U}$ ,  $S^u$ ,  $\bar{V}$ ,  $S^v$ ,  $\bar{A}$ ,  $S^a$ ,  $\bar{B}$ ,  $S^b$ .
2:  $R = X - \bar{U}^\top \bar{V}$ ,
3:  $R^u = \bar{U} - \bar{A}^\top F$ ,  $R^v = \bar{V} - \bar{B}^\top G$ 
4: for  $t = 1 \dots T$  do
5:   /****** Update  $q(U)$  *****/
6:   parallel for  $i = 1, \dots, I$  do
7:     for  $k = 1, \dots, K$  do
8:        $\xi \leftarrow \bar{u}_{ki}$ 
9:        $s_{ki}^u \leftarrow (\alpha_k + \tau \sum_{j \in \Omega_i} (\bar{v}_{kj}^2 + s_{kj}^v))^{-1}$ 
10:       $\theta \leftarrow (\alpha_k \bar{a}_k^\top \mathbf{f}_i + \tau \sum_{j \in \Omega_i} (R_{ij} + \bar{u}_{ki} \bar{v}_{kj}) \bar{v}_{kj})$ 
11:       $\bar{u}_{ki} \leftarrow s_{ki}^u \theta$ 
12:      for  $j \in \Omega_i$  do
13:         $R_{ij} \leftarrow R_{ij} - (\bar{u}_{ki} - \xi) \bar{v}_{kj}$ 
14:      end for
15:    end for
16:  end parallel for
17:  /****** Update  $q(V)$  *****/
18:  parallel for  $j = 1, \dots, J$  do
19:    for  $k = 1, \dots, K$  do
20:       $\xi \leftarrow \bar{v}_{kj}$ 
21:       $s_{kj}^v \leftarrow (\beta_k + \tau \sum_{i \in \Omega_j} (\bar{u}_{ki}^2 + s_{ki}^u))^{-1}$ 
22:       $\theta \leftarrow (\beta_k \bar{b}_k^\top \mathbf{g}_j + \tau \sum_{i \in \Omega_j} (R_{ij} + \bar{u}_{ki} \bar{v}_{kj}) \bar{u}_{ki})$ 
23:       $\bar{v}_{kj} \leftarrow s_{kj}^v \theta$ 
24:      for  $i \in \Omega_j$  do
25:         $R_{ij} \leftarrow R_{ij} - \bar{u}_{ki} (\bar{v}_{kj} - \xi)$ 
26:      end for
27:    end for
28:  end parallel for

```

Our linear time complexity learning algorithm for VBMF with side information, which is referred to as VBMFSI-CA in this paper, is summarized in Algorithm 1. The space complexity of VBMFSI-CA is also linear with respect to K because we only have to save the variational posterior means and variances instead of covariances. VBMFSI-CA is a coordinate ascent method, where a single variational posterior parameter or hyperparameter is updated at a time while keeping others fixed. The convergence of coordinate ascent method is guaranteed when the solution for each coordinate is uniquely attained [19] and VBMFSI-CA fulfills this condition. Note that parallelization can be easily done in a column-by-column manner because each column of variational parameters $\{S^u, \bar{U}, S^v, \bar{V}, S^b, \bar{A}, S^a, \bar{B}\}$ can be updated independently from the updates of other columns.

Our approach bases on the element-wise independence assumption for variational distributions which is much

```

29:  /****** Update  $q(A)$  *****/
30:  parallel for  $k = 1, \dots, K$  do
31:    for  $m = 1, \dots, M$  do
32:       $\xi \leftarrow \bar{a}_{mk}$ 
33:       $s_{mk}^a \leftarrow (\phi_k + \alpha_k \|F_m\|^2)^{-1}$ 
34:       $\bar{a}_{mk} \leftarrow s_{mk}^a \alpha_k (R_k^u + \bar{a}_{mk} F_m) F_m^\top$ 
35:       $R_k^u \leftarrow R_k^u - (\bar{a}_{mk} - \xi) F_m$ 
36:    end for
37:  end parallel for
38:  /****** Update  $q(B)$  *****/
39:  parallel for  $k = 1, \dots, K$  do
40:    for  $n = 1, \dots, N$  do
41:       $\xi \leftarrow \bar{b}_{nk}$ 
42:       $s_{nk}^b \leftarrow (\varphi_k + \beta_k \|G_n\|^2)^{-1}$ 
43:       $\bar{b}_{nk} \leftarrow s_{nk}^b \beta_k (R_k^v + \bar{b}_{nk} G_n) G_n^\top$ 
44:       $R_k^v \leftarrow R_k^v - (\bar{b}_{nk} - \xi) G_n$ 
45:    end for
46:  end parallel for
47:  /****** Update hyperparameters *****/
48:   $\tau \leftarrow |\Omega| / \sum_{(i,j) \in \Omega} (R_{ij}^2 + W_{ij})$ 
49:  for  $k = 1, \dots, K$  do
50:     $\alpha_k \leftarrow \frac{I}{\sum_{i=1}^I \left[ (\bar{u}_{ki} - \bar{a}_k^\top \mathbf{f}_i)^2 + s_{ki}^u + \sum_{m=1}^M s_{mk}^a f_{mi}^2 \right]}$ 
51:     $\beta_k \leftarrow \frac{J}{\sum_{j=1}^J \left[ (\bar{v}_{kj} - \bar{b}_k^\top \mathbf{g}_j)^2 + s_{kj}^v + \sum_{n=1}^N s_{nk}^b g_{nj}^2 \right]}$ 
52:     $\phi_k \leftarrow M / \left( \sum_{m=1}^M \bar{a}_{mk}^2 + s_{mk}^a \right)$ 
53:     $\varphi_k \leftarrow N / \left( \sum_{n=1}^N \bar{b}_{nk}^2 + s_{nk}^b \right)$ 
54:  end for
55: end for

```

stricter than matrix-wise independence assumption, hence there will be an loss in prediction accuracy. However our experimental results confirmed that the element-wise independence assumption does not sacrifice the prediction accuracy. We empirically observed that variational posterior covariances in (13) are almost diagonal, so it looks like that information losses due to the element-wise independent assumption is not critical, although we cannot theoretically guarantee it. In fully observed case, there is an interesting theoretical results [20] but it cannot be applied in the presence of missing values.

3 RELATED WORK

3.1 Matrix Factorization with Side Information

We briefly review several matrix factorization models which incorporates side information. *Matrix co-factorization* is one of an efficient ways to exploit the side information. Matrix co-factorization jointly decomposes

multiple data matrices, where each decomposition is coupled by sharing some factor matrices [21, 22]. 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 [23] and the variational method [24, 25].

There are two different approaches to incorporate the side information with regression in Bayesian matrix factorization framework [26, 17]. In the first approach [26], 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:

$$\begin{aligned} \mathbf{X} &= \left[\mathbf{U}^\top \mathbf{A}^\top \mathbf{F}^\top \right] \left[\mathbf{V}^\top \mathbf{G}^\top \mathbf{B}^\top \right]^\top + \mathbf{E}^x \\ &= \mathbf{U}^\top \mathbf{V} + \mathbf{A}^\top \mathbf{G} + \mathbf{F}^\top \mathbf{B} + \mathbf{E}^x. \end{aligned}$$

On the other hand, *Regression-based Latent Factor Model* (RLFM) [17] assumes that the user and item factor matrices are generated from the side information via the linear regression,

$$\mathbf{U} = \mathbf{A}^\top \mathbf{F} + \mathbf{E}^u, \quad \mathbf{V} = \mathbf{B}^\top \mathbf{G} + \mathbf{E}^v, \quad (23)$$

and then rating matrix \mathbf{X} is generated by (1). Recently RLMF is generalized to hierarchical Bayesian matrix factorization with side information [18].

It was hard to directly compare the effectiveness of previous three different modelings for side information, because each model was trained by different inference algorithm such as Markov Chain Monte Carlo [26], Monte Carlo Expectation Maximization [17], and variational method [24, 25, 18]. The experimental results in [18], where all discussed models are trained by variational method for fair comparisons, showed that the third approach is the most promising way to incorporate side information in both cold and warm start situation. Note that our work also belongs to this third approach. All discussed methods assumed that the side information is given in form of the predefined feature vector. Recently a general functional matrix factorization which automatically constructs the feature function has been proposed [27].

3.2 Optimization

In the perspective of optimization techniques, our work is closely related with [8], where similar idea is used to reduce the time complexity of ALS algorithm. The trick that

reuses the residuals on the observation and efficient updating of them by (18), is a key point of reducing the time complexity from quadratic to linear and it was presented in [8]. Note that computational complexity for optimizing the objective function of weighed trace norm regularized matrix factorization (2) and the variational lower-bound from the fully factorized variational distribution (16) via coordinate ascent method is similar:

$$O(4|\Omega|K) \text{ vs. } O(6|\Omega|K).$$

However the benefit from the extra small effort is crucial: more accurate model is obtainable without tuning of regularization parameter.

Similar to our VBMFSI-CA, the learning algorithm for VBMF in [13], which is referred to as VBMF-GA in this paper, also considers the fully factorized variational distribution. VBMF-GA updates variational posterior variances $\{\mathbf{S}^u, \mathbf{S}^v\}$ and hyperparameters $\{\tau, \alpha, \beta\}$ by fixed-point rules as in VBMFSI-CA, however updates variational posterior means $\{\bar{\mathbf{U}}, \bar{\mathbf{V}}\}$ by scaled gradient ascent method,

$$\begin{aligned} \bar{u}_{ki} &\leftarrow \bar{u}_{ki} + \eta \left(\frac{\partial^2 \mathcal{F}}{\partial^2 \bar{u}_{ki}} \right)^{-\gamma} \frac{\partial \mathcal{C}}{\partial \bar{u}_{ki}}, \\ \bar{v}_{kj} &\leftarrow \bar{v}_{kj} + \eta \left(\frac{\partial^2 \mathcal{F}}{\partial^2 \bar{v}_{kj}} \right)^{-\gamma} \frac{\partial \mathcal{C}}{\partial \bar{v}_{kj}}, \end{aligned} \quad (24)$$

where \mathcal{F} is the variational lower-bound, η is a learning rate, and diagonal part of the Hessian matrix is used for scaling with control parameter γ . The parameter γ allows the learning algorithm to vary from standard gradient ascent ($\gamma = 0$) to diagonal Newton's method ($\gamma = 1$). VBMF-GA also takes linear space and time complexity but it requires tuning of hyperparameters η and γ . In addition, empirical results showed that learning speed is slower than our VBMFSI-CA.

4 NUMERICAL EXPERIMENTS

4.1 Datasets and Side Information

We performed experiments on four large-scale datasets (MovieLens-10M¹, Douban², Netflix³, and Yahoo! Music⁴) with various side information (binary implicit feedback, neighborhood, social network, movie contents, and music taxonomy). Table 2 presents the statistics of these four datasets.

Binary implicit feedback: One of the lessons learnt from the Netflix Prize competition is the importance of integrating implicit feedback into models [28, 29]. For a dataset

¹<http://www.grouplens.org/datasets/movielens>

²<https://www.cse.cuhk.edu.hk/irwin.king/pub/data/douban>

³<http://www.netflixprize.com>

⁴<http://kddcup.yahoo.com>

Table 2: Statistics of datasets.

Dataset	$ \Omega $	I	J	$\frac{ \Omega }{IJ}$ (%)
MovieLens	10M	69,878	10,677	1.34
Douban	17M	129,490	58,541	0.22
Netflix	100M	480,189	17,770	1.18
Yahoo!	263M	1,000,990	624,961	0.04

such as the Netflix, the most natural choice for implicit feedback would be movie rental history, which tells us about user preferences without requiring them to explicitly provide their ratings. For other datasets, browsing or purchase history could be used as implicit feedback. Although such data is not available to us for experiments, we alternatively used the binary indicator matrix which indicates whether the user rates the movie as implicit feedback. In our experiments, we used the normalized binary implicit feedback as the user side information $\mathbf{F} \in \mathbb{R}^{M \times I}$ defined as

$$f_{mi} = \begin{cases} 1/\sqrt{|\Omega_i|} & \text{if } (i, m) \in \Omega, \\ 0 & \text{otherwise,} \end{cases} \quad (25)$$

where the dimension of feature M is equal to the number of item N , and each feature vector \mathbf{f}_i is normalized to unit length. If we incorporate the normalized binary implicit feedback feature (25) with our user generative model (5), then user i is modeled as

$$\mathbf{u}_i = \mathbf{A}^\top \mathbf{f}_i + \mathbf{e}_i^u = \frac{\sum_{j \in \Omega_i} A_j^\top}{\sqrt{|\Omega_i|}} + \mathbf{e}_i^u,$$

where $A_j \in \mathbb{R}^{1 \times K}$ is the j^{th} row vector of \mathbf{A} . Note that this user modeling is equivalent to the one used in SVD++.

Neighborhood: The neighborhood-based approach is one of the primary method for recommendation system. In [29], the item-based neighborhood model was integrated with matrix factorization and showed good prediction accuracy. In our experiments, the item-based neighborhood information is incorporated by defining the item feature matrix $\mathbf{G} \in \mathbb{R}^{N \times J}$ as

$$g_{nj} = \begin{cases} c_{nj}/Z_j & \text{if } c_{nj} \geq \theta, \\ 0 & \text{otherwise,} \end{cases} \quad (26)$$

where dimension of feature N is equal to the number of item J , c_{nj} is the similarity between two items n and j , Z_j is the appropriate normalizer which makes unit length \mathbf{g}_j , θ is the threshold for determining the neighborhood, and Pearson correlation coefficient is used as the similarity criterion [30]. With neighborhood feature (26), item j is modeled as

$$\mathbf{v}_j = \mathbf{B}^\top \mathbf{g}_j + \mathbf{e}_j^v = \sum_{n \in \Upsilon_j} g_{nj} \mathbf{B}_n^\top + \mathbf{e}_j^v,$$

where Υ_j is the set of indices n for which $g_{nj} > 0$. Note g_{nj} has the positive value only if item n is the neighborhood of item j .

Social network: Based on the intuition that users' social relations can be employed to enhance the recommender system, several social recommendation methods have been proposed [31, 32, 33, 34, 35]. In our experiments, the social network information is incorporated by defining the user feature matrix $\mathbf{F} \in \mathbb{R}^{M \times I}$ as

$$f_{mi} = \begin{cases} c_{mi}/Z_i & \text{if } (m, i) \in \Gamma \text{ and } c_{mi} \geq \theta, \\ 0 & \text{otherwise,} \end{cases} \quad (27)$$

where dimension of feature M is equal to the number of user I , $(m, i) \in \Gamma$ denotes that user m is the friend of user i , c_{mi} is the similarity between two users m and i , Z_i is the appropriate normalizer which makes \mathbf{f}_i to unit length.

Movie content: MovieLens-10M dataset provides the genre information of the movies. We extend the content information by collecting the director, actor, and tag information from the IMDB. The extended content information consists of 19 genres, 1710 directors, 35779 actors, and 5294 tags. This movie content information is encoded in binary valued vector with dimension $N = 42802$ and then normalized to unit length.

Music taxonomy: A distinctive feature of the Yahoo! Music dataset is that user ratings are given to entities of four different types: *tracks*, *albums*, *artists*, and *genres*. The majority of items are tracks (81.15%), followed by albums (14.23%), artists (4.46%), and genres (0.16%) [36]. In addition, items are tied together within a taxonomy. That is, for a track we know the identity of its album, performing artist and associated genres. Similarity we have artist and genre annotation for the albums. There is no genre information for artists, because artists may switch between many genres in their career. This taxonomy information is encoded in normalized binary valued vector with dimension $N = 117789$, which equals to the sum of the number of albums, artists, and genres.

4.2 Learning Speed

We compare our learning algorithm for VBMF with side information, which is referred to as VBMFSI-CA in this paper, with three learning algorithms for VBMF.

VBMF-CA: This is a reduced version of VBMFSI-CA, where side information is not incorporated.

VBMF-BCA: The variational lower bound from the matrix-wisely factorized variational distributions is optimized by block coordinate ascent method [2]. The time complexity of VBMF-BCA is cubic with respect to K .

VBMF-GA: The variational lower bound (16) is optimized by gradient ascent method [13]. The time complexity of VBMF-GA is linear with respect to K . We followed the same adaptation rule for learning rate η and setting of the scaling parameter $\gamma = 2/3$ presented in [13] for fair comparison.

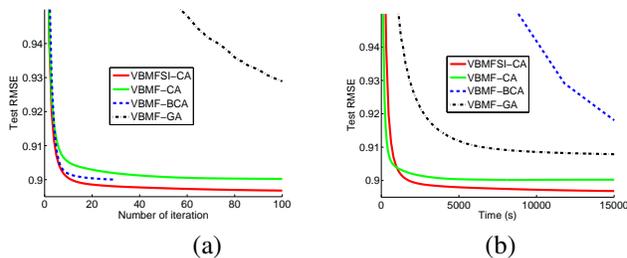


Figure 2: Comparison between VBMFSI-CA, VBMF-CA, VBMF-BCA, VBMF-GA on the Netflix dataset with $K = 50$: (a) Number of iteration vs. RMSE, (b) Time vs. RMSE.

We do not compare the existing learning algorithm for VBMF with side information proposed in [18], which uses similar updating rules described in Table 1, because it can't handle the side information in the form of high dimensional feature vector (e.g. $M, N \gg 10000$). We used a quad-core Intel(R) core(TM) i7-3820 @ 3.6GHz processor and 64GB memory for the comparison. All algorithms were implemented in Matlab2011a, where main computational modules are implemented in C++ as mex files and parallelized with the OpenMP library.⁵

For the Netflix dataset, we trained the VBMF model with four algorithms, where rank K was set to 50. In the case of VBMFSI-CA, the binary implicit feedback was used for side information. Fig. 2-(a) shows the comparison of the number of iteration versus test RMSE for the four learning algorithms. In the perspective of the number of iteration, VBMFSI-CA is the best, VBMF-GA is the worst and VBMF-BCA (RMSE=0.9000) is slightly better than VBMF-CA (RMSE=0.9002). However the time per iteration of each algorithm is quite different: 158 seconds, 70 seconds, 66 minutes, and 19 seconds for VBMFSI-CA, VBMF-CD, VBMF-BCA and VBMF-GA respectively. Hence we also compare the four learning algorithm in the perspective of training time as shown in Fig. 2-(b). After factoring the time per iteration, VBMF-GA converges faster than VBMF-BCA but it usually finds bad local solution. Our VBMFSI-CA and VBMF-CA clearly outperforms VBMF-GA and VBMF-BCA both in the learning speed and the prediction accuracy. Especially VBMFSI-CA finds the most accurate model by incorporating the binary implicit feedback information and its additional computational cost is acceptable because it is linear in the number of nonzero elements in feature matrix, that is almost equal to the number of observed ratings ($\|\mathbf{F}\|_0 \approx |\Omega|$).

4.3 Effect of Side Information

We trained VBMFSI-CA for each dataset, where $K = 50$ and the maximum number of training iteration was set to 100. We performed 5-fold cross-validation for Douban

Table 3: Test RMSE and time per iteration on different datasets with various kinds of side information.

MovieLens 10M, r_a			
	RMSE	% Improvement	Time (s)
none	0.8582	0.00%	5.3
C	0.8553	0.34%	5.5
B	0.8455	1.48%	8.2
B, C	0.8444	1.61%	8.4
Douban, 5-fold cross-validation			
	RMSE	% Improvement	Time (s)
none	0.6913	0.00%	17
S	0.6894	0.30%	18
B	0.6856	0.84%	24
Netflix, probe10			
	RMSE	% Improvement	Time (s)
none	0.9002	0.00%	70
N	0.9000	0.02%	80
B	0.8968	0.38%	158
B, N	0.8944	0.64%	167
Yahoo! Music, validation (not test)			
	RMSE	% Improvement	Time (s)
none	22.0313	0.00%	237
T	21.8808	0.68%	242
B	21.7045	1.48%	524
B, T	21.5673	1.87%	531

dataset and used pre-defined train/test split for others. As shown in Table 3, significant improvement is achieved when side information is included, where B, N, S, C, and T denote the binary implicit feedback, neighborhood, social network, movie contents, and music taxonomy respectively. Especially, the binary implicit feedback gave the most effective result.

5 CONCLUSIONS

We have presented a scalable algorithm for VBMF with side information, the time complexity of which is linear in K to update variational parameters in element-wise fashion. Incorporating side information also requires only linear time complexity with respect to the number of nonzero elements in feature vectors provided by side information. In addition, the proposed algorithm is easily parallelized on multi-core systems. Experimental results on large-scale datasets with various side information confirmed the useful behavior of our algorithm such as scalability, fast learning, and prediction accuracy.

Acknowledgements: This work was supported by National Research Foundation (NRF) of Korea (NRF-2013R1A2A2A01067464) and POSTECH Rising Star Program.

⁵Source code: <http://mlg.postech.ac.kr/~karma13>

References

- [1] R. M. Bell and Y. Koren. Scalable collaborative filtering with jointly derived neighborhood interpolation weights. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, Omaha, NE, 2007.
- [2] 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.
- [3] R. Salakhutdinov and A. Mnih. Bayesian probabilistic matrix factorization using MCMC. In *Proceedings of the International Conference on Machine Learning (ICML)*, Helsinki, Finland, 2008.
- [4] Y.-D. Kim and S. Choi. Weighted nonnegative matrix factorization. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Taipei, Taiwan, 2009.
- [5] 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.
- [6] 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.
- [7] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [8] I. Pilászy, D. Zibriczky, and D. Tikk. Fast ALS-based matrix factorization for explicit and implicit feedback datasets. In *Proceedings of the ACM International Conference on Recommender Systems (RecSys)*, Barcelona, Spain, 2010.
- [9] H.-F. Yu, C.-J. Hsieh, S. Si, and I. Dhillon. Scalable coordinate descent approaches to parallel matrix factorization for recommender systems. In *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, Brussels, Belgium, 2012.
- [10] R. Gemulla, P. J. Hass, E. Nijkamp, and Y. Sismanis. Large-scale matrix factorization with distributed stochastic gradient descent. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, San Diego, CA, USA, 2011.
- [11] Y. Zhuang, W. Chin, Y. Juan, and C. Lin. A fast parallel SGD for matrix factorization in shared memory systems. In *Proceedings of the ACM International Conference on Recommender Systems (RecSys)*, Hong Kong, China, 2013.
- [12] S. Schelter, C. Boden, M. Schenck, A. Alexandrov, and V. Markl. Distributed matrix factorization with mapreduce using a series of broadcast-joins. In *Proceedings of the ACM International Conference on Recommender Systems (RecSys)*, Hong Kong, China, 2013.
- [13] 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.
- [14] D. Stern, R. Herbrich, and R. Herbrich. Matchbox: Large scale online Bayesian recommendations. In *Proceedings of the International Conference on World Wide Web (WWW)*, pages 111–120, 2009.
- [15] 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.
- [16] N. Koenigstein and U. Paquet. Xbox movies recommendations: Variational Bayes matrix factorization with embedded feature selection. In *Proceedings of the ACM International Conference on Recommender Systems (RecSys)*, Hong Kong, China, 2013.
- [17] 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.
- [18] S. Park, Y.-D. Kim, and S. Choi. Hierarchical Bayesian matrix factorization with side information. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, Beijing, China, 2013.
- [19] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, 1999.
- [20] S. Nakajima, M. Sugiyama, and D. Babacan. Global solution of fully-observed variational bayesian matrix factorization is column-wise independent. In *Advances in Neural Information Processing Systems (NIPS)*, volume 24, 2011.
- [21] 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.
- [22] G. Bouchard, S. Guo, and D. Yin. Convex collective matrix factorization. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, Scottsdale, AZ, USA, 2013.

- [23] 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.
- [24] 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.
- [25] 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.
- [26] 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.
- [27] T. Chen, H. Li, Q. Yang, and Y. Yu. General functional matrix factorization using gradient boosting. In *Proceedings of the International Conference on Machine Learning (ICML)*, Atlanta, Georgia, USA, 2013.
- [28] R. Salakhutdinov, A. Mnih, and G. Hinton. Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the International Conference on Machine Learning (ICML)*, Corvallis, OR, USA, 2007.
- [29] Y. Koren. Factorization meets the neighborhood: A multifaceted collaborative filtering model. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, Las Vegas, Nevada, USA, 2008.
- [30] J. S. Breese, D. Heckerman, and C. Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, Madison, WI, 1998.
- [31] H. Ma, H. Yang, M. R. Lyu, and I. King. SoRec: Social recommendation using probabilistic matrix factorization. In *Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM)*, Napa Valley, CA, 2008.
- [32] H. Ma, I. King, and M. R. Lyu. Learning to recommend with social trust ensemble. In *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, Boston, MA, 2009.
- [33] H. Ma, D. Zhou, C. Liu, M. R. Lyu, and I. King. Recommender systems with social regularization. In *Proceedings of the ACM International Conference on Web Search and Data Mining (WSDM)*, Hong Kong, 2010.
- [34] M. Jamali and M. Ester. TrustWalker: A random walk model for combining trust-based and item-based recommendation. In *Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, Paris, France, 2009.
- [35] M. Jamali and M. Ester. A matrix factorization technique with trust propagation for recommendation in social networks. In *Proceedings of the ACM International Conference on Recommender Systems (RecSys)*, Barcelona, Spain, 2010.
- [36] G. Dror, N. Koenigstein, Y. Koren, and M. Weimer. The yahoo! music dataset and KDD-Cup'11. In *Proceedings of KDD Cup and Workshop*, 2011.