

Attribute-aware Collaborative Filtering: Survey and Classification

Recommendation system

Goal: discover user preferences over unrated items based on current user preferences

Current user preferences can be interpreted in three ways:

- Pointwise: preserve exact numerical rating that a user gives to an item
- Pairwise: preserve any pairs of item preferences of a particular user
- Listwise: preserve the relative order of all item ratings of a user as a list

Pointwise approach seeks to preserve item **rating** while pairwise and listwise seek to preserve item **ranking**, which does not require the recommendation system to produce ratings for each item

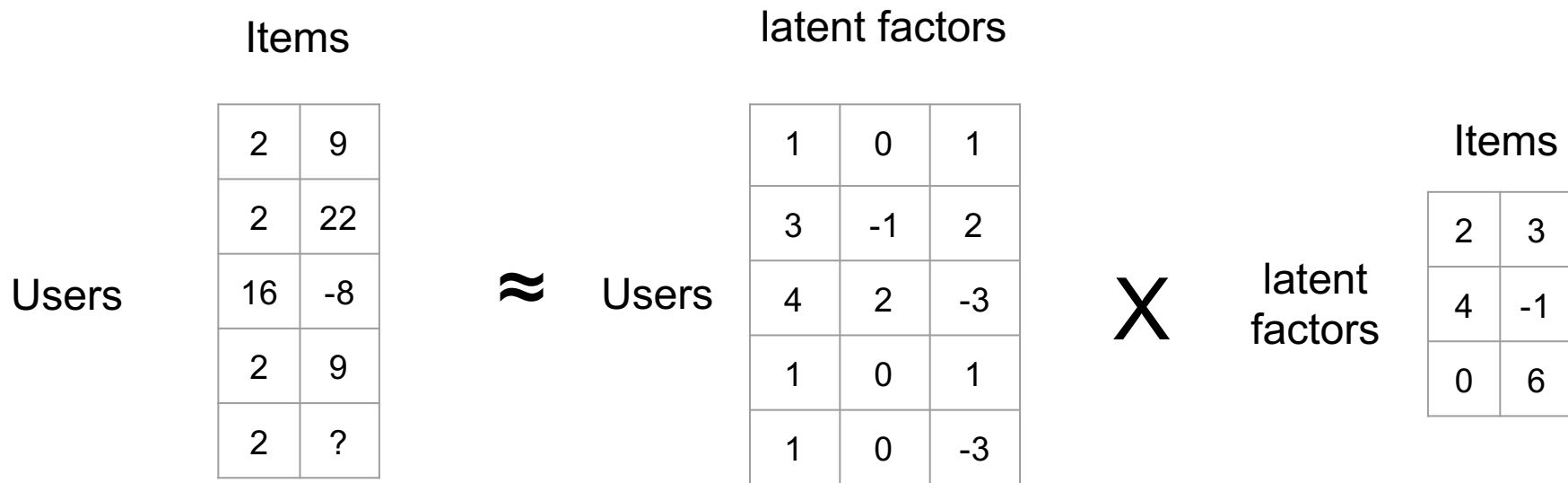
Rating type

- Explicit feedback: real value ratings e.g. $[0, 5]$, $[0, 100]$
- Implicit feedback: binary ratings e.g. $\{0, 1\}$, $\{-1, 1\}$
 - e.g. whether a user clicks a webpage, where 0 indicates not clicked and 1 indicates clicked
 - Sigmoid function can be used to constrain real value predictions into $[0, 1]$

$$\text{sig}(x) = \frac{1}{1+\exp(-x)}$$

- Negative sampling can be applied when there is only positive examples in training data

Matrix Factorization - Toy example



1. Train user and item latent factors with observed ratings
 2. Predict unobserved ratings with trained user and item factors
- e.g. $? = \langle 1, 0, -3 \rangle * \langle 3, -1, 6 \rangle = -15$

Matrix factorization - Formal definition

The goal of MF can be seen as finding user and item latent factors which minimize the square error between predicted ratings and true ratings:

$$\mathbf{W}^*, \mathbf{H}^* = \operatorname{argmin}_{\mathbf{W}, \mathbf{H}} \sum_{(u,i) \in \delta(R)} \frac{1}{2} (r_{ui} - \mathbf{w}_u^\top \mathbf{h}_i)^2 + \frac{\lambda_W}{2} \sum_{u=1}^{N_u} \|\mathbf{w}_u\|_2^2 + \frac{\lambda_H}{2} \sum_{i=1}^{N_i} \|\mathbf{h}_i\|_2^2$$

Regularizations are added to prevent overfit of latent factors

Relation of MF to CF

- MF gives each user and item a latent factor which represents its underlying pattern
- Users with similar latent factors would have similar ratings on items
- Users with similar ratings on a same set of items are trained to have similar latent factors
- MF assumes rating matrix to be at most rank K
 - Low rank implies some rows or columns in rating matrix are linear dependent to others
 - Ratings of a user might be obtained by linear combinations of ratings from other users

Biased matrix factorization

The drawback of MF is it only considers the interaction between user and item, which neglects the bias of individual user or item. Bias MF improves it by considering them:

$$\hat{r}_{ui} = \mu + c_u + d_i + \mathbf{w}_u^\top \mathbf{h}_i$$

The formal definition then becomes:

$$\begin{aligned} \mathbf{W}^*, \mathbf{H}^*, \mathbf{c}^*, \mathbf{d}^*, \mu^* = \operatorname{argmin}_{\mathbf{W}, \mathbf{H}, \mathbf{c}, \mathbf{d}, \mu} & \sum_{(u,i) \in \delta(R)} (r_{ui} - \mu - c_u - d_i - \mathbf{w}_u^\top \mathbf{h}_i)^2 \\ & + \lambda (\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2 + \|\mathbf{c}\|_2^2 + \|\mathbf{d}\|_2^2) \end{aligned}$$

Probabilistic Matrix Factorization (1/2)

- An equivalent probabilistic form of MF
- Assume each user and item follows a zero-mean spherical multivariate Gaussian distribution

$$p(\mathbf{W} \mid \sigma_W^2) = \prod_{u=1}^{N_u} \mathcal{N}(\mathbf{w}_u \mid \mathbf{0}, \sigma_W^2 \mathbf{I}), \quad p(\mathbf{H} \mid \sigma_H^2) = \prod_{i=1}^{N_i} \mathcal{N}(\mathbf{h}_i \mid \mathbf{0}, \sigma_H^2 \mathbf{I})$$

Conditional probability over observed ratings:

$$p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \sigma^2) = \prod_{(i,j) \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} \mid \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2)$$

Probabilistic Matrix Factorization (2/2)

Maximum a posteriori (MAP):

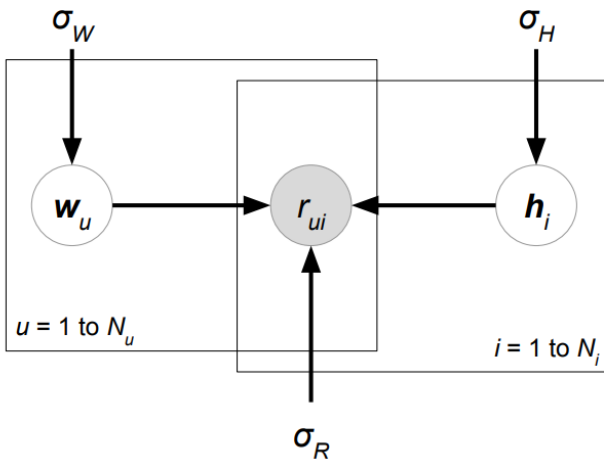
$$\begin{aligned}\log p(\mathbf{W}, \mathbf{H} \mid \mathbf{R}, \sigma_R^2, \sigma_W^2, \sigma_H^2) &= \log p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \sigma_R^2) + \log p(\mathbf{W} \mid \sigma_W^2) + \log p(\mathbf{H} \mid \sigma_H^2) + C \\ &= -\frac{1}{2\sigma_R^2} \sum_{(u,i) \in \delta(\mathbf{R})} (r_{ui} - \mathbf{w}_u^\top \mathbf{h}_i)^2 - \frac{1}{2\sigma_W^2} \sum_{u=1}^{N_u} \mathbf{w}_u^\top \mathbf{w}_u - \frac{1}{2\sigma_H^2} \sum_{i=1}^{N_i} \mathbf{h}_i^\top \mathbf{h}_i \\ &\quad - \frac{1}{2} \left(|\delta(\mathbf{R})| \log \sigma_R^2 + N_u K \log \sigma_W^2 + N_i K \log \sigma_H^2 \right) + C\end{aligned}$$

With Gaussian noise observed, it is equivalent to minimize the following objective function:

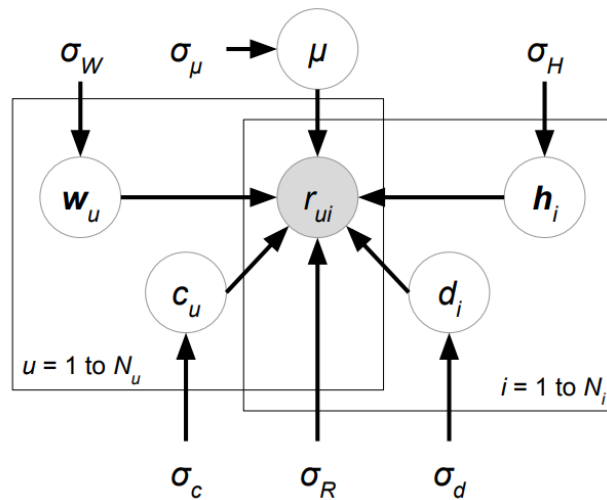
$$\sum_{(u,i) \in \delta(\mathbf{R})} \frac{1}{2} (r_{ui} - \mathbf{w}_u^\top \mathbf{h}_i)^2 + \frac{\lambda_W}{2} \sum_{u=1}^{N_u} \|\mathbf{w}_u\|_2^2 + \frac{\lambda_H}{2} \sum_{i=1}^{N_i} \|\mathbf{h}_i\|_2^2$$

which is the same as MF

Probabilistic MF - Graphical model



(a) PMF



(b) Biased PMF

Recommendation goal - Rating prediction (1/3)

- Pointwise learning
- The objective function of rating prediction recommendation system is often related to Root Mean Square Error (RMSE):

$$\text{RMSE} = \sqrt{\frac{1}{|\delta(\mathbf{R})|} \sum_{(u,i) | r_{ui} \in \delta(\mathbf{R})} (\hat{r}_{ui} - r_{ui})^2}$$

Recommendation goal - Rating prediction (2/3)

Minimize MSE (root is ignored for convenience) or maximize normal likelihood

Explicit feedback:

$$\operatorname{argmin}_{\hat{r}} \sum_{(u,i) | r_{ui} \in \delta(R)} (\hat{r}_{ui} - r_{ui})^2$$

$$\operatorname{argmax}_{\hat{r}} \prod_{(u,i) | r_{ui} \in \delta(R)} \mathcal{N}(r_{ui} \mid \mu = \hat{r}_{ui}, \sigma^2)$$

Implicit feedback:

$$\operatorname{argmin}_{\hat{r}} \sum_{(u,i) | r_{ui} \in \delta(R)} (\operatorname{sig}(\hat{r}_{ui}) - r_{ui})^2$$

$$\operatorname{argmax}_{\hat{r}} \prod_{(u,i) | r_{ui} \in \delta(R)} \mathcal{N}(r_{ui} \mid \mu = \operatorname{sig}(\hat{r}_{ui}), \sigma^2)$$

Recommendation goal - Rating prediction (3/3)

If ratings are $\{0, 1\}$, the problem can be formulated as a binary classification problem, where logistic regression (Bernoulli likelihood) can be applied:

$$\begin{aligned} \operatorname{argmax}_{\hat{r}} & \prod_{(u,i) | 1=r_{ui} \in \delta(R)} \Pr(\hat{r}_{ui} = 1) \prod_{(u,i) | 0=r_{ui} \in \delta(R)} \Pr(\hat{r}_{ui} = 0) \\ = & \underbrace{\prod_{(u,i) | 1=r_{ui} \in \delta(R)} \operatorname{sig}(\hat{r}_{ui})}_{\text{Positive set}} \underbrace{\prod_{(u,i) | 0=r_{ui} \in \delta(R)} (1 - \operatorname{sig}(\hat{r}_{ui}))}_{\text{Negative set}}. \end{aligned}$$

Recommendation goal - Item ranking (1/2)

- Pairwise, listwise learning
- The objective function is to preserve user preferences over items

$$\begin{aligned} \operatorname{argmax}_{\hat{r}} \quad & \prod_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(R), \\ r_{ui} > r_{uj}}} \Pr(\hat{r}_{ui} > \hat{r}_{uj}) \\ = \quad & \prod_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(R), \\ r_{ui} > r_{uj}}} \operatorname{sig}(\hat{r}_{ui} - \hat{r}_{uj}) \end{aligned}$$

Recommendation goal - Item ranking (2/2)

Most item ranking recommendation systems have objective function related to:

- AUC

$$\text{AUC} = \frac{1}{T} \sum_{(u,i,j) | \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(R), \\ r_{ui} > r_{uj}}} \mathbb{I}(\hat{r}_{ui} > \hat{r}_{uj})$$

- Hinge loss

$$\underset{\hat{r}}{\operatorname{argmin}} \sum_{(u,i,j) | \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(R), \\ r_{ui} > r_{uj}}} \max \{0, \hat{r}_{uj} - \hat{r}_{ui}\}$$

Item ranking - Bayesian Personalized Ranking

Bayesian Personalized Ranking (BPR)

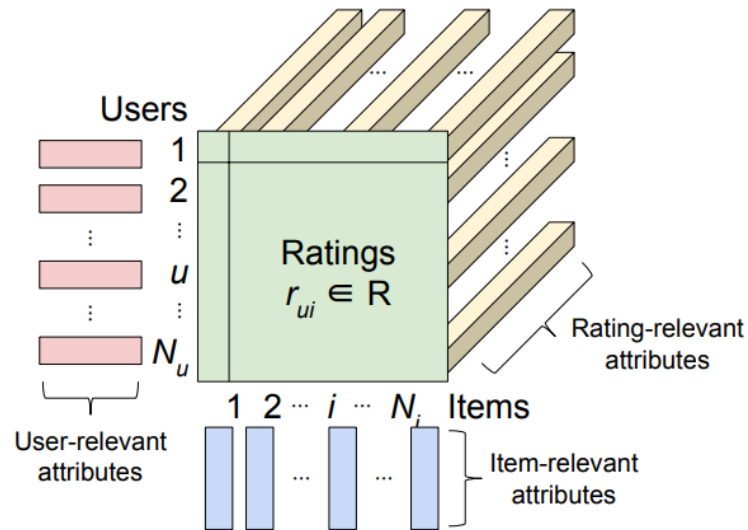
- Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence (UAI '09). AUA Press, Arlington, Virginia, United States, 452-461.
- The objective function of BPR is to maximize AUC

$$\begin{aligned}\operatorname{argmax}_{\hat{r}} \text{AUC} &= \sum_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(\mathbf{R}), \\ r_{ui} > r_{uj}}} \mathbb{I}(\hat{r}_{ui} - \hat{r}_{uj} > 0) \\ &\approx \sum_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(\mathbf{R}), \\ r_{ui} - r_{uj} > 0}} \operatorname{sig}(\hat{r}_{ui} - \hat{r}_{uj}) \\ &= \sum_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(\mathbf{R}), \\ r_{ui} - r_{uj} > 0}} \log \operatorname{sig}(\hat{r}_{ui} - \hat{r}_{uj}) \\ &= \log \prod_{(u,i,j) \mid \substack{\{r_{ui}, r_{uj}\} \subseteq \delta(\mathbf{R}), \\ r_{ui} - r_{uj} > 0}} \operatorname{sig}(\hat{r}_{ui} - \hat{r}_{uj})\end{aligned}$$

What are attributes?

Attribute source:

- User: age, job, gender etc.
- Item: location, price etc.
- Rating: upvotes, downvotes etc.



We refer to user/item attributes as **side information** and rating attributes as **context**

Attribute type

- Numerical
 - Real value, integer
 - e.g. Age, # of likes
- Categorical
 - Hard to be digitalized
 - e.g. Types of food, words

Attribute type - Numerical

Three common uses of numerical attributes

- Attributes fit latent factors

$$\operatorname{argmin}_{\theta, W} \|f_{\theta}(X) - W\| \text{ or } W = f_{\theta}(X) \text{ for user-relevant attributes}$$

- Latent factors fit attributes

$$\operatorname{argmin}_{\theta, W} \|f_{\theta}(W) - X\| \text{ or } X = f_{\theta}(W) \text{ for user-relevant attributes}$$

- Independent of latent factors

$$\operatorname{argmin}_{\theta, W, H} \|f_{\theta}(X) + W^{\top} H - R\|$$

Attribute type - Categorical

Categorical attributes are inappropriate to be directly transformed into numerical attributes because there is no ordering

=> Each dimension of attributes can be transformed into an **one-hot encoding vector** first and then be treated the same as numerical attributes

Attribute transformation (1/2)

Despite the sources of attributes are different (side information and context), they are actually interchangeable

- Side information can be convert to context by simply concatenate the original context with the user and item side information

$$\mathbf{z}'_{\pi(u,i)} = \begin{bmatrix} \mathbf{z}_{\pi(u,i)} \\ \mathbf{x}_u \\ \mathbf{y}_i \end{bmatrix} \in \mathbb{R}^{K_Z + K_X + K_Y}$$

Attribute transformation (2/2)

- Context can be converted to side information by appending all context of users/items to its original side information
 - If a particular feature is absent, the value should be all zeros. This is to ensure the dimension of user/item attributes to remain the same, as this is a requirement of most recommendation systems

$$\mathbf{x}'_u = \left[\mathbf{x}_u^\top \quad \mathbf{z}_{\pi(u,1)}^\top \quad \mathbf{z}_{\pi(u,2)}^\top \quad \cdots \quad \mathbf{z}_{\pi(u,i)}^\top \quad \cdots \quad \mathbf{z}_{\pi(u,N_i)}^\top \right]^\top \in \mathbb{R}^{K_X + K_Z N_i}$$

Why using attributes?

- More constraints can be added
 - Users with similar attributes should have similar latent factors
- Similarity between users or items can be pre-computed and further utilized
- Cold-start problem
 - Sometimes users or items have few observed ratings
 - Additional attributes are used as auxiliary information

How to use attributes?

- Discriminative matrix factorization (DMF)
- Generative matrix factorization
- Generalized factorization
- Heterogeneous graph

DMF - framework

Main idea: attributes model latent factors

$$\begin{aligned}\arg\max_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{W}, \mathbf{H} \mid \mathbf{R}, \mathbf{X})}_{\text{Posterior}} &= \arg\max_{\mathbf{W}, \mathbf{H}} \frac{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \mathbf{X}) p(\mathbf{W}, \mathbf{H} \mid \mathbf{X})}{p(\mathbf{R} \mid \mathbf{X})} \\ &= \arg\max_{\mathbf{W}, \mathbf{H}} p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \mathbf{X}) p(\mathbf{W}, \mathbf{H} \mid \mathbf{X}) \\ &= \arg\max_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \mathbf{X})}_{\text{Likelihood}} \underbrace{p(\mathbf{W} \mid \mathbf{X}) p(\mathbf{H} \mid \mathbf{X})}_{\text{Prior}}\end{aligned}$$

DMF - Linear model

Weight vectors are applied on attributes (like linear regression)

$$\operatorname{argmax}_{\mathbf{W}, \mathbf{H}, \theta} \prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} | \mu_R = \mathbf{w}_u^\top \mathbf{h}_i + \mathbf{a}_u^\top \mathbf{x}_u + \mathbf{b}_i^\top \mathbf{y}_i + \mathbf{c}^\top \mathbf{z}_{\pi(u,i)}, \sigma_R^2) p(\mathbf{W} | \mathbf{X}) p(\mathbf{H} | \mathbf{Y})$$

- Pros: flexible; can be applied when some resources of attributes are absent, e.g. only user attributes are given
- Cons: user attributes and item attributes are considered separately

Example: Bayesian Matrix Factorization with Side Information (BMFSI)

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{W}, \mathbf{H}, \theta} \underbrace{p(\mathbf{R} | \mathbf{W}, \mathbf{H}, \theta)}_{\text{Likelihood}} \underbrace{p(\mathbf{W}) p(\mathbf{H})}_{\text{Priors}} \\ &= \operatorname{argmax}_{\mathbf{W}, \mathbf{H}, \theta} \underbrace{\prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} | \mathbf{w}_u^\top \mathbf{h}_i + \mathbf{a}_u^\top \mathbf{x}_u + \mathbf{b}_i^\top \mathbf{y}_i, \sigma_R^2)}_{\text{Matrix factorization using attributes}} \underbrace{\prod_u \mathcal{N}(\mathbf{w}_u | \boldsymbol{\mu}_u, \boldsymbol{\Sigma}_u) \prod_i \mathcal{N}(\mathbf{h}_i | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}_{\text{Regularization}} \end{aligned}$$

DMF - Bilinear model (1/2)

In addition to weight vectors, a matrix is presented to model attributes interactions (usually between user and item)

$$\operatorname{argmax}_{\mathbf{W}, \mathbf{H}, \theta} \prod_{(u, i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} | \mu_R = \boxed{\mathbf{x}_u^\top \mathbf{A} \mathbf{y}_i} + \mathbf{c}_u^\top \mathbf{x}_u + \mathbf{d}_i^\top \mathbf{y}_i + \mathbf{b}, \sigma_R^2) p(\mathbf{W} | \mathbf{X}) p(\mathbf{H} | \mathbf{Y})$$

The likelihood can be reformed to the following:

$$\mu_R = \mathbf{x}_u^\top \mathbf{A} \mathbf{y}_i + \mathbf{c}_u^\top \mathbf{x}_u + \mathbf{d}_i^\top \mathbf{y}_i + \mathbf{b} = \tilde{\mathbf{x}}_u^\top \tilde{\mathbf{A}} \tilde{\mathbf{y}}_i$$

which can be seen as learning the matrix only

DMF - Bilinear model (2/2)

Some bilinear models are implicit, for example:

learning $w_u^\top h_i$ where $w_u = Sx_u$ and $h_i = Ty_i$ is the same as learning

$$x_u^\top (S^\top T) y_i \text{ where } A = S^\top T$$

- Pros: interactions between user attributes and item attributes are captured
- Cons: requires at least two sources of attributes, otherwise it is reduced to linear model

DMF - Bilinear model (3/3)

Example: Regression-based Latent Factor Model (RLFM)

$$\begin{aligned}
 & \underset{\mathbf{W}, \mathbf{H}, \theta}{\operatorname{argmax}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \mathbf{c}, \mathbf{d}, \gamma, \mathbf{Z})}_{\text{Likelihood}} \underbrace{p(\mathbf{W} \mid \mathbf{A}, \mathbf{X})p(\mathbf{H} \mid \mathbf{B}, \mathbf{Y})p(\mathbf{c} \mid \boldsymbol{\alpha}, \mathbf{X})p(\mathbf{d} \mid \boldsymbol{\beta}, \mathbf{Y})}_{\text{Prior}} \\
 &= \underset{\mathbf{W}, \mathbf{H}, \theta}{\operatorname{argmax}} \underbrace{\prod_{(u,i) \mid r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} \mid \mathbf{w}_u^\top \mathbf{h}_i + c_u + d_i + \gamma^\top \mathbf{z}_{\pi(u,i)}, \sigma_R^2)}_{\text{Matrix factorization using attributes}} \\
 & \quad \underbrace{\prod_u \mathcal{N}(\mathbf{w}_u \mid \mathbf{A} \mathbf{x}_u, \boldsymbol{\Sigma}_W) \mathcal{N}(c_u \mid \boldsymbol{\alpha}^\top \mathbf{x}_u, \sigma_c^2) \prod_i \mathcal{N}(\mathbf{h}_i \mid \mathbf{B} \mathbf{y}_i, \boldsymbol{\Sigma}_H) \mathcal{N}(d_i \mid \boldsymbol{\beta}^\top \mathbf{y}_i, \sigma_d^2)}_{\text{Regularization using attributes}}
 \end{aligned}$$

Deepak Agarwal and Bee-Chung Chen. 2009. Regression-based latent factor models. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining* (KDD '09)

DMF - Similarity Matrix (1/2)

A matrix is pre-computed to capture similarities between users or items.

The usage of the matrix varies. For example, $S_{ij} ||W_i - W_j||_2^2$ regularizes user/item latent factors based on their attribute similarity

- Pros:
 - Human knowledge involves in collaborative filtering (explicitly defines how similarity should be measured)
 - Various choices of similarity measure, e.g. kernel methods
- Cons:
 - High space complexity for similarity matrix (quadratic to number of users/items)

DMF - Similarity Matrix (2/2)

Example: Kernelized Probabilistic Matrix Factorization (KPMF)

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H})}_{\text{Likelihood}} \underbrace{p(\mathbf{W} \mid \mathbf{X})p(\mathbf{H} \mid \mathbf{Y})}_{\text{Prior}} \\ &= \operatorname{argmax}_{\mathbf{W}, \mathbf{H}} \underbrace{\prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} \mid \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2)}_{\text{Matrix factorization}} \underbrace{\prod_k \mathcal{N}(\mathbf{w}^k \mid \mathbf{0}, \mathbf{S}_\mathbf{X}) \prod_l \mathcal{N}(\mathbf{h}^l \mid \mathbf{0}, \mathbf{S}_\mathbf{Y})}_{\text{Regularization using attributes}}. \end{aligned}$$

Kernelized Probabilistic Matrix Factorization: Exploiting Graphs and Side Information Tinghui Zhou, Hanhuai Shan, Arindam Banerjee, and Guillermo Sapiro Proceedings of the 2012 SIAM International Conference on Data Mining

How to use attributes?

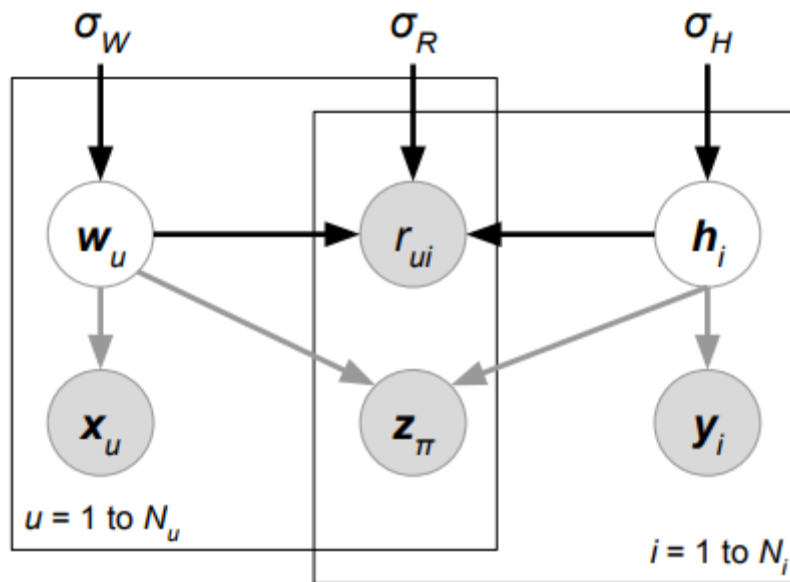
- Discriminative matrix factorization
- Generative matrix factorization (GMF)
- Generalized factorization
- Heterogeneous graph

GMF - framework

Main idea: distribution of attributes are modeled by latent factors

$$\begin{aligned}\arg\max_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{W}, \mathbf{H} \mid \mathbf{R}, \mathbf{X})}_{\text{Posterior}} &= \arg\max_{\mathbf{W}, \mathbf{H}} \frac{p(\mathbf{R}, \mathbf{X} \mid \mathbf{W}, \mathbf{H}) p(\mathbf{W}, \mathbf{H})}{p(\mathbf{R}, \mathbf{X})} \\ &= \arg\max_{\mathbf{W}, \mathbf{H}} p(\mathbf{R}, \mathbf{X} \mid \mathbf{W}, \mathbf{H}) p(\mathbf{W}, \mathbf{H}) \\ &= \arg\max_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}) p(\mathbf{X} \mid \mathbf{W}, \mathbf{H})}_{\text{Likelihood}} \underbrace{p(\mathbf{W}) p(\mathbf{H})}_{\text{Prior}}.\end{aligned}$$

GMF - graphical model



GMF - Matrix Factorization (1/2)

Attribute matrices are factorized by latent matrices

$$\mathbf{R} \approx \mathbf{W}^\top \mathbf{H} \qquad \mathbf{X} \approx \mathbf{A}^\top \mathbf{W}, \mathbf{Y} \approx \mathbf{B}^\top \mathbf{H}$$

$$\begin{aligned} \operatorname{argmax}_{\mathbf{W}, \mathbf{H}, \mathbf{A}, \mathbf{B}} \quad & \prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} | \mu_R = \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2) \prod_{(j,u)} \mathcal{N}(x_{ju} | \mathbf{a}_j^\top \mathbf{w}_u, \sigma_X^2) \prod_{(v,i)} \mathcal{N}(y_{vi} | \mathbf{b}_v^\top \mathbf{h}_i, \sigma_Y^2) \\ & \prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(z_{ui} | \mathbf{w}_u^\top \mathbf{C} \mathbf{h}_i, \sigma_Z^2) p(\mathbf{W}) p(\mathbf{H}), \end{aligned}$$

- Pros: can be easily extended when multiple attribute sources are available
- Cons: difficult to train (jointly learn rating and attribute distribution)

GMF - Matrix Factorization (2/2)

Example: Collective Matrix Factorization (CMF)

$$\begin{aligned}
 & \underset{\mathbf{W}, \mathbf{H}, \mathbf{A}, \mathbf{B}}{\operatorname{argmax}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}) p(\mathbf{X} \mid \mathbf{W}, \mathbf{A}) p(\mathbf{Y} \mid \mathbf{H}, \mathbf{B})}_{\text{Likelihood}} \underbrace{p(\mathbf{W}) p(\mathbf{H}) p(\mathbf{A}) p(\mathbf{B})}_{\text{Prior}} \\
 &= \underset{\mathbf{W}, \mathbf{H}, \mathbf{A}, \mathbf{B}}{\operatorname{argmax}} \underbrace{\prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} \mid \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2)}_{\text{Matrix factorization of } R} \underbrace{\prod_{(j,u)} \mathcal{N}(x_{ju} \mid \mathbf{a}_j^\top \mathbf{w}_u, \sigma_X^2)}_{\text{Matrix factorization of } X} \underbrace{\prod_{(v,i)} \mathcal{N}(y_{vi} \mid \mathbf{b}_v^\top \mathbf{h}_i, \sigma_Y^2)}_{\text{Matrix factorization of } Y} \\
 & \quad \underbrace{\prod_u \mathcal{N}(\mathbf{w}_u \mid \mathbf{0}, \Sigma_W) \prod_i \mathcal{N}(\mathbf{h}_i \mid \mathbf{0}, \Sigma_H) \prod_j \mathcal{N}(\mathbf{a}_j \mid \mathbf{0}, \Sigma_A) \prod_v \mathcal{N}(\mathbf{b}_v \mid \mathbf{0}, \Sigma_B)}_{\text{Regularization}}
 \end{aligned}$$

Ajit P. Singh and Geoffrey J. Gordon. 2008. Relational learning via collective matrix factorization. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining* (KDD '08)

GMF - DNN (1/2)

Use autoencoder to model attribute distribution and generate latent factors

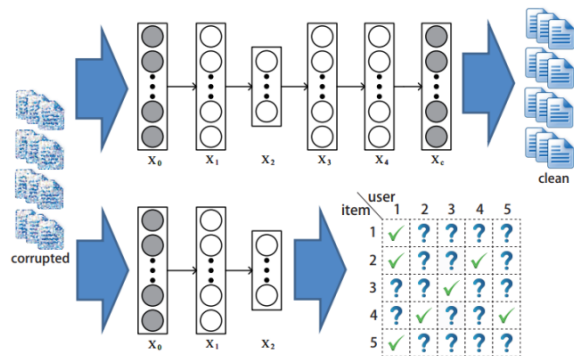
$$\begin{aligned}\operatorname{argmax}_{\mathbf{W}, \mathbf{H}} p(\mathbf{W}, \mathbf{H} \mid \mathbf{R}, \mathbf{X}, \tilde{\mathbf{X}}) &= \operatorname{argmax}_{\mathbf{W}, \mathbf{H}} \frac{p(\mathbf{R}, \mathbf{X} \mid \mathbf{W}, \mathbf{H}, \tilde{\mathbf{X}}) p(\mathbf{W}, \mathbf{H} \mid \tilde{\mathbf{X}})}{p(\mathbf{R}, \mathbf{X} \mid \tilde{\mathbf{X}})} \\ &= \operatorname{argmax}_{\mathbf{W}, \mathbf{H}} p(\mathbf{R}, \mathbf{X} \mid \mathbf{W}, \mathbf{H}, \tilde{\mathbf{X}}) p(\mathbf{W}, \mathbf{H} \mid \tilde{\mathbf{X}}) \\ &= \operatorname{argmax}_{\mathbf{W}, \mathbf{H}} \underbrace{p(\mathbf{R} \mid \mathbf{W}, \mathbf{H}, \tilde{\mathbf{X}}) p(\mathbf{X} \mid \mathbf{W}, \mathbf{H}, \tilde{\mathbf{X}})}_{\text{Likelihood}} \underbrace{p(\mathbf{H} \mid \tilde{\mathbf{X}}) p(\mathbf{W} \mid \tilde{\mathbf{X}})}_{\text{Priors, i.e., Encoder } \mathcal{E}}\end{aligned}$$

- Pros:
 - Non-linear mapping is available (with appropriate activation function)
 - Stacked denoising autoencoder makes it more robust to attribute corruption
- Cons:
 - Much more parameters to be tuned

GMF - DNN (2/2)

Example: Collaborative Deep Learning (CDL)

$$\begin{aligned}
 & \underset{\mathbf{W}, \mathbf{H}, \boldsymbol{\theta}, \phi}{\operatorname{argmax}} \underbrace{p(\mathbf{R}, | \mathbf{W}, \mathbf{H})}_{\text{Likelihood}} \underbrace{p(\mathbf{Y} | \mathbf{H}, \tilde{\mathbf{Y}}) p(\mathbf{H} | \tilde{\mathbf{Y}}) p(\mathbf{W})}_{\text{Prior}} \\
 &= \underset{\mathbf{W}, \mathbf{H}, \boldsymbol{\theta}, \phi}{\operatorname{argmax}} \underbrace{\prod_{(u,i) | r_{ui} \in \delta(\mathbf{R})} \mathcal{N}(r_{ui} | \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2)}_{\text{Matrix factorization}} \underbrace{\prod_i \mathcal{N}(\mathbf{y}_i | \mathcal{D}_\phi(\mathcal{E}_\theta(\tilde{\mathbf{y}}_i)), \boldsymbol{\Sigma}_Y) \mathcal{N}(\mathbf{h}_i | \mathcal{E}_\theta(\tilde{\mathbf{y}}_i), \boldsymbol{\Sigma}_H)}_{\text{Stacked denoising auto-encoder for } \mathbf{Y}} \underbrace{\prod_u \mathcal{N}(\mathbf{w}_u | \mathbf{0}, \boldsymbol{\Sigma}_W)}_{\text{Regularization}}.
 \end{aligned}$$



- The embeddings of autoencoder regularize item latent factors, i.e. $\mathcal{N}(\mathbf{h}_i | \mathcal{E}_\theta(\tilde{\mathbf{y}}_i), \boldsymbol{\Sigma}_H)$
- User latent factors are generated from Gaussian distribution
- Rating prediction is generated by product of user and item latent factors, i.e. $\mathcal{N}(r_{ui} | \mathbf{w}_u^\top \mathbf{h}_i, \sigma_R^2)$
- The model can be easily extended if user attributes are given

Comparison of DMF and GMF

	DMF	GMF
Learning target	Rating	Rating, Attributes
Training difficulty	Less difficult	More difficult
Number of models	Many	Few

How to use attributes?

- Discriminative matrix factorization
- Generative matrix factorization
- Generalized factorization (GF)
- Heterogeneous graph

GF - framework (1/2)

Main idea: a more general form of DMF

1. Consider user/item as a kind of categorical attribute (one-hot encoding vector)
2. All attributes (including user, item) can be concatenated as a vector x
 - e.g. given a $\langle u, i, r \rangle$ pair and let E_u, E_i be the one-hot encoding of user u and item i respectively, the concatenated vector x is $\langle E_u, E_i, X_u, Y_i, Z_r \rangle$
3. The vector is then a typical training example of traditional classification or regression problem, with rating as the label

$$\operatorname{argmax}_w \prod_{r \in \delta(\mathbf{R})} \mathcal{N} \left(r \mid \mu_R = \sum_{d=D_m}^{D_M} \sum_{j_1=1}^N \sum_{j_2=j_1+1}^N \cdots \sum_{j_d=j_{d-1}+1}^N w_{j_1 j_2 \dots j_d} (x_{j_1} x_{j_2} \dots x_{j_d}), \sigma_R^2 \right)$$

GF - framework (2/2)

Due to large number of parameters in w , which often results in overfitting, w is often a function of latent factors

$$w_{j_1 j_2 \dots j_d} = f_d(\mathbf{v}_{j_1}, \mathbf{v}_{j_2}, \dots, \mathbf{v}_{j_d})$$

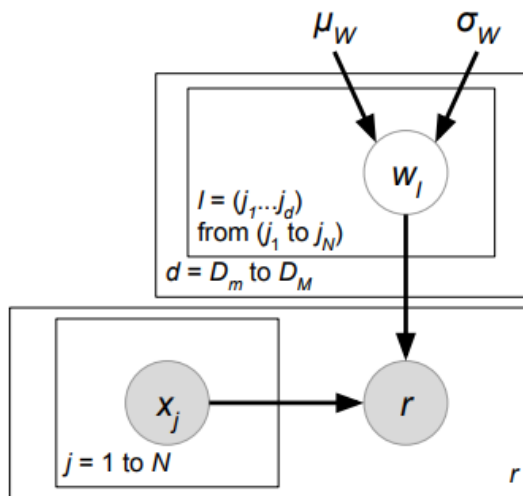
where \mathbf{v}_j represents the latent factor for attribute x_j

Matrix factorization is a special case of GF:

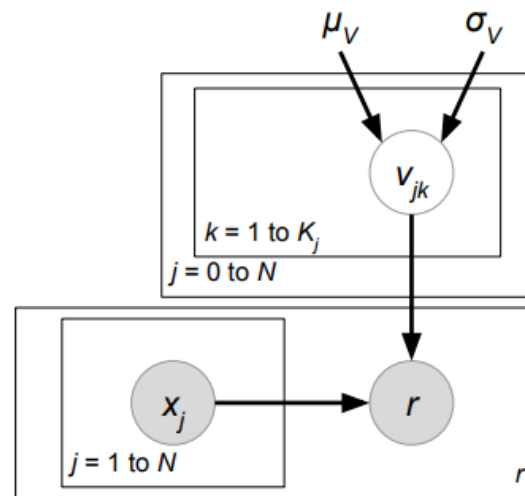
$$\operatorname{argmax}_{\mathbf{v}} \prod_{r_{ui} \in \delta(\mathbf{R})} \mathcal{N} \left(\hat{r}_{ui} \mid \mu_R = \sum_{j_1=1}^N \sum_{j_2=j_1+1}^N \mathbf{v}_{j_1}^\top \mathbf{v}_{j_2} (x_{j_1} x_{j_2}) = \mathbf{v}_u^\top \mathbf{v}_{N_u+i}, \sigma_R^2 \right)$$

where \mathbf{x} is the concatenation of user and item one-hot encoding vector

GF - graphical model



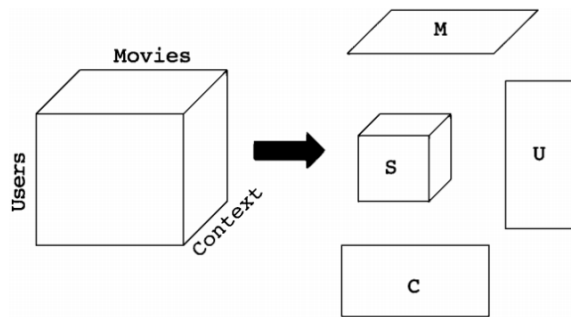
(a) w -weighted generalization



(b) v -approximate generalization

GF - Tensor Factorization

An N-dimensional rating tensor is factorized into a core tensor and several matrices, where each dimension is (user, item, attribute1, attribute2, ...)



e.g. 3-dimension tensor (user, item, attr1)

$$U \in \mathbb{R}^{n \times d_U}, M \in \mathbb{R}^{m \times d_M} \text{ and } C \in \mathbb{R}^{c \times d_C}$$

$$S \in \mathbb{R}^{d_U \times d_M \times d_C}$$

Prediction: $F_{ijk} = S \times_U U_{i*} \times_M M_{j*} \times_C C_{k*}$

- Pros: highly extendable
- Cons:
 - High space complexity for core tensor (exponential to number of attributes)
 - Attributes are limited to be categorical

Tensor multiplication:

$$T = Y \times_U U$$

$$\Rightarrow T_{ljk} = \sum_{i=1}^n Y_{ijk} U_{ij}$$

Alexandros Karatzoglou, Xavier Amatriain, Linas Baltrunas, and Nuria Oliver. 2010. Multiverse recommendation: n-dimensional tensor factorization for context-aware collaborative filtering. In *Proceedings of the fourth ACM conference on Recommender systems* (RecSys '10)

GF - Factorization Machine

The objective function of FM can be simplified from GF framework, which is

$$\begin{aligned}\mu_R &= \underbrace{w_0}_{d=0} + \underbrace{\sum_{l=1}^N w_l x_l}_{d=1} + \underbrace{\sum_{j_1=1}^N \sum_{j_2=j_1+1}^N w_{j_1 j_2} (x_{j_1} x_{j_2})}_{d=2} \\ &= w_0 + \sum_{l=1}^N w_l x_l + \sum_{j_1=1}^N \sum_{j_2=j_1+1}^N \mathbf{v}_{j_1}^\top \mathbf{v}_{j_2} (x_{j_1} x_{j_2})\end{aligned}$$

- Pros: attributes can be numerical
- Cons: higher-order attribute interactions (more than 2) is not considered
- Extension:
 - Multi-view machines (consider higher-order interactions)
 - Neural Factorization Machines (consider non-linear attribute interactions)

More on FM

- Lets start with a simple linear regression solution: X = user features \mathbf{f}_u and item features \mathbf{g}_v , Y = ratings:

$$\min_{\mathbf{w}} \sum_{u,v \in R} \left(R_{u,v} - \mathbf{w}^T \begin{bmatrix} \mathbf{f}_u \\ \mathbf{g}_v \end{bmatrix} \right)^2$$

- Let's try degree-2 polynomial mappings to solve:

$$\min_{w_{t,s}, \forall t,s} \sum_{u,v \in R} \left(r_{u,v} - \sum_{t'=1}^U \sum_{s'=1}^V w_{t',s'} (f_u)_{t'} (g_v)_{s'} \right)^2$$

- It is likely to overfit since we have many parameters to learn ($\#Para > \#Ratings$)

More on FM

- Why not using $P^T Q$ to model W , where P and Q are low-rank matrices

$$\min_{u,v \in R} (R_{u,v} - \mathbf{f}_u^T P^T Q \mathbf{g}_v)^2$$

How to use attributes?

- Discriminative matrix factorization
- Generative matrix factorization
- Generalized factorization
- Heterogeneous graph (HG)

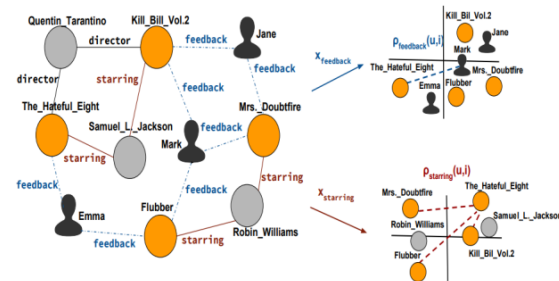
HG

User, item and attributes are represented as nodes in a heterogeneous graph, and the recommendation problem becomes a link prediction problem

- Each user, item, attribute is a node
 - a. since nodes are discrete, attributes must be categorical
- Ratings are weights for links ($\{0, 1\}$ for implicit feedback while real value for explicit feedback)
- Random walk or meta-path algorithms extract similarity from the graph
- MF or other supervised ML algorithms extracts features from similarity information

HG (2/2)

Example: entity2rec



- A knowledge graph with different types of nodes and edges are given
- For each edge type, a subgraph which consists of nodes and edges associated with this edge type only is formed
- For each subgraph, node2vec is applied to extract embeddings
- Similarity can be computed for each $\langle u, i \rangle$ pair, e.g. cosine similarity
- Similarities belonging to a user from different subgraphs are features for any supervised learning-to-rank algorithm (e.g. Adarank) where labels are observed ratings, i.e. $\langle u, i, r \rangle$ pair