Piyush Rai IIT Kanpur

#### Probabilistic Machine Learning (CS772A)

Feb 8, 2016

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



3

<ロ> (日) (日) (日) (日) (日)

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



• Each entry of **X** can be written as  $X_{nm} \approx \boldsymbol{u}_n^\top \boldsymbol{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$ 

3

・ロト ・個ト ・ヨト ・ヨト

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



• Each entry of **X** can be written as  $X_{nm} \approx \boldsymbol{u}_n^\top \boldsymbol{v}_m = \sum u_{nk} v_{mk}$ 

• Note:  $K \ll \min\{M, N\}$ 

▲口> ▲圖> ▲国> ▲国>

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



- Each entry of **X** can be written as  $X_{nm} \approx \boldsymbol{u}_n^\top \boldsymbol{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$
- Note:  $K \ll \min\{M, N\}$
- U:  $N \times K$  row latent factor matrix,  $u_n$ :  $K \times 1$  latent factors of row n

・ロト ・回ト ・ヨト ・ヨト

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



- Each entry of **X** can be written as  $X_{nm} \approx \boldsymbol{u}_n^\top \boldsymbol{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$
- Note:  $K \ll \min\{M, N\}$
- U:  $N \times K$  row latent factor matrix,  $\boldsymbol{u}_n$ :  $K \times 1$  latent factors of row n
- V:  $M \times K$  column latent factor matrix,  $v_m$ :  $K \times 1$  latent factors of column m

・ロン ・四 と ・ ヨ と ・ ヨ と

• Given a matrix **X** of size  $N \times M$ , approximate it via a low-rank decomposition



- Each entry of **X** can be written as  $X_{nm} \approx \boldsymbol{u}_n^\top \boldsymbol{v}_m = \sum_{k=1}^K u_{nk} v_{mk}$
- Note:  $K \ll \min\{M, N\}$
- U:  $N \times K$  row latent factor matrix,  $\boldsymbol{u}_n$ :  $K \times 1$  latent factors of row n
- V:  $M \times K$  column latent factor matrix,  $v_m$ :  $K \times 1$  latent factors of column m
- X may have missing entries

イロン イヨン イヨン イヨン



Some applications:

• Learning embeddings from dyadic/relational data (each matrix entry is a dyad, e.g., user-item rating, document-word count, user-user link, etc.). Thus it also performs dimensionality reduction.



Some applications:

• Learning embeddings from dyadic/relational data (each matrix entry is a dyad, e.g., user-item rating, document-word count, user-user link, etc.). Thus it also performs dimensionality reduction.



Some applications:

• Learning embeddings from dyadic/relational data (each matrix entry is a dyad, e.g., user-item rating, document-word count, user-user link, etc.). Thus it also performs dimensionality reduction.



Some applications:

- Learning embeddings from dyadic/relational data (each matrix entry is a dyad, e.g., user-item rating, document-word count, user-user link, etc.). Thus it also performs dimensionality reduction.
- Matrix Completion, i.e., predicting missing entries in X via the learned embeddings (useful in recommender systems/collaborative filtering - Netflix Prize competition, link prediction in social networks, etc.): X<sub>nm</sub> ≈ u<sub>n</sub><sup>T</sup> v<sub>m</sub>

э

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

# Interpreting the Embeddings

- The embeddings/latent factors/latent features can be given interpretations (e.g., as genres if the matrix **X** represents a user-movie rating matrix case)
- A cartoon illustation of matrix factorization based embeddings (or "generes") learned from a user-movie rating data set using embedding dimension K = 2



Picture courtesy: Matrix Factorization Techniques for Recommender Systems: Koren et al, 2009 🗆 🕨 🖉 🕨 👌 🖉 🔌 😒 🦿 🖉 🔷 🔍

# Interpreting the Embeddings

- The embeddings/latent factors/latent features can be given interpretations (e.g., as genres if the matrix **X** represents a user-movie rating matrix case)
- A cartoon illustation of matrix factorization based embeddings (or "generes") learned from a user-movie rating data set using embedding dimension K = 2



• Similar things (users/movies) get embedded nearby in the embedding space (two things will be deemed similar if their embeddings are similar). Thus useful for computing similarities and/or making recommendations

Picture courtesy: Matrix Factorization Techniques for Recommender Systems: Koren et al, 2009 🗆 🕨 🔄 🖉 🕂 🚊 👘 🚊 👘 🖓 🔍 (^{>}

# Interpreting the Embeddings

• Another illustation of two-dimensional embeddings of movies only



Embedding dimension 1 (or latent factor 1)

• Similar movies get embedded nearby in the embedding space

Picture courtesy: Matrix Factorization Techniques for Recommender Systems: Koren et al, 2009 🗆 🕨 👍 👘 🖌 🚖 🖢 👋 🚊 🖉 🔍 🔇

- Recall our model  $\bm{X}\approx \bm{U}\bm{V}^{\top}$  or  $\bm{X}=\bm{U}\bm{V}^{\top}+\bm{E}$  where  $\bm{E}$  is the noise matrix
- Goal: learn **U** and **V**, given a subset  $\Omega$  of **X** (let's call it **X**<sub> $\Omega$ </sub>)
- Some notations:
  - $\Omega = \{(n, m)\}$ :  $X_{nm}$  is observed
  - $\Omega_{u_n}$ : column indices of observed entries in rows n
  - $\Omega_{\mathbf{v}_m}$ : row indices of observed entries in column m



• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

$$\boldsymbol{u}_n \sim \mathcal{N}(\boldsymbol{u}_n|\boldsymbol{0},\lambda_U^{-1}\boldsymbol{I}_K), \qquad n=1,\ldots,N$$

• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

$$\begin{array}{lll} \boldsymbol{u}_n & \sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & & n = 1, \dots, N \\ \boldsymbol{v}_m & \sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & & m = 1, \dots, M \end{array}$$

• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

$$\begin{array}{lll} \boldsymbol{u}_n & \sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m & \sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} & \sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^{\top} \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

$$\begin{array}{lll} \boldsymbol{u}_n &\sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m &\sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} &\sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^{\top} \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$



• Assuming latent factors  $\boldsymbol{u}_n$ ,  $\boldsymbol{v}_m$  and each matrix entry  $X_{nm}$  to be real-valued

$$\begin{array}{lll} \boldsymbol{u}_n &\sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m &\sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} &\sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$



• This is also equivalent to  $X_{nm} = \boldsymbol{u}_n^\top \boldsymbol{v}_m + \epsilon_{nm}$  where the noise/residual

$$\epsilon_{nm} \sim \mathcal{N}(0,\sigma^2)$$

Probabilistic Machine Learning (CS772A)

イロト イヨト イヨト イヨト

• Our basic model

$$\begin{array}{lll} \boldsymbol{u}_n &\sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m &\sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} &\sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^{\top} \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

2

イロン イ団 とくほと くほとう

• Our basic model

$$\begin{array}{lll} \boldsymbol{u}_n & \sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m & \sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} & \sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

• Note: Many variations possible, e.g., adding row/column biases  $(a_n, b_m)$ , rows/column features  $(\mathbf{X}^U, \mathbf{X}^V)$ ; will not consider those here

$$X_{nm} = \mathcal{N}(X_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m + \boldsymbol{a}_n + \boldsymbol{b}_m + \boldsymbol{\beta}_U^\top \boldsymbol{x}_n^U + \boldsymbol{\beta}_V^\top \boldsymbol{x}_m^V, \sigma^2)$$

< ロ > < 同 > < 回 > < 回 > < 回 > <

• Our basic model

$$\begin{array}{ll} \boldsymbol{u}_n & \sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m & \sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} & \sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

• Note: Many variations possible, e.g., adding row/column biases  $(a_n, b_m)$ , rows/column features  $(\mathbf{X}^U, \mathbf{X}^V)$ ; will not consider those here

$$X_{nm} = \mathcal{N}(X_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m + \boldsymbol{a}_n + \boldsymbol{b}_m + \boldsymbol{\beta}_U^\top \boldsymbol{x}_n^U + \boldsymbol{\beta}_V^\top \boldsymbol{x}_m^V, \sigma^2)$$

• Note: Gaussian assumption on X<sub>nm</sub> may not be appropriate if data is not real-valued, e.g., is binary/counts/ordinal (but it still works well nevertheless)

< ロ > < 同 > < 回 > < 回 > < 回 > <

• Our basic model

$$\begin{array}{lll} \boldsymbol{u}_n &\sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m &\sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{\chi}_{nm} &\sim & \mathcal{N}(\boldsymbol{\chi}_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

• Note: Many variations possible, e.g., adding row/column biases  $(a_n, b_m)$ , rows/column features  $(\mathbf{X}^U, \mathbf{X}^V)$ ; will not consider those here

$$X_{nm} = \mathcal{N}(X_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m + \boldsymbol{a}_n + \boldsymbol{b}_m + \boldsymbol{\beta}_U^\top \boldsymbol{x}_n^U + \boldsymbol{\beta}_V^\top \boldsymbol{x}_m^V, \sigma^2)$$

- Note: Gaussian assumption on X<sub>nm</sub> may not be appropriate if data is not real-valued, e.g., is binary/counts/ordinal (but it still works well nevertheless)
- Likewise, if we want to impose specific constraints on the latent factors (e.g., non-negativity, sparsity, etc.) then Gaussians on  $u_n$ ,  $v_m$  are not appropriate

< ロ > < 同 > < 回 > < 回 > < □ > <

• Our basic model

$$\begin{array}{ll} \boldsymbol{u}_n & \sim & \mathcal{N}(\boldsymbol{u}_n | \boldsymbol{0}, \lambda_U^{-1} \boldsymbol{I}_K), & n = 1, \dots, N \\ \boldsymbol{v}_m & \sim & \mathcal{N}(\boldsymbol{v}_n | \boldsymbol{0}, \lambda_V^{-1} \boldsymbol{I}_K), & m = 1, \dots, M \\ \boldsymbol{X}_{nm} & \sim & \mathcal{N}(\boldsymbol{X}_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m, \sigma^2), & \forall (n, m) \in \Omega \end{array}$$

• Note: Many variations possible, e.g., adding row/column biases  $(a_n, b_m)$ , rows/column features  $(\mathbf{X}^U, \mathbf{X}^V)$ ; will not consider those here

$$X_{nm} = \mathcal{N}(X_{nm} | \boldsymbol{u}_n^\top \boldsymbol{v}_m + \boldsymbol{a}_n + \boldsymbol{b}_m + \boldsymbol{\beta}_U^\top \boldsymbol{x}_n^U + \boldsymbol{\beta}_V^\top \boldsymbol{x}_m^V, \sigma^2)$$

- Note: Gaussian assumption on X<sub>nm</sub> may not be appropriate if data is not real-valued, e.g., is binary/counts/ordinal (but it still works well nevertheless)
- Likewise, if we want to impose specific constraints on the latent factors (e.g., non-negativity, sparsity, etc.) then Gaussians on  $u_n$ ,  $v_m$  are not appropriate
- Here, we will only focus on the Gaussian case (leads to a simple algorithm)

< ロ > < 同 > < 回 > < 回 > < 回 > <

• Let's do MAP estimation (recall, we have priors on the latent factors)

э

イロト イヨト イヨト イヨト

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior log  $p(X_{\Omega}, U, V) = \log p(X_{\Omega}|U, V)p(U)p(V)$  is given by

 $\mathcal{L} = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$ 

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior  $\log p(\mathbf{X}_{\Omega}, \mathbf{U}, \mathbf{V}) = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) p(\mathbf{U}) p(\mathbf{V})$  is given by

$$\mathcal{L} = \log p(\mathbf{X}_{\Omega}|\mathbf{U},\mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$$

$$= \log \prod_{(n,m)\in\Omega} p(X_{nm}|\boldsymbol{u}_n,\boldsymbol{v}_m)$$

イロン イロン イヨン イヨン 三日

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior  $\log p(\mathbf{X}_{\Omega}, \mathbf{U}, \mathbf{V}) = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) p(\mathbf{U}) p(\mathbf{V})$  is given by

$$\mathcal{L} = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$$
$$= \log \prod_{(n,m)\in\Omega} p(X_{nm} | \mathbf{u}_n, \mathbf{v}_m) + \log \prod_{n=1}^{N} p(\mathbf{u}_n)$$

イロン イロン イヨン イヨン 三日

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior log  $p(\mathbf{X}_{\Omega}, \mathbf{U}, \mathbf{V}) = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) p(\mathbf{U}) p(\mathbf{V})$  is given by

$$\mathcal{L} = \log p(\mathbf{X}_{\Omega}|\mathbf{U},\mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$$
  
= 
$$\log \prod_{(n,m)\in\Omega} p(X_{nm}|\boldsymbol{u}_n,\boldsymbol{v}_m) + \log \prod_{n=1}^N p(\boldsymbol{u}_n) + \log \prod_{m=1}^M p(\boldsymbol{v}_m)$$

イロン イロン イヨン イヨン 三日

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior log  $p(\mathbf{X}_{\Omega}, \mathbf{U}, \mathbf{V}) = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) p(\mathbf{U}) p(\mathbf{V})$  is given by

$$\mathcal{L} = \log p(\mathbf{X}_{\Omega}|\mathbf{U},\mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$$
  
= 
$$\log \prod_{(n,m)\in\Omega} p(X_{nm}|\mathbf{u}_n,\mathbf{v}_m) + \log \prod_{n=1}^N p(\mathbf{u}_n) + \log \prod_{m=1}^M p(\mathbf{v}_m)$$

• With Gaussian likelihood and priors, ignoring the constants, we have

$$\mathcal{L} = \sum_{(n,m)\in\Omega} -\frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 - \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 - \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

- Let's do MAP estimation (recall, we have priors on the latent factors)
- Log-posterior log  $p(\mathbf{X}_{\Omega}, \mathbf{U}, \mathbf{V}) = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) p(\mathbf{U}) p(\mathbf{V})$  is given by

$$\mathcal{L} = \log p(\mathbf{X}_{\Omega} | \mathbf{U}, \mathbf{V}) + \log p(\mathbf{U}) + \log p(\mathbf{V})$$
  
= 
$$\log \prod_{(n,m)\in\Omega} p(X_{nm} | \mathbf{u}_n, \mathbf{v}_m) + \log \prod_{n=1}^N p(\mathbf{u}_n) + \log \prod_{m=1}^M p(\mathbf{v}_m)$$

• With Gaussian likelihood and priors, ignoring the constants, we have

$$\mathcal{L} = \sum_{(n,m)\in\Omega} -\frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^{\top} \boldsymbol{v}_m)^2 - \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 - \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• Can solve for row and column latent factors  $\boldsymbol{u}_n, \boldsymbol{v}_m$  in an alternating fashion

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

э

・ロト ・回ト ・ヨト ・ヨト

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For row latent factors  $\boldsymbol{u}_n$  (with all column factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{u}_n} = \sum_{m \in \Omega_{\boldsymbol{u}_n}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_U}{2} \boldsymbol{u}_n^\top \boldsymbol{u}_n$$

イロト イポト イヨト イヨト

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^{\top} \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For row latent factors  $\boldsymbol{u}_n$  (with all column factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{u}_n} = \sum_{m \in \Omega_{\boldsymbol{u}_n}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_U}{2} \boldsymbol{u}_n^\top \boldsymbol{u}_n$$

• Taking derivative w.r.t.  $\boldsymbol{u}_n$  and setting to zero, we get

$$\boldsymbol{u}_n = \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} \boldsymbol{v}_m \boldsymbol{v}_m^\top + \lambda_U \sigma^2 \boldsymbol{I}_K\right)^{-1} \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} X_{nm} \boldsymbol{v}_m\right)$$

イロト イポト イヨト イヨト

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For row latent factors  $\boldsymbol{u}_n$  (with all column factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{u}_n} = \sum_{m \in \Omega_{\boldsymbol{u}_n}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_U}{2} \boldsymbol{u}_n^\top \boldsymbol{u}_n$$

• Taking derivative w.r.t.  $\boldsymbol{u}_n$  and setting to zero, we get

$$\boldsymbol{u}_n = \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} \boldsymbol{v}_m \boldsymbol{v}_m^\top + \lambda_U \sigma^2 \boldsymbol{I}_K\right)^{-1} \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} X_{nm} \boldsymbol{v}_m\right)$$

• Note: with **V** fixed, we can solve for all  $\boldsymbol{u}_n$  (n = 1, ..., N) in parallel

(日) (周) (王) (王)

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

э

イロト イヨト イヨト イヨト

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^{\top} \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For column latent factors  $\boldsymbol{v}_m$  (with all row factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{v}_m} = \sum_{n \in \Omega_{\boldsymbol{v}_m}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_V}{2} \boldsymbol{v}_m^\top \boldsymbol{v}_m$$

(日) (周) (王) (王)

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For column latent factors  $\boldsymbol{v}_m$  (with all row factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{v}_m} = \sum_{n \in \Omega_{\boldsymbol{v}_m}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_V}{2} \boldsymbol{v}_m^\top \boldsymbol{v}_m$$

• Taking derivative w.r.t.  $\boldsymbol{v}_m$  and setting to zero, we get

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{\mathbf{v}_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \sigma^2 \mathbf{I}_K\right)^{-1} \left(\sum_{n \in \Omega_{\mathbf{u}_m}} X_{nm} \mathbf{u}_n\right)$$

(日) (周) (王) (王)

• The (negative) log-posterior

$$\mathcal{L} = \sum_{(n,m)\in\Omega} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^{\top} \boldsymbol{v}_m)^2 + \sum_{n=1}^N \frac{\lambda_U}{2} ||\boldsymbol{u}_n||^2 + \sum_{m=1}^M \frac{\lambda_V}{2} ||\boldsymbol{v}_m||^2$$

• For column latent factors  $\boldsymbol{v}_m$  (with all row factors fixed), the objective will be

$$\mathcal{L}_{\boldsymbol{v}_m} = \sum_{n \in \Omega_{\boldsymbol{v}_m}} \frac{1}{2\sigma^2} (X_{nm} - \boldsymbol{u}_n^\top \boldsymbol{v}_m)^2 + \frac{\lambda_V}{2} \boldsymbol{v}_m^\top \boldsymbol{v}_m$$

• Taking derivative w.r.t.  $\boldsymbol{v}_m$  and setting to zero, we get

$$\boldsymbol{v}_{m} = \left(\sum_{n \in \Omega_{\boldsymbol{v}_{m}}} \boldsymbol{u}_{n} \boldsymbol{u}_{n}^{\top} + \lambda_{V} \sigma^{2} \boldsymbol{I}_{K}\right)^{-1} \left(\sum_{n \in \Omega_{\boldsymbol{u}_{m}}} X_{nm} \boldsymbol{u}_{n}\right)$$

• Note: with **U** fixed, we can solve for all  $\mathbf{v}_m$   $(m = 1, \dots, M)$  in parallel

• Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$ 

イロト イポト イヨト イヨト

- Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$
- Initialize the column latent factors  $\mathbf{v}_1, \ldots, \mathbf{v}_M$  randomly, e.g., from the prior, i.e.,  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \lambda_V^{-1} \mathbf{I}_K)$

イロト 不得 トイヨト イヨト

- Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$
- Initialize the column latent factors  $\mathbf{v}_1, \ldots, \mathbf{v}_M$  randomly, e.g., from the prior, i.e.,  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \lambda_V^{-1} \mathbf{I}_K)$
- Iterate until converge

< ロ > < 同 > < 回 > < 回 > < 回 > <

- Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$
- Initialize the column latent factors  $\mathbf{v}_1, \ldots, \mathbf{v}_M$  randomly, e.g., from the prior, i.e.,  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \lambda_V^{-1} \mathbf{I}_K)$
- Iterate until converge
  - Update each row latent factor  $u_n$ , n = 1, ..., N (can be in parallel)

$$\boldsymbol{u}_n = \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} \boldsymbol{v}_m \boldsymbol{v}_m^\top + \lambda_U \sigma^2 \mathbf{I}_K\right)^{-1} \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} X_{nm} \boldsymbol{v}_m\right)$$

< ロ > < 同 > < 回 > < 回 > < □ > <

- Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$
- Initialize the column latent factors  $\mathbf{v}_1, \ldots, \mathbf{v}_M$  randomly, e.g., from the prior, i.e.,  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \lambda_V^{-1} \mathbf{I}_K)$
- Iterate until converge
  - Update each row latent factor  $u_n$ , n = 1, ..., N (can be in parallel)

$$\boldsymbol{u}_n = \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} \boldsymbol{v}_m \boldsymbol{v}_m^\top + \lambda_U \sigma^2 \boldsymbol{I}_K\right)^{-1} \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} X_{nm} \boldsymbol{v}_m\right)$$

• Update each column latent factor  $v_m$ , m = 1, ..., M (can be in parallel)

$$\mathbf{v}_m = \left(\sum_{n \in \Omega_{\mathbf{v}_m}} \mathbf{u}_n \mathbf{u}_n^\top + \lambda_V \sigma^2 \mathbf{I}_K\right)^{-1} \left(\sum_{n \in \Omega_{\mathbf{u}_m}} X_{nm} \mathbf{u}_n\right)$$

イロト イポト イヨト イヨト 二日

- Input: Partially complete matrix  $\boldsymbol{X}_{\Omega}$
- Initialize the column latent factors  $\mathbf{v}_1, \ldots, \mathbf{v}_M$  randomly, e.g., from the prior, i.e.,  $\mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \lambda_V^{-1} \mathbf{I}_K)$
- Iterate until converge
  - Update each row latent factor  $u_n$ , n = 1, ..., N (can be in parallel)

$$\boldsymbol{u}_n = \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} \boldsymbol{v}_m \boldsymbol{v}_m^\top + \lambda_U \sigma^2 \boldsymbol{I}_K\right)^{-1} \left(\sum_{m \in \Omega_{\boldsymbol{u}_n}} X_{nm} \boldsymbol{v}_m\right)$$

• Update each column latent factor  $v_m$ ,  $m = 1, \dots, M$  (can be in parallel)

$$\boldsymbol{v}_m = \left(\sum_{n \in \Omega_{\boldsymbol{v}_m}} \boldsymbol{u}_n \boldsymbol{u}_n^\top + \lambda_V \sigma^2 \boldsymbol{I}_K\right)^{-1} \left(\sum_{n \in \Omega_{\boldsymbol{u}_m}} X_{nm} \boldsymbol{u}_n\right)$$

• Final prediction for any entry:  $X_{nm} = \boldsymbol{u}_n^\top \boldsymbol{v}_m$ 

Suppose we are solving for the column latent factor  $\mathbf{v}_m$  (with **U** fixed)



...1

・ロト ・回ト ・ヨト ・ヨト

3

Suppose we are solving for the column latent factor  $\mathbf{v}_m$  (with **U** fixed)



イロト イポト イヨト イヨト

Suppose we are solving for the column latent factor  $\mathbf{v}_m$  (with **U** fixed)



Likewise, solving for each row latent factor  $\boldsymbol{u}_n$  is a least-squares regression problem

• A very useful way to think about matrix factorization

イロト イポト イヨト イヨト

- A very useful way to think about matrix factorization
- Can modify the regularized least-squares like objective

$$\arg\min_{\boldsymbol{u}_n}\sum_{m\in\Omega_{\boldsymbol{u}_n}}\frac{1}{2\sigma^2}(\boldsymbol{X}_{nm}-\boldsymbol{u}_n^{\top}\boldsymbol{v}_m)^2+\frac{\lambda_U}{2}\boldsymbol{u}_n^{\top}\boldsymbol{u}_n$$

.. and replace it by other loss functions and regularizers

- A very useful way to think about matrix factorization
- Can modify the regularized least-squares like objective

$$\arg\min_{\boldsymbol{u}_n}\sum_{m\in\Omega_{\boldsymbol{u}_n}}\frac{1}{2\sigma^2}(X_{nm}-\boldsymbol{u}_n^{\top}\boldsymbol{v}_m)^2+\frac{\lambda_U}{2}\boldsymbol{u}_n^{\top}\boldsymbol{u}_n$$

.. and replace it by other loss functions and regularizers

- Can easily extend the model in various ways, e.g.
  - Handle other types of entries in the matrix **X**, e.g., binary, counts, etc. (by changing the loss function or the likelihood function term)
  - Impose constraints on the latent factors (by changing the regularizer or prior on latent factors)

< ロ > < 同 > < 回 > < 回 > < 回 > <