

## Case Study 4: Collaborative Filtering

### Review: Probabilistic Matrix Factorization

Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

February 20<sup>th</sup>, 2014

©Emily Fox 2014

1

## Probabilistic Matrix Factorization (PMF)

### ■ A generative process:

- Pick user factors

$$L_{u_1}, \dots, L_{u_k}$$

$$L_{u_i} \stackrel{\text{iid}}{\sim} N(0, \sigma_u^2)$$

$P(L)$   
prior on  
user factors

- Pick movie factors

$$R_{v_1}, \dots, R_{v_k}$$

$$R_{v_i} \stackrel{\text{iid}}{\sim} N(0, \sigma_v^2)$$

$P(R)$

- For each (user, movie) pair observed:

- Pick rating as  $L_u^T R_v + \text{noise}$

$$\Rightarrow r_{uv} | L_u, R_v \sim N(L_u \cdot R_v, \sigma_r^2)$$

$N(0, \sigma_r^2)$

$P(X | L, R)$   
likelihood

### ■ Joint probability:

$$P(L, R, X) = P(L) P(R) P(X | L, R)$$

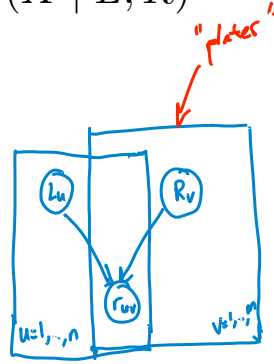
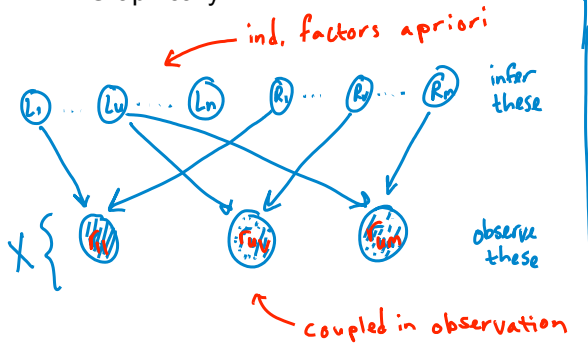
©Emily Fox 2014

2

# PMF Graphical Model

$$P(L, R | X) \propto P(L)P(R)P(X | L, R)$$

- Graphically: *posterior & joint prob.*



©Emily Fox 2014

3

# MAP versus Regularized Least-Squares for Matrix Completion

- MAP under Gaussian Model:

$$\max_{L,R} \log P(L, R | X) =$$

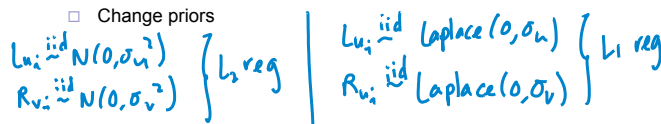
$$-\frac{1}{2\sigma_u^2} \sum_u \sum_i L_{ui}^2 - \frac{1}{2\sigma_v^2} \sum_v \sum_i R_{vi}^2 - \frac{1}{2\sigma_r^2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \text{const}$$

- Least-squares matrix completion with  $L_2$  regularization:

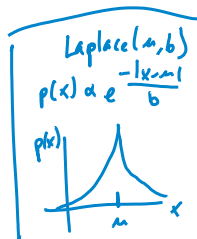
$$\min_{L,R} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2$$

*objectives are equivalent!*

- Understanding as a probabilistic model is very useful! E.g.,



- Incorporate other sources of information or dependencies *very key! many extensions of PMF*



©Emily Fox 2014

4

# Posterior Computations

- MAP estimation focuses on point estimation:

$$\hat{\theta}^{MAP} = \arg \max_{\theta} p(\theta | x)$$

← data  
← parameters

- What if we want a full characterization of the posterior?

- Maintain a measure of uncertainty
- Estimators other than posterior mode (different loss functions)
- Predictive distributions for future observations



$$p(x^{N+1} | x^1, \dots, x^N) = \int p(x^{N+1} | \theta) p(\theta | x^1, \dots, x^N) d\theta$$

← integrate over uncertainty in model params  
← belief about  $\theta$  having seen obs.  $x^1, \dots, x^N$   
← assuming  $x^i$  iid given  $\theta$  (exch.)

Contrast with:

$$p(x^{N+1} | \hat{\theta}^{MAP}(x^1, \dots, x^N)) \leftarrow \text{make pred w/ } \hat{\theta}^{MAP} \text{ after } N \text{ obs.}$$

- Often ~~no closed form characterization~~ (e.g., mixture models, PMF, etc.)
- "plug-in estimator"

# Bayesian PMF Example

- Latent user and movie factors:

$$L_u \sim N(\mu_u, \Sigma_u) \quad u=1, \dots, n$$

$$R_v \sim N(\mu_v, \Sigma_v) \quad v=1, \dots, m$$

- Observations  $r_{uv} \sim N(L_u^T R_v, \sigma_r^2)$

- Hyperparameters:

$$\phi = \{ \mu_u, \Sigma_u, \mu_v, \Sigma_v, \sigma_r^2 \}$$

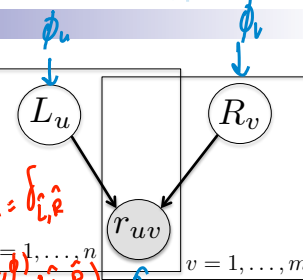
←  $\phi_u$ 
←  $\phi_v$ 
←  $\phi_r$

- Want to predict new movie rating:

$$p(r_{uv}^* | x, \phi) = \int p(r_{uv}^* | L_u, R_v) p(L, R | x, \phi) dL dR$$

↑ new rating
↑ obs. ratings
← integrate over posterior uncertainty in user/movie factors

Full Bayesian approach  
place priors on  $\phi$  as well!



(more general than before)

if post. of  $L, R = \int L, R \Rightarrow p(r_{uv} | x, \phi) = P(r_{uv} | L, R)$

new user/movie combo same as plug-in estimator  $\phi_r$

posterior given obs. so far



# Bayesian PMF Example

$$p(r_{uv}^* | X, \phi) = \int p(r_{uv}^* | L_u, R_v) p(L, R | X, \phi) dL dR$$

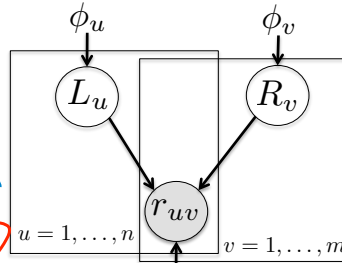
ANALYTICALLY INTRACTABLE!

- Monte Carlo methods:

Approx. as

$$p(r_{uv}^* | X, \phi) \approx \frac{1}{N} \sum_{k=1}^N p(r_{uv}^* | L_u^{(k)}, R_v^{(k)})$$

← samples from posterior  
... how?



- Ideally:  $(L^{(k)}, R^{(k)}) \stackrel{iid}{\sim} p(L, R | X, \phi)$  ← ind samples from posterior

$$p(L, R | X) = \frac{p(X | L, R) p(L) p(R)}{p(X) \int p(X | L, R) p(L) p(R) dL dR}$$

← Again, intractable !!  
← issue!

©Emily Fox 2014

7

# Bayesian PMF Gibbs Sampler

← return dependent

- Outline of Bayesian PMF sampler

1. Init  $L^{(1)}, R^{(1)}$

2. For  $k=1, \dots, N_{iter}$

(i) Sample hyperparams  $\phi^{(k)} = \{\phi_u^{(k)}, \phi_v^{(k)}, \phi_r^{(k)}\}$  ← from cond. posterior

(ii) For each user  $u=1, \dots, n$  sample in parallel

$$L_u^{(k+1)} \sim P(L_u | X, R^{(k)}, \phi^{(k)})$$

(iii) For each movie  $v=1, \dots, m$  sample in parallel

$$R_v^{(k+1)} \sim P(R_v | X, L^{(k+1)}, \phi^{(k)})$$

← just Gaussian dist.

Very similar to ideas of ALS (systematically)

Samples from the posterior eventually

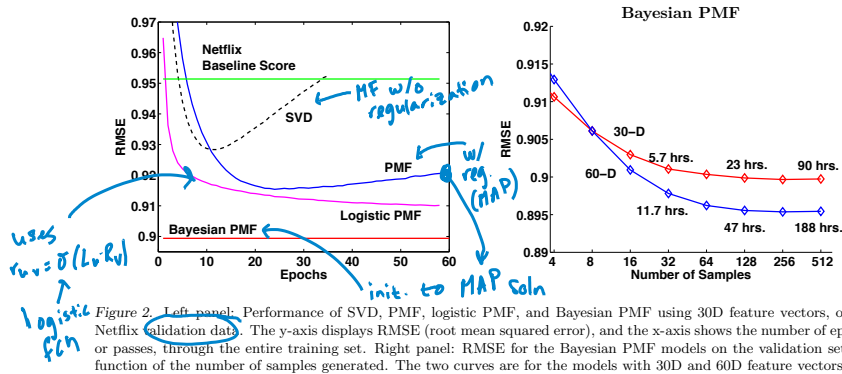
©Emily Fox 2014

8

# Bayesian PMF Results

From Salakhutdinov and Mnih, ICML 2008

- Netflix data with:
  - Training set = 100,480,507 ratings from 480,189 users on 17,770 movie titles
  - Validation set = 1,408,395 ratings.
  - Test set = 2,817,131 user/movie pairs with the ratings withheld.



©Emily Fox 2014

9

# Bayesian PMF Results

From Salakhutdinov and Mnih, ICML 2008

- Bayesian model better controls for overfitting by averaging over possible parameters (instead of committing to one)

D	Valid. RMSE			Test RMSE		
	PMF	BPMF	% Inc.	PMF	BPMF	% Inc.
30	0.9154	0.8994	1.74	0.9188	0.9029	1.73
40	0.9135	0.8968	1.83	0.9170	0.9002	1.83
60	0.9150	0.8954	2.14	0.9185	0.8989	2.13
150	0.9178	0.8931	2.69	0.9211	0.8965	2.67
300	0.9231	0.8920	3.37	0.9265	0.8954	3.36

dim of user/movie factors

Bayesian model improves

Table 1. Performance of Bayesian PMF (BPMF) and linear PMF on Netflix validation and test sets.

Note: Each sampling stage of BPMF requires an  $O(D^3)$  operation, so not for free

©Emily Fox 2014

10

## What you need to know...

- Idea of full posterior inference vs. MAP estimation
- Gibbs sampling as an MCMC approach
- Example of inference in Bayesian probabilistic matrix factorization model

## Case Study 4: Collaborative Filtering

### Matrix Factorization and Probabilistic LFMs for Network Modeling

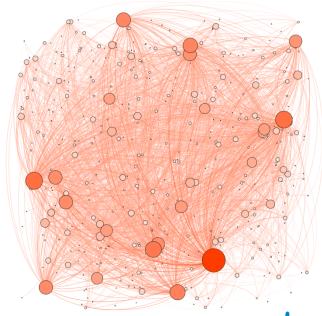
Machine Learning for Big Data  
CSE547/STAT548, University of Washington

Emily Fox

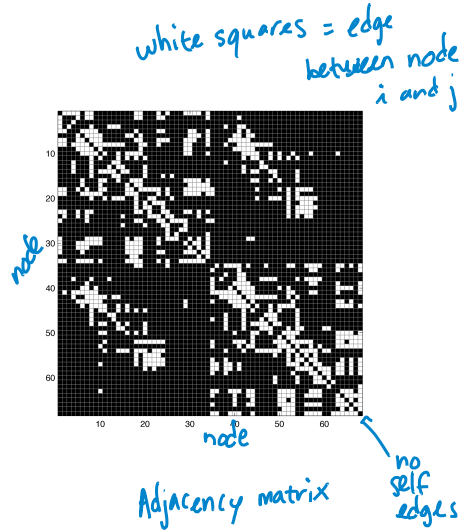
February 20<sup>th</sup>, 2014

# Network Data

- Structure of network data



nodes in a network w/ undirected edges



©Emily Fox 2014

13

# Properties of Data Source

- Similarities to Netflix data:

- Matrix-valued data (adj. matrix)
- High-dimensional many nodes
- Sparse few links between nodes (eg ppl in social network)

- Differences

- Square ← same indices for rows + columns
- Binary ← yes/no for link, though you could have multigraph (multiple link instances bt nodes)

If undirected, then matrix is symmetric



©Emily Fox 2014

14

# Matrix Factorization for Network Data

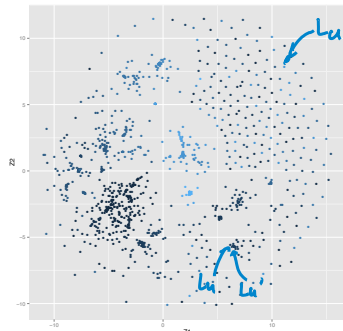
- Vanilla matrix factorization approach:
  - In undirected case, just introduce node (eg user) factors  $L_u$
  - $r_{uv} \approx L_u \cdot L_v \leftarrow$  edge bt users  $u + v$
  - In directed (or asymm) case, introduce sender factors  $L_u$  and receiver factors  $\tilde{L}_v$  (every node/user has both  $L_u$  and  $\tilde{L}_v$ )
  - $r_{uv} \approx L_u \cdot \tilde{L}_v \leftarrow$  edge from user  $u$  to user  $v$
- What to return for link prediction?
  - Is  $r_{uv}$  binary?  $L_u \in \mathbb{R}^k \rightarrow$  no.
  - Many options, but can return top  $\tilde{k}$
  - $r_{uv_1}, \dots, r_{uv_{\tilde{k}}}$  (just use threshold rule)
- Slightly fancier:
  - More appropriate to have  $r_{uv} \in [0, 1]$
  - Use  $r_{uv} \approx \sigma(L_u \cdot L_v)$   $\sigma =$  logistic fcn

©Emily Fox 2014

15

# Probabilistic Latent Space Models

- Assume features (covariates) of the user  $X_u$  or relationship  $X_{uv}$
- Each user has a "position" in a  $k$ -dimensional latent space
  - $L_u \in \mathbb{R}^k$ , just as in matrix factorization
    - unobserved, learn
    - $\leftarrow$  could be  $X_{uv}$
- Probability of link:
 
$$\begin{aligned} & \text{log odds } p(r_{uv}=1 \mid L_u, L_v, X_{uv}, \beta) \\ &= \log \frac{p(r_{uv}=1 \mid L_u, L_v, X_{uv}, \beta)}{p(r_{uv}=0 \mid L_u, L_v, X_{uv}, \beta)} \\ &= \beta_0 + \beta^T X_{uv} - \|L_u - L_v\| \\ &= \beta_0 + \beta^T X_{uv} \quad \text{OR} \\ &= \beta_0 + \beta^T X_{uv} - \|L_u - L_v\| \end{aligned}$$



2D example

©Emily Fox 2014

16



# Probabilistic Latent Space Models

- Probability of link:

$$\text{log odds } p(r_{uv} = 1 \mid L_u, L_v, x_{uv}, \beta) = \beta_0 + \beta^T x_{uv} - |L_u - L_v|$$

$$\text{log odds } p(r_{uv} = 1 \mid L_u, L_v, x_{uv}, \beta) = \beta_0 + \beta^T x_{uv} + |L_u^T L_v|$$

prob. of link is high for  $L_u$  close to  $L_v$

can modify as  $\frac{|L_u^T L_v|}{|L_v|}$

prob. of link is high if  $\angle$  bt  $L_u + L_v$  is small

- Bayesian approach:

- Place prior on user factors and regression coefficients
- Place hyperprior on user factor hyperparameters

- Many other options and extensions (e.g., can use GMM for  $L_u \rightarrow$  clustering of users in the latent space)

©Emily Fox 2014

17

# What you need to know...

- Representation of network data as a matrix
  - Adjacency matrix
- Similarities and differences between adjacency matrices and general matrix-valued data
- Matrix factorization approaches for network data
  - Just use standard MF and threshold output
  - Introduce link functions to constrain predicted values
- Probabilistic latent space models
  - Model link probabilities using distance between latent factors

©Emily Fox 2014

18