

Case Study 4: Collaborative Filtering

Review: Cold Start Problem

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

Emily Fox

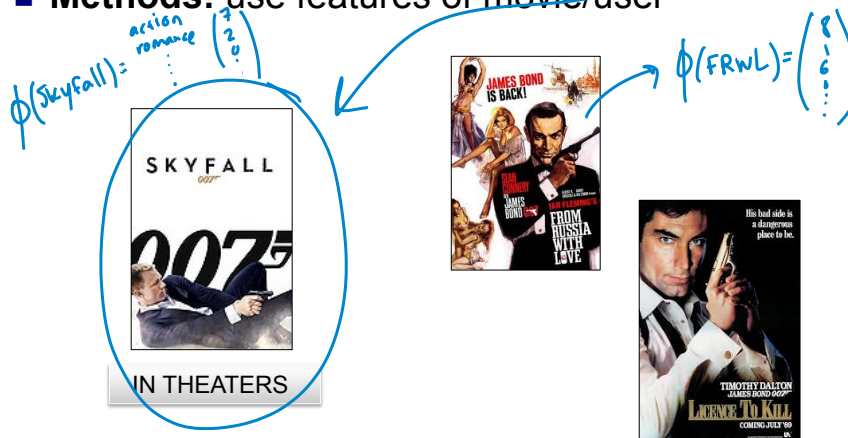
February 18th, 2014

©Emily Fox 2014

1

Cold-Start Problem

- **Challenge:** Cold-start problem (new movie or user)
- **Methods:** use features of movie/user



©Emily Fox 2014

2

Cold-Start Problem More Formally

- Consider a new user u' and predicting that user's ratings

- No previous observations

$r_{u'v} = ? \quad \forall v$ want $L_{u'}$ to predict ratings for this user

- Objective considered so far:

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

- Optimal user factor:

$L_{u'} = 0$ only penalty term in opt. L_u (min $\|L_u\|_F^2$)

- Predicted user ratings: always predict $r_{u'v} = 0 \quad \forall v$... problem

©Emily Fox 2014

3

An Alternative Formulation

- A simpler model for collaborative filtering

- We would not have this issue if we assumed all users were identical

- i.e. all users shared a feature vector w
 - w informed by all ratings + can use it for new user u'

- What about for new movies? What if we had side information?

Create a movie feature vector: \rightarrow synopsis/info

$$\phi(v) = (\text{action}, 1994, \text{Tarantino}, \dots)$$

genre year director \leftarrow dim k

- What dimension should w be?

- Fit linear model: \leftarrow same length as movie feature vector

For all users, u , $r_{uv} \approx w \cdot \phi(v)$ \leftarrow want to learn this

- Minimize: $\min_w \sum_{r_{uv}} (w \cdot \phi(v) - r_{uv})^2 + \lambda_w \|w\|$ \leftarrow LS Lasso

©Emily Fox 2014

4

Personalization

- If we don't have any observations about a user, use wisdom of the crowd
 - Address cold-start problem

$$\text{For user } u', \text{ predict } r_{u'v} \approx w \cdot \phi(v)$$

- Clearly, not all users are the same ... shared w is strong assumption
- Just as in personalized click prediction, consider model with global and user-specific parameters

Consider user-specific deviations w_u from the crowd w

$$r_{uv} \approx (w + w_u) \cdot \phi(v)$$

init. to 0

- As we gain more information about the user, forget the crowd
 w_u more informed

©Emily Fox 2014

5

User Features...

- In addition to movie features, may have information about the user:

$$\phi(u) = (25, F, MS, A^+, \dots)$$

age, gender, education, grade in Big Data

- Combine with features of movie:

$$\phi(u, v) = (\dots, \phi(u), \dots, \dots, \phi(v), \dots, \dots)$$

cross features...

- Unified linear model:

$$r_{uv} \approx (w + w_u) \cdot \phi(u, v)$$

©Emily Fox 2014

6

Feature-based Approach versus Matrix Factorization

- Feature-based approach:
 - Feature representation of user and movies fixed
 - Can address cold-start problem

- Matrix factorization approach:
 - Suffers from cold-start problem
 - User & movie features are learned from data

- A unified model:

©Emily Fox 2014

7

Unified Collaborative Filtering via SGD

$$\min_{L,R,w,\{w_u\}_u} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v + (w + w_u) \cdot \phi(u, v) - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2 + \frac{\lambda_w}{2} \|w\|_2^2 + \frac{\lambda_{wu}}{2} \sum_u \|w_u\|_2^2$$

- Gradient step observing r_{uv}
 - For L,R
$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$
 - For w and w_u :

©Emily Fox 2014

8

What you need to know...

- Cold-start problem
- Feature-based methods for collaborative filtering
 - Help address cold-start problem
- Unified approach

Case Study 4: Collaborative Filtering

Connections with Probabilistic Matrix Factorization

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

Emily Fox

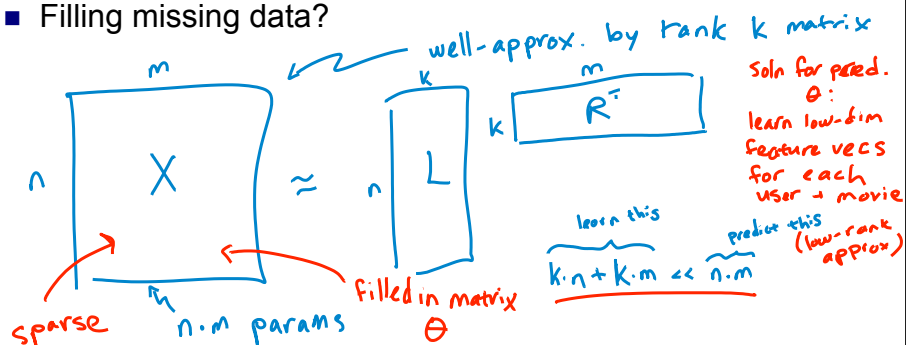
February 18th, 2014

Matrix Completion Problem



X_{ij} known for black cells
 X_{ij} unknown for white cells
 Rows index users
 Columns index movies

- Filling missing data?



©Emily Fox 2014

11

Coordinate Descent for Matrix Factorization: Alternating Least-Squares

$$\min_{L,R} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$

- Fix movie factors, optimize for user factors

□ Independent least-squares over users

$$\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|$$

- Fix user factors, optimize for movie factors

□ Independent least-squares over movies

$$\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda_v \|R\|$$

- System may be underdetermined: use regularization

- Converges to local optima

©Emily Fox 2014

12

Probabilistic Matrix Factorization (PMF)

- A generative process:
 - Pick user factors
 - Pick movie factors
 - For each (user,movie) pair observed:
 - Pick rating as $L_u R_v + \text{noise}$
- Joint probability:

©Emily Fox 2014

13

PMF Graphical Model

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

- Graphically:

©Emily Fox 2014

14

Maximum A Posteriori for Matrix Completion

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

$$\propto e^{-\frac{1}{2\sigma_u^2} \sum_{u=1}^n \sum_{i=1}^k L_{ui}^2} e^{-\frac{1}{2\sigma_v^2} \sum_{v=1}^m \sum_{i=1}^k R_{vi}^2} e^{-\frac{1}{2\sigma_r^2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2}$$

©Emily Fox 2014

15

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:

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

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

- Incorporate other sources of information or dependencies

©Emily Fox 2014

16

What you need to know...

- Probabilistic model for collaborative filtering
 - Models, choice of priors
 - MAP equivalent to optimization for matrix completion

Case Study 4: Collaborative Filtering

Gibbs Sampling for Bayesian Inference

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

Emily Fox

February 18th, 2014

Posterior Computations

- MAP estimation focuses on point estimation:

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

- 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

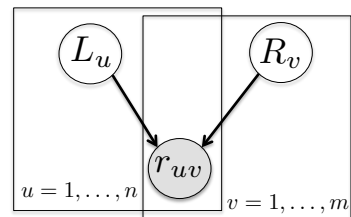
- Often no closed-form characterization (e.g., mixture models, PMF, etc.)

©Emily Fox 2014

19

Bayesian PMF Example

- Latent user and movie factors:



- Observations
- Hyperparameters:

- Want to predict new movie rating:

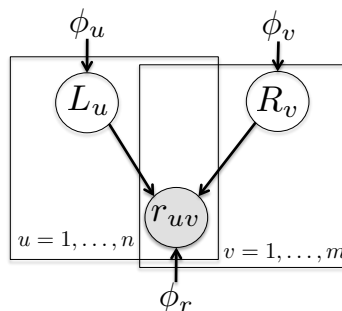
©Emily Fox 2014

20

Bayesian PMF Example

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

- Monte Carlo methods:



- Ideally:

©Emily Fox 2014

21

Bayesian PMF Example

- Want posterior samples $(L^{(k)}, R^{(k)}) \sim p(L, R | X, \phi)$
- What can we sample from?
 - Hint: Same reasoning as behind ALS, but sampling rather than maximization

©Emily Fox 2014

22

Bayesian PMF Example

- For user u :

$$p(L_u | X, R, \phi_u) \propto p(L_u | \phi_u) \prod_{v \in V_u} p(r_{uv} | L_u, R_v, \phi_r)$$

- Symmetrically for R_v conditioned on L (breaks down over movies)
- Luckily, we can use this to get our desired posterior samples

©Emily Fox 2014

23

Gibb Sampling

- Want draws:

- Construct Markov chain whose steady state distribution is
- Then, asymptotically correct
- Simplest case:

©Emily Fox 2014

24

Bayesian PMF Gibbs Sampler

- Outline of Bayesian PMF sampler

©Emily Fox 2014

25

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.

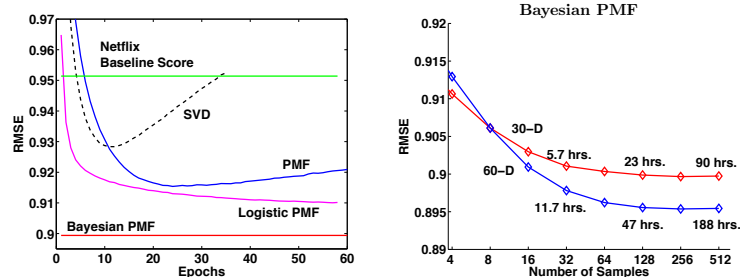


Figure 2. Left panel: Performance of SVD, PMF, logistic PMF, and Bayesian PMF using 30D feature vectors, on the Netflix validation data. The y-axis displays RMSE (root mean squared error), and the x-axis shows the number of epochs, or passes, through the entire training set. Right panel: RMSE for the Bayesian PMF models on the validation set as a function of the number of samples generated. The two curves are for the models with 30D and 60D feature vectors.

©Emily Fox 2014

26

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

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

©Emily Fox 2014

27

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

©Emily Fox 2014

28

Case Study 4: Collaborative Filtering

Matrix Factorization and Probabilistic LFM for Network Modeling

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

Emily Fox

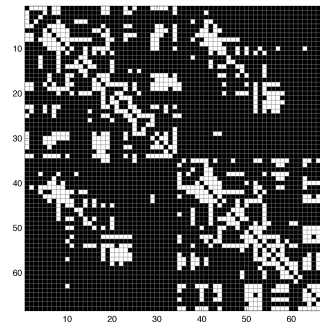
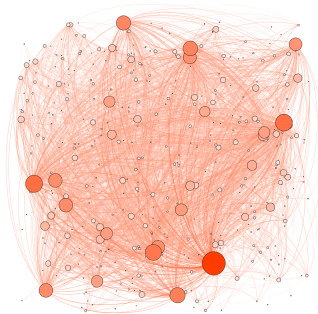
February 18th, 2014

©Emily Fox 2014

29

Network Data

- Structure of network data



©Emily Fox 2014

30

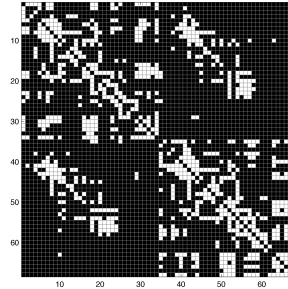
Properties of Data Source

- Similarities to Netflix data:

- Matrix
- High-dimensional
- Sparse

- Differences

- Square
- Binary



©Emily Fox 2014

31

Matrix Factorization for Network Data

- Vanilla matrix factorization approach:

- What to return for link prediction?

- Slightly fancier:

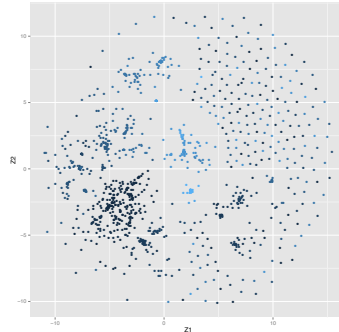
©Emily Fox 2014

32

Probabilistic Latent Space Models

- Assume features (covariates) of the user or relationship
- Each user has a “position” in a k -dimensional latent space

- Probability of link:



©Emily Fox 2014

33

Probabilistic Latent Space Models

- Probability of link:

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

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

- 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

34

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