

## Case Study 4: Collaborative Filtering

### Probabilistic Models for Matrix Factorization

### Cold-Start Problem

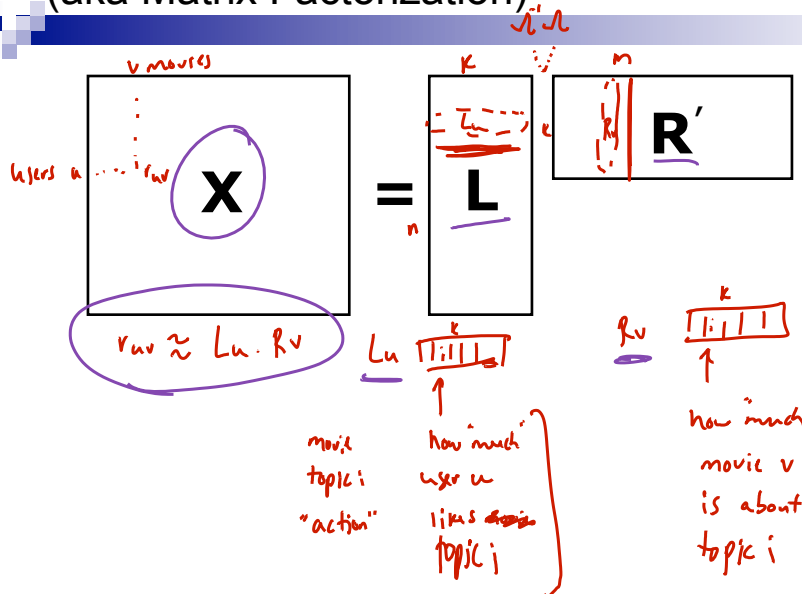
Machine Learning/Statistics for Big Data  
CSE599C1/STAT592, University of Washington

Carlos Guestrin  
March 12<sup>th</sup>, 2013

©Carlos Guestrin 2013

1

## Interpreting Low-Rank Matrix Completion (aka Matrix Factorization)



©Carlos Guestrin 2013

2

# Stochastic Gradient Descent

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

$\sum_u L_u \cdot L_u$

- Observe one rating at a time  $r_{uv}^t$   $\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} - r_{uv}$

- Gradient observing  $r_{uv}$ :

$$\frac{\partial F}{\partial L_u} = \epsilon_t R_v + \lambda_u L_u \quad \frac{\partial F}{\partial R_v} = \epsilon_t L_u + \lambda_v R_v$$

- Updates: step size  $\eta_t$ ,  $\begin{bmatrix} L \\ R \end{bmatrix} \leftarrow \begin{bmatrix} L \\ R \end{bmatrix} - \eta_t \nabla F_t$

$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$

fast & easy to implement

©Carlos Guestrin 2013

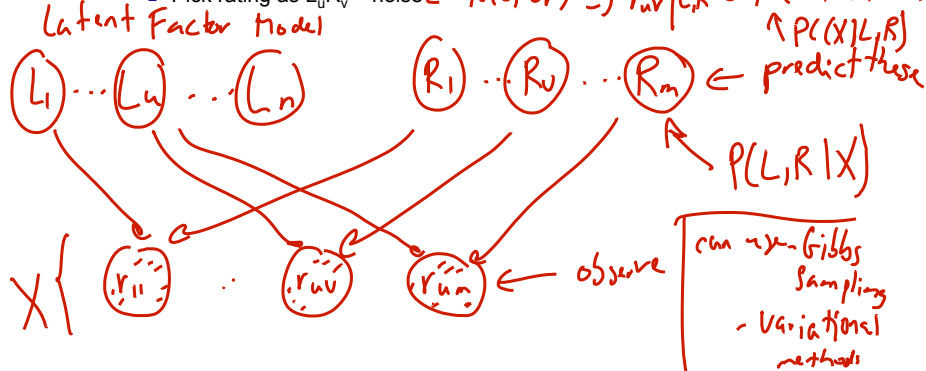
3

# What's Matrix Factorization Optimizing???

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

- A generative process:

- Pick user factors  $L_{u1}, \dots, L_{un} \Rightarrow L_{ui} \overset{iid}{\sim} N(0, \sigma_u^2) \leftarrow P(L)$
- Pick movie factors  $R_{v1}, \dots, R_{vm} \Rightarrow R_{vi} \overset{iid}{\sim} N(0, \sigma_v^2) \leftarrow P(R)$
- For each (user, movie) pair observed:
  - Pick rating as  $L_u R_v + \text{noise} \leftarrow N(0, \sigma_r^2) \Rightarrow r_{uv} | L, R \sim N(L_u \cdot R_v, \sigma_r^2) \leftarrow P(X|L,R)$



©Carlos Guestrin 2013

4

# Maximum A Posteriori for Matrix Completion

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

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

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

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

$\|L\|_F^2 \quad \|R\|_F^2 \quad \text{prediction error on train data / obs}$

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

$\lambda_u = \frac{\sigma_r^2}{\sigma_u^2} \quad \lambda_v = \frac{\sigma_r^2}{\sigma_v^2} \quad \text{yay!!}$

# MAP versus Regularized Least-Squares for Matrix Completion

- MAP under Gaussian Model:

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

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

- 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 models is very useful! E.g.

- Change priors
  - $L_{ui} \sim N(0, \sigma_u^2)$  }  $L_2$
  - $R_{vi} \sim N(0, \sigma_v^2)$  } regularization
- Incorporate other sources of information or dependencies

$\uparrow$  equivalent

Laplace  $(\mu, b)$   
 $P(x) \propto e^{-\frac{1}{b}(\mu - x)}$

# Cold-Start Problem

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

$\phi(\text{Skyfall}) = \begin{pmatrix} \text{action} \\ \text{romance} \\ 7 \\ 2 \\ 0 \\ \vdots \end{pmatrix}$

$\phi(\text{FRWL}) = \begin{pmatrix} 8 \\ 1 \\ 0 \\ \vdots \end{pmatrix}$

©Carlos Guestrin 2013

7

# Cold-Start More Formally

- No observations about a particular user:  $u$ 's ratings:  $\forall v, r_{u,v} = ?$

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

$\Rightarrow L_u = 0$ , always predict  $r_{u,v} = 0$  (constant)

- A simpler model for collaborative filtering:

□ Observe ratings:  $r_{uv} \in X$

□ Given features of a movie:

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

genre
year
director

□ Fit linear model:

For all users  $u$ ,  $r_{uv} \sim N(w \cdot \phi(v), \sigma_r^2)$

□ Minimize: for movie  $v$

$$\min_v \sum_{r_{uv}} (w \cdot \phi(v) - r_{uv})^2 + \lambda_w \|w\| \leftarrow \text{least squares / Lasso}$$

©Carlos Guestrin 2013

8

# Personalization

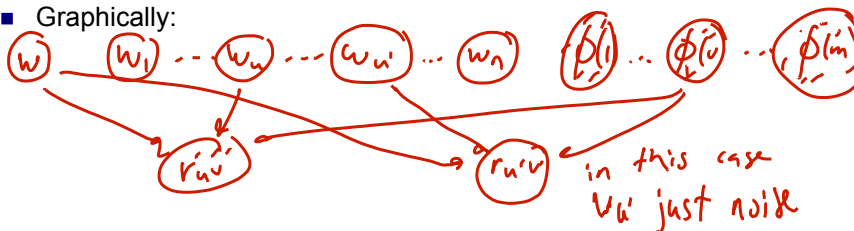


- If we don't have any observations about a user, use wisdom of the crowd
  - Address cold-start problem

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

- But, as we gain more information about the user, forget the crowd:  
 parameter  $w_u$  per user, represents deviation from crowd prob  
 $r_{uv} \sim N((w + w_u) \cdot \phi(v), \sigma_r^2)$

- Graphically:



©Carlos Guestrin 2013

9

# User Features...

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

$$\phi(u) = (25, \text{ } \epsilon, \text{ MSc}, \text{ A+}, \dots)$$

age      gender      education      grade in Big Data (class)

- Combine with features of movie:

$$\phi(u, v) = (\dots \phi(u) \dots, \dots \phi(v) \dots, \dots \text{cross features} \dots)$$

- Unified linear model:

$$r_{uv} \sim N((w + w_u) \cdot \phi(u, v), \sigma_r^2)$$

©Carlos Guestrin 2013

10

# Feature-based Approach versus Matrix Factorization

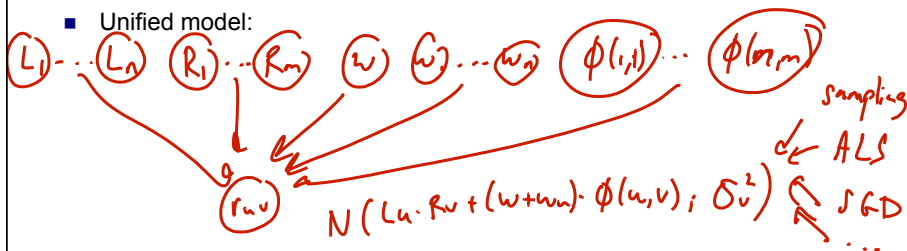
- Feature-based approach:
  - Feature representation of user and movies fixed
  - Can address cold-start

*important side info*

$\phi(u,v)$

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

- Unified model:



©Carlos Guestrin 2013

11

## MAP for Unified Collaborative Filtering via SGD

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

*hierarchical model*  
 $\|w_u - w\|_2^2$

- Gradient step observing  $r_{uv}^{(t)}$ :  $\epsilon_t = L_u^{(t)} \cdot R_v^{(t)} + (w_u^{(t)} + w^{(t)}) \cdot \phi(u,v) - r_{uv}^{(t)}$

- 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  $w_u$ :  $\partial_w F^{(t)} = \epsilon_t \phi(u,v) + \lambda_w w^{(t)}$

$$\begin{bmatrix} w^{(t+1)} \\ w_u^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t + \lambda_w) w^{(t)} - \eta_t \epsilon_t \phi(u,v) \\ (1 - \eta_t + \lambda_{w_u}) w_u^{(t)} - \eta_t \epsilon_t \phi(u,v) \end{bmatrix}$$

*only update  $w_u$  for user  $u$  in  $r_{uv}^{(t)}$*

©Carlos Guestrin 2013

12

## What you need to know...

- Probabilistic model for collaborative filtering
  - Models, choice of priors
  - MAP equivalent to optimization for matrix completion
  
- Cold-start problem
  
- Feature-based methods for collaborative filtering
  - Help address cold-start problem
  
- Unified approach