

# Collaborative Filtering

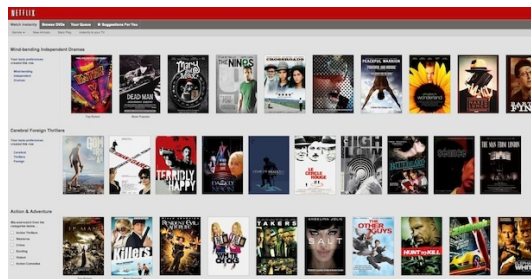
## Matrix Completion

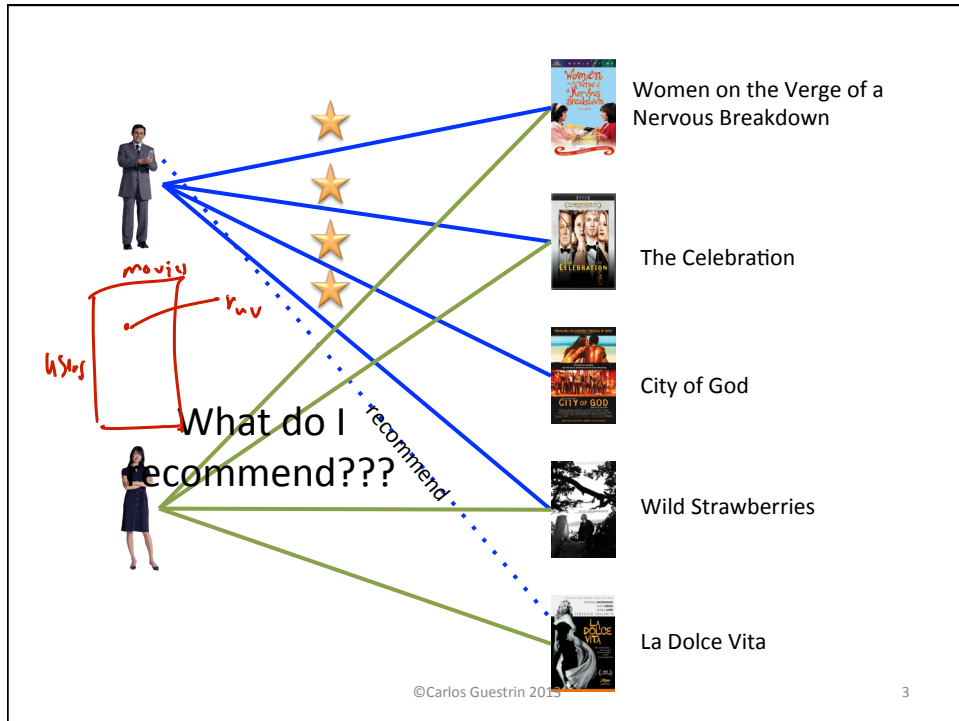
### Alternating Least Squares

Carlos Guestrin  
February 28<sup>th</sup>, 2013

1

- **Goal:** Find movies of interest to a user based on movies watched by the user and others
- **Methods:** matrix factorization, GraphLab





## Cold-Start Problem

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

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



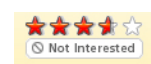
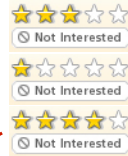
$$\phi(\text{FRWL}) = \begin{pmatrix} 8 \\ 1 \\ 6 \\ \vdots \end{pmatrix}$$

# Netflix Prize



- Given 100 million ratings on a scale of 1 to 5, predict 3 million ratings to highest accuracy

Skyfall  
WGB  
FRWL



PoC

- 17770 total movies
- 480189 total users
- Over 8 billion total ratings
- How to fill in the blanks?

mass  
users  
8B params

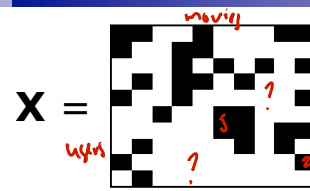
100M obs

Figures from Ben Recht

©Carlos Guestrin 2013

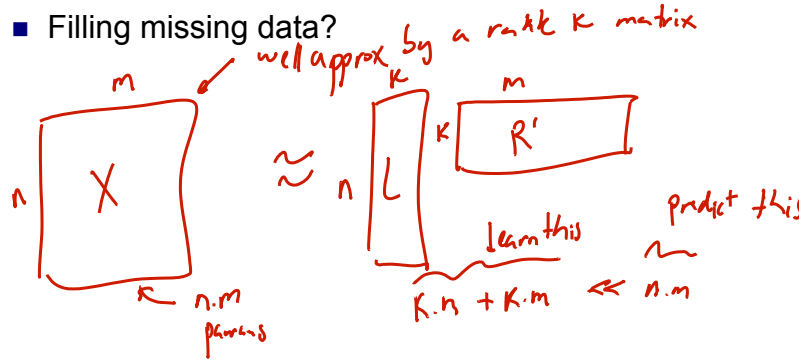
5

# Matrix Completion Problem



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

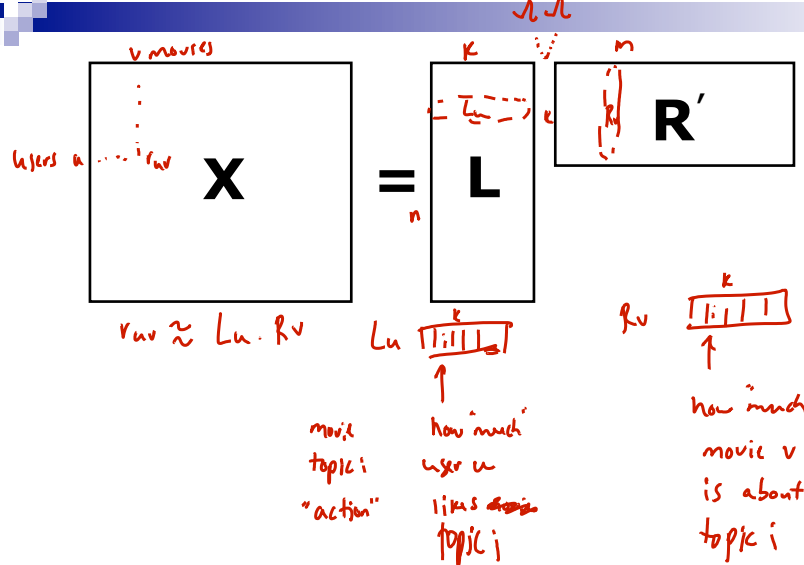
- Filling missing data?



©Carlos Guestrin 2013

6

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



©Carlos Guestrin 2013

7

## Matrix Completion via Rank Minimization

- Given observed values:  $(u, v, r_{uv}) \in X$  some  $r_{uv} = ?$
- Find matrix  $\Theta$
- Such that:  $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$   
fit  $r_{uv} \neq ?$  perfectly
- But...
- Introduce bias:  $\min_{\Theta} \text{rank}(\Theta)$   
 $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$
- Two issues:   
 (up-hard) you can't hope to get exact matchings

©Carlos Guestrin 2013

8

# Approximate Matrix Completion

- Minimize squared error:
  - (Other loss functions are possible)
- Choose rank  $k$ :
- Optimization problem:

©Carlos Guestrin 2013

9

# Coordinate Descent for Matrix Factorization

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

- Fix movie factors, optimize for user factors
- First Observation:

©Carlos Guestrin 2013

10

# Minimizing Over User Factors

- For each user  $u$ :  $\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$

- In matrix form:

- Second observation: Solve by

©Carlos Guestrin 2013

11

# Coordinate Descent for Matrix Factorization: Alternating Least-Squares

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

- Fix movie factors, optimize for user factors

- Independent least-squares over users

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

- Fix user factors, optimize for movie factors

- Independent least-squares over movies

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

- System may be underdetermined:

- Converges to

©Carlos Guestrin 2013

12

## Effect of Regularization

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

A diagram illustrating the matrix factorization equation  $X = LR'$ . It consists of three rectangular boxes. The first box on the left is a square and contains the letter **X**. To its right is an equals sign. The second box is a tall, narrow rectangle and contains the letter **L**. To its right is a third box, which is a horizontal rectangle and contains the letter **R'**.

©Carlos Guestrin 2013

13

## What you need to know...

- Matrix completion problem for collaborative filtering
- Over-determined  $\rightarrow$  low-rank approximation
- Rank minimization is NP-hard
- Minimize least-squares prediction for known values for given rank of matrix
  - Must use regularization
- Coordinate descent algorithm = “Alternating Least Squares”

©Carlos Guestrin 2013

14

## Case Study 4: Collaborative Filtering

### SGD for Matrix Completion Matrix-norm Minimization

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

Carlos Guestrin

March 7<sup>th</sup>, 2013

©Carlos Guestrin 2013

15

## Stochastic Gradient Descent

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

- Observe one rating at a time  $r_{uv}$

- Gradient observing  $r_{uv}$ :

- Updates:

©Carlos Guestrin 2013

16



## Local Optima v. Global Optima

- We are solving:

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

- We (kind of) wanted to solve:

- Which is NP-hard...

- ☐ How do these things relate???

©Carlos Guestrin 2013

17

## Eigenvalue Decompositions for PSD Matrices

- Given a (square) symmetric positive semidefinite matrix:

- ☐ Eigenvalues:

- Thus rank is:

- Approximation:

- Property of trace:

- Thus, approximate rank minimization by:

©Carlos Guestrin 2013

18

# Generalizing the Trace Trick

- Non-square matrices ain't got no trace
- For (square) positive definite matrices, matrix factorization:
- For rectangular matrices, singular value decomposition:
- Nuclear norm:

©Carlos Guestrin 2013

19

# Nuclear Norm Minimization

- Optimization problem:
- Possible to relax equality constraints:
- Both are convex problems!  
(solved by semidefinite programming)

©Carlos Guestrin 2013

20

# Analysis of Nuclear Norm

- Nuclear norm minimization is a convex relaxation of rank minimization problem:

$$\min_{\Theta} \text{rank}(\Theta)$$

$$\min_{\Theta} \|\Theta\|_*$$

$$r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq ?$$

$$r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq ?$$

- Theorem [Candes, Recht '08]:

- If there is a true matrix of rank  $k$ ,
- And, we observe at least

$$C k n^{1.2} \log n$$

random entries of true matrix

- Then true matrix is recovered exactly with high probability with convex nuclear norm minimization!
- Under certain conditions

©Carlos Guestrin 2013

21

## Nuclear Norm Minimization versus Direct (Bilinear) Low Rank Solutions

- Nuclear norm minimization:  $\min_{\Theta} \sum_{r_{uv}} (\Theta_{uv} - r_{uv})^2 + \lambda \|\Theta\|_*$

- Annoying because:

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

- Annoying because:

- But  $\|\Theta\|_* = \inf \left\{ \min_{L,R} \frac{1}{2} \|L\|_F^2 + \frac{1}{2} \|R\|_F^2 : \Theta = LR' \right\}$

- So
- And

- Under certain conditions [Burer, Monteiro '04]

©Carlos Guestrin 2013

22

## What you need to know...

- Stochastic gradient descent for matrix factorization
- Norm minimization as convex relaxation of rank minimization
  - Trace norm for PSD matrices
  - Nuclear norm in general
- Intuitive relationship between nuclear norm minimization and direct (bilinear) minimization

©Carlos Guestrin 2013

23

## Case Study 4: Collaborative Filtering

Nonnegative Matrix Factorization  
Projected Gradient

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

Carlos Guestrin

March 7<sup>th</sup>, 2013

©Carlos Guestrin 2013

24

## Matrix factorization solutions can be unintuitive...

- Many, many, many applications of matrix factorization
- E.g., in text data, can do topic modeling (alternative to LDA):

$$\mathbf{X} = \mathbf{L} \mathbf{R}'$$

- Would like:
- But...

©Carlos Guestrin 2013

25

## Nonnegative Matrix Factorization

$$\mathbf{X} = \mathbf{L} \mathbf{R}'$$

- Just like before, but

$$\min_{L \geq 0, R \geq 0} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u ||L||_F^2 + \lambda_v ||R||_F^2$$

- Constrained optimization problem
  - Many, many, many, many solution methods... we'll check out a simple one

©Carlos Guestrin 2013

26

# Projected Gradient

- Standard optimization:
  - Want to minimize:  $\min_{\Theta} f(\Theta)$
  - Use gradient updates:
 
$$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$
- Constrained optimization:
  - Given convex set  $\mathcal{C}$  of feasible solutions
  - Want to find minima within  $\mathcal{C}$ :  $\min_{\substack{\Theta \\ \Theta \in \mathcal{C}}} f(\Theta)$
- Projected gradient:
  - Take a gradient step (ignoring constraints):
  - Projection into feasible set:

©Carlos Guestrin 2013

27

# Projected Stochastic Gradient Descent for Nonnegative Matrix Factorization

$$\min_{L \geq 0, R \geq 0} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2} \|L\|_F^2 + \frac{\lambda_v}{2} \|R\|_F^2$$

- Gradient step observing  $r_{uv}$  ignoring constraints:
 
$$\begin{bmatrix} \tilde{L}_u^{(t+1)} \\ \tilde{R}_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$
- Convex set:
- Projection step:

©Carlos Guestrin 2013

28

## What you need to know...

- In many applications, want factors to be nonnegative
- Corresponds to constrained optimization problem
- Many possible approaches to solve, e.g., projected gradient