

# Recommender Systems

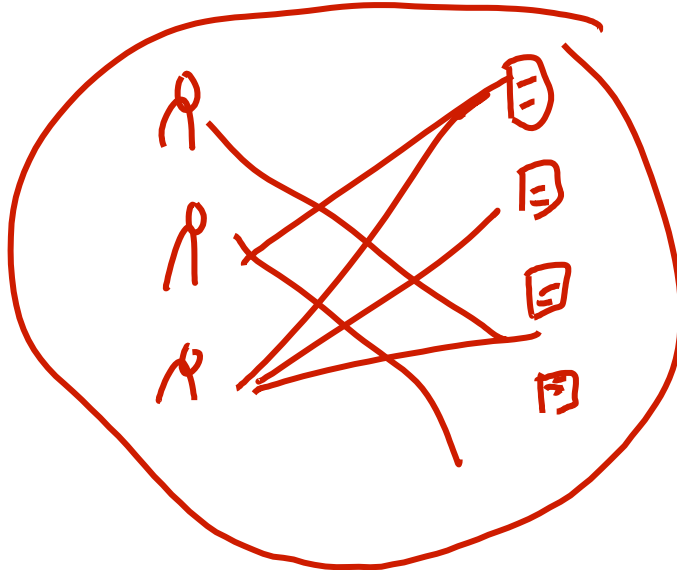
Machine Learning – CSEP546  
Carlos Guestrin  
University of Washington  
February 10, 2014

# Personalization is transforming our experience of the world



100 Hours a Minute  
*What do I care about?*

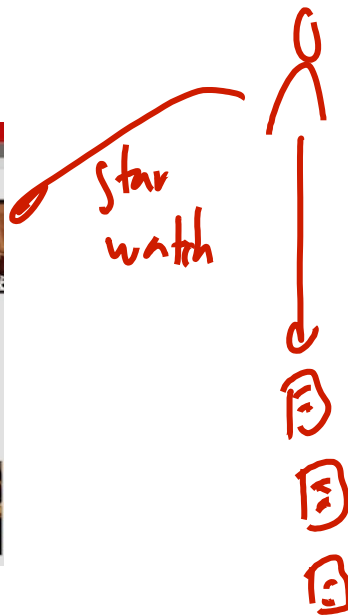
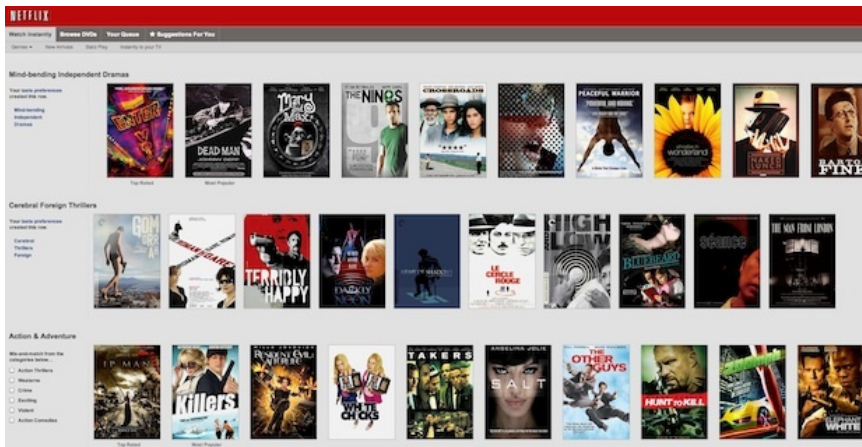
- Information overload →  
"Browsing" is history
  - Need fundamentally new ways to discover content
- Personalization: Connects *users & items*



⊗ recommendations



# Movie recommendations



# Product recommendations

Recommendations combine  
global & session interests

what I bought in the past

intention  
- - - -

The screenshot shows the Amazon product page for the book "Processing: A Programming Handbook for Visual Designers and Artists (Hardcover)". The book is by Casey Reas (Author), Ben Fry (Author), and John Mase (Foreword). It has a 4.3-star rating from 43 customer reviews. The price is \$47.95 for new and \$43.56 for used. The page includes a "Get Free Two-Day Shipping" offer, a "Share your own customer images" section, a "Please tell the publisher" section, and a "Related Education & Training Services in Pittsburgh" section. At the bottom, there is a "Customers Who Bought This Item Also Bought" section with five recommended books: "Processing: Creative Coding and Computational A...", "Visualizing Data: Exploring and Explaining Data...", "Making Things Talk: Practical Methods for Control...", "Physical Computing: Sensing and Controlling the...", and "Learning Processing: A Beginner's Guide to...".

**Processing: A Programming Handbook for Visual Designers and Artists (Hardcover)**  
by Casey Reas (Author), Ben Fry (Author), John Mase (Foreword)  
★★★★★ (43 customer reviews)  
Available from these sellers.  
31 new from \$47.95    8 used from \$43.56  
Get Free Two-Day Shipping  
Get Free Two-Day Shipping for three months with a special extended free trial of Amazon Prime™. Add this eligible textbook to your cart to qualify. Sign up at checkout. [See details.](#)

[Share your own customer images](#)  
Publisher: learn how customers can search inside this book.

Please tell the publisher:  
I'd like to read this book on Kindle.  
Don't have a Kindle? [Get yours here.](#)

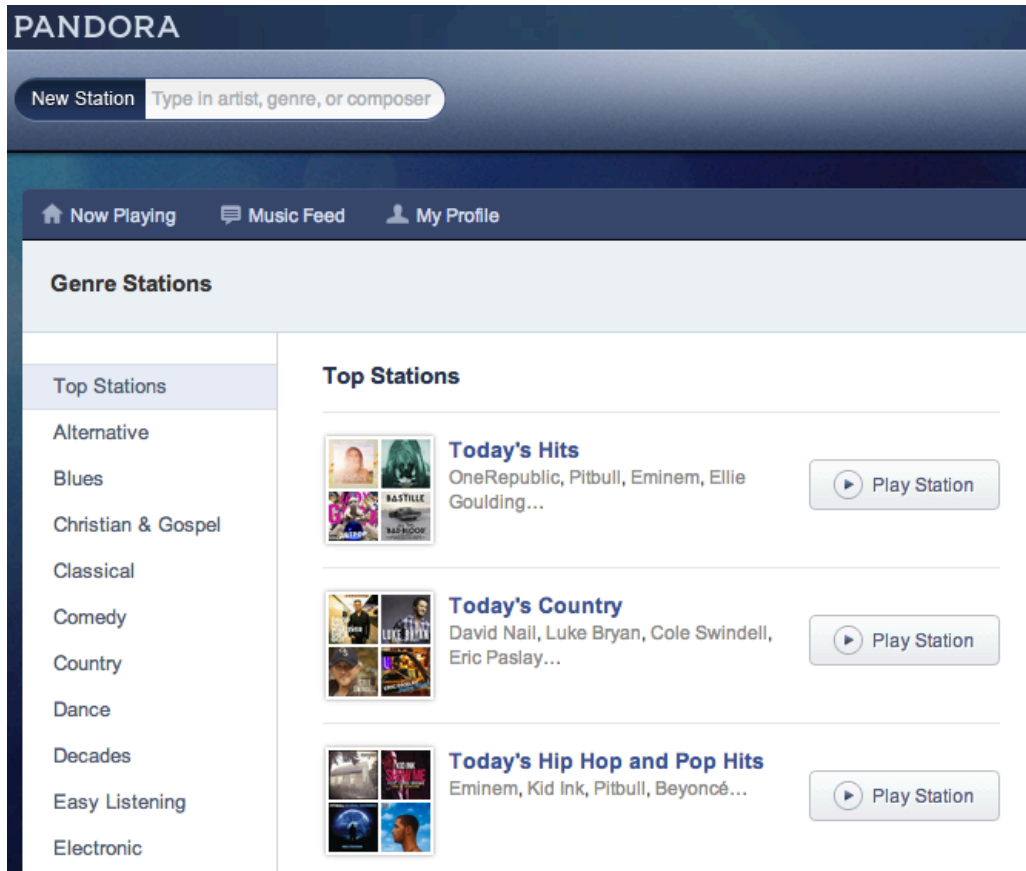
**Related Education & Training Services in Pittsburgh** [Compare map](#) | [Change location](#)  
[Learn HTML Coding](#)  
www.FullSelf.edu • Earn Your Bachelor's Degree in Web Design and Development.  
[Create Websites with HTML](#)  
http://www.unix.Berkeley.edu • Learn HTML Online, Start Anytime! with UC Berkeley Extension  
[Intensive XSLT Training](#)  
www.objectdatalabs.com/course10.asp • OnSite or in NYC, LA, SFO, ORD, DC Will customize & train as few as 3

**Customers Who Bought This Item Also Bought**

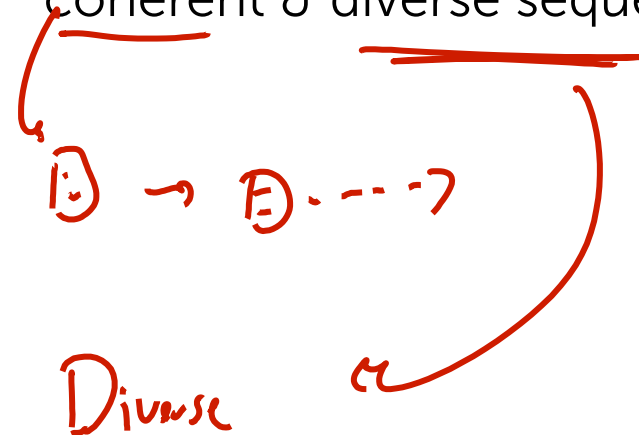
Product	Price
<a href="#">Processing: Creative Coding and Computational A...</a> by Ira Greenberg	★★★★★ (7) \$43.99
<a href="#">Visualizing Data: Exploring and Explaining Data...</a> by Ben Fry	★★★★★ (13) \$26.30
<a href="#">Making Things Talk: Practical Methods for Control...</a> by Tom Igoe	★★★★★ (13) \$19.79
<a href="#">Physical Computing: Sensing and Controlling the...</a> by Tom Igoe	★★★★★ (20) \$19.00
<a href="#">Learning Processing: A Beginner's Guide to...</a> by Daniel Shiffman	★★★★★ (7) \$44.95



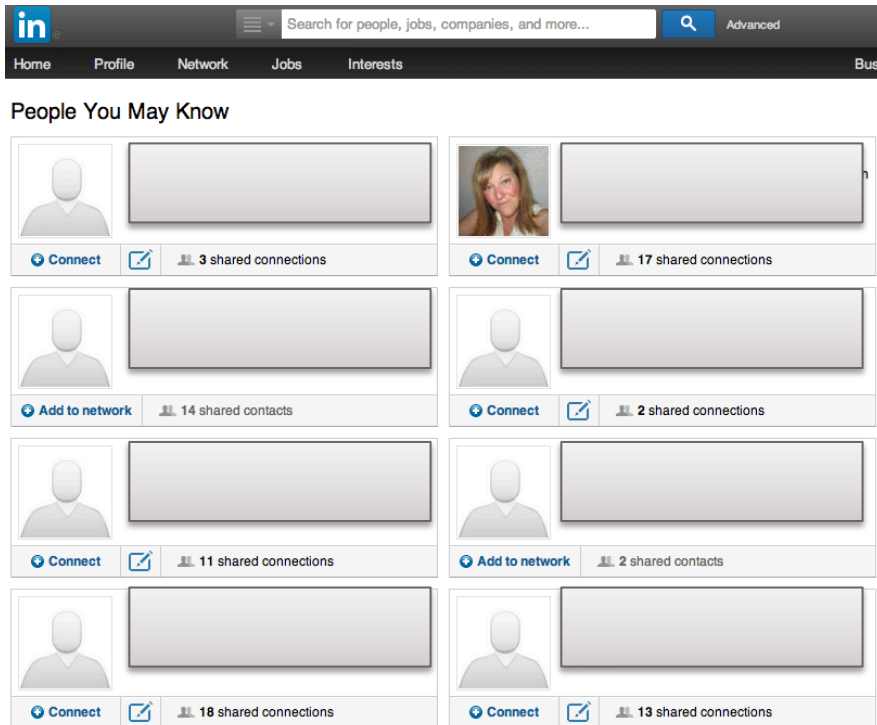
# Playlist recommendations



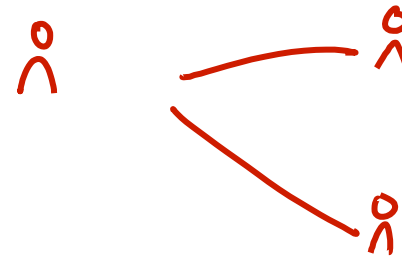
Recommendations form  
coherent & diverse sequence



# Friend recommendations

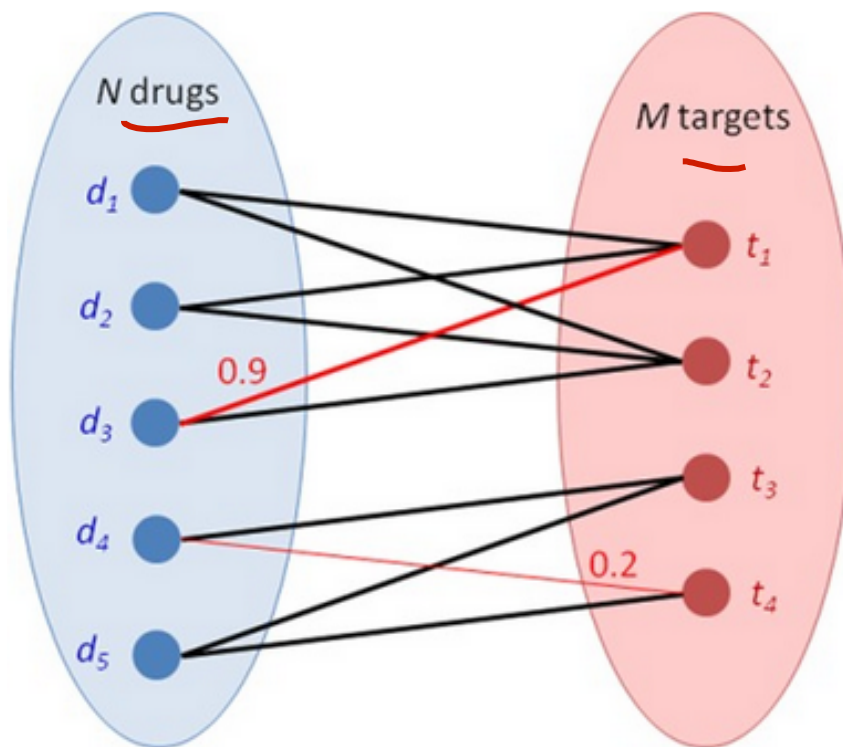


Users and "items" are of the same type



# Drug-target interactions

What drug should we  
“repurpose” for some disease?



Cobanoglu et al. '13



# Challenges of developing recommender systems





# Type of feedback

- Explicit – user tells us what she likes



- Implicit – we try to infer what she likes from usage data



clicked  
time on site  
...



# Top K versus diverse outputs

- Top K recommendations may be very redundant
  - *People who liked Rocky 1 also enjoyed Rocky 2, Rocky 3, Rocky 4, Rocky 5,...*
- Diverse recommendations
  - Users are multi-faceted & want to hedge our bets
  - *Rocky 2, It never rains in Philadelphia, Gandhi*



# A new movies walks into a bar...



IN THEATERS



- Cold-start problem: recommendations for new users or new movies

- Need side information about user/movie
  - A.K.A. features!

*action, actors, sequel,...*

- Could also play 20-questions game...



# That's so last year...

- Interests change over time...
  - Is it 1967?
  - Or 1977?
  - Or 1988?
  - Or 1998?
  - Or 2011?
- Models need flexibility to adapt to users
  - Macro scale
  - Micro scale *intention now*
- And keep checking that system still accurate



macys.com



# Scalability

- For  $N$  users and  $M$  movies, some approaches take  $O(N^3 + M^3)$ 
  - Not so good for billions of users...
- GraphLab can help...
  - Efficient implementations
  - Fast exact & approximate methods as needed



# Building a recommender system (easily with GraphLab)



# Solution 0: Popularity



# Simplest approach: popularity

- What people are viewing now
  - Super popular



- Limitations:
  - No context (what's my intention now)
  - No personalization

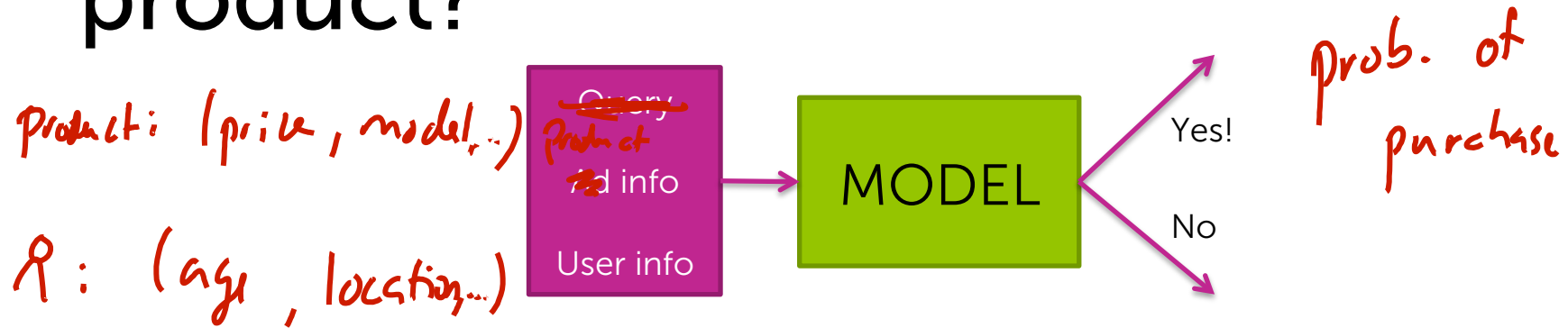




# Solution 1: Click prediction



# What's the probability I'll buy this product?



- Features capture context
  - Time of the day, what I just saw, user info, what I bought in the past
- Helps mitigate cold-start problem
  - Rate new movie from features of other movies user liked
- Limitation:
  - May not have context available
  - Often doesn't perform as well as collaborative filtering methods (next)



Solution 2: People who bought this  
also bought...



# Co-occurrence matrix

- Matrix C: item by item

$x \rightarrow$  Beer, diapers ...

$C = \text{item}$

Beer  $\rightarrow$  +1

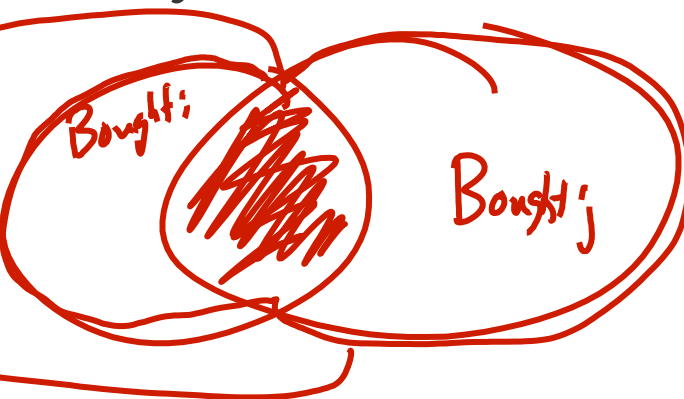
diapers  $\rightarrow$  +1

- $C_{ij} = C_{ji}$  number of users who bought both items  $i$  &  $j$



# Normalizing co-occurrences: similarity matrix

- $C_{ij}$  very large if either  $i$  or  $j$  are very popular movies → drowns out other effects
  - just recommends by popularity
- Jaccard similarity: normalizes by popularity
  - Who watched  $i$  and  $j$  divided by who watched  $i$  or  $j$

$$S_{ij} = \frac{\text{bought } i \text{ and } j}{\text{bought } i \text{ or } j}$$


- Many other similarity metrics possible, e.g., cosine similarity



# Using similarity matrix to recommend

- People who bought diapers also bought beer
- For  $i=\text{diapers}$ , sort  $S_{ij}$  and find  $j$  with highest similarity  
– Beer, milk, baby food,...  
*Handwritten notes: Bud Light, O.G. milk, Gerber, Pampers, Pampers & Pampers mix*
- Limitation:
  - Only current page matters, no history



## Solution 3: Average item-item similarity



# (Weighted) average over items user bought


- User  $u$  bought items  $B_u$ 
  - Define user specific similarity (score) for each item  $j$ 
    - Average similarity for items in  $B_u$

$$\text{Score}(R_u, \text{Beer}) = \frac{1}{2} (S_{\text{Beer}, \text{Milk}} + S_{\text{Beer}, \text{Diapers}})$$

- Could also weight recent purchases more
- For  $B_u = \{\text{diapers}, \text{beer}\}$  sort  $\text{Score}(u, j)$  and find  $j$  with highest similarity





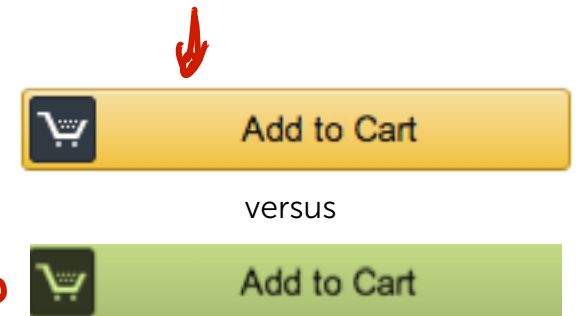


# Demo: Item-based collaborative filtering with GraphLab

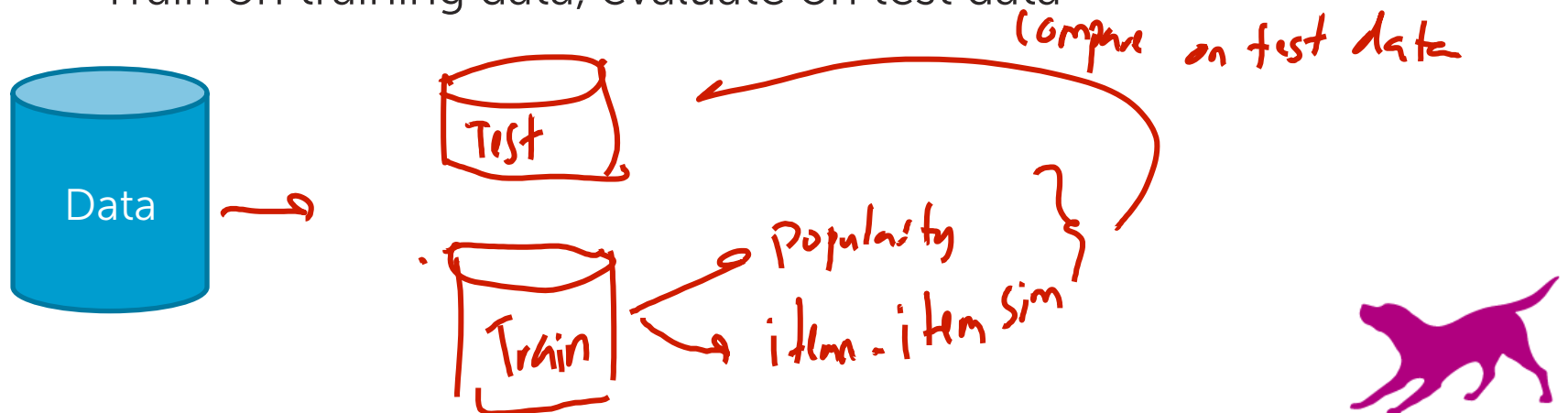


# Training versus testing data

- A/B testing standard in industry:
  - Randomly split users into groups A & B
  - Show different websites
  - Compare outcomes

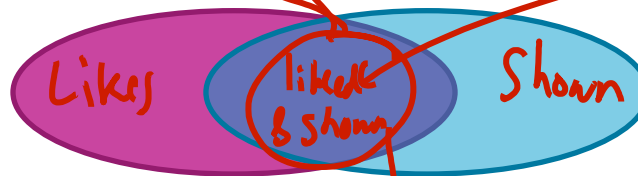


- Same idea fundamental in ML
  - Randomly split data into train and test sets
  - Train on training data, evaluate on test data

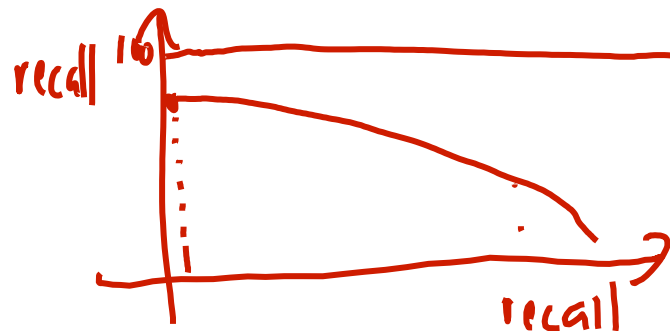



# Example Performance metric for recommenders

- User  $u$  liked  $m$  movies, we showed her  $k$  movies *want to maximize*



- Recall*: what fraction of the liked movies we found  
$$\frac{\text{Liked \& Shown}}{\text{Liked}}$$
- Precision*: what fraction of the movies we showed she liked  
$$\frac{\text{Liked \& Shown}}{\text{Shown}}$$
- Precision-Recall curve:





# Demo for reals: Item-based collaborative filtering with GraphLab



# Limitations of item-based similarity

- Scalability – similarity matrix  $M^2$  size
- Cold-start problem



## Solution 4: Discovering hidden structure by matrix factorization

# Suppose we had $d$ features of movies and users

- Describe movie  $v$  with features  $R_v$ 
  - How much is it action, romance, drama,...

$$R_v = (\overset{\text{action}}{0.9}, \overset{\text{romance}}{0.2}, \overset{\text{drama}}{0.5})$$

- Describe user  $u$  with features  $L_u$ 
  - How much she likes action, romance, drama,...

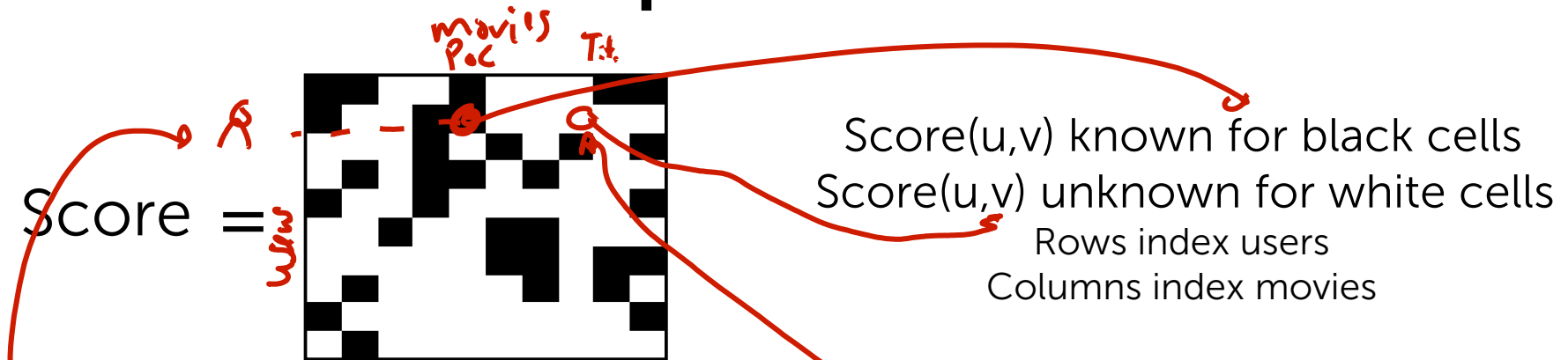
$$L_u = (0.01, 0, 0.9)$$

- $\text{Score}(u, v)$  is the product of the two vectors

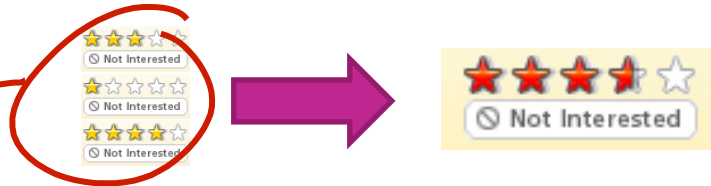
$$\text{Score}(u, v) = 0.9 \times 0.01 + 0.2 \times 0 + 0.5 \times 0.9$$

But we don't know features of users and movies...

# Matrix Completion Problem



- Users score some movies



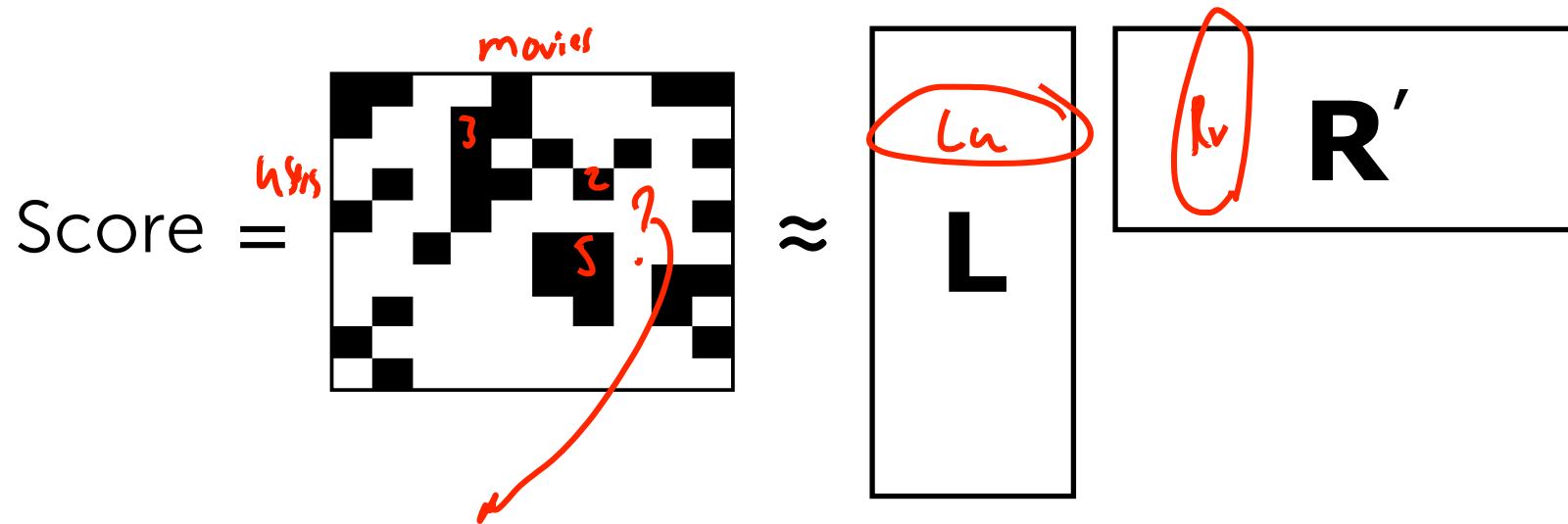
- Filling missing data?

If I had user & movie features

$$\text{Score}(u, \text{Titanic}) = u \times \boxed{\text{Titanic}}$$



# Matrix Factorization: discovering topics for users and movies



$$Score(R_{17}, P_{OC}) = L_{u_{17}} \cdot R_{P_{OC}}$$

Many efficient algorithms for  
matrix factorization implemented in GraphLab



The diagram shows a document matrix (a grid of black and white squares) labeled "document" and "words". This matrix is approximately equal to the product of matrix L and matrix R'.

Matrix L is a vertical rectangle with a circled element. A red arrow points from this circled element to the text "L\_u = for each doc how much is it about each of these".

Matrix R' is a horizontal rectangle with a circled element. A red arrow points from this circled element to the text "K topics each topic is associated with words".

7



# Using the results of matrix factorization

- Discover "features"  $R_v$  for each movie  $v$
- Discover "features"  $L_u$  for each user  $u$
- $\text{Score}(u,v)$  is the product of the two vectors  $\rightarrow$  predict how much a user will like a movie

$$\text{Score}(u, v) = L_u \cdot R_v$$

- Recommendations: sort movies user hasn't watched by  $\text{Score}(u,v)$

$\hat{v}$  ?  $\rightarrow$  Tianyi

$\text{Score}(\text{Tianyi}, v)$   $\forall v$   
sort

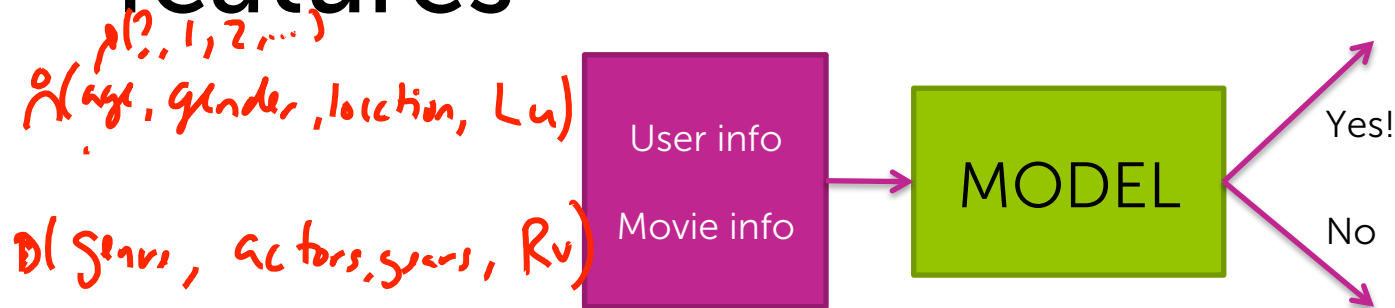
# Limitations of matrix factorization

- Cold-start problem



Bringing it all together:  
Featurized matrix factorization

# Combining real and discovered features



- Real features capture context
  - Time of the day, what I just saw, user info, what I bought in the past
- Discovered features from matrix factorization capture groups of users who behave similarly
  - Hipster wannabes from Seattle who teach and have a startup
- Mitigates cold-start problem
  - Ratings for a new user from real features only
  - As more information about user is discovered, matrix factorization “features” become more relevant

# Matrix Factorization

## Alternating Least Squares

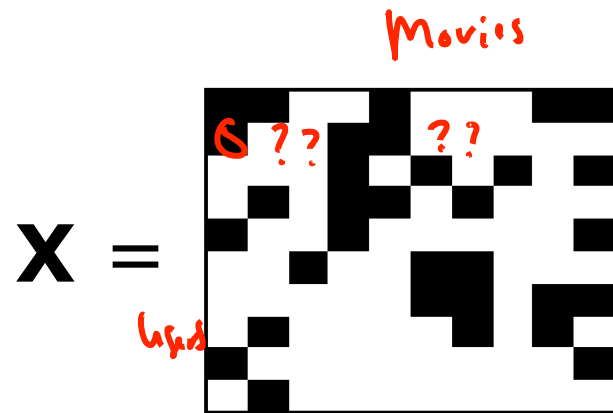
Machine Learning – CSEP546

Carlos Guestrin

University of Washington

February 10, 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?

Handwritten diagram illustrating matrix factorization:

$$X \approx L R$$

where  $X$  is  $n \times m$ ,  $L$  is  $n \times k$ , and  $R$  is  $k \times m$ .

if  $k = 100$

want to estimate  $n \cdot m$  numbers

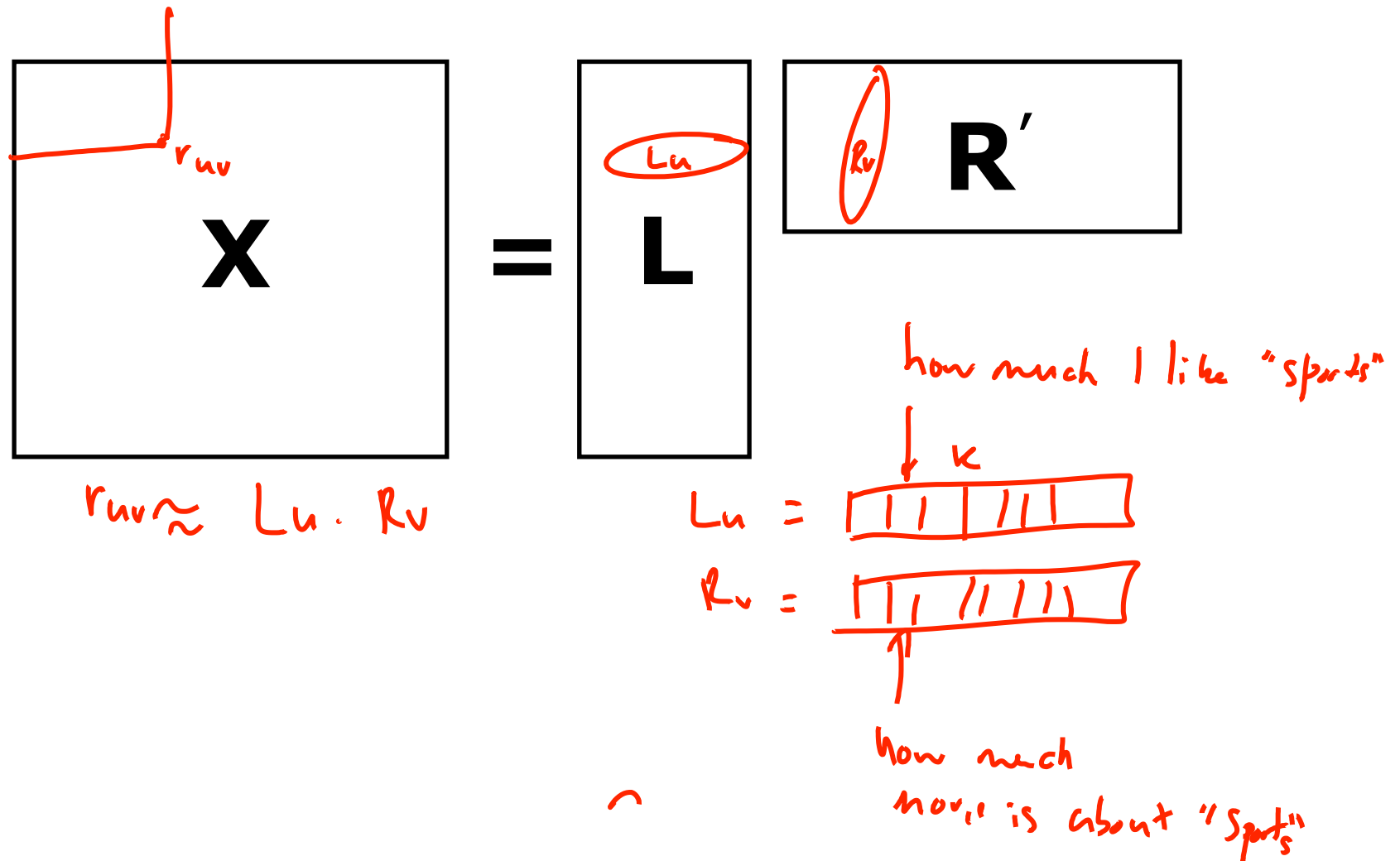
$n \cdot k + m \cdot k$

$5064 + 1.7M \approx 51.7M$  params

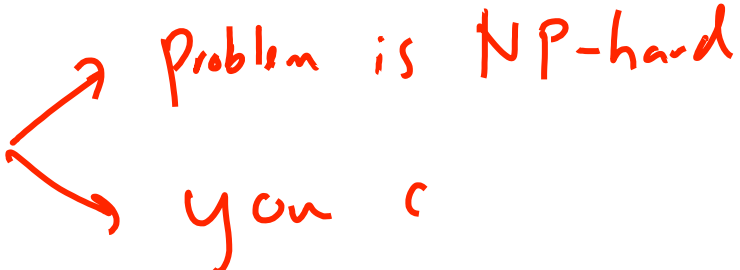
$m = 17k$     $n = 5064$     $n \cdot m \approx 8B$



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



# Matrix Completion via Rank Minimization

- Given observed values:  $(u, v, r_{uv}) \in X$   $r_{uv} \neq ?$
- Find matrix  $\Theta$
- Such that:  $\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$
- But... predictions for  $r_{uv} = ?$
- Introduce bias:  $\text{rank}(\Theta) = k$
- Two issues: 
  - Problem is NP-hard
  - you c

# Approximate Matrix Completion

- Minimize squared error:
  - (Other loss functions are possible)

$$\min_{\Theta} \sum_{u,v: r_{uv} \neq ?} (\Theta_{uv} - r_{uv})^2$$

- Choose rank  $k$ :

$$n \times m \Theta = n \times k \begin{matrix} \text{ } \\ \text{ } \end{matrix} \begin{matrix} k \\ m \end{matrix} R^T$$

- Optimization problem:

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

Non-convex problem  $\rightarrow$  local opt only.

# Coordinate Descent for Matrix Factorization

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

OPT ↓ fixed ↓

- Fix movie factors  $R$ , optimize for user factors  $L$
- First Observation:

$$\begin{aligned}
 & \min_{L_1, \dots, L_n} \sum_{(u,v): r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 && V_u \equiv \text{set of movies user } u \text{ rated} \\
 & \equiv \min_{L_1, \dots, L_n} \sum_{u=1}^n \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 && \leftarrow \text{indep opt probs per user} \\
 & \equiv \sum_{u=1}^n \underbrace{\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2}_{\text{Local regression, only depends on movies user } u \text{ rated}}
 \end{aligned}$$

# 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:

Hand-drawn matrix diagram illustrating the minimization over user factors  $L_u$ . The diagram shows a large matrix  $R_v$  with dimensions  $|V_u|$  (rows) by  $K$  (columns). A vector  $L_u$  is shown with dimensions  $K$  (rows) by  $1$  (column). The product  $L_u \cdot R_v$  is shown as a vector  $r_u$  with dimensions  $1$  (row) by  $|V_u|$  (column). The equation is  $L_u \cdot R_v \approx r_u$ .

become basis functions/features in regression

- Second observation: Solve by  
matrix inversion

# Coordinate Descent for Matrix Factorization: Alternating Least-Squares

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

Don't initialize  $R_v$  to 0!! usually init to random

- Fix movie factors, optimize for user factors
  - Independent least-squares over users  $L_u^{(t+1)}$   $\min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v^{(t)} - r_{uv})^2 + \lambda_u \|L_u\|$

- Fix user factors, optimize for movie factors
  - Independent least-squares over movies  $R_v^{(t+1)}$   $\min_{R_v} \sum_{u \in U_v} (L_u \cdot R_v^{(t+1)} - r_{uv})^2 + \lambda_v \|R_v\|$

$U_v \leftarrow$  users who rated movie  $v$

- System may be underdetermined:

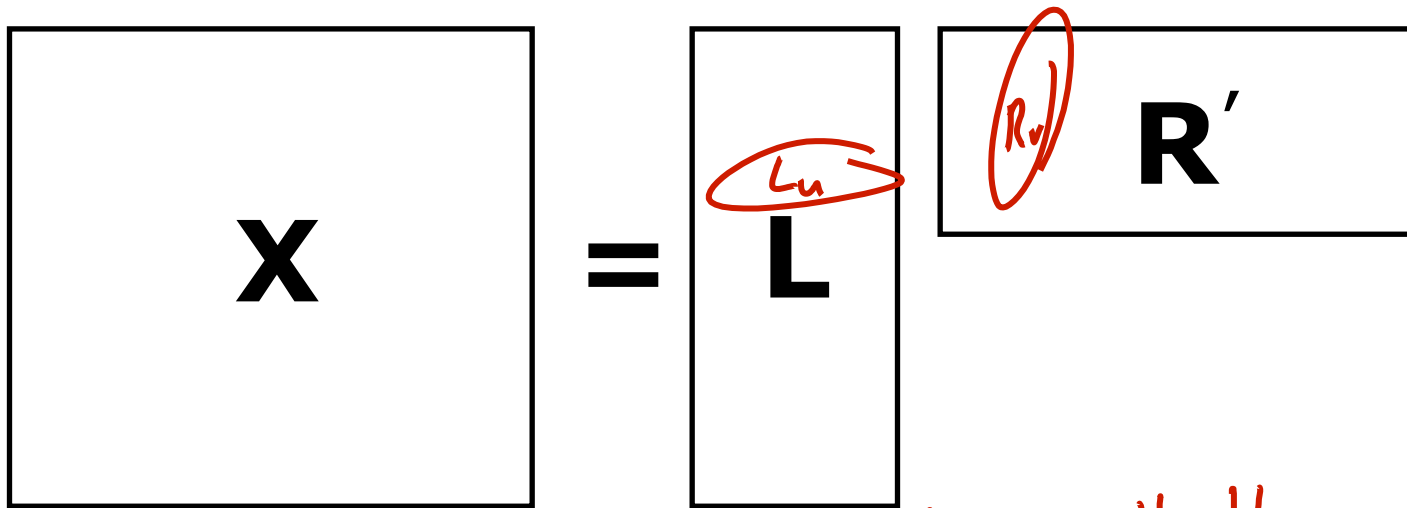
use regularization

- Converges to local opt

# Effect of Regularization

Frobenius  $\|L\|_F = \sqrt{\sum_{u,v} L_{uv}^2}$

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



if  $\|\cdot\| = \|\cdot\|_F$   
 each sub problem  
 uses  $L_2$  regularization  
 $\rightarrow$  ridge regression

$\|\cdot\| = \|\cdot\|_1$   
 $\rightarrow$  each sub problem LASSO  
 $\wedge$  sparse latent factors

# Stochastic Gradient Descent

$F$ :

$$\min_{L, R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\|_2^2 + \lambda_v \|R\|_2^2$$

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

- Gradient:

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

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

SGD observe ratings  $r_{uv}$

$$\begin{bmatrix} L^{(t+1)} \\ R^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} L^{(t)} \\ R^{(t)} \end{bmatrix} - \eta_t \nabla F_t$$

- Updates:

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

very simple to implement



# What you need to know...

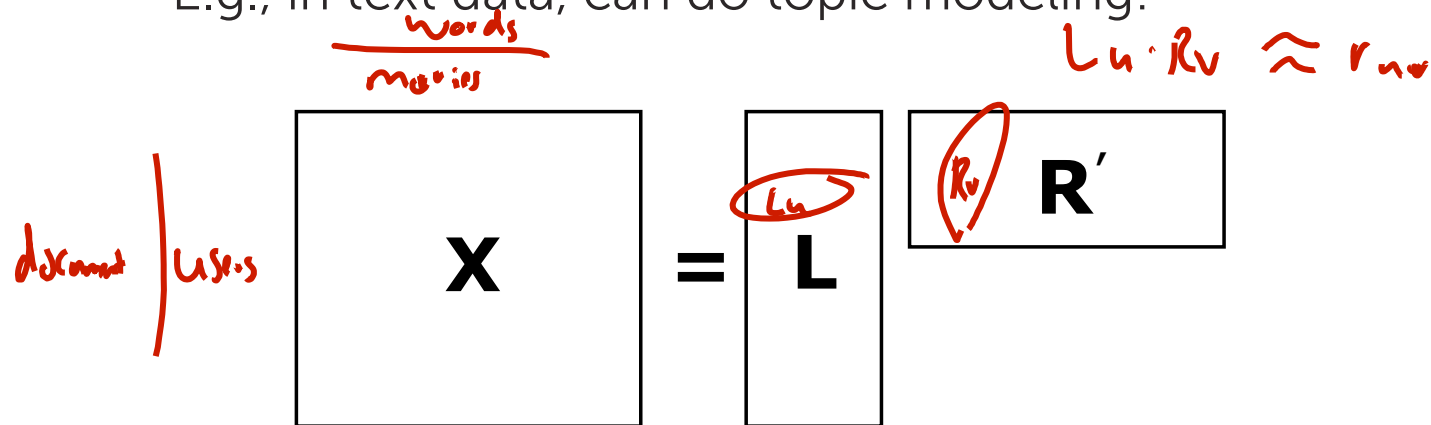
- Matrix completion problem for collaborative filtering
- Over-determined  $\rightarrow$  low-rank approximation
- Rank minimization is NP-hard
- Minimize least-squares prediction for known values for given rank of matrix
  - Must use regularization
- Coordinate descent algorithm = “Alternating Least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

# Non-Negative Matrix Factorization

Machine Learning – CSEP546  
Carlos Guestrin  
University of Washington  
February 10, 2014

# Matrix factorization solutions can be unintuitive...

- Many, many, many applications of matrix factorization
- E.g., in text data, can do topic modeling:



- Would like:  $\mathbf{L}_u \geq 0$   
 $\mathbf{R}_v \geq 0$
- But... possible

# 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

# 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:

# 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  $C$  of feasible solutions

- Want to find minima within  $C$ :  $\min_{\substack{\Theta \\ \Theta \in C}} f(\Theta)$

- Projected gradient:

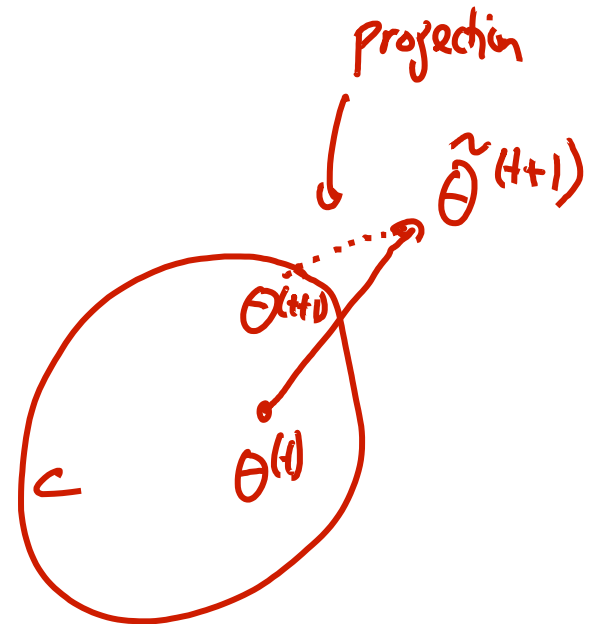
- Take a gradient step (ignoring constraints):

$$\tilde{\Theta}^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$

- Projection into feasible set:

$$\Pi_C(\Theta) \equiv \operatorname{argmin}_{\beta \in C} \|\Theta - \beta\|_2^2 \quad \left. \begin{array}{l} \text{often easy} \\ \text{to compute} \\ \text{(always convex)} \end{array} \right\}$$

$$\Theta^{(t+1)} = \Pi_C(\tilde{\Theta}^{(t+1)})$$



# 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:  $L_u \geq 0 ; R_v \geq 0 \quad \forall u, v$

- Projection step:

$\Pi_C(\theta) = \underset{\beta \in C}{\operatorname{argmin}} \|\theta - \beta\|_2^2 \leftarrow$  totally indep prob per dimension  
 Single dime  
 $= \underset{\beta \geq 0}{\operatorname{argmin}} (\theta - \beta)^2$   
 $= \begin{cases} \theta, & \text{if } \theta \geq 0 \\ 0, & \text{if } \theta < 0 \end{cases} = (\theta)_+$

$\begin{bmatrix} L^{(t+1)} \\ R^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} \tilde{L}^{(t+1)} \\ \tilde{R}^{(t+1)} \end{bmatrix}_+$ 
 set all neg coords to zero, universe on our side!!

# 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



# The Cold-Start Problem

Machine Learning – CSEP546  
Carlos Guestrin  
University of Washington  
February 10, 2014

# Cold-Start Problem

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



IN THEATERS



# Cold-Start More Formally

- No observations about a particular user:

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

*no ratings:  $\forall v \ r_{u,v} = ?$*

*doesn't depend on  $L_u$*

*$\Rightarrow L_u = 0$*

*always predict  $r_{u,v} = 0$  (constant)*

- A simpler model for collaborative filtering:

- Observe ratings:  $r_{uv} \in \mathcal{X}$

- 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 model

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

# Personalization

- If we don't have any observations about a user, use wisdom of the crowd
  - Address cold-start problem
- But, as we gain more information about the user, forget the crowd:

# User Features...

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

$$\phi(u) = ( \underset{\text{age}}{25}, \underset{\text{gender}}{F}, \underset{\text{education}}{MSc}, \underset{\text{grade in Big Data class}}{A+}, \dots )$$

- Combine with features of movie:

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

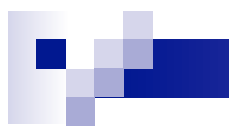
- Unified linear model:

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

# Feature-based Approach versus Matrix Factorization

- Feature-based approach:
  - Feature representation of user and movies fixed
  - Can address cold-start
- Matrix factorization approach:
  - Suffers from cold-start problem
  - User & movie features are learned from data
- Unified model:

# MAP for Unified Collaborative Filtering via SGD



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

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

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

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

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

# What you need to know...

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