

Say hi or ask questions in chat
before/during/after class

CSE/STAT 416

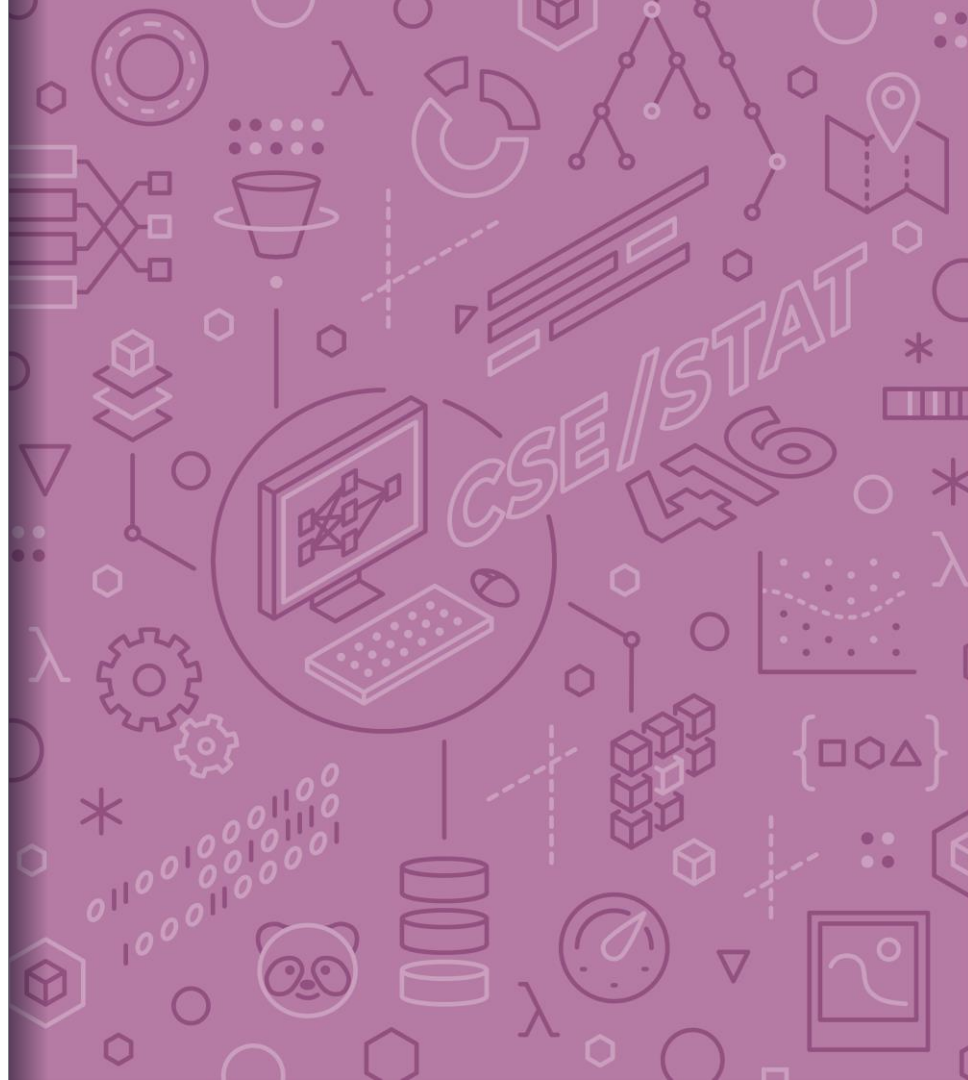
Recommender Systems

Hunter Schafer

Paul G. Allen School of Computer Science & Engineering
University of Washington

May 17, 2021

Music: HAIM



Last Time...

Co-occurrence Matrix

Solution 2

Co-occurrence Matrix

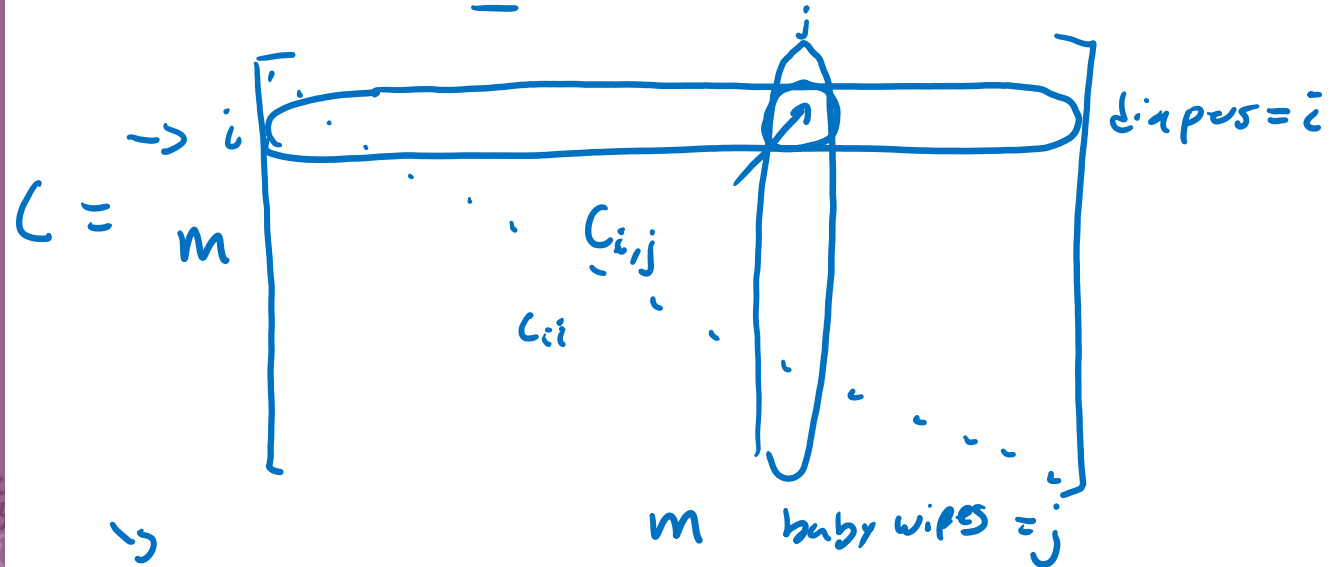
$m = \# \text{ items}$

$n = \# \text{ users}$

Idea: People who bought this, also bought ...

- E.g. people who buy diapers also buy baby wipes

Make co-occurrence matrix $C \in \mathbb{R}^{m \times m}$ (m is the number of items) of item purchases. $C_{ij} = \#$ of users who bought both item i and j



C will be symmetric ($C_{ij} = C_{ji}$)

Normalization

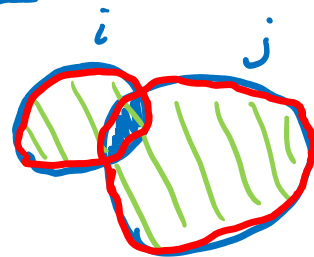
$$C \rightarrow S$$

The count matrix C needs to be normalized, otherwise popular items will drown out others (will just reduce to popularity).

Normalize the counts by using the Jaccard similarity instead

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

$$= \frac{C_{ij}}{T_i + T_j - C_{ij}}$$



Where $T_k = \#$ of people who bought item k

Could also use something like Cosine similarity, but Jaccard is popular

Analysis

Pros:

- It personalizes to the user

Cons

- Does **not** utilize
 - Context (e.g. time of day)
 - User features (e.g. age)
 - Product features (e.g. baby vs electronics)
- Scalability
 - Similarity is size m^2 where m is the number of items
- Cold start problem

Matrix Factorization

Solution 4


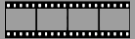
















Matrix Completion


Want to recommend movies based on user ratings for movies.

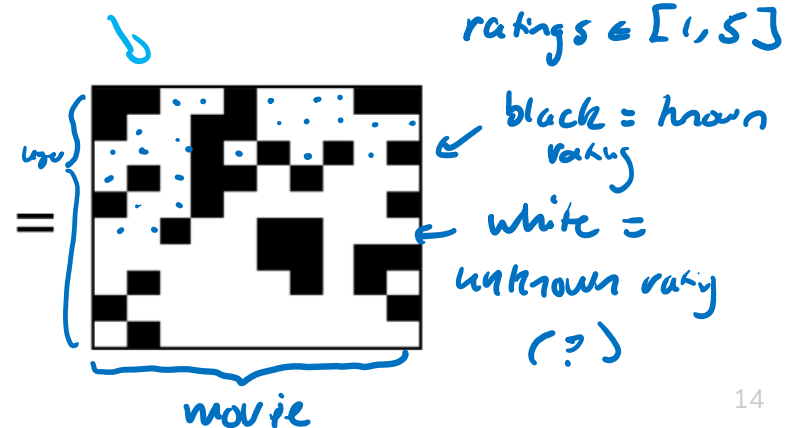
Challenge: Users have rated relatively few of the entire catalog

Can think of this as a matrix of users and ratings with missing data!

Input Data

User	Movie	Rating
		★★★★☆
		★★★★★
		★★★☆☆
		★★★★☆
		★★★☆☆
		★★★★★
		★★★★★
		★★★★★
		★★★★☆

					
User 1	5				3
User 2		2		4	
User 3			3		
User 4	1				
User 5			4		
User 6		5			2



Assumption

Matrix completion is an impossible task without some assumptions on data (unknowns could be anything otherwise).

Assume: There are k types of movies (e.g. action, romance, etc.) which users have various interests in.

This means we can describe a movie v with feature vector \underline{R}_v

- How much is the movie action, romance, drama, ...
- Example: $\underline{R}_v = [\underline{0.3}, \underline{0.01}, \underline{1.5}, \dots] \in \mathbb{R}^k$
action mystery romance

We can describe each user u with a feature vector \underline{L}_u

- How much she likes action, romance, drama,
- Example: $\underline{L}_u = [\underline{2.3}, \underline{0}, \underline{0.7}, \dots] \in \mathbb{R}^k$
action mystery romance

Estimate rating for user u and movie v as

$$\widehat{Rating}(u, v) = \underline{L}_u \cdot \underline{R}_v = 2.3 \cdot 0.3 + 0 \cdot 0.01 + 0.7 \cdot 1.5 + \dots$$

Example

"factors"

Suppose we have learned the following user and movie features using 2 features $k=2$

u_1
 u_2
 u_3
 u_4

User ID	Feature
1	(2, 0)
2	(1, 1)
3	(0, 1)
4	(2, 1)

L_u

v_1
 v_2
 v_3

Movie ID	Feature vector
1	(3, 1)
2	(1, 2)
3	(2, 1)

R_v

$L: 4 \times 2$ ($n \times k$)

$R: 3 \times 2$ ($m \times k$)

Then we can predict what each user would rate each movie

n

2	0
1	1
0	1
2	1

L

m

3	1	2
1	2	1

R^T

k

$=$

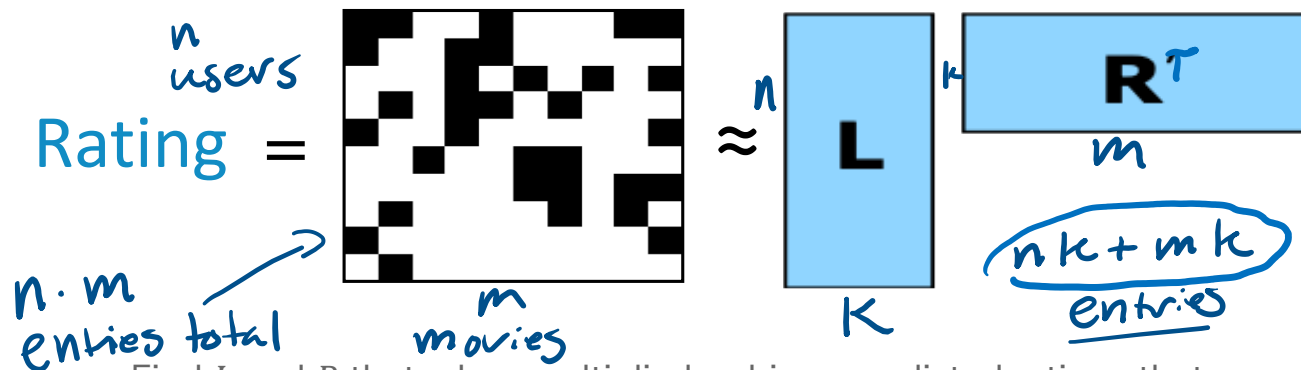
n

6	2	4
4	3	3
1	2	1
7	4	5

m

Rating (u_i, v_j)

Matrix Factorization



Find L and R that when multiplied, achieve predicted ratings that are close to the values that we have data for.

Our quality metric will be (could use others)

$$\hat{L}, \hat{R} = \underset{L, R}{\operatorname{argmin}} \sum_{u, v: r_{uv} \neq ?} \underbrace{(L_u \cdot R_v)}_{\text{predicted rating}} - \underbrace{r_{uv}}_{\text{actual rating}}^2$$

entries we have ratings for

Unique Solution?

Is this problem well posed? Unfortunately, there is not a unique solution ☹️

For example, assume we had a solution

$$\begin{array}{|c|c|c|} \hline 6 & 2 & 4 \\ \hline 4 & 3 & 3 \\ \hline 1 & 2 & 1 \\ \hline 7 & 4 & 5 \\ \hline \end{array} = \begin{array}{|c|c|} \hline L & \\ \hline 2 & 0 \\ \hline 1 & 1 \\ \hline 0 & 1 \\ \hline 2 & 1 \\ \hline \end{array} \begin{array}{|c|c|c|} \hline R^T & \\ \hline 3 & 1 & 2 \\ \hline 1 & 2 & 1 \\ \hline \end{array}$$

Then doubling everything in L and halving everything in R is also a valid solution. The same is true for all constant multiples.

$$\begin{array}{|c|c|c|} \hline 6 & 2 & 4 \\ \hline 4 & 3 & 3 \\ \hline 1 & 2 & 1 \\ \hline 7 & 4 & 5 \\ \hline \end{array} = \begin{array}{|c|c|} \hline L & \\ \hline 4 & 0 \\ \hline 2 & 2 \\ \hline 0 & 2 \\ \hline 4 & 2 \\ \hline \end{array} \begin{array}{|c|c|c|} \hline R^T & \\ \hline 1.5 & 0.5 & 1.0 \\ \hline 0.5 & 1.0 & 0.5 \\ \hline \end{array}$$

Coordinate Descent

Find \hat{L} & \hat{R}

Remember, our quality metric is

$$\hat{L}, \hat{R} = \operatorname{argmin}_{L, R} \sum_{u, v: r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

Gradient descent is not used much in practice to optimize this, since it is much easier to implement **coordinate descent** (i.e. Alternating Least Squares) to find \hat{L} and \hat{R}



Coordinate Descent

Goal: Minimize some function $g(w) = g(w_0, w_1, \dots, w_D)$

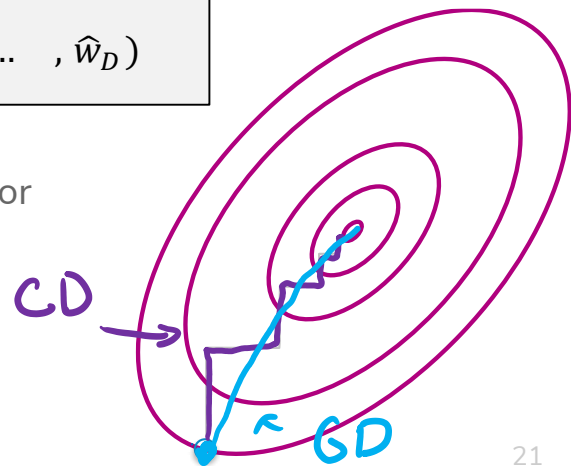
Instead of finding optima for all coordinates, do it for one coordinate at a time!

Initialize $\hat{w} = 0$ (or smartly)
while not converged:
 pick a coordinate j
 $\hat{w}_j = \underset{w}{\operatorname{argmin}} g(\hat{w}_0, \dots, \hat{w}_{j-1}, \underline{w}, \hat{w}_{j+1}, \dots, \hat{w}_D)$

To pick coordinate, can do round robin or pick at random each time.

Guaranteed to find an optimal solution under some constraints

Strongly convex, smooth



Coordinate Descent for Matrix Factorization

$V_u =$ set of all movies rated by user u

$$\hat{L}, \hat{R} = \operatorname{argmin}_{L, R} \sum_{u, v: r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

Fix movie factors R and optimize for L_u

First key insight:

$$\begin{aligned} & \operatorname{argmin}_{L_1, \dots, L_n} \sum_{u, v: r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 \\ &= \operatorname{argmin}_{L_1, \dots, L_n} \sum_u \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 \end{aligned}$$

Independent optimization for each user u !

$$\hat{L}_u = \operatorname{argmin}_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$$

Coordinate Descent for Matrix Factorization

Holding movies fixed, we can solve for each user separately!

For each user u

$$\hat{L}_u = \operatorname{argmin}_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$$

Handwritten annotations:
- An arrow points from "true data" to r_{uv} .
- An arrow points from "parameters" to L_u .
- An arrow points from "input (fixed)" to R_v .

Second key insight:

Looks like linear regression!

$$\operatorname{argmin}_w \sum_{i=1}^n (w^T h(x_i) - y_i)^2$$



Overall Algorithm

Want to optimize

$$\hat{L}, \hat{R} = \operatorname{argmin}_{L, R} \sum_{u, v: r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

Fix movie factors R , and optimize for user factors separately

- Independent least squares for each user

$$\hat{L}_u = \operatorname{argmin}_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L_u\|$$

Fix user factors, and optimize for movie factors separately

- Independent least squares for each movie

$$\hat{R}_v = \operatorname{argmin}_{R_v} \sum_{u \in U_v} (L_u \cdot R_v - r_{uv})^2 + \lambda_v \|R_v\|$$

↑
Users who rated movie v

Repeatedly do these steps until convergence (to local optima)

System might be underdetermined: Use regularization

Consider we had the ratings matrix

	Movie 1	Movie 2
User 1	4	?
User 2	?	2

During one step of optimization, user and movie factors are

$k=3$

	User Factors
User 1	[1, 2, 1]
User 2	[1, 1, 0]

	Movie Factors
Movie 1	[2, 1, 0]
Movie 2	[0, 0, 2]

Two questions:

What is the current residual sum of squares loss? (number)

If the next step of coordinate descent updates the user factors, which factors would change?

- User 1
- User 2
- User 1 and 2
- None



Think 

3 minutes

pollev.com/cs416

Consider we had the ratings matrix

	Movie 1	Movie 2
User 1	4	2
User 2	?	0

$$\widehat{\text{Rating}}(U_i, V_j) = L_i \cdot R_j = 1 \cdot 2 + 2 \cdot 1 + 1 \cdot 0 = 4$$

During one step of optimization, user and movie factors are

	User Factors
User 1	[1, 2, 1]
User 2	[1, 1, 0]



	Movie Factors
Movie 1	[2, 1, 0]
Movie 2	[0, 0, 2]

Two questions:

$$RSS(L, R) = (4 - 4)^2 + (2 - 0)^2 = 4$$

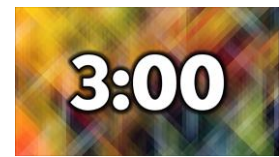
What is the current residual sum of squares loss? (number)

If the next step of coordinate descent updates the user factors, which factors would change?

- User 1 
- User 2 
- User 1 and 2
- None

$$\hat{L}_u = \underset{L_u}{\operatorname{argmin}} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2$$

Error is 0 for user 0, no update



Using Results

Use movie factors \hat{R} to discover “topics” for movie v : \hat{R}_v

Use user factors \hat{L} to discover “topic preferences” for user u : \hat{L}_u

Predict how much a user u will like a movie v

$$\widehat{Rating}(u, v) = \hat{L}_u \cdot \hat{R}_v$$

Recommendations: Sort movies user hasn't watched by

$\widehat{Rating}(u, v)$

- Recommend movies with highest predicted rating





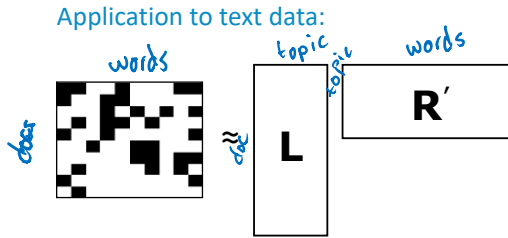
Brain Break



Topics

The “features” found by matrix factorization don’t always correspond to something meaningful (like film genre), but sometimes they do!

- Remember, the exact values are meaningless since we can scale them an infinite number of ways, but directions might mean something



partylaw government election court president elected council general minister political national members committee united office federal member house parliament vote public elections democratic vote

son died his childhood years family campaign married family king daughter john death william father born wife royal ireland irish henry house lord charles sir prince brother children england queen duke thomas years marriage george east edward english succeeded robert i member castle

school students university high college schools education year program student campus community programs training member members science national years public academic association officers arts educational inside class middle department teachers colleges classes offers activities companies district engineering learning founded faculty gifts sports children boy international board teaching academy secondary established

york county american united city washington john texas served virginia pennsylvania war moved ohio chicago william george james died army centuries dynasty royal yorkshire led byzantine defeated ruled great throne others digital castle military late tv middle center

season team game league games birds small long large animals bird plants genus giant natural habitats tree fern tropical white black order leaves brown summer keeps three arches flowers eggs worldwide feed coop sculptures wild length male breeding habitats range food female full short means endemic forest group including include moist threatened tall

album band radio station news television channel broadcast stations network media tv broadcasting time format local program bbc programming live tv morning host began sports live tv radio call hosted overnight music pre weekday daily channels digital also aired changed current broadcasted communications programme day broadcasts moved obs week assembly late night

century king roman empire greek production built engines speed vehicles designed produced power front system version type series motor rear standard gun company introduced range ford sold five wheel tank fleet factory machine developed based released wheels time powered small high weight electric body moving early

art museum work works artists collection design arts painting artist gallery paintings exhibition style painter 16 early created sculpture artists history contemporary collections prints museum worked images time photographs figures academies exhibitions modern portrait photographs began studio drawing include exhibited produced designed period visual

age 18 population income average years median living 65 males females households 100 family people families older town size city household miles density american township later area sparsity issues census 2000 issues 42 28 84 children 24 44 white females land including girls housing times individuals located poverty united village

war army military forces battle force british command general navy ship division ships troops corps commander infantry attack men officer fleet soldiers units officers operations unit pine august brigade jay the sailing master sailmaker ship operations captain september three enemy united soldier sea royal german marine major

white red black blue called color will head green gold side small hand long arms top flag horse wear silver common light dog wood body type large yellow farm worn dogs cut popular left generally traditional ball front horses phase ball field colors line used three specialty modern face cross

music musical opera festival orchestra dance performance works concert symphony composer original performances instruments musicians classical including work composed major stage songs folk instrument ball competition composers play performing ensemble playing stage years inside popular choir ensemble sound style time wide ball three chamber recordings string



Poll Everywhere

Think 

1 min

~~pollev.com/cs416~~

Which of the following are true about matrix factorization for recommendation systems?

- A. Provides personalization ✓
- B. Captures context (e.g. time of day) ✗
- C. Solves the cold start problem ✗

1:00

Blending Models: Featurized Matrix Factorization

Final Solution

Cold Start Again

$\|L\|_F^2$ is the L2 norm equivalent for matrices

Consider a new user u' and we want to predict their ratings

No previous ratings for them so: for all v , $r_{u'v} = ?$

Objective

$$\hat{L}, \hat{R} = \operatorname{argmin}_{L, R} \sum_{u, v: r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_U \|L\|_F^2 + \lambda_V \|R\|_F^2$$

Optimal user factor: $L_{u'} = 0$ because there is only penalty

Therefore, for all v , $\hat{r}_{u'v} = 0$ which seems like a problem

Blend Models

Idea: Learn a model to supplement the matrix factorization model!

Create a feature vector for each movie

$$h(v) = \left(\begin{array}{c} \text{genre} \quad \text{year} \quad \text{director} \\ \text{'action'} \quad 1994 \quad \text{'...'} \quad , \dots \end{array} \right)$$

Define weights on these features for **all users** (global)

$$w_G \in \mathbb{R}^d$$

Fit linear model : $\hat{r}_{uv} = w_G^T h(v)$

$$\hat{w}_G = \underset{w}{\operatorname{argmin}} \sum_{u,v: r_{uv} \neq ?} (w_G^T h(v) - r_{uv})^2 + \lambda \|w\|$$

Add Personalization

Of course, not all users have same preferences.

Include a user-specific deviation from global model

$$\hat{r}_{uv} = (\hat{w}_G + \hat{w}_u)^T h(v)$$

New user $\Rightarrow \hat{w}_u = 0$ but can
be updated over time

Can also add user specific features to model

$$h(u, v) = \left(\begin{array}{cc} \text{genre} & \text{year} & & \text{gender} & \text{age} \\ \text{action, 1994, \dots,} & & & \text{F,} & \text{25, \dots} \end{array} \right)$$

Learn \hat{w}_G, \hat{w}_u for $h(u, v)$

Featurized Matrix Factorization

Feature-based approach

- Feature representation of user and movie fixed
- Can address cold start problem

Matrix factorization approach

- Suffers from cold start problem
- User & Movie features are learned from data

A unified model

$$\hat{r}_{uv} = f \left(\underbrace{\hat{L}_u \cdot \hat{R}_v}_{\text{MF}}, \underbrace{(\hat{w}_G + \hat{w}_u)^T h(u,v)}_{\text{FB}} \right)$$

f is some arbitrary method to combine results

Evaluating Recommendations

Accuracy?

Could we use classification accuracy to identify which recommender system is performing best?

- We don't really care to predict what a person **does not** like
- Instead, we want to find the relatively few items from the catalog that they will like
- Similar to a class imbalance

Instead, we want to look at our set of recommendations and ask:

- How many of our recommendations did the user like? *Precision*
- How many of the items that the user liked did we recommend? *Recall*

Sound familiar?



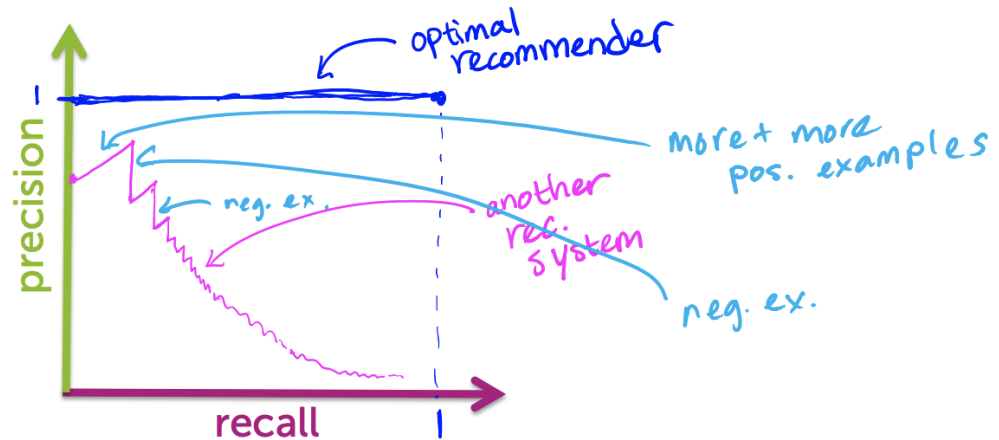
Precision - Recall

Precision and recall for recommender systems

$$\text{precision} = \frac{\# \text{ liked \& shown}}{\# \text{ shown}}$$

$$\text{recall} = \frac{\# \text{ liked \& shown}}{\# \text{ liked}}$$

For a given recommender system, plot precision and recall for different number of recommended items



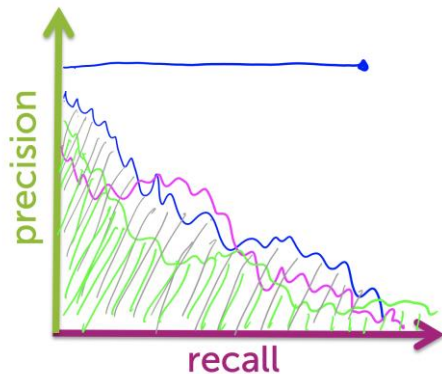
Comparing Algorithms

In general, it depends

- What is true always is that for a given precision, we want recall to be as large as possible (and vice versa)
- What target precision/recall depends on your application

One metric: area under the curve (**AUC**)

Another metric: Set desired recall and maximize precision (**precision at k**)



Optional: Offline Replay

Other approach: A/B Testing

Challenge: Deploying a new recommender system in-the-wild can be scary! Might not have confidence in how it will do.

One clever idea: use a simple model to start and record it's predictions and if the user interacted with the content (e.g., watched movie). Use this data of the model in deployment to test new models without needing to deploy them!

Test your new model on this logged data, if your model makes the same prediction as the simple model did originally, use that as a test example!

Deployed : (U_1, A, X) (U_2, B, \checkmark) (U_3, B, \checkmark)
Experimental : ~~(U_1, B, \checkmark)~~ (U_2, B, \checkmark) (U_3, B, X)

user pred. label

Recap

Now you can:

- Describe the goal of a recommender system
- Provide examples of applications where recommender systems are useful
- Implement a co-occurrence based recommender system
- Describe the input (observations, number of “topics”) and output (“topic” vectors, predicted values) of a matrix factorization model
- Implement a coordinate descent algorithm for optimizing the matrix factorization objective presented
- Exploit estimated “topic” vectors to make recommendations
- Describe the cold-start problem and ways to handle it (e.g., incorporating features)
- Analyze performance of various recommender systems in terms of precision and recall
- Use AUC or precision-at-k to select amongst candidate algorithms