

# **15-388/688 - Practical Data Science: Recommender systems**

J. Zico Kolter  
Carnegie Mellon University  
Spring 2018

# Outline

Recommender systems

Collaborative filtering

User-user and item-item approaches

Matrix factorization

# Outline

Recommender systems

Collaborative filtering

User-user and item-item approaches

Matrix factorization

# Recommender systems

amazon Prime

All ▾

COUNTDOWN TO BLACK FRIDAY DEAL

Hello, Zico Your Account ▾ Prime ▾ Lists ▾

Departments ▾ Browsing History ▾ Zico's Amazon.com Today's Deals Gift Cards & Registry Sell Help

Recommendations for you in Toys & Games

The screenshot shows the Amazon website interface. At the top, there's a search bar and a 'COUNTDOWN TO BLACK FRIDAY DEAL' banner. Below the header, user account information is displayed. The main content area features a heading 'Recommendations for you in Toys & Games'. It displays several product images: a 'Pocket Flash Cards' book titled 'Alphabet' with an apple illustration; two 'Scooby-Doo' activity books; a 'MY 1ST BATH BOOK' with a cow and other animals; a 'Teach My TODDLER' kit with flashcards for shapes, numbers, and colors; and a large collection of 'PLAYSKOOL' educational flash cards for Math, Language, and Basic Skills.

NETFLIX

Browse ▾ DVD

Search

Because you watched Curious George: A Halloween Boo Fest

The screenshot shows the Netflix website interface. At the top, there's a search bar and a 'Search' button. The main content area features a recommendation message 'Because you watched Curious George: A Halloween Boo Fest'. Below this, a row of six movie and TV show covers is displayed: 'Curious George: A Very Monkey Christmas', 'Clifford the Big Red Dog', 'The Tiger Who Came to Tea', 'Pooh's Grand Adventure', 'Alvin & the Chipmunks Meet Frankenstein', 'Little Einsteins Children's Favorites: Halloween Treats', and 'Masha and the Bear'.

# Information we can use to make predictions

“Pure” user information:

- Age
- Location
- Profession

“Pure” item information:

- Movie budget
- Main actors
- (Whether it is a Netflix release)

User-item information:

- Which items are most similar to those I have bought before?
- What items have users most similar to me bought?

# Supervised or unsupervised?

Do recommender systems fit more within the “supervised” or “unsupervised” setting?

Like supervised learning, there are known outputs (items that the user purchases), but like unsupervised learning, we want to find structure/similarity between users/items

We won’t worry about classifying this as just one or the other, but we will again formulate the problem within the three elements of a machine learning algorithm: 1) hypothesis function, 2) loss function, 3) optimization

# Challenges in recommender systems

There are many challenges beyond what we will consider here in recommender systems:

1. Lack of user ratings / only “presence” data
2. Balancing personalization with generic “good” items
3. Privacy concerns

# Historical note: Netflix Prize



Public competition ran from 2006 to 2009, goal was to produce a recommender system with 10% improvement in RMSE over existing Netflix system (based upon user-user Pearson correlation), \$1M prize

Sparked a great deal of research in collaborative filtering, especially matrix factorization techniques

Larger impacts: put “data science competitions” in the public eye, emphasized practical importance of ensemble methods (though winning solution was never fielded)

# Outline

Recommender systems

Collaborative filtering

User-user and item-item approaches

Matrix factorization

# Collaborative filtering

Collaborative filtering refers to recommender systems that make recommendations based solely upon the preferences that other users have indicated for these item (e.g., past ratings)

The mathematical setting to have in mind is that of a matrix with mostly unknown entries

$$X = \begin{bmatrix} 1 & 2 & 5 \\ 3 & 4 & 5 \\ 4 & 4 & 3 \end{bmatrix}$$

columns correspond to different items

rows correspond to different users

entries correspond to known (given by user) scores for that user, for that items

# Matrix view of collaborative filtering

Collaborative filtering  $X$  matrix is sparse, but unknown entries do not correspond to zero, are just missing

Goal is to “fill in” the missing entries of the matrix

$$X = \begin{bmatrix} 1 & ? & ? & 3 \\ ? & 2 & 5 & ? \\ ? & 3 & ? & 5 \\ 4 & ? & 4 & ? \end{bmatrix}$$

# Approaches to collaborative filtering

**User – user approaches:** find the users that are most similar to myself (based upon only those items that are rated for *both* of us), and predict scores for other items based upon the average

**Item – item approaches:** find the items most similar to a given item (based upon all users rated both items), and predict scores for other users based upon the average

**Matrix factorization approaches:** find some low-rank decomposition of the  $X$  matrix that agrees at observed values

# Outline

Recommender systems

Collaborative filtering

User-user and item-item approaches

Matrix factorization

# User-user and item-item approaches

**Basic intuition of user-user approach:** find other users who are similar to me, e.g. by correlation coefficient or cosine similarity, look at how they ranked other items that I did not rank

**One difference:** correlation coefficient, etc, are only defined for vectors of the same size, so we only typically compute correlation across items that both users ranked

$$X = \begin{bmatrix} 1 & ? & ? & 3 \\ ? & 2 & 5 & ? \\ ? & 3 & ? & 5 \\ 4 & ? & 4 & ? \end{bmatrix}$$

Item-item approaches do the same thing but by row instead of column

## User-user approach: formally

To match with our previous notation as much as possible, we will our prediction of  $X_{ij}$  as  $\hat{X}_{ij}$  (we will later also refer to this as  $h_\theta(i, j)$ , our hypothesis evaluated on point  $i, j$ )

User-user methods typically make predictions:

$$\hat{X}_{ij} = \bar{x}_i + \frac{\sum_{k:X_{kj} \neq 0} w_{ik} (X_{kj} - \bar{x}_k)}{\sum_{k:X_{kj} \neq 0} w_{ik}}$$

- $\bar{x}_i$  - mean of user  $i$ 's ratings
- $w_{ik}$  - similarity function between users  $i$  and  $k$

Common modification: restrict sum to only  $K$  users “most similar” to  $i$

# Similarity measures

How do we measure similarity between two users?

Two example approaches:

1. Pearson correlation ( $\mathcal{I}_{ik}$  denotes items ranked by users  $i$  and  $k$ ):

$$w_{ik} = \frac{\sum_{j \in \mathcal{I}_{ik}} (X_{ij} - \bar{x}_i)(X_{kj} - \bar{x}_k)}{\left( \sum_{j \in \mathcal{I}_{ik}} (X_{ij} - \bar{x}_i)^2 \cdot \sum_{j \in \mathcal{I}_{ik}} (X_{kj} - \bar{x}_k)^2 \right)^{1/2}}$$

2. Raw cosine similarity (treating missing as zero):

$$w_{ik} = \frac{\sum_j X_{ij} \cdot X_{kj}}{\left( \sum_j X_{ij}^2 \cdot \sum_j X_{kj}^2 \right)^{1/2}}$$

# Poll: predictions of user-user method

Given the following matrix, what would a user-user method based upon Pearson correlation predict for the following point:

$$X = \begin{bmatrix} 1 & ? & ? & 3 \\ ? & 2 & 5 & ? \\ ? & 3 & ? & 5 \\ 4 & ? & 4 & ? \end{bmatrix}$$

1. 2
2. 2.5
3. 3
4. 3.5
5. 4

$$\hat{X}_{ij} = \bar{x}_i + \frac{\sum_{k:X_{kj} \neq 0} w_{ik} (X_{kj} - \bar{x}_k)}{\sum_{k:X_{kj} \neq 0} w_{ik}}$$

$$w_{ik} = \frac{\sum_{j \in \mathcal{I}_{ik}} (X_{ij} - \bar{x}_i)(X_{kj} - \bar{x}_k)}{\left( \sum_{j \in \mathcal{I}_{ik}} (X_{ij} - \bar{x}_i)^2 \cdot \sum_{j \in \mathcal{I}_{ik}} (X_{kj} - \bar{x}_k)^2 \right)^{1/2}}$$

# Item-item approaches

Item-item approaches just do the same process flipping rows/columns

Make predictions:

$$\hat{X}_{ij} = \bar{x}_j + \frac{\sum_{k:X_{ik} \neq 0} w_{jk} (X_{ik} - \bar{x}_k)}{\sum_{k:X_{ik} \neq 0} w_{jk}}$$

Similarity function, e.g.:

$$w_{jk} = \frac{\sum_{i \in \mathcal{I}_{jk}} (X_{ij} - \bar{x}_j)(X_{ik} - \bar{x}_k)}{\left( \sum_{i \in \mathcal{I}_{jk}} (X_{ij} - \bar{x}_j)^2 \cdot \sum_{i \in \mathcal{I}_{jk}} (X_{ik} - \bar{x}_k)^2 \right)^{1/2}}$$

## Poll: efficiency of user and item based method

Suppose we have many more users than items. Assuming we use dense matrix operations for everything, which method would be more efficient for computing *all* the predictions  $\hat{X}_{ij}$  for all missing elements?

1. The user-user approach will be more efficient
2. The item-item approach will be more efficient
3. They will both have the same complexity

# Outline

Recommender systems

Collaborative filtering

User-user and item-item approaches

Matrix factorization

# Matrix factorization approach

Approximate the  $i, j$  entry of  $X \in \mathbb{R}^{m \times n}$  as  $\hat{X}_{ij} = u_i^T v_j$  where  $u_i \in \mathbb{R}^k$  denotes user-specific weights and  $v_j \in \mathbb{R}^k$  denotes item-specific weights

1. Hypothesis function

$$\hat{X}_{ij} = h_{\theta}(i, j) = u_i^T v_j, \quad \theta = \{u_{1:m}, v_{1:n}\}$$

2. Loss function: squared error (on observed entries)

$$\ell(h_{\theta}(i, j), X_{ij}) = (h_{\theta}(i, j) - X_{ij})^2$$

leads to optimization problem ( $S$  denotes set of observed entries)

$$\underset{\theta}{\text{minimize}} \sum_{i,j \in S} \ell(h_{\theta}(i, j), X_{ij})$$

# Optimization approaches

3. How do we optimize the matrix factorization objective? (Like k-means, EM, possibility of local optima)

Consider the objective with respect to a *single*  $u_i$  term:

$$\underset{u_i}{\text{minimize}} \sum_{j:(i,j) \in S} (v_j^T u_i - X_{ij})^2$$

This is just a least-squares problem, can solve analytically:

$$u_i = \left( \sum_{j:(i,j) \in S} v_j v_j^T \right)^{-1} \left( \sum_{j:(i,j) \in S} v_j X_{ij} \right)$$

**Alternating minimization algorithm:** Repeatedly solve for all  $u_i$  for each user,  $v_j$  for each item (may not give global optimum)

# Matrix factorization interpretation

What we are effectively doing here is factorizing  $X$  as a low rank matrix

$$X \approx UV, \quad U \in \mathbb{R}^{m \times k}, V \in \mathbb{R}^{k \times n}$$

where

$$U = \begin{bmatrix} - u_1^T - \\ \vdots \\ - u_m^T - \end{bmatrix}, \quad V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_n \\ | & & | \end{bmatrix}$$

However, we are *only* requiring the  $X$  match the factorization at the observed entries of  $X$

## Relationship to PCA

PCA also performs a factorization of  $X \approx UV$  (if you want to follow the precise notation of the PCA slides, it would actually be  $X^T = UV$  where  $V$  contains the columns  $Wx^{(i)}$ )

But unlike collaborative filtering, in PCA, *all* the entries of  $X$  are observed

Though we won't get into the details: this difference is what lets us solve PCA exactly, while we can only solve matrix factorization for collaborative filtering locally