# Collaborative Filtering Implementation and Evaluation

## Recommender System Research for Ecommerce Data

Jun Yang

School of Electronic and Computer Engineering, Peking University

June 12, 2014

## Overview

# Content

**Introduction**   Item-based Collaborative Filtering   Implementation   Evaluation
○●○○○○○                              ○○○○○○○○○○○○○○          ○○○○○○○○○○○○○○○○○○○○○○○○○

Recommender System

# Content

# What is Recommender System?

Identify a set of items that will be of interest to a certain user.

- Similarity between items/users will be analyzed
- Historical data will be used



$1,000,000 Prize

# Categories of Recommeder System

- Collaborative Filtering (CF)
  - User-based: Facebook, LinkedIn
  - Item-based: Amazon
- Content-based
  - Information-Retrieval: tf-idf, LSM
  - Machine-Learning: clusters, classifiers, neural networks

# Content

**Introduction**
○○○○●○○

Item-based Collaborative Filtering
○○○○○○○○○○○○○

Implementation
○○○○○○○○○○○○○○

Evaluation
○○○○○○○○○○○○○○○○○○○○○○○○

Collaborative Filtering

# User-based Collaborative Filtering

1. Identify a neighborhood of people with similar behavior
2. Analyze this neighborhood to find out recommends

$$P_{u,i} = \sum_{Rank(s_{u,v}) > k} s_{u,v} \cdot R_{v,i} \tag{1}$$

**Introduction**
ooooooeo

Item-based Collaborative Filtering

Implementation
oooooooooooooo

Evaluation
oooooooooooooooooooooooo

Collaborative Filtering

# Problems with User-based CF

### Sparsity

Limited information for a certain user caused inaccuracy when identifying neiborhood, thus poor recommendations.

### Scalability

Computation grows linearly with the number of users.

### New User Problem

We have nothing to recommend to new users.

# Item-based Collaborative Filtering

1. Identify a neighborhood of items
2. Analyze this neighborhood to find out recommends

$$P_{u,i} = \sum_{Rank(s_{i,j}) > k} s_{i,j} \cdot R_{u,j} \qquad (2)$$

# Content

Introduction
0000000

Item-based Collaborative Filtering

Implementation
0000000000000

Evaluation
0000000000000000000000000

# Terminology

- Users $U$
  The collection of users.
- Items $I$
  The collection of items.
- Rating $R$
  $R_{i,j}$ represents user $i$'s rating for item $j$.
- Similarity $S$
  $S_{i,j}$ represents the similarity between user/item $i$ and $j$.
- Prediction $P$
  $P_{i,j}$ represents user $i$'s predicted rating for item $j$.

Introduction
0000000

Item-based Collaborative Filtering

Implementation
0000000000000

Evaluation
00000000000000000000000

## Algorithm I

1. Get input ratings: $R$
2. Compute similarities: $S$

$$S_{i,j} = sim(\vec{R_i}, \vec{R_j}) \tag{3}$$

where:
$\vec{R_i} = (R_{1,i}, R_{2,i}, \cdots, R_{N,i})$
$\vec{R_j} = (R_{1,j}, R_{2,j}, \cdots, R_{N,j})$

3. Identify $k$ neighbors $\{j_1, j_2, \cdots, j_k\}$ for each item

Introduction
○○○○○○○

**Item-based Collaborative Filtering**

Implementation
○○○○○○○○○○○○○

Evaluation
○○○○○○○○○○○○○○○○○○○○○○○○○○

# Algorithm II

4. For each user who bought a set of items $C$:

   1. Neighbors of $c \in C$ forms candidates: $N$, remove $n \in N$ that is already in $C$
   2. For each $n \in N$, compute its similarity with $C$ as the sum of similarities with $c \in C$
   3. Sort $N$ respect to that similarities.

5. The sorted $N$ for each user forms recommendation set.

# Content

Introduction    Item-based Collaborative Filtering    **Implementation**    Evaluation
○○○○○○○    ○○○○○○○○○○○○○○○○    ●○○○○○○○○○○○○○    ○○○○○○○○○○○○○○○○○○○○○○○○○○

Rating

# Content

# Data Format

- Formart:

```
1 #user_id      item_id type      date
  847750        2235    0         4.15
3 847750        2215    1         4.16
  6694750       14020   0         6.16
```

- Statistic:

| No. Users | 860 |
|-----------|-----|
| No. Items | 7842 |
| No. Lines | 42531 |
| Begin Date | 15th, April |
| End Date | 15th, August |

Introduction    Item-based Collaborative Filtering    **Implementation**    Evaluation
○○○○○○○    ○○○○○○○○○○○○○○○○○○    ○○●○○○○○○○○○○○○    ○○○○○○○○○○○○○○○○○○○○○○○○○○

Rating

# Behavior Statistics

Ratings matrix $R$ is required to do collaborative filtering, how to get $R$ from the dataset?

- If user $u$ bought item $i$, then $R_{u,i} = 1$

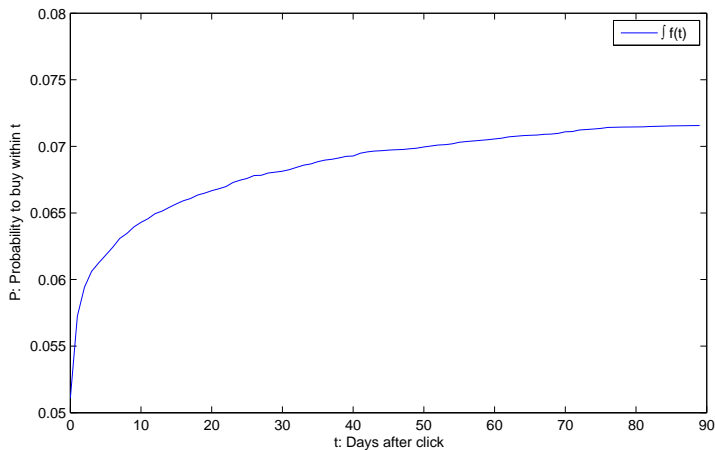## What about user $u$ clicked item $i$ ?

- When a user clicked an item, it's likely that he'll buy it. As time goes, the probability buy decreases.
- Make statistics about probability with respective to time:

$$p = f(t) \tag{4}$$

- The rating:

$$R_{u,i} = \int_{prediction\ timespan} f(t)dt \tag{5}$$

Introduction          Item-based Collaborative Filtering          Implementation          Evaluation
○○○○○○○                                                  ○○○●○○○○○○○○○○○          ○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Rating

# Integrated Probability for Click Behavior

# Content

1. Introduction
   - Recommender System
   - Collaborative Filtering

2. Item-based Collaborative Filtering

3. Implementation
   - Rating
   - Item Similarity
   - Prediction

4. Evaluation
   - Effect of Pre and Post Processing
   - Similarity Method Comparison
   - Summing and Normalization
   - Comparison with User-based CF

Introduction | Item-based Collaborative Filtering | **Implementation** | Evaluation
0000000 | | 000000000000000 | 000000000000000000000000

Item Similarity

# Similarity Methods

Similarity between items are usually computed in the space of *users*, treating each item as a vector. Classes of similarity methods:

- Cosine similarity
- Adjusted cosine similarity
- Pearson Correlation coefficient

## Cosine Similarity

$$sim(\vec{u}, \vec{v}) = \cos(\vec{u}, \vec{v}) = \frac{\vec{u} \cdot \vec{v}}{||\vec{u}|| \cdot ||\vec{v}||} \qquad (6)$$

- Similarity tends to be high when if each user who purchases u also purchases v as well.
- Frequently purchased items are de-emphasized by the denominator.
- The result range: $[0, 1]$

Introduction    Item-based Collaborative Filtering    **Implementation**    Evaluation
○○○○○○○    ○○○○○○○○○○○○○○○    ○○○○○○○○●○○○○○    ○○○○○○○○○○○○○○○○○○○○○○○○○○○○

Item Similarity

## Adjusted Similarity

$$sim(i,j) = \frac{\sum_{u \in U}(R_{u,i} - \overline{R_u})(R_{u,j} - \overline{R_u})}{\sqrt{\sum_{u \in U}(R_{u,i} - \overline{R_u})^2}\sqrt{\sum_{u \in U}(R_{u,j} - \overline{R_u})^2}} \quad (7)$$

- Different rating scales between users are taken into account.
- The result range: $[-1, 1]$

# Pearson Correlation Coefficient

$$sim(i,j) = \frac{\sum_{u \in U}(R_{u,i} - \overline{R_i})(R_{u,j} - \overline{R_j})}{\sqrt{\sum_{u \in U}(R_{u,i} - \overline{R_i})^2}\sqrt{\sum_{u \in U}(R_{u,j} - \overline{R_j})^2}} \tag{8}$$

- *Pearson − r* is a measure of the linear dependence between two variables.
- We need isolate co-rated cases in advance.
- The result range: $[-1, 1]$

Introduction          Item-based Collaborative Filtering          **Implementation**          Evaluation
○○○○○○○          ○○○○○○○○○○○○○○○○○○○○          ○○○○○○○○○○●○○○○          ○○○○○○○○○○○○○○○○○○○○○○○○○

Prediction

# Content

| Introduction | Item-based Collaborative Filtering | **Implementation** | Evaluation |
|---|---|---|---|
| ○○○○○○○ | | ○○○○○○○○○○○●○○○○ | ○○○○○○○○○○○○○○○○○○○○○○ |

Prediction

# Similarity Summing

As mentioned above, $P_{u,i}$ is computed by summing the similarities between item $i$ and items in user $u$'s basket:

$$P_{u,i} = \sum_{Rank(s_{i,j})>k} s_{i,j} \cdot R_{u,j} \tag{9}$$

### Better summing method?

- Treat recommendations from each item $j$ as independent events.
- Only when all recommendation fails, $R_{u,i} = 0$. Thus:

$$P_{u,i} = 1 - \prod_{Rank(s_{i,j})>k} (1 - s_{i,j} \cdot R_{u,j}) \tag{10}$$

# Similarity Normalization

Yet our summing method don't take into account the difference in density of the $k$ neighbors.

$$P_{u,i} = \sum_{Rank(s_{i,j}) > k} s_{i,j} \cdot R_{u,j} \tag{11}$$

We can normalize the k similarities so that they add-up to 1.

### Problem Case

An infrequently purchased item $i$ have a moderate overlap with another infrequently purchased item $j$, this will cause high similarity and lead to wrong recommendation.

# Row Normalization

For users who purchases a lot, each item reflects his appetite less.
We'll normalize $R$ before doing prediction.

# Post Procession

- If user $u$ once purchased item $i$, whether he'll purchase another depends on the lifetime of item $i$.
- Lifetime of item $i$:

$$\tau_i \approx \frac{\sum_{u \in U, R_{u,i}=1} R_{u,i}}{\sum_{u \in U, R_{u,i}=1} 1 \cdot TimeSpan} \quad (12)$$

- Thus the modified $R_{u,i}$ will be:

$$R'_{u,i} = \begin{cases} R_{u,i} & \text{if } R_{u,i} \neq 1 \\ \frac{Timespan_{prediction}}{\tau_i} \cdot R_{u,i} & \text{otherwise} \end{cases} \quad (13)$$

# Content

## Quality Measure

The quality of recommender system is measured by *recall* and *precision*.

- *recall* is the percentage of hits in the test set.
- *precision* is the percentage of hits in the prediction set.
- *F*1 is used to combine these two parameters:

$$F1 = \frac{2}{1/recall + 1/precision} \qquad (14)$$

## Train Set & Test Set

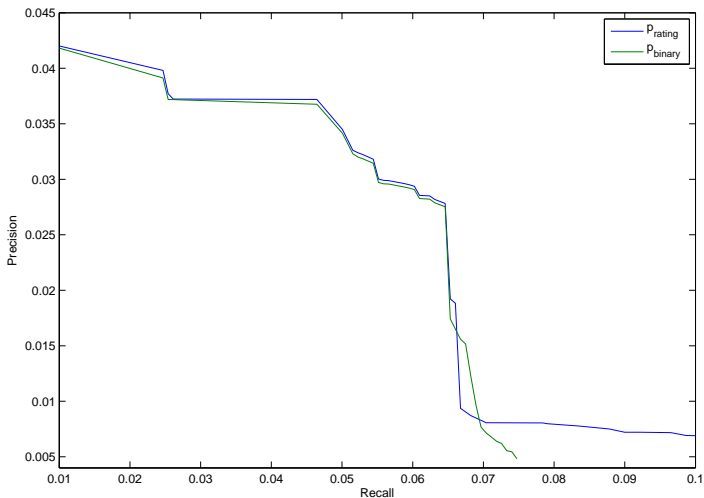In order to measure recommendation quality, we split the dataset into train set and test set:

|                 | Train Set  | Test Set    |
|-----------------|------------|-------------|
| Begin Date      | 15th April | 15th July   |
| End Date        | 14th July  | 14th August |
| No. Transaction | 4971       | 2013        |

# Content

Introduction
0000000

Item-based Collaborative Filtering

Implementation
00000000000000

Evaluation
0●000000000000000000000000

Effect of Pre and Post Processing

# Effect of Pre-processing

Introduction    Item-based Collaborative Filtering    Implementation    **Evaluation**
0000000         0000000000000                          0000000000000       00●00000000000000000000

Effect of Pre and Post Processing

# Effect of Pre-processing

Introduction
○○○○○○○

Item-based Collaborative Filtering

Implementation
○○○○○○○○○○○○○○

Evaluation
○○○●○○○○○○○○○○○○○○○○○○○○○○○○○

Effect of Pre and Post Processing

# Effect of Post-processing

Introduction        Item-based Collaborative Filtering        Implementation        **Evaluation**
0000000            0000000000000                              00000000000000      00000●000000000000000000000

Effect of Pre and Post Processing

# Effect of Post-processing

# Content

# Similarity Method Comparison

# Similarity Method Comparison

# Similarity Method Comparison

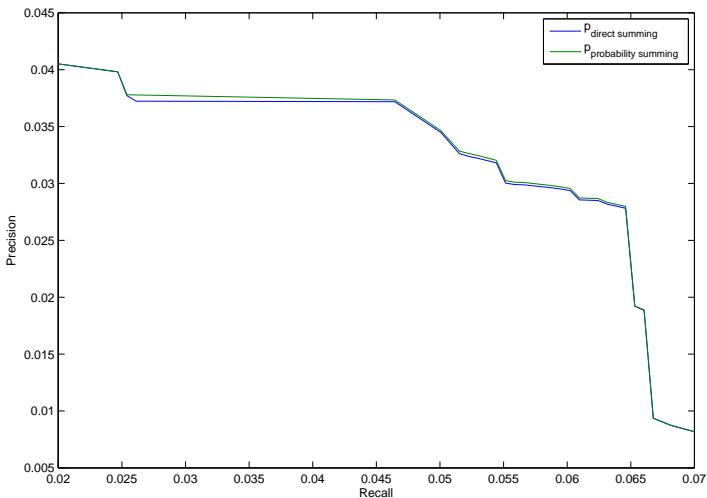- Correlation and Cosine are approximate.
- Adjusted cosine is bad.

### Why adjusted cosine not suitable?

Adjusted cosine removes the difference in rating scales, which means positive rating could contain negative information.
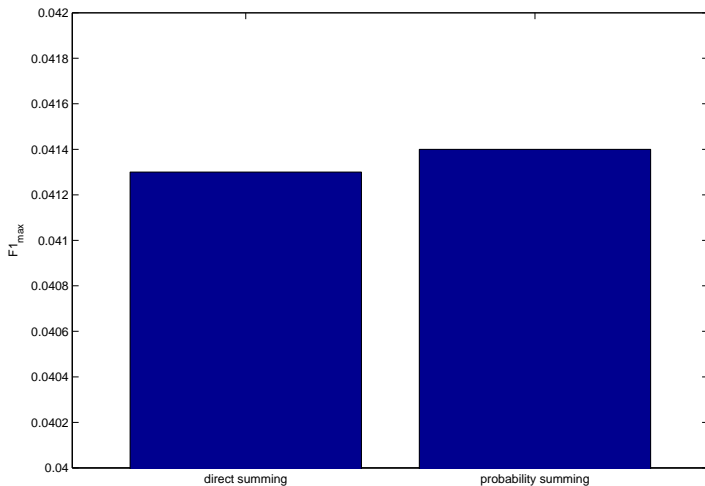While for ecommerce data, positive rating is always positive.

Introduction    Item-based Collaborative Filtering    Implementation    **Evaluation**
0000000         00000000000000                        000000000000000    0000000000●000000000000000

Summing and Normalization

# Content

Introduction
0000000

Item-based Collaborative Filtering

Implementation
00000000000000

Evaluation
0000000000●000000000000

Summing and Normalization

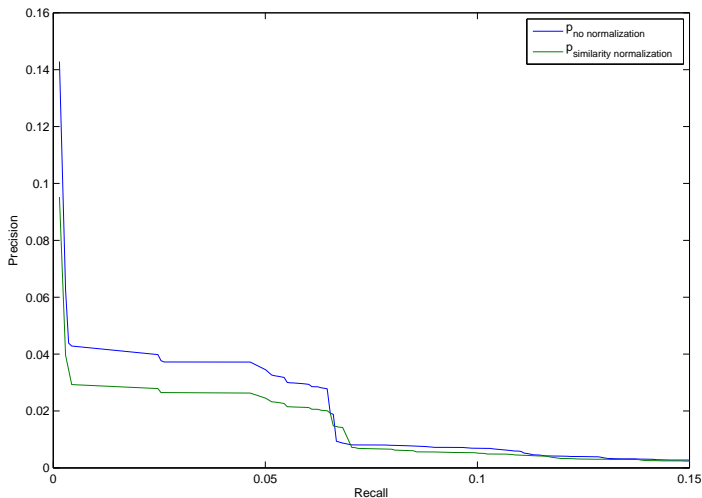# Effect of Similarity Summing Method

# Effect of Similarity Summing Method

# Similarity Normalization

Introduction
0000000

Item-based Collaborative Filtering

Implementation
00000000000000

Evaluation
000000000000000●000000000000

Summing and Normalization

# Similarity Normalization

Introduction    Item-based Collaborative Filtering    Implementation    **Evaluation**
0000000        0000000000000              000000000000000    000000000000000000000
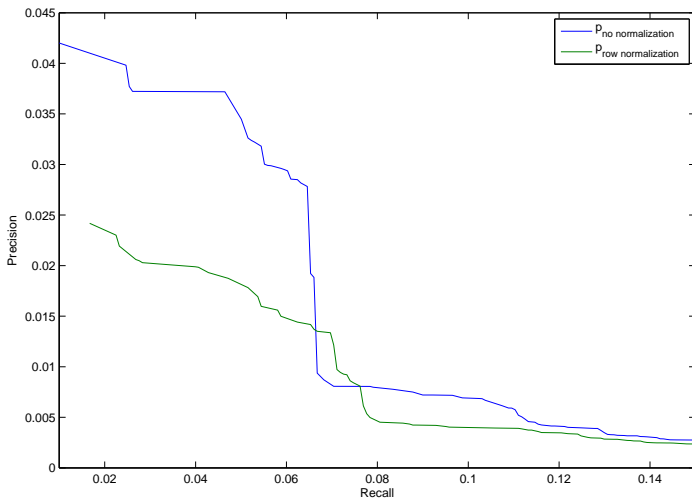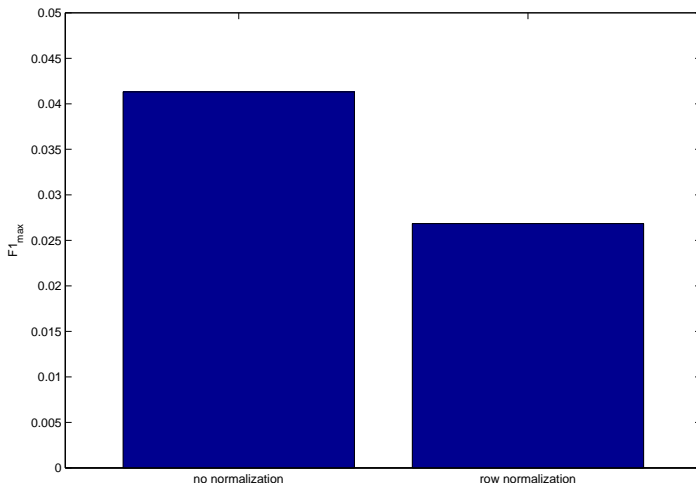
Summing and Normalization

# Similarity Normalization

- Dataset is so small, that many items do not have $k$ neighbors.
- The sparse item will divide a smaller denominator, rather than larger(expected).

Introduction          Item-based Collaborative Filtering          Implementation          Evaluation
○○○○○○○          ○○○○○○○○○○○○○○○          ○○○○○○○○○○○○○○          ○○○○○○○○○○○○○○○○●○○○○○○○

Summing and Normalization

# Row Normalization

# Row Normalization

Introduction
0000000

Item-based Collaborative Filtering

Implementation
0000000000000000

Evaluation
00000000000000000000000000000

Summing and Normalization

# Row Normalization

- Row normalization results in equal purchasing power for each user.
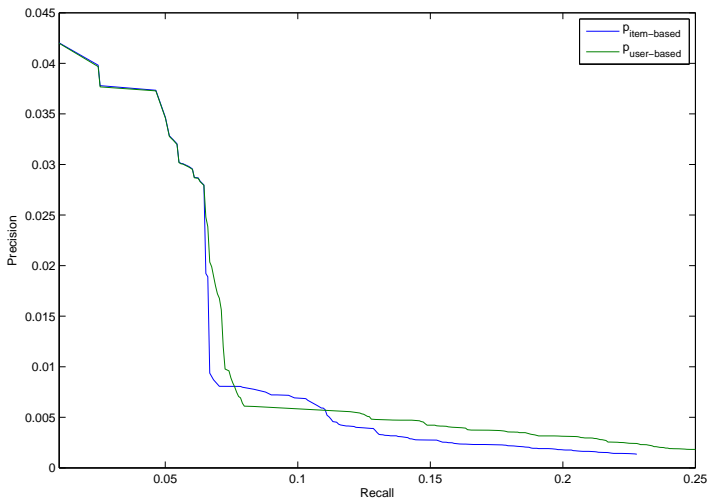- This side-effect is considerable. Users who purchase alot certainly will purchase a lot.

## Case for Top-N Recommendation

- Top-N recommendation system will recommend N item for each user.
- Normalizing purchasing power has no side-effect.
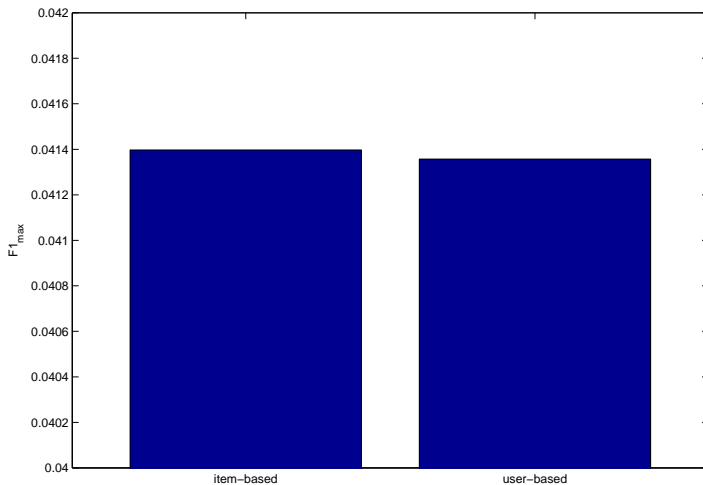- Row Normalization is useful now.

Introduction   Item-based Collaborative Filtering   Implementation   **Evaluation**
○○○○○○○   ○○○○○○○○○○○○○   ○○○○○○○○○○○○○   ○○○○○○○○○○○○○○○○○○●○○○○

Comparison with User-based CF

# Content

# P-R Relationship

Introduction
○○○○○○○

Item-based Collaborative Filtering

Implementation
○○○○○○○○○○○○○○

Evaluation
○○○○○○○○○○○○○○○○○○○○○○○●○○

Comparison with User-based CF

# Max F1

Introduction  Item-based Collaborative Filtering  Implementation  **Evaluation**
○○○○○○○  ○○○○○○○○○○○○○  ○○○○○○○○○○○○○○○○○○○○○○○●○

Thanks

# References

📄 Karypis G.
*Evaluation of item-based top-n recommendation algorithms[C].*
Proceedings of the 10th international conference on Information and knowledge management. ACM, 2001: 247-254.

📄 Sarwar B, Karypis G, Konstan J, et al.
*Item-based collaborative filtering recommendation algorithms[C].*
Proceedings of the 10th international conference on World Wide Web. ACM, 2001: 285-295.

# Thanks

Thank you!