# Data Intelligence

## Recommendation-1
## Content based & Collaborative Filtering

**U Kang**
**Seoul National University**

# In This Lecture

- Understand the motivation and the problem of recommendation

- Compare the content-based vs. collaborative filtering approaches for recommender system

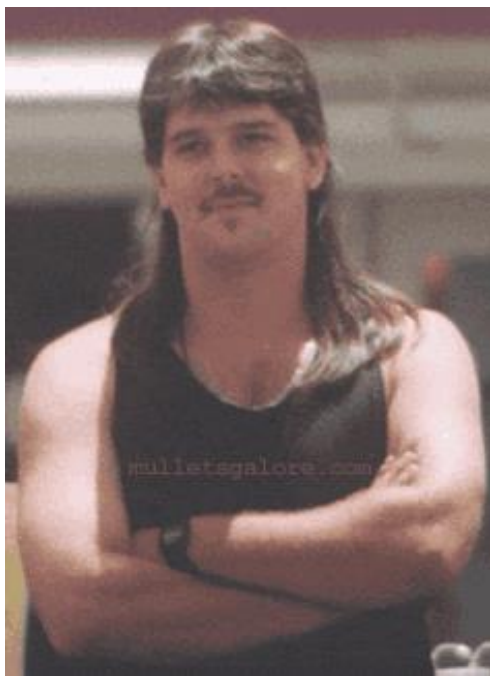- Learn how to evaluate methods for recommendation

# Outline

➡ ☐ **Overview**

☐ Content-based Recommender System

☐ Collaborative Filtering

☐ Evaluation & Complexity

# Example: Recommender Systems

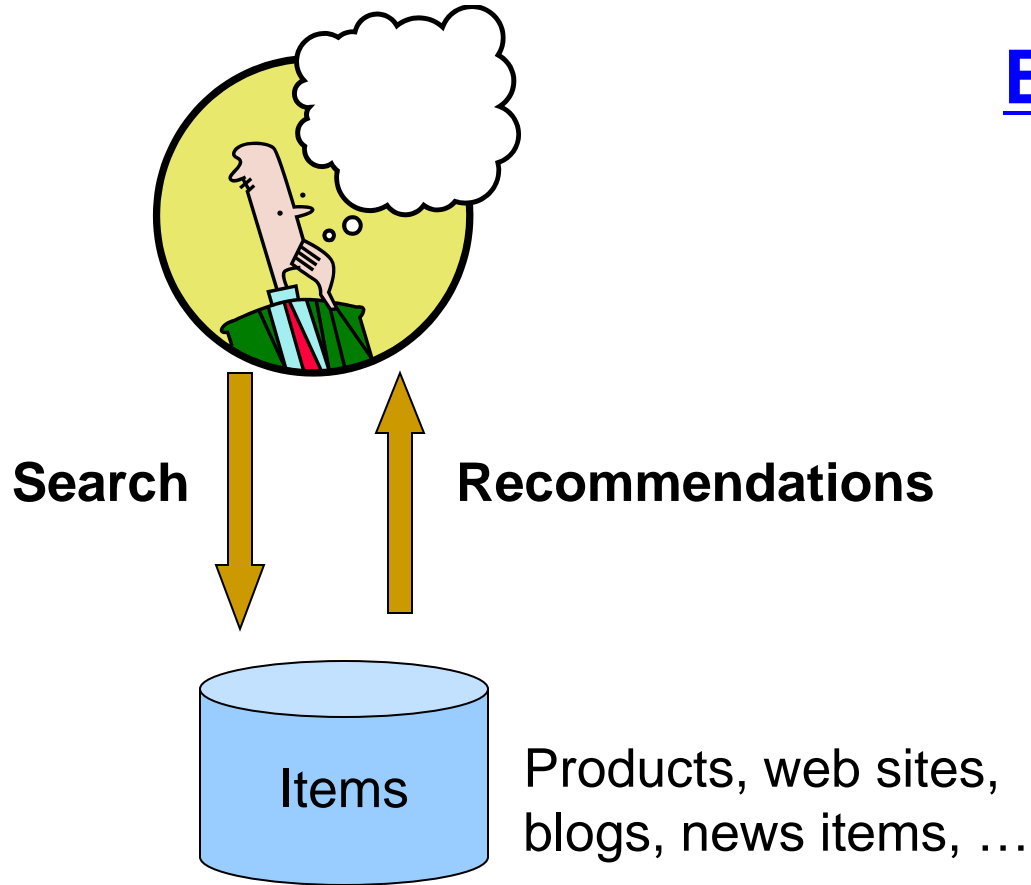

■ **Customer X**

❑ Buys Metallica CD

❑ Buys Megadeth CD

■ **Customer Y**

❑ Does search on Metallica

❑ Recommender system suggests Megadeth from data collected about customer **X**

U Kang                                                                 4

# Recommendations



**Search**  **Recommendations**

Items

Products, web sites, blogs, news items, …

## Examples:

amazon.com.

PANDORA

StumbleUpon

del.icio.us

NETFLIX

m o v i e l e n s
helping you find the *right* movies

last·fm
the social music revolution

Google News

You Tube

XBOX LIVE

# Offline vs. Online Recommendation

- **Offline recommendation: popular item**
  - Wall-mart: shelf space contains only 'popular' items
  - Also: TV networks, movie theaters,…

- **Web enables near-zero-cost dissemination of information about products**
  - Can recommend scarce items, too

- **More choice necessitates better filters**
  - Recommendation engines
  - How **Into Thin Air (1998)** made **Touching the Void (1988)** a bestseller: http://www.wired.com/wired/archive/12.10/tail.html

# Sidenote: The Long Tail



Sources: Erik Brynjolfsson and Jeffrey Hu, MIT, and Michael Smith, Carnegie Mellon; Barnes & Noble; Netflix; RealNetworks

Source: Chris Anderson (2004)

# Types of Recommendations

- **Editorial and hand curated**
  - ❑ List of favorite cities
  - ❑ List of "essential" items for travel

- **Simple aggregates**
  - ❑ Top 10, Most Popular, Recent Uploads

- **Tailored to individual users**
  - ❑ Amazon, Netflix, …

# Formal Model

- **$X$** = set of **Customers**
- **$S$** = set of **Items**

- **Utility function $u$: $X \times S \rightarrow R$**
  - ❑ **$R$** = set of ratings
  - ❑ **$R$** is a totally ordered set
  - ❑ e.g., **0-5** stars, real number in **[0,1]**

# Utility Matrix

|       | Avatar | LOTR | Matrix | Pirates |
|-------|--------|------|--------|---------|
| Alice | 1      |      | 0.2    |         |
| Bob   |        | 0.5  |        | 0.3     |
| Carol | 0.2    |      | 1      |         |
| David |        |      |        | 0.4     |

# Key Problems

- **(1) Gathering "known" ratings for matrix**
  - ❏ How to collect the data in the utility matrix

- **(2) Extrapolate unknown ratings from the known ones**
  - ❏ Mainly interested in high unknown ratings
    - We are not interested in knowing what you don't like but what you like

- **(3) Evaluating extrapolation methods**
  - ❏ How to measure success/performance of recommendation methods

# (1) Gathering Ratings

- **Explicit**
  - Ask people to rate items
  - Doesn't work well in practice – people can't be bothered

- **Implicit**
  - Learn ratings from user actions
    - E.g., purchase implies high rating
  - What about low ratings?
    - "not buying an item" = "don't like the item"   ?

# (2) Extrapolating Utilities

- **Key problem:** Utility matrix $U$ is **sparse**
    - Most people have not rated most items
    - **Cold start:**
        - New items have no ratings
        - New users have no history

- **Three approaches to recommender systems:**
    - **1)** Content-based
    - **2)** Collaborative
    - **3)** Latent factor based

# Outline

☑ Overview

➡ ☐ **Content-based Recommender System**

☐ Collaborative Filtering

☐ Evaluation & Complexity

# Content-based Recommendations

- **Main idea:** Recommend items to customer *x* similar to previous items rated highly by *x*
  - John enjoyed watching "Avengers Infinity War". John will also like "Avengers End Game" as well since they are similar in content

*Example:*

- **Movie recommendations**
  - Recommend movies with same actor(s), genre, …
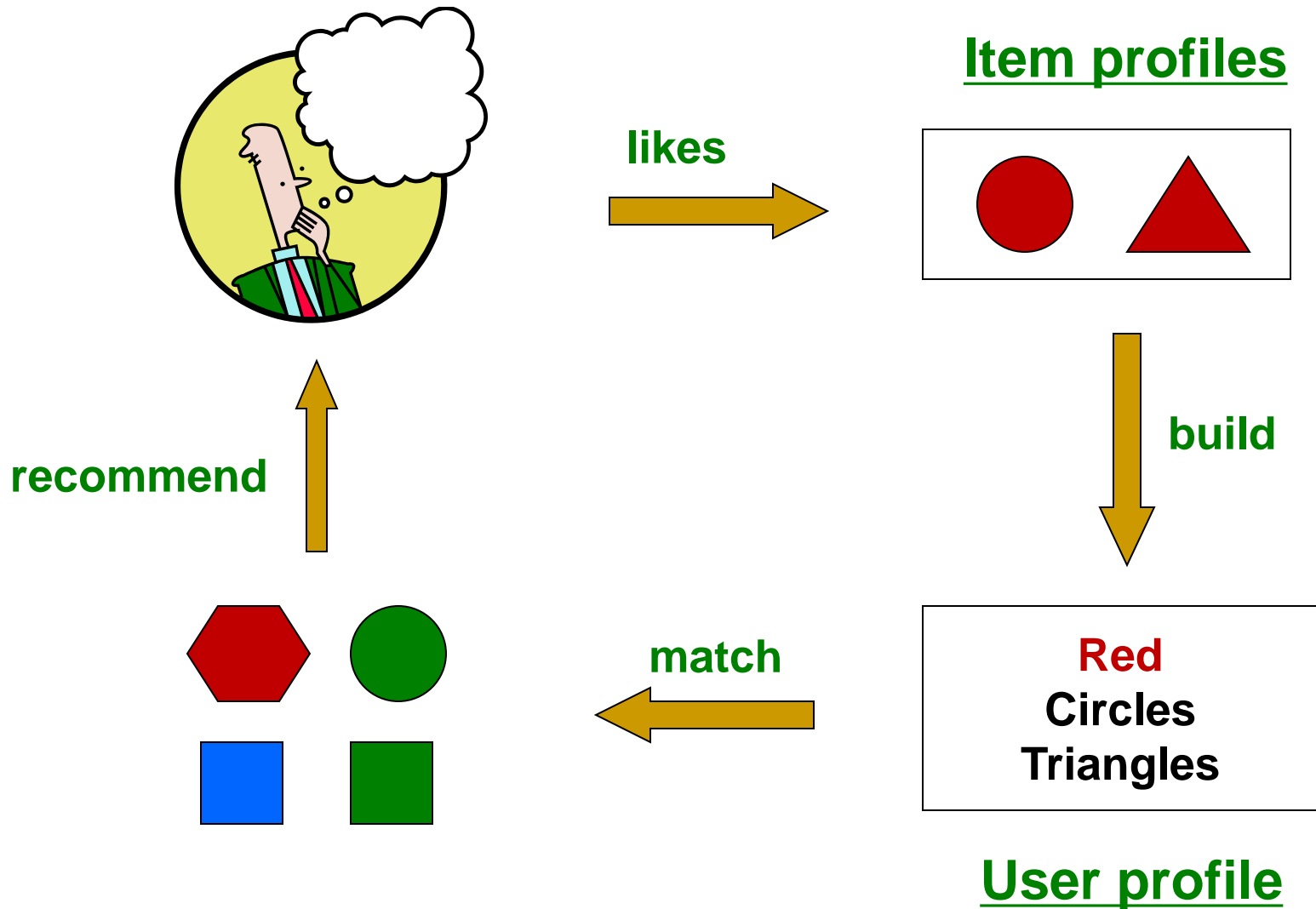
- **Websites, blogs, news**
  - Recommend other sites with "similar" content

# Plan of Action



**Item profiles**

**likes**

**build**

**recommend**

**match**

Red
**Circles**
**Triangles**

**User profile**

# Item Profiles

- For each item, create an **item profile**

- **Profile is a set (vector) of features**
  - **Movies:** author, title, actor, …
  - **Text:** set of "important" words in document

- **How to pick important features?**
  - Usual heuristic from text mining is **TF-IDF** (Term frequency * Inverse Doc Frequency)
    - **Term … Feature**
    - **Document … Item**

# Sidenote: TF-IDF

$f_{ij}$ = frequency of term (feature) **i** in doc (item) **j**

$$TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}}$$

**Note:** we normalize TF to discount for "longer" documents

$n_i$ = number of docs that mention term **i**

**N** = total number of docs

$$IDF_i = \log \frac{N}{n_i}$$

**TF-IDF score:** $w_{ij} = TF_{ij} \times IDF_i$

**Doc profile =** set of words with highest **TF-IDF** scores, together with their scores

# User Profiles and Prediction

- **User profile possibilities:**
  - Weighted average of rated item profiles
  - **Variation:** weight by difference from average rating for item
  - …

- **Prediction heuristic:**
  - Given user profile $x$ and item profile $i$, estimate
  $$u(x, i) = \cos(x, i) = \frac{x \cdot i}{||x|| \cdot ||i||}$$

# Pros: Content-based Approach

- **+: No need for data on other users**
  - No cold-start or sparsity problems
- **+: Able to recommend to users with unique tastes**
- **+: Able to recommend new & unpopular items**
  - No first-rater problem
- **+: Able to provide explanations**
  - Can provide explanations of recommended items by listing content-features that caused an item to be recommended

# Cons: Content-based Approach

- **–: Finding the appropriate features is hard**
  - E.g., images, movies, music
- **–: Recommendations for new users**
  - **How to build a user profile?**
- **–: Overspecialization**
  - Never recommends items outside user's content profile
    - Users want to be surprised sometimes
  - People might have multiple interests
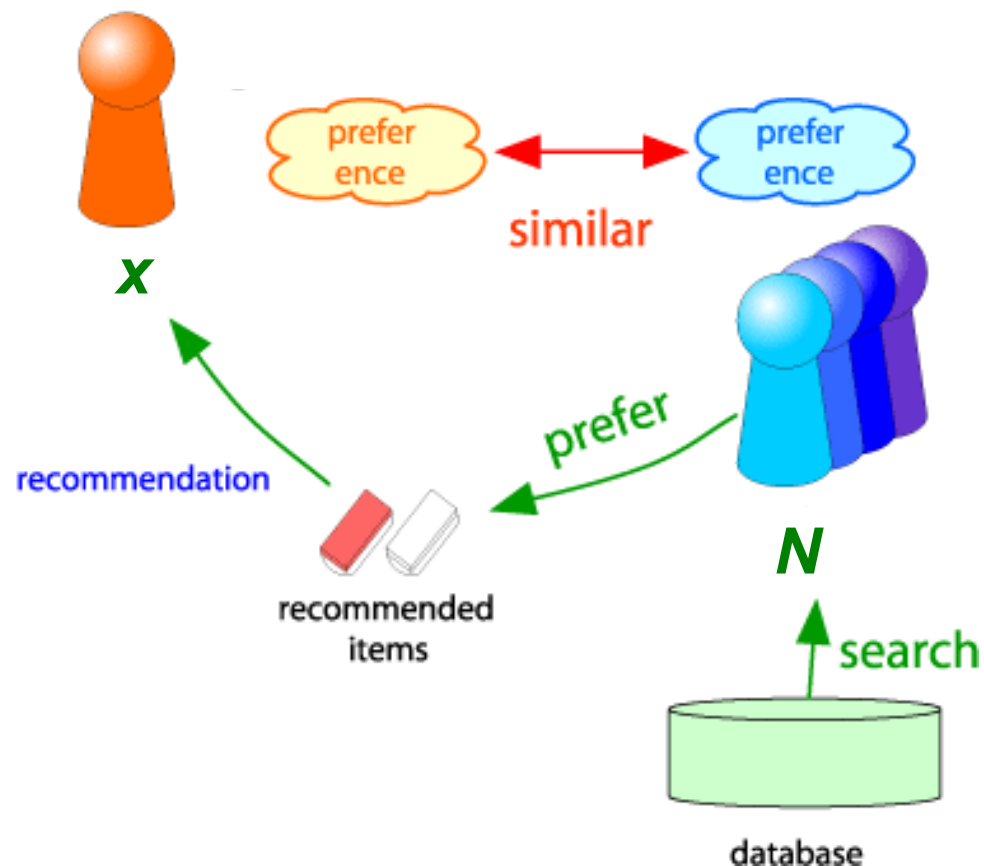  - **Unable to exploit quality judgments of other users**

# Outline

☑ Overview

☑ Content-based Recommender System

➡ **Collaborative Filtering**

☐ Evaluation & Complexity

# Collaborative Filtering

- Consider user $x$

- Find set $N$ of other users whose ratings are "**similar**" to $x$'s ratings

- Estimate $x$'s ratings based on ratings of users in $N$

Note that contents of items are not used here.

# Finding "Similar" Users

$r_x = [*, \_, \_, *, ***]$
$r_y = [*, \_, **, **, \_]$

- Let $r_x$ be the vector of user $x$'s ratings

- **Jaccard similarity measure**
  - **Problem:** Ignores the value of the rating

- **Cosine similarity measure**
  - sim($x$, $y$) = cos($r_x$, $r_y$) = $\dfrac{r_x \cdot r_y}{||r_x|| \cdot ||r_y||}$
  - **Problem:** low rating is not penalized much

- **Pearson correlation coefficient**
  - $S_{xy}$ = items rated by both users $x$ and $y$

$$sim(x, y) = \frac{\sum_{s \in S_{xy}}(r_{xs} - \overline{r_x})(r_{ys} - \overline{r_y})}{\sqrt{\sum_{s \in S_{xy}}(r_{xs} - \overline{r_x})^2}\sqrt{\sum_{s \in S_{xy}}(r_{ys} - \overline{r_y})^2}}$$

*$r_x$, $r_y$ as sets:*
$r_x = \{1, 4, 5\}$
$r_y = \{1, 3, 4\}$

*$r_x$, $r_y$ as points:*
$r_x = \{1, 0, 0, 1, 3\}$
$r_y = \{1, 0, 2, 2, 0\}$

$\overline{r}_x, \overline{r}_y$ … avg. rating of **x**, **y**

# Similarity Metric

|   | HP1 | HP2 | HP3 | TW | SW1 | SW2 | SW3 |
|---|-----|-----|-----|-----|-----|-----|-----|
| $A$ | 4 | | | 5 | 1 | | |
| $B$ | 5 | 5 | 4 | | | | |
| $C$ | | | | 2 | 4 | 5 | |
| $D$ | | 3 | | | | | 3 |

- **Intuitively we want: sim($A$, $B$) > sim($A$, $C$)**

- **Jaccard similarity: 1/5 < 2/4**

- **Cosine similarity: 0.380 > 0.322**

  - ❑ Problem: low rating is not penalized much

  - ❑ **Solution: subtract the (row) mean**

**sim A,B vs. A,C:**
0.092 **>** -0.559

|   | HP1 | HP2 | HP3 | TW | SW1 | SW2 | SW3 |
|---|-----|-----|-----|-----|-----|-----|-----|
| $A$ | 2/3 | | | 5/3 | −7/3 | | |
| $B$ | 1/3 | 1/3 | −2/3 | | | | |
| $C$ | | | | −5/3 | 1/3 | 4/3 | |
| $D$ | | 0 | | | | | 0 |

# Rating Predictions

**From similarity metric to recommendations:**

- Let $r_x$ be the vector of user $x$'s ratings

- Let $N$ (called 'k-nearest neighbors') be the set of $k$ users most similar to $x$ who have rated item $i$

- **Prediction $r_{xi}$ for item i of *user x*:**

  - $r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi}$

  **Shorthand:**
  $s_{xy} = sim(x, y)$

  - $r_{xi} = \frac{\sum_{y \in N} s_{xy} \cdot r_{yi}}{\sum_{y \in N} s_{xy}}$

- **Many other tricks possible…**

# Item-Item Collaborative Filtering

- **So far: User-user collaborative filtering**

- **Another view: Item-item**

  - For item $i$, find other similar items rated by user x
    - Use the utility matrix for computing similarity

  - Estimate rating for item $i$ based on ratings for similar items

  - Can use same similarity metrics and prediction functions as in user-user model

$$r_{xi} = \frac{\sum_{j \in N(i;x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i;x)} s_{ij}}$$

$s_{ij}$… similarity of items $i$ and $j$
$r_{xj}$…rating of user $x$ on item $j$
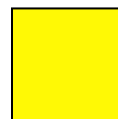$N(i;x)$… set items rated by $x$ similar to $i$

# Item-Item CF (|N|=2)

| movies \ users | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |   | 3 |   |   | 5 |   |   | 5 |   | 4 |   |
| 2 |   |   | 5 | 4 |   |   | 4 |   |   | 2 | 1 | 3 |
| 3 | 2 | 4 |   | 1 | 2 |   | 3 |   | 4 | 3 | 5 |   |
| 4 |   | 2 | 4 |   | 5 |   |   | 4 |   |   | 2 |   |
| 5 |   |   | 4 | 3 | 4 | 2 |   |   |   |   | 2 | 5 |
| 6 | 1 |   | 3 |   | 3 |   | 2 |   |   |   | 4 |   |

☐ - unknown rating  ▨ - rating between 1 to 5

# Item-Item CF (|N|=2)

|        | 1 | 2 | 3 | 4 | **5** | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|--------|---|---|---|---|-------|---|---|---|---|----|----|----|
| **1**  | 1 |   | 3 |   | ?     | 5 |   |   | 5 |    | 4  |    |
| **2**  |   |   | 5 | 4 |       |   | 4 |   |   | 2  | 1  | 3  |
| **3**  | 2 | 4 |   | 1 | 2     |   | 3 |   | 4 | 3  | 5  |    |
| **4**  |   | 2 | 4 |   | 5     |   |   | 4 |   |    | 2  |    |
| **5**  |   |   | 4 | 3 | 4     | 2 |   |   |   |    | 2  | 5  |
| **6**  | 1 |   | 3 |   | 3     |   |   | 2 |   |    | 4  |    |

movies

- estimate rating of movie **1** by user **5**

# Item-Item CF (|N|=2)



**users**

**movies**

| | 1 | 2 | 3 | 4 | **5** | 6 | 7 | 8 | 9 | 10 | 11 | 12 | sim(1,m) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | | 3 | | ? | 5 | | | 5 | | 4 | | 1.00 |
| 2 | | | 5 | 4 | | | 4 | | | 2 | 1 | 3 | -0.18 |
| **3** | 2 | 4 | | 1 | 2 | | 3 | | 4 | 3 | 5 | | **0.41** |
| 4 | | 2 | 4 | | 5 | | | 4 | | | 2 | | -0.10 |
| 5 | | | 4 | 3 | 4 | 2 | | | | | 2 | 5 | -0.31 |
| **6** | 1 | | 3 | | 3 | | | 2 | | | 4 | | **0.59** |

**Neighbor selection:**
Identify movies similar to movie **1**, **rated by user 5**

**Similarity computation:**
1) Subtract mean rating $m_i$ from each movie $i$
   $m_1 = (1+3+5+5+4)/5 = 3.6$
   *row 1: [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]*
2) Compute cosine similarities between rows

U Kang

30

# Item-Item CF (|N|=2)

**users**

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | sim(1,m) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 |  | 3 |  | ? | 5 |  |  | 5 |  | 4 |  | 1.00 |
| 2 |  |  | 5 | 4 |  |  | 4 |  |  | 2 | 1 | 3 | -0.18 |
| **3** | 2 | 4 |  | 1 | 2 |  | 3 |  | 4 | 3 | 5 |  | **0.41** |
| 4 |  | 2 | 4 |  | 5 |  |  | 4 |  |  | 2 |  | -0.10 |
| 5 |  |  | 4 | 3 | 4 | 2 |  |  |  |  | 2 | 5 | -0.31 |
| **6** | 1 |  | 3 |  | 3 |  |  | 2 |  |  | 4 |  | **0.59** |

**movies**

**Compute similarity weights:**
$s_{1,3}=0.41, s_{1,6}=0.59$

# Item-Item CF (|N|=2)

**users**

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 1 | | 3 | | **2.6** | 5 | | | 5 | | 4 | |
| **2** | | | 5 | 4 | | | 4 | | | 2 | 1 | 3 |
| **3** | 2 | 4 | | 1 | 2 | | 3 | | 4 | 3 | 5 | |
| **4** | | 2 | 4 | | 5 | | | 4 | | | 2 | |
| **5** | | | 4 | 3 | 4 | 2 | | | | | 2 | 5 |
| **6** | 1 | | 3 | | 3 | | | 2 | | | 4 | |

**movies**

**Predict by taking weighted average:**

$r_{1.5}$ = **(0.41*2 + 0.59*3) / (0.41+0.59) = 2.6**

$$r_{ix} = \frac{\sum_{j \in N(i;x)} s_{ij} \cdot r_{jx}}{\sum s_{ij}}$$

# CF: Common Practice

- Define **similarity** $s_{ij}$ of items $i$ and $j$
- Select $k$ nearest neighbors $N(i; x)$
  - Items most similar to $i$, that were rated by $x$
- Estimate rating $r_{xi}$ as the weighted average:

$$r_{xi} = b_{xi} + \frac{\sum_{j \in N(i;x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i;x)} s_{ij}}$$

**baseline estimate for $r_{xi}$**

$$b_{xi} = \mu + b_x + b_i$$

- $\mu$ = overall mean movie rating
- $b_x$ = rating deviation of user $x$
  = (avg. rating of user $x$) − $\mu$
- $b_i$ = rating deviation of movie $i$
  = (avg. rating of movie i) − $\mu$

# CF: Baseline Predictor

- Mean movie rating: **3.7 stars**

- *The Sixth Sense* is **0.5** stars above avg.

- Joe rates **0.2** stars below avg.

  $\Rightarrow$ **Baseline estimation:**
  *Joe* **will rate** *The Sixth Sense* **4 stars**

# Item-Item vs. User-User

| | Avatar | LOTR | Matrix | Pirates |
|---|---|---|---|---|
| **Alice** | 1 | | 0.8 | |
| **Bob** | | 0.5 | | 0.3 |
| **Carol** | 0.9 | | 1 | 0.8 |
| **David** | | | 1 | 0.4 |

- **In practice, it has been observed that <u>item-item</u> often works better than user-user**
- **Why?** Items are simpler, users have multiple tastes

# Pros/Cons of Collaborative Filtering

- **+ Works for any kind of item**
    - No feature selection needed
- **+ Can use other people's suggestions**

# Pros/Cons of Collaborative Filtering

- **- Cold Start:**
  - ❑ Needs enough users in the system to find a match
- **- Sparsity:**
  - ❑ The user/ratings matrix is sparse
  - ❑ Hard to find users that have rated the same items
- **- First rater:**
  - ❑ Cannot recommend an item that has not been previously rated (e.g., new items, esoteric items)
- **- Popularity bias:**
  - ❑ Cannot recommend items to someone with unique taste
  - ❑ Tends to recommend popular items

# Hybrid Methods

- **Implement two or more different recommenders and combine predictions**
  - ❑ Perhaps using a linear model
- **Add content-based methods to collaborative filtering**
  - ❑ Item profiles for new item problem
  - ❑ Demographics to deal with new user problem

# Outline

☑ Overview

☑ Content-based Recommender System

☑ Collaborative Filtering

➡ ☐ **Evaluation & Complexity**

# Evaluation



movies

users

# Evaluation

# Evaluating Predictions

- **Compare predictions with known ratings**
  - **Root-mean-square error** (RMSE)
    - $\sqrt{\sum_{xi}(r_{xi} - r_{xi}^*)^2}$ where $\boldsymbol{r_{xi}}$ is predicted, $\boldsymbol{r_{xi}^*}$ is the true rating of $\boldsymbol{x}$ on $\boldsymbol{i}$
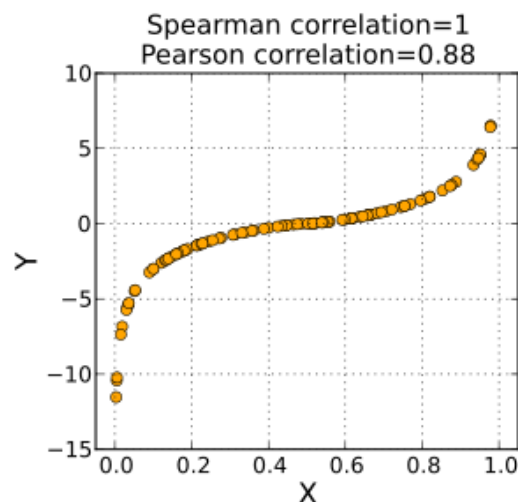  - **Precision at top 10**: error in top 10 highest predictions
  - **Rank Correlation**:
    - Spearman's *correlation* between system's and user's complete rankings

(From Wikipedia)

Pearson correlation coefficient

Spearman correlation=1
Pearson correlation=0.88

Rank correlation coefficient=1

U Kang

# Problems with Error Measures

- **Narrow focus on accuracy sometimes misses the point**
  - ❑ E.g., prediction diversity

- **In practice, we care only to predict high ratings:**
  - ❑ RMSE might penalize a method that does well for high ratings and badly for others

# Collaborative Filtering: Complexity

- Expensive step is finding $k$ most similar customers: **O(|X|)**
  - X … set of customers
- **Too expensive to do at runtime**

  - Could pre-compute
- Pre-compute finding similar customers
  - Near-neighbor search in high dimensions (**LSH**)
  - Clustering
  - Dimensionality reduction

# Tip: Add Data

- **Simple method on large data is better than complex method on small data**
    - Leverage all the data
    - Don't try to reduce data size in an effort to make fancy algorithms work

- **Add more data**
    - e.g., add IMDB data on genres

- **More data beats better algorithms**
  http://anand.typepad.com/datawocky/2008/03/more-data-usual.html

# What You Need to Know

- Motivation and the problem of recommendation
- Compare the content-based vs. collaborative filtering approaches for recommender system
  - Content-based: less cold-start problem
  - Collaborative filtering: works for any item
- Evaluation methods for recommendation
  - Training set and test set

# Questions?