

Mining Massive Datasets

Lecture 3

Artur Andrzejak

<http://pvs.ifi.uni-heidelberg.de>



RUPRECHT-KARLS-
UNIVERSITÄT
HEIDELBERG



Note on Slides

A substantial part of these slides come (either verbatim or in a modified form) from the book *Mining of Massive Datasets* by Jure Leskovec, Anand Rajaraman, Jeff Ullman (Stanford University).

For more information, see the website accompanying the book: <http://www.mmds.org>.

Current Topic

High dim. data

Locality
sensitive
hashing

Clustering

Dimensio-
nality
reduction

Graph data

PageRank,
SimRank

Community
Detection

Spam
Detection

Infinite data

Filtering
data
streams

Web
advertising

Queries on
streams

Machine learning

SVM

Decision
Trees

Perceptron,
kNN

Apps

Recommender
systems

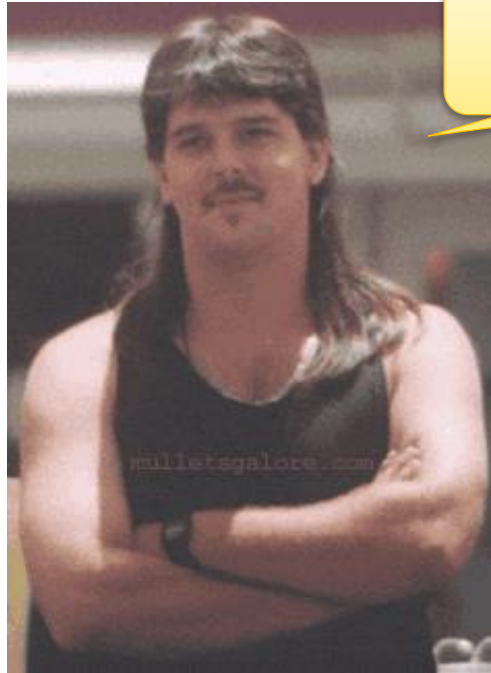
Association
Rules

Duplicate
document
detection

Programming in Spark & MapReduce

Recommender Systems: Motivation and Overview

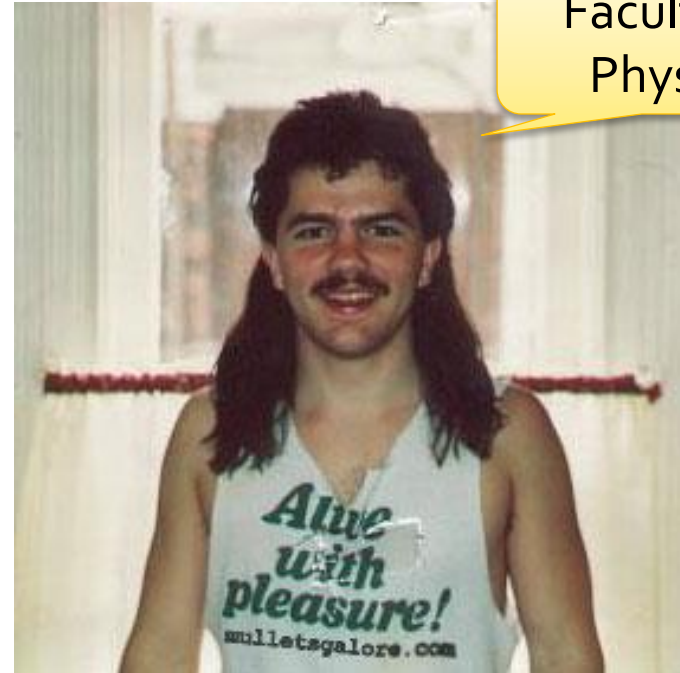
Recommender Systems - Example



Faculty of
Math & CS

■ Scholar X

- Buys Metallica CD
- Buys Megadeth CD

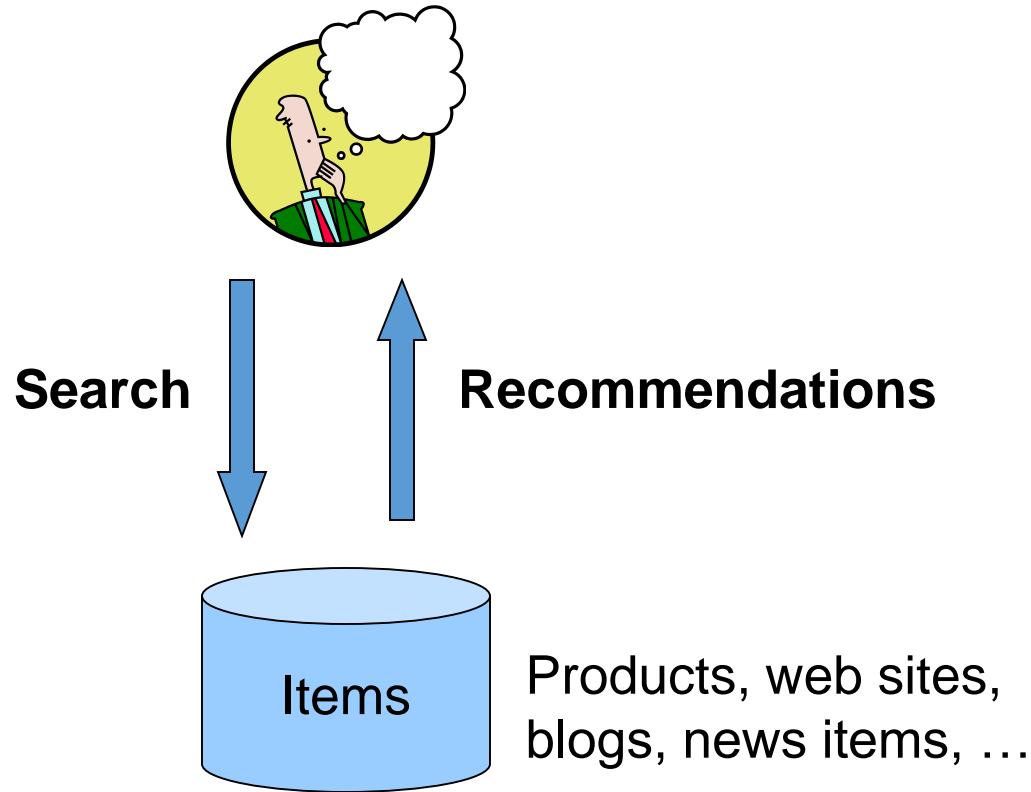


Faculty of
Physics

■ Scholar Y

- Does search on Metallica
- Recommender system suggests Megadeth
 - From data about X

Recommender Systems



Examples:

amazon.com



movielens
helping you find the *right* movies



Why using Recommender Systems?

- Value for users
 - Find things that are interesting
 - Narrow down the set of choices
 - Discover new things ...
- Value for providers
 - Increase trust and customer loyalty
 - Increase sales, click rates, conversion etc.
 - Opportunities for promotion


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

		users			
movies		Alice	Bob	Carol	David
	Star Wars	1		0.2	
	Matrix		0.5		0.3
	Avatar	0.2		1	
	Pirates				0.4

Utility Matrix

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

users

items (movies)

	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	

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 what people like (high scores)
 - Key problem: Utility matrix \mathbf{U} is sparse
 - **Cold start**: New items have no ratings / New users have no history
- (3) **Evaluating extrapolation methods**
 - How to measure success/performance of recommendation methods

Overview of the Approaches

A. Collaborative filtering (CF)

- Two Versions: user-user and item-item

B. Content-based recommenders

- Difficult in practice due to manually drafted features

C. Latent factor models (LF)

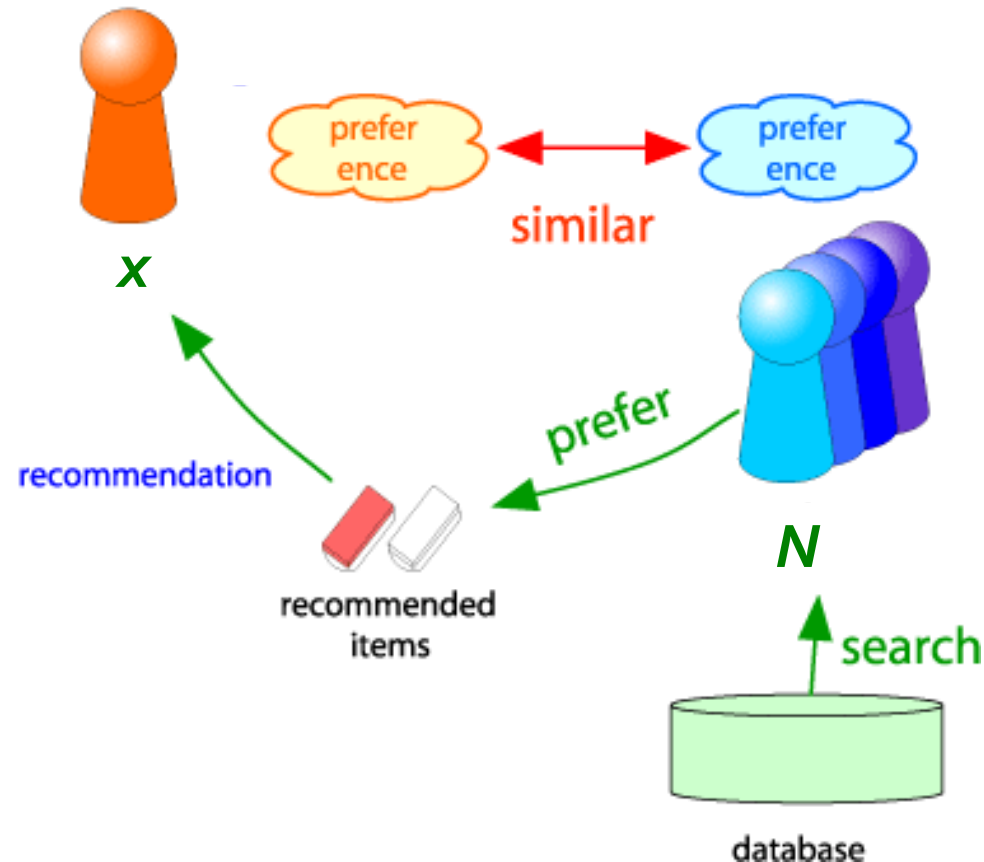
- Improves on B by finding latent features *automagically*

A1. Collaborative Filtering

Version: User-User

User-User 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



Algorithm for User-User CF

- Let r_x be the vector of user x 's ratings
- Let N be the set of k “neighbors” of x
 - = users most similar to x (by metric $\text{sim}(x, y)$, later)
 - ... who have already rated item i
- The **predicted rating of user x for item i** is:

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

Average rating of k
“neighbors” of x for item i

Rating Predictions - Improved

- Let r_x be the vector of user x 's ratings
- Let N be the set of k “neighbors” of x
 - =users most similar to x (by $\text{sim}(x, y)$)
 - ... who have already rated item i
- The predicted rating of user x for item i is:

$$r_{xi} = \frac{\sum_{y \in N} \text{sim}(x, y) \cdot r_{yi}}{\sum_{y \in N} |\text{sim}(x, y)|}$$

Improvement here: we scale the recommendations of other user y by similarity $\text{sim}(x, y)$ of x to this “neighbor” y

Note: Expression $\text{sim}(x, y)$ for $y \in N$ should be positive since y and x are by definition similar; but for general case, use $|\text{sim}(x, y)|$ in denominator!

Example: User-User CF

Computing N = the set of k users most similar to x

- According to $\text{sim}(x, y)$, y 's have already rated item i

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

- 1. Find all users y who have already rated item 1

users

items

	1	2	3	4	5	6
1	1		3		?	5
2			5	4		
3	2	4		1	2	
4		2	4		5	
5			4	3	4	2
6	1		3		3	

Example: User-User CF

Computing N = the set of k users most similar to x

- According to $\text{sim}(x, y)$, y 's have already rated item i

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

users

	1	2	3	4	5	6
1	1		3		?	5
2			5	4		
3	2	4		1	2	
4		2	4		5	
5			4	3	4	2
6	1		3		3	

items

- 1. Find all users y who have already rated item 1
 - Solution: $y \in \{1, 3, 6\}$
- 2. Compute $\text{sim}(x, y)$ for users y
 - Let $s_{5,1} = \text{sim}(5,1)$, $s_{5,3} = \dots$

Example: User-User CF

Computing N = the set of k users most similar to x

- According to $\text{sim}(x, y)$, y 's have already rated item i

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

users


	1	2	3	4	5	6
1	1		3		?	5
2			5	4		
3	2	4		1	2	
4		2	4		5	
5			4	3	4	2
6	1		3		3	

items

- 1. Find all users who have already rated item 1
- 2. For such users y (here $y \in \{1, 3, 6\}$) compute $\text{sim}(x, y)$
 - Let $s_{5,1} = \text{sim}(5,1)$, $s_{5,3} = \dots$
- 3. Among $\{s_{5,1}, s_{5,3}, s_{5,6}\}$, find top k values (say $k=2$)
 - E.g. $s_{5,3}, s_{5,6}$ largest $\Rightarrow N = \{3,6\}$

Finding “Similar” Users by Cosine

- Let r_z be the column of user z 's ratings:

$\begin{aligned} r_x &= [* , _ , _ , * , ***] \\ r_y &= [* , _ , ** , ** , _] \end{aligned}$		$\begin{aligned} r_x, r_y \text{ as vectors, NaN's } \rightarrow 0\text{'s:} \\ r_x &= \{1, 0, 0, 1, 3\} \\ r_y &= \{1, 0, 2, 2, 0\} \end{aligned}$
--	--	---

- Turn stars to numbers, fill missing values with 0's
- => We can consider r_x, r_y as vectors
- Use **Cosine similarity measure**
 - $sim(x, y) = \cos(r_x, r_y) = \frac{r_x \cdot r_y}{||r_x|| \cdot ||r_y||}$
 - $\cos(r_x, r_y)$ is just the angle bw. r_x, r_y , from 1 (=similar) to -1 (“opposite”)

Problems with Cosine

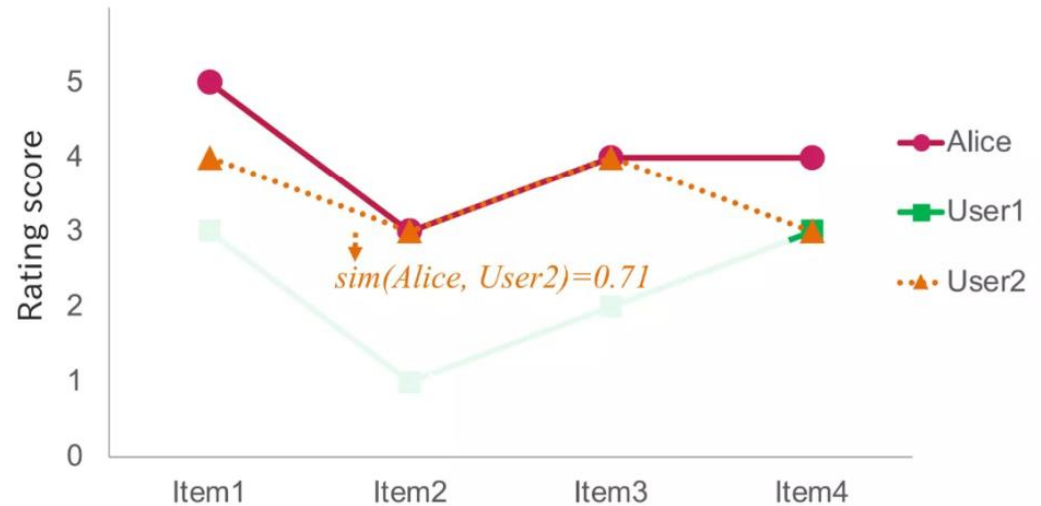
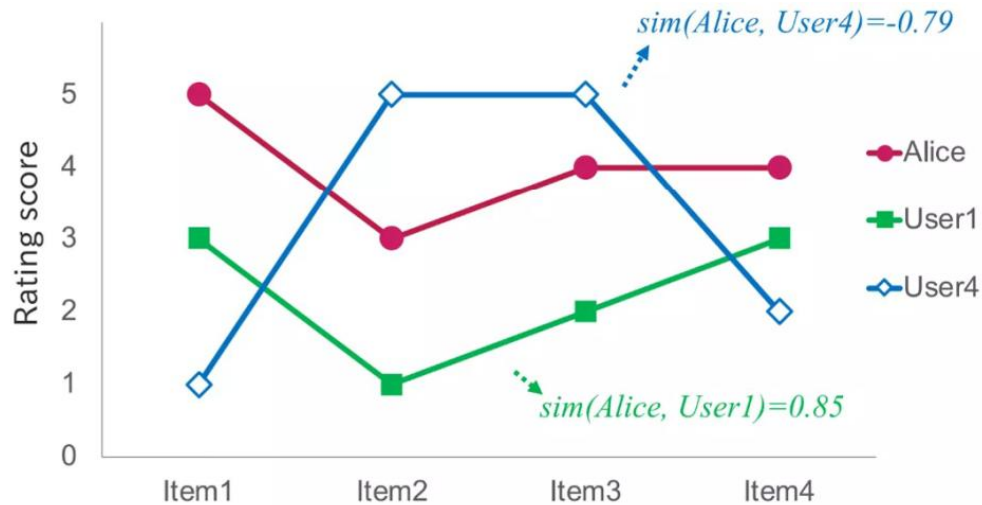
- **Cosine similarity measure**

- $sim(\mathbf{x}, \mathbf{y}) = \cos(\mathbf{r}_x, \mathbf{r}_y) = \frac{\mathbf{r}_x \cdot \mathbf{r}_y}{\|\mathbf{r}_x\| \cdot \|\mathbf{r}_y\|}$

$\mathbf{r}_x, \mathbf{r}_y$ as vectors:
 $\mathbf{r}_x = \{1, 0, 0, 1, 3\}$
 $\mathbf{r}_y = \{1, 0, 2, 2, 0\}$

- Problem 1: Missing ratings become low ratings!
 - 0's are interpreted as (very) low ratings
 - Items not rated by a user are considered as disliked!
- Problem 2: Users have different bias
 - I.e. some give higher ratings on average, others low ratings on average

Example: User Bias vs. Similarity



From Cosine to Pearson

$$\begin{aligned} r_x &= [* , _ , _ , * , ***] \\ r_y &= [* , _ , ** , ** , _] \end{aligned}$$

- Solution for problems 1 & 2
 - For P1, *consider only parts of the vectors* r_x, r_y
 - Vector entries for items rated by both users x and y
 - For P2, *normalize or center* each vector: subtract average user's rating
- => Use **Pearson correlation coefficient** on partial vectors
 - Pearson CC is computed as cosine similarity between centered vectors
- Definition and formula
 - S_{xy} = items rated by both users x and y
 - $\bar{r}_x, \bar{r}_y, \dots$ = average ratings of users x, y, \dots

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

Computing Most Similar Users

- For recommendations, we use only most similar users, or “neighbors” $N(x)$ of x
- A. Set a threshold for user similarity
 - If a user has higher similarity than a threshold, he/she can be regarded as a “similar” user
- B. Focus on top k similar users (kNN method)
 - If a user ranks at the top k similarity, he/she can be regarded as a similar user
 - k is often set to between 50 ~ 200
 - In worst cases, a system uses rating information of users with low similarity

Computing Predictions

- The predicted rating for item i and user x is

$$r_{xi} = \frac{\sum_{y \in N(x)} \text{sim}(x, y) \cdot r_{yi}}{\sum_{y \in N(x)} \text{sim}(x, y)}$$

- But it is good consider the bias or baseline of user x :

$$r_{xi} = \overline{r_x} + \frac{\sum_{y \in N(x)} \text{sim}(x, y) \cdot (r_{yi} - \overline{r_x})}{\sum_{y \in N(x)} \text{sim}(x, y)}$$

A2. Collaborative Filtering

Version: Item-Item

Problems with User-Based CF

- Lack of data: If users haven't rate the same items yet, user similarity cannot be computed
- Instability
 - It is rare that two users rated the same item
 - => User sim. drastically changes with few new ratings
 - User preferences (user features) often change, while item features do not often change
- Computational cost
 - In general, $\#Users \gg \#Items$ => high cost of finding nearest neighbors (similar users)

Item-Item Collaborative Filtering

- For item i , find **similar items** rated by user x
 - Let $N(i; x)$ be the set of such items
 - We can use similarity metrics $sim()$ as before
- Estimate rating of user x for item i *based on her/his previous ratings* for the items in $N(i; x)$
- The **predicted rating of user x for item i** is:

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

Example: Item-Item CF

		users											
		1	2	3	4	5	6	7	8	9	10	11	12
items	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	

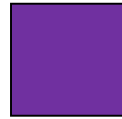


- unknown rating



- rating between 1 to 5

Item-Item CF



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

items	users											
	1	2	3	4	5	6	7	8	9	10	11	12
	1		3		?	5			5		4	
	2		5	4			4			2	1	3
	3	2	4		2		3		4	3	5	
	4		2	4	5			4			2	
	5		4	3	4	2					2	5
6	1		3		3			2			4	

Item-Item CF (Set $|N|=2$)

Neighbor selection:

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

We use Pearson correlation as similarity:

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

		users												
		1	2	3	4	5	6	7	8	9	10	11	12	sim(1,m)
items	1	1		3		?	5			5		4		1.00
	2			5	4			4			2	1	3	
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

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

Compute similarity weights:

$$s_{1,3} = \mathbf{0.41}, s_{1,6} = \mathbf{0.59}$$

		users												sim(1,m)
		1	2	3	4	5	6	7	8	9	10	11	12	
items	1	1		3		?	5			5		4		1.00
	2			5	4			4			2	1	3	-
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

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

$$r_{xi} = \frac{\sum_{j \in N(i;x)} \text{sim}(i,j) \cdot r_{xj}}{\sum_{j \in N(i;x)} \text{sim}(i,j)}$$

Predict by taking the weighted average:

$$r_{15} = (0.41 \cdot 2 + 0.59 \cdot 3) / (0.41 + 0.59) = 2.6$$

		users												
		1	2	3	4	5	6	7	8	9	10	11	12	sim(1,m)
items	1	1		3		2.6	5			5		4		1.00
	2			5	4			4			2	1	3	
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

CF: Using Baseline Estimates

- Define similarity $sim(i, j)$ 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:

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

baseline estimate for (x, i) :

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

baseline estimate for (x, j)

- μ = overall mean item rating
- b_x = rating deviation of user x
= (avg. rating of user x) - μ
- b_i = rating deviation of item i

Pros/Cons of Collaborative Filtering

- **+ Works for any kind of item**
 - No feature selection needed
- **- Cold Start:**
 - Need 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

Collaborative Filtering

Remarks & Practical Tips

Item-Item vs. User-User

	Alice	Bob	Carol	David
Star Wars	1		0.2	
Matrix		0.5		0.3
Avatar	0.2		1	
Pirates				0.4

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

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

Evaluation

The diagram shows a 10x6 grid representing a user-movie rating matrix. The grid is labeled "users" on the left and "movies" on the top. The grid contains numerical ratings from 1 to 5 in some cells, while others are empty.

1	3	4			
	3	5			5
		4	5		5
		3			
		3			
2			2		2
				5	
	2	1			1
	3			3	
1					

Evaluation

movies

users

1	3	4			
	3	5			5
		4	5		5
		3			
		3			
2			?		?
				?	
	2	1			?
	3			?	
1					

Test Data Set

Collaborative Filtering: Complexity

- Expensive step is **finding k most similar customers**: $O(|X|)$
 - X ... set of customers
- Too expensive to do at runtime
- But we could pre-compute for all customers
 - Naïve pre-computation takes time $O(|X|^2)$
- We will learn how to do this faster!
 - Near-neighbor search in high dimensions via **Locality-Sensitive Hashing**
 - Other means: clustering, dimensionality reduction

Tip: Add Data

- **Leverage all the data**

- Don't try to reduce data size in an effort to make fancy algorithms work
- Simple methods on large data do best

- **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>

Computing Pearson Efficiently

A Side-Note

Computing Pearson Efficiently

- „Our“ formula for Pearson skips missing values (NaN's) in the computation
 - => We use only items in \mathbf{S}_{xy} (= items rated by both users \mathbf{x} and \mathbf{y}) in each of the 3 sums
 - E.g. numerator: $\sum_{s \in \mathbf{S}_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)$
- Better: A. pre-process each vector as follows:
 - 1. Normalize non-missing values in each vector
 - I.e. compute mean over (only) non-missing values, and subtract m from each non-missing value
 - 2. Replace each missing value by 0
- ... and B. use standard Pearson formula (no skipping)

Computing Pearson - Example

- Example phase A:
 - 1. Normalize non-missing values in each vector
 - I.e. compute mean over (only) non-missing values, and subtract m from each non-missing value
 - 2. Replace each missing value by 0's
 - $xr = [1, \text{nan}, 3, \text{nan}, \text{nan}, 5, \text{nan}, \text{nan}, 5, \text{nan}, 4, \text{nan}]$
 - 1a. Compute mean: $m = (1+3+5+5+4)/5 = 3.6$
 - 1b. Subtract mean m from each non-missing value:
 $xr-m = [-2.6, \text{nan}, -.6, \text{nan}, \text{nan}, 1.4, \text{nan}, \text{nan}, 1.4, \text{nan}, .4, \text{nan}]$
 - 2. Replace NaN's by 0's:
- ⇒ $x = [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]$

Computing Pearson - Example

- $xr = [1, \text{nan}, 3, \text{nan}, \text{nan}, 5, \text{nan}, \text{nan}, 5, \text{nan}, 4, \text{nan}]$
⇒ $x = [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]$
 $yr = [2, 4, \text{nan}, 1, 2, \text{nan}, 3, \text{nan}, 4, 3, 5, \text{nan}]$
⇒ $y = [-1, 1, 0, -2, -1, 0, 0, 0, 1, 0, 2, 0]$
- B. Compute „normal“ Pearson of x , y : 0.4140393356
- Why is the result same as using S_{xy} and “skipping”?
 - If a value (rating) is missing, it becomes 0 after pre-processing => same as “skipping” it in each sum
 - E.g. numerator: $\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)$
 - $s \notin S_{xy} \Rightarrow (r_{xs} - \bar{r}_x) = 0$ or $(r_{ys} - \bar{r}_y) = 0$
 - => No contribution to the sum, i.e. as if skipped!

Computing Pearson - Code

```
from scipy import nan
import numpy as np
from scipy.stats import pearsonr
```

```
def normalize (input: list):
    mean = np.nanmean(input)
    return input - mean
```

```
def preprocess_vec (input: list):
    "Normalize and remove NaN's"
    return np.nan_to_num(normalize(input))
```

Computing Pearson - Code

```
# Rows (items) from slide "Item-Item CF (Set |N|=2)"
# x is item 1, y is item 3
xraw = [1, nan, 3, nan, nan, 5, nan, nan, 5, nan, 4, nan]
x = preprocess_vec(xraw)
yraw = [2, 4, nan, 1, 2, nan, 3, nan, 4, 3, 5, nan]
y = preprocess_vec(yraw)

# Correlation of x and y
corr_xy, p_value = pearsonr(x, y)
print (f"Pearson correlation of x and y is {corr_xy}")
```

=> Pearson correlation of x and y is 0.4140393356054126

B. Content-based Recommender Systems

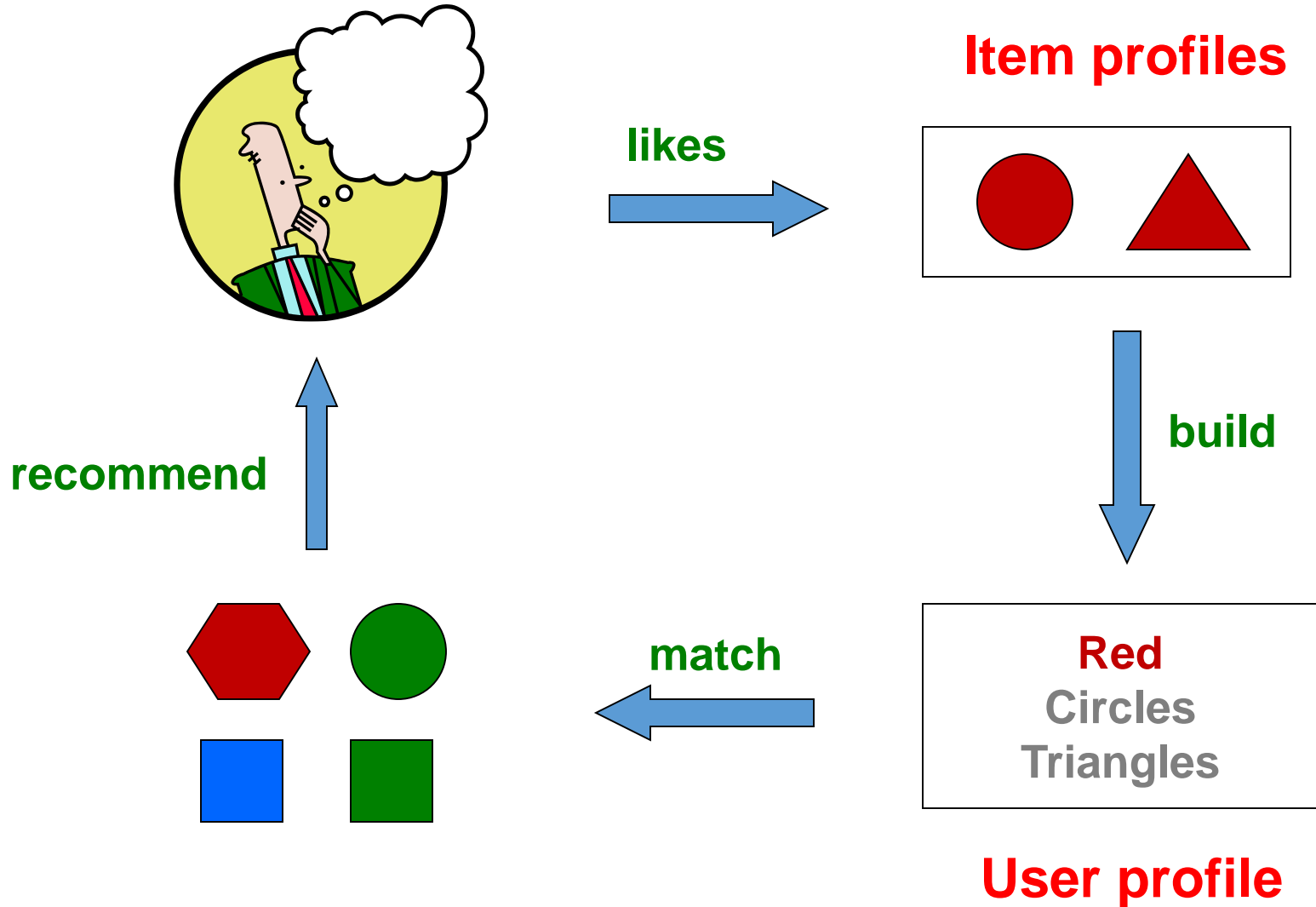
Content-based Recommendations

- **Main idea:** Recommend to customer x items similar to previous items rated highly by x

Example:

- **Movie recommendations**
 - Recommend movies with same actor(s), genre, director,...
- **Websites, blogs, news**
 - Recommend other sites with “similar” content

Plan of Action



Item Profiles

- For each item, create an **item profile**
- Profile is a **set (vector) v of features**
 - **Movies:** author, title, actor, director,...
 - **Text:** Set of “important” words in document
- Binary encoding (0/1) is used frequently
 - Each (famous) actor gets a fixed vector position k
 - If present in the movie, v_k is set to 1, else to 0

Example: Items are Movies

- Representing item profile – a “mixed” vector

$$i_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \\ 1 \\ \vdots \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 3.2 \\ 2.9 \\ \vdots \end{bmatrix} \begin{array}{l} \left. \begin{array}{c} \\ \\ \\ \\ \end{array} \right\} \text{Set of actors} \\ \\ \left. \begin{array}{c} \\ \\ \\ \\ \end{array} \right\} \text{Director(s)} \\ \\ \left. \begin{array}{c} \\ \end{array} \right\} \text{Ratings from various} \\ \text{movie DBs (as numbers)} \end{array}$$

User Profiles – A Simple Approach

- Idea: average the profiles of all items rated by a user
 - Weight each such item profile by the rating of this user
- Let use **u** give items 1, 2, ..., n ratings r_1, r_2, \dots, r_n
- **User profile x** = weighted average of rated **item profiles**

$$\mathbf{x} = (r_1 * i_1 + r_2 * i_2 + \dots + r_n * i_n) / n$$

Weights = ratings

Profiles of items 1, 2, .., n (those rated by user u) - vectors

Example: Star-Based Ratings

$$i_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, i_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, i_3 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, i_4 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, i_5 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

← actor A present
← actor B present

- Items are movies, only features are “actors”
 - Item profile has 2 components (for actor A and actor B)
- User ratings are 1 to 5 stars (per movie)
- User watched 5 movies
 - Actor A – movies got 3 and 5 stars (movies 1 & 2)
 - Actor B – movies got 1, 2 and 4 stars (movies 3, 4, 5)
- Ratings are $r_1=3, r_2=5, r_3=1, r_4=2, r_5=4$
- Item profiles are as above

Example: Star-Based Ratings

- The user profile becomes:

$$(r_1 i_1 + r_2 i_2 + r_3 i_3 + r_4 i_4 + r_5 i_5) / 5 =$$

$$\left(3 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 5 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 4 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right) / 5 = \begin{bmatrix} 8/5 \\ 7/5 \end{bmatrix}$$

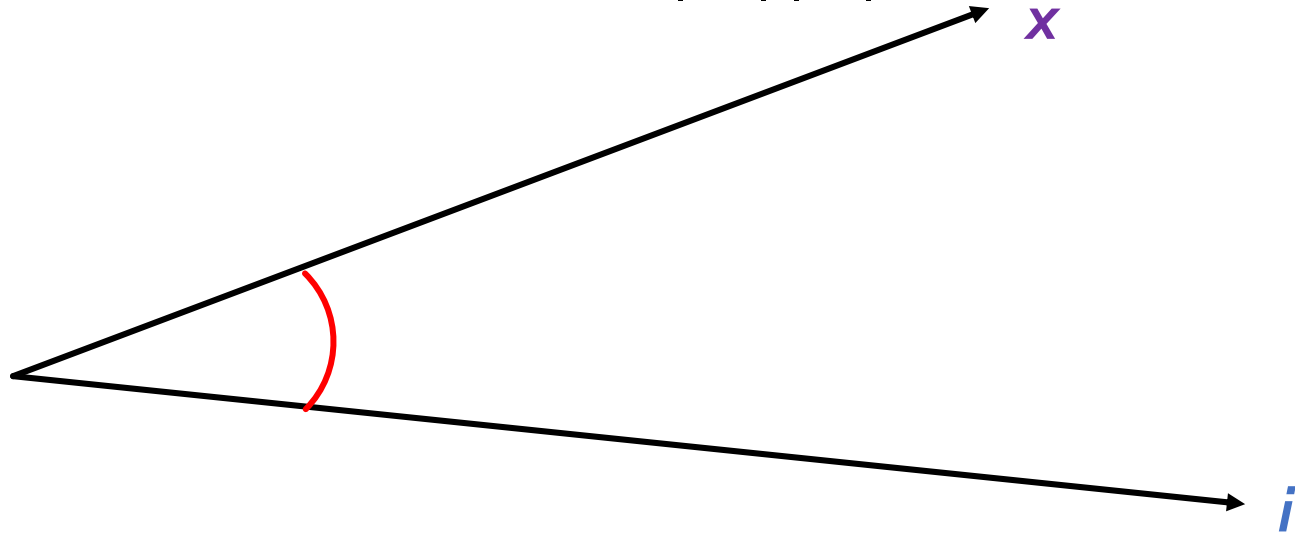
Such simple user profiles can exhibit low quality
(esp. if a user rated only few items) =>

See additional slides for problems & solutions:
Slides from “Problems with Simple User Profiles”

Matching User and Item Profiles

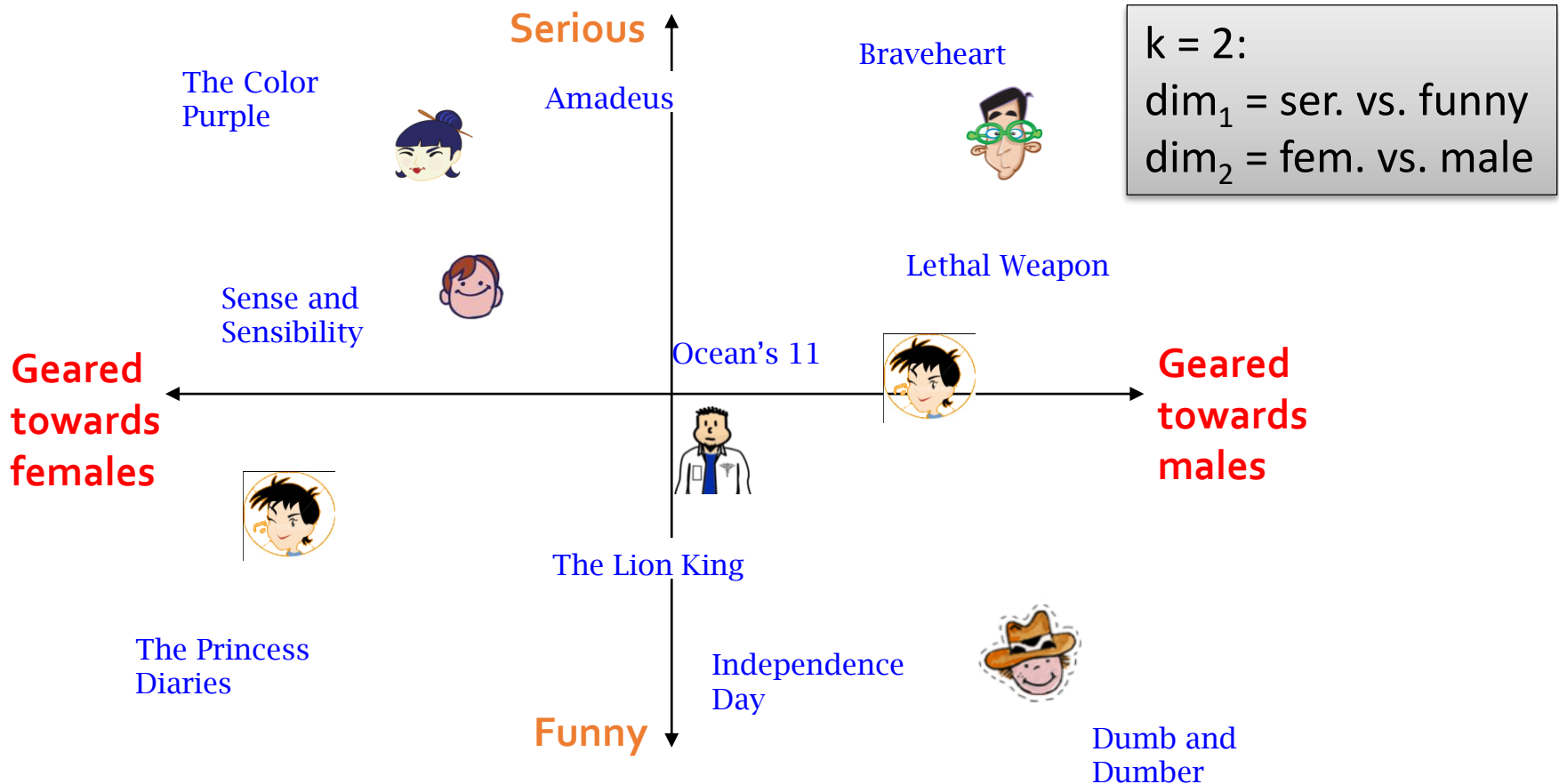
- To compute similarity of user profile and item profile, use a **prediction heuristic**:
 - Given **user profile x** and **item profile i** , estimate

$$u(x, i) = \cos(x, i) = \frac{x \cdot i}{||x|| \cdot ||i||}$$



Content-Based R.: Interpretation

- In content-based recommendation, we represented each **item** and each **user** as a vector in a k-dimensional space
- => Item **i** close to a user **x** gets a high recommendation rating



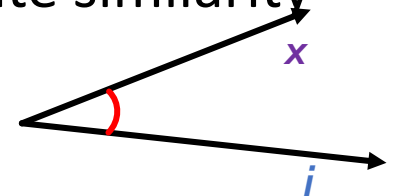
Summary: Content-Based R.

- We construct a vector i **for each item** (“**item profile**”) and a vector x (of size s) **for each user**
 - Item profile i : “natural” attributes of an item
 - User vector x : combination of item profiles rated by this user (“synthetic” profile)

- **Prediction heuristic:**

- Given a user vector x and item vector i , estimate similarity

$$u(x, i) = \cos(x, i) = \frac{x \cdot i}{||x|| \cdot ||i||}$$



- For a user with vector x , recommend by various criteria:
 - E.g. all items with $u(x, i) > \text{threshold}$
 - Rank items by $u(x, i)$, recommend top k (e.g. $k=5$)

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
 - People might have multiple interests
 - **Unable to exploit quality judgments of other users**

Thank you.

Questions?

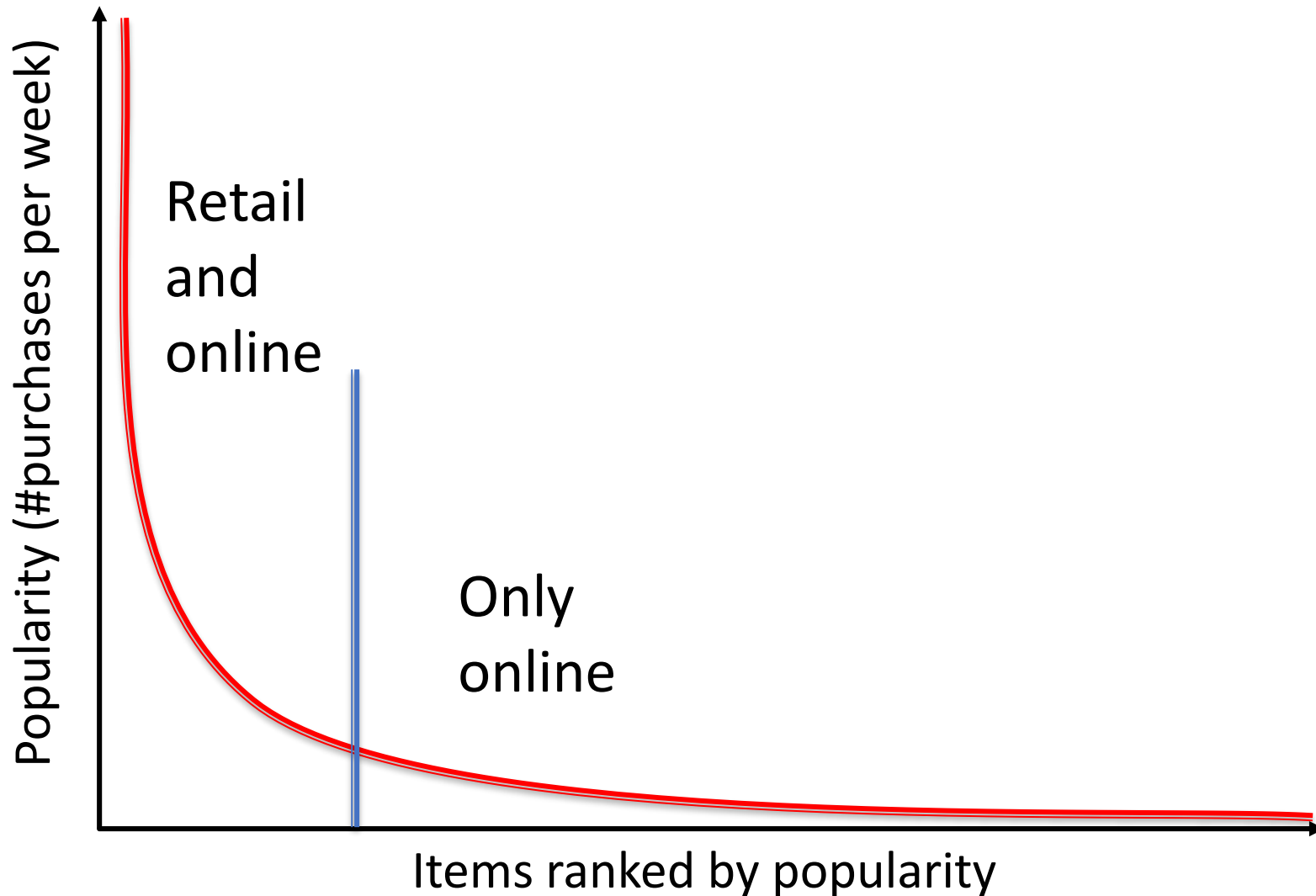
Additional Slides

Overview

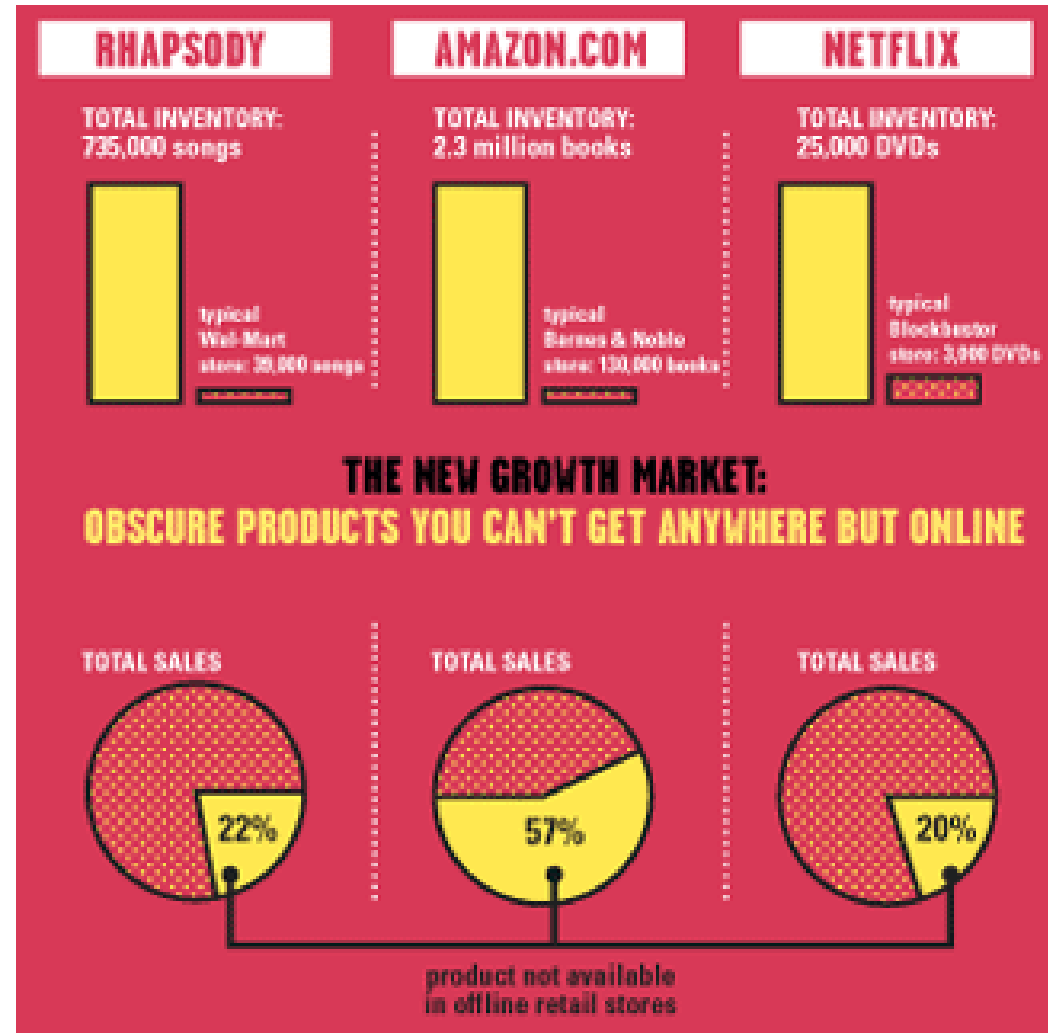
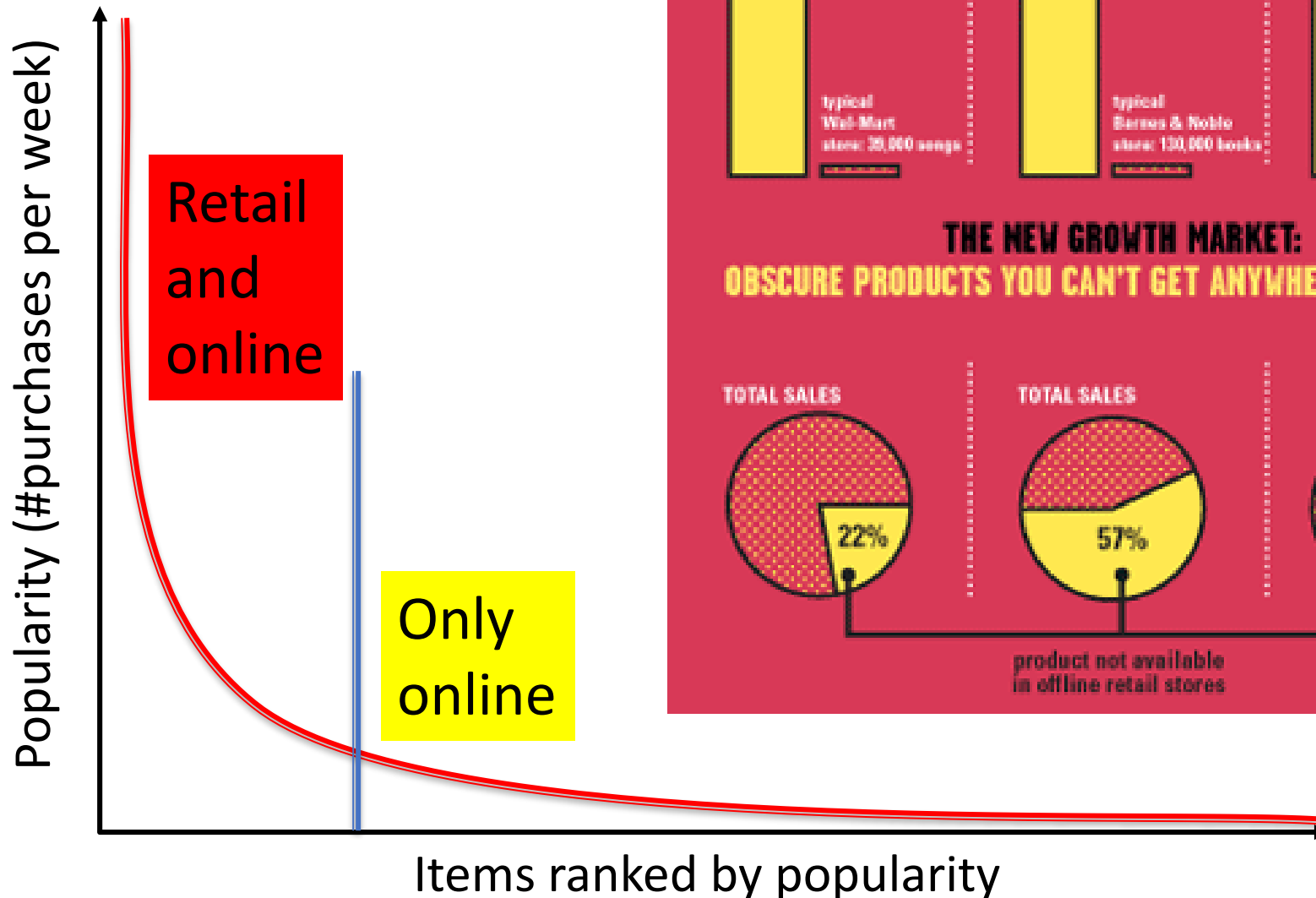
From Scarcity to Abundance

- Shelf space is a scarce commodity for traditional retailers
 - Also: TV networks, movie theaters,...
- Web enables near-zero-cost dissemination of information about products
 - From scarcity to abundance
- More choice necessitates better filters
 - Recommendation engines
 - How **Into Thin Air** made **Touching the Void** a bestseller: <http://www.wired.com/wired/archive/12.10/tail.html>

The Long Tail



The Long Tail



(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?

Additional Slides

Collaborative Filtering
Details on Similarity and Evaluation

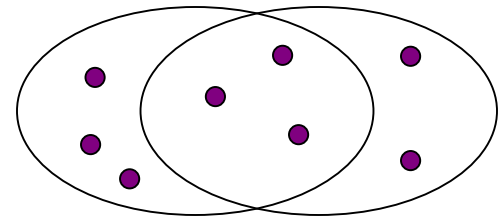
Jaccard Measures

- The **Jaccard similarity** of two **sets** is the size of their intersection divided by the size of their union:

$$\text{sim}(\mathbf{C}_1, \mathbf{C}_2) = |\mathbf{C}_1 \cap \mathbf{C}_2| / |\mathbf{C}_1 \cup \mathbf{C}_2|$$

- **Jaccard distance:** $d(\mathbf{C}_1, \mathbf{C}_2) = 1 - |\mathbf{C}_1 \cap \mathbf{C}_2| / |\mathbf{C}_1 \cup \mathbf{C}_2|$

- For measuring similarity of users, we consider only sets of items for which users voted
- Problem? Values of ratings are ignored!



3 in intersection
8 in union
Jaccard similarity = $3/8$
Jaccard distance = $5/8$

$$\begin{aligned} r_x &= [* , _ , _ , * , ***] \\ r_y &= [* , _ , ** , ** , _] \end{aligned}$$

$$\begin{aligned} &r_x, r_y \text{ as sets:} \\ r_x &= \{1, 4, 5\} \\ r_y &= \{1, 3, 4\} \end{aligned}$$

A Good 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: $\text{sim}(A, B) > \text{sim}(A, C)$
 - Jaccard similarity: $1/5 < 2/4 \Rightarrow$ bad
 - Cosine similarity: $0.386 > 0.322 \Rightarrow$ not good
- Solution: subtract the (row) mean**

	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

Notice: cos similarity
is = **correlation**
when data is
centered at 0!

$\text{sim } A, B$ vs. A, C :
 $0.092 > -0.559$

$$\Rightarrow \text{sim}(x, y) = \frac{\sum_i r_{xi} \cdot r_{yi}}{\sqrt{\sum_i r_{xi}^2} \cdot \sqrt{\sum_i r_{yi}^2}}$$

Evaluating Predictions

How to compare predictions with known ratings?

- **Root-mean-square error (RMSE)**, details: [link](#)

- $\sqrt{\frac{1}{N} \sum_{xi} (r_{xi} - r_{xi}^*)^2}$
- where r_{xi} is predicted, r_{xi}^* is the true rating of x on i , and N is the number of ratings (= number of used (x,i) combinations)

- **0/1 model**

- Coverage: Number of items/users for which system can make predictions
- Precision: Accuracy of predictions
- Receiver operating characteristic (ROC): Tradeoff curve between false positives and false negatives

Additional Slides

Content-based
Recommender Systems

Problems with Simple User Profiles

- The user profile becomes:

$$(r_1 i_1 + r_2 i_2 + r_3 i_3 + r_4 i_4 + r_5 i_5) / 5 =$$

$$(3 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 5 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 4 \begin{bmatrix} 0 \\ 1 \end{bmatrix}) / 5 = \begin{bmatrix} 8/5 \\ 7/5 \end{bmatrix}$$

- Problem 1: user likes actor A more than actor B, but this shows only weakly in his profile!
- Problem 2: with more ratings, each component becomes smaller (as **n** gets larger)
 - Because components with value 0 disturb the average, but should be treated as “don’t care about corresp. rating”

For 1: Normalizing Ratings

- Solution for 1: Normalize ratings by subtracting user's mean rating (which is $3 = (3+5+1+2+4)/5$)
 - Normalized ratings for actor A movies $\Rightarrow 0, +2$
 - Normalized ratings for actor A movies $\Rightarrow -2, -1, +1$
- Then the user profile is:
 - With $r_1=0, r_2=+2, r_3=-2, r_4=-1, r_5=+1$

$$\left(0 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) / 5 = \begin{bmatrix} 2/5 \\ -2/5 \end{bmatrix}$$

Now better: clear distinction
for actor A and actor B

For 2: Per-component Weights

Only ratings for these 2 items should be counted for actor A

$$\left(0 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix} + 1 \begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) / 5 = \begin{bmatrix} 2/5 \\ -2/5 \end{bmatrix}$$

Only the ratings for these 3 items
should be counted for actor B

- Essence of problem 2: a 0 in an item's component (=attribute) k should mean “don't care”, but now mean “one more neutral rating for attribute k”
- => Use “individual” n for each vector component
 - For actor A: $n_A = 2$, for actor B: $n_B = 3$

For 2: Per-component Weights

- => Use “individual” n for each vector component
 - For actor A: $n_A = 2$, for actor B: $n_B = 3$
 - Recall: Normalized ratings are $r_1=0$, $r_2=+2$, $r_3=-2$, $r_4=-1$, $r_5=+1$
- Then the user profile becomes:

$$\begin{bmatrix} r_1 / n_a \\ 0 \end{bmatrix} + \begin{bmatrix} r_2 / n_a \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ r_3 / n_b \end{bmatrix} + \begin{bmatrix} 0 \\ r_4 / n_b \end{bmatrix} + \begin{bmatrix} 0 \\ r_5 / n_b \end{bmatrix} =$$

$$\begin{bmatrix} 0/2 \\ 0 \end{bmatrix} + \begin{bmatrix} 2/2 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ -2/3 \end{bmatrix} + \begin{bmatrix} 0 \\ -1/3 \end{bmatrix} + \begin{bmatrix} 0 \\ 1/3 \end{bmatrix} = \begin{bmatrix} 1 \\ -2/3 \end{bmatrix}$$

Sidenote: Text Features

- How to pick important text features?
- Usual heuristic from text mining is **TF-IDF**
(Term-frequency * Inverse-Doc-Frequency)
 - Term == Feature
 - Doc(ument) == Item
- Now popular: word embeddings like **word2vec**
 - Each word is represented as a (sub)vector
 - Such vectors represent typical contexts (other words) in which this one occur
 - This **distributed representation** is learned from data

Sidenote: **TF-IDF**

- For term = i and doc = j , the **TF-IDF score** w_{ij} is:

$$w_{ij} = TF_{ij} * IDF_i$$

TF (term frequency):

- Let f_{ij} be frequency of term i in document j
- Then:

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

Note: we normalize TF
to discount for “longer” docs

IDF (inverse document frequency (of a term)):

- Let N = total number of docs, n_i = number of docs that mention term i
- Then:

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