

COMP9313: Big Data Management

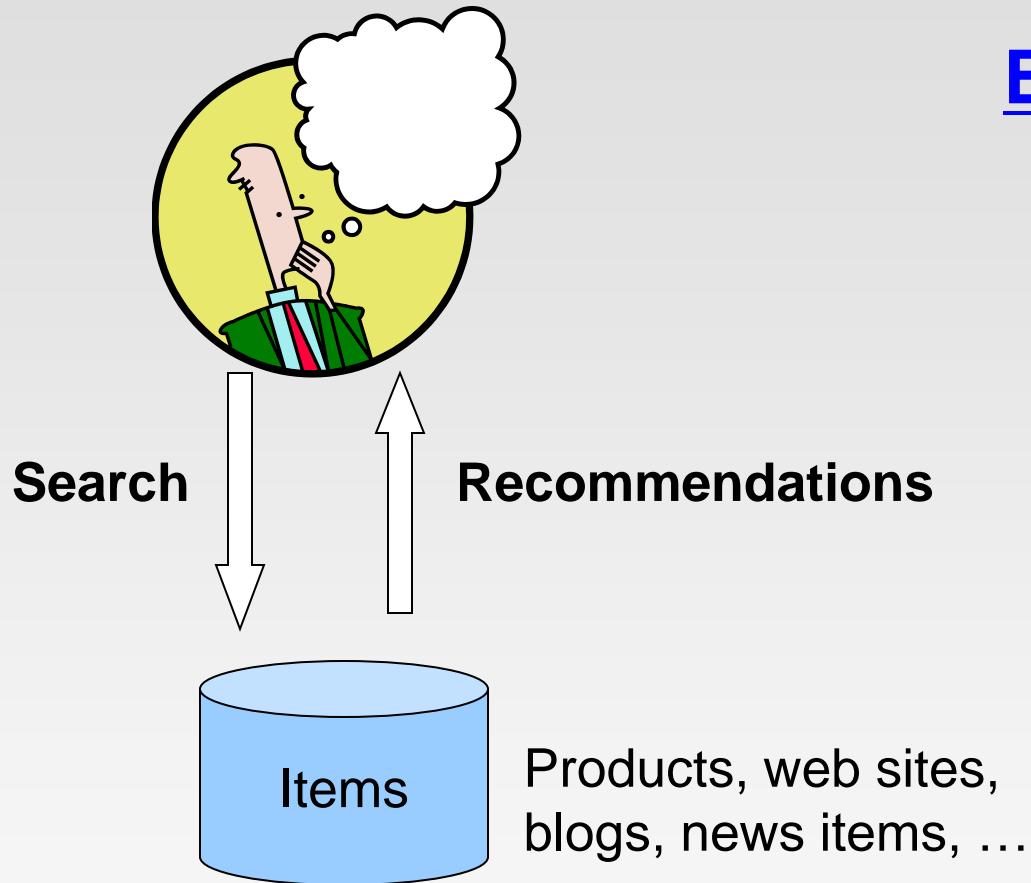


Lecturer: Xin Cao

Course web site: <http://www.cse.unsw.edu.au/~cs9313/>

Chapter 9: Recommender Systems

Recommendations



Examples:

amazon.com.
PANDORA

StumbleUpon

del.icio.us

P
PANDORA

NETFLIX

movielens
helping you find the *right* movies

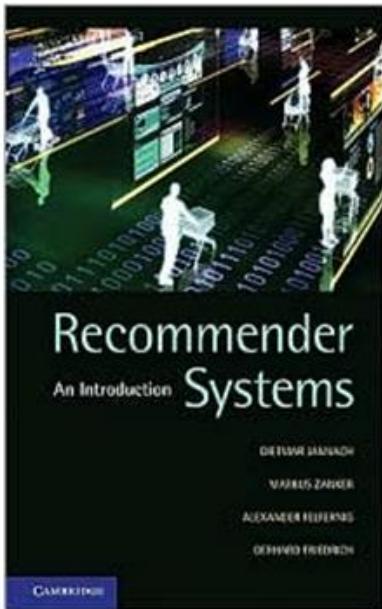
last.fm™
the social music revolution

Google™
News

You Tube

XBOX
LIVE

Recommender Systems



Recommender Systems: An Introduction

by Dietmar Jannach, Markus Zanker, Alexander Felfernig, Gerhard Friedrich

AVERAGE CUSTOMER RATING:

★★★★★ (Be the first to review)



Registrieren, um sehen zu können, was
deinen Freunden gefällt.

FORMAT:

Hardcover

NOOKbook (eBook) - not available

[Tell the publisher you want this in NOOKbook format](#)

NEW FROM BN.COM

\$65.00 List Price

\$52.00 Online Price
(You Save 20%)

Add to Cart

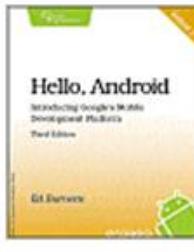
NEW & USED FROM OUR

New starting at **\$56.46** (You Save 20%)
Used starting at **\$51.98** (You Save 20%)

See All Prices

Table of Contents

Customers who bought this also bought



Recommender Systems

❖ Application areas

- Movie recommendation (Netflix)
- Related product recommendation (Amazon)
- Web page ranking (Google)
- Social recommendation (Facebook)
-

Related hotels...



Hotel 41
1,170 Reviews
London, England

Show Prices

Jobs you may be interested in Beta

Email Alerts | See More »

 **Technical Sales Manager - Europe**
Thermal Transfer Products - Home office

 **Senior Program Manager (f/m)**
Johnson Controls - Germany-NW-Burscheid

You may also like



★★★★★ (109)



★★★★★ (53)



★★★★★ (33)

Picasa™ -Webalben

Startseite Meine Fotos Erkunden Hochladen

Empfohlene Fotos Alle anzeigen



Netflix Movie Recommendation

Training data

| user | movie | date | score |
|------|-------|----------|-------|
| 1 | 21 | 5/7/02 | 1 |
| 1 | 213 | 8/2/04 | 5 |
| 2 | 345 | 3/6/01 | 4 |
| 2 | 123 | 5/1/05 | 4 |
| 2 | 768 | 7/15/02 | 3 |
| 3 | 76 | 1/22/01 | 5 |
| 4 | 45 | 8/3/00 | 4 |
| 5 | 568 | 9/10/05 | 1 |
| 5 | 342 | 3/5/03 | 2 |
| 5 | 234 | 12/28/00 | 2 |
| 6 | 76 | 8/11/02 | 5 |
| 6 | 56 | 6/15/03 | 4 |

Test data

| user | movie | date | score |
|------|-------|----------|-------|
| 1 | 62 | 1/6/05 | ? |
| 1 | 96 | 9/13/04 | ? |
| 2 | 7 | 8/18/05 | ? |
| 2 | 3 | 11/22/05 | ? |
| 3 | 47 | 6/13/02 | ? |
| 3 | 15 | 8/12/01 | ? |
| 4 | 41 | 9/1/00 | ? |
| 4 | 28 | 8/27/05 | ? |
| 5 | 93 | 4/4/05 | ? |
| 5 | 74 | 7/16/03 | ? |
| 6 | 69 | 2/14/04 | ? |
| 6 | 83 | 10/3/03 | ? |

Why using Recommender Systems?

- ❖ Value for the customer
 - Find things that are interesting
 - Narrow down the set of choices
 - Help me explore the space of options
 - Discover new things
 - Entertainment
 - ...
- ❖ Value for the provider
 - Additional and probably unique personalized service for the customer
 - Increase trust and customer loyalty
 - Increase sales, click trough rates, conversion etc.
 - Opportunities for promotion, persuasion
 - Obtain more knowledge about customers
 - ...

The Impact of Recommender Systems

- ❖ 35% of the purchases on Amazon are the result of their recommender system, according to McKinsey.
- ❖ During the Chinese global shopping festival of November 11, 2016, Alibaba achieved growth of up to 20% of their conversion rate using personalized landing pages, according to Alizila.
- ❖ Recommendations are responsible for 70% of the time people spend watching videos on YouTube.
- ❖ 75% of what people are watching on Netflix comes from recommendations, according to McKinsey
- ❖ Employing a recommender system enables Netflix to save around \$1 billion each year, according to the paper "[The Netflix Recommender System: Algorithms, Business Value, and Innovation](#)"

<https://tryolabs.com/blog/introduction-to-recommender-systems/>

Recommender systems

- ❖ RS seen as a function
- ❖ Given:
 - User model (e.g. ratings, preferences, demographics, situational context)
 - Items (with or without description of item characteristics)
- ❖ Find:
 - Relevance score. Used for ranking.
- ❖ Finally:
 - Recommend items that are assumed to be relevant
- ❖ But:
 - Remember that relevance might be context-dependent
 - Characteristics of the list itself might be important (diversity)

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 | 用户给物品的评分 |
|-------|--------|------|--------|---------|----------|
| Alice | 1 | | 0.2 | | |
| Bob | | 0.5 | | 0.3 | |
| Carol | 0.2 | | 1 | | |
| David | | | | 0.4 | |

Key Problems

获取一些已知的评分喜好

- ❖ Gathering “known” ratings for matrix
 - How to collect the data in the utility matrix

基于已知的评分喜好推断未知的喜好

- ❖ 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
- ❖ Evaluating extrapolation methods
 - How to measure success/performance of recommendation methods

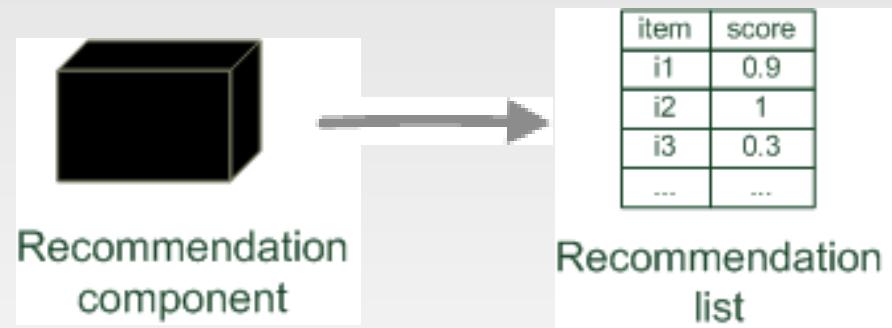
Gathering Ratings

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

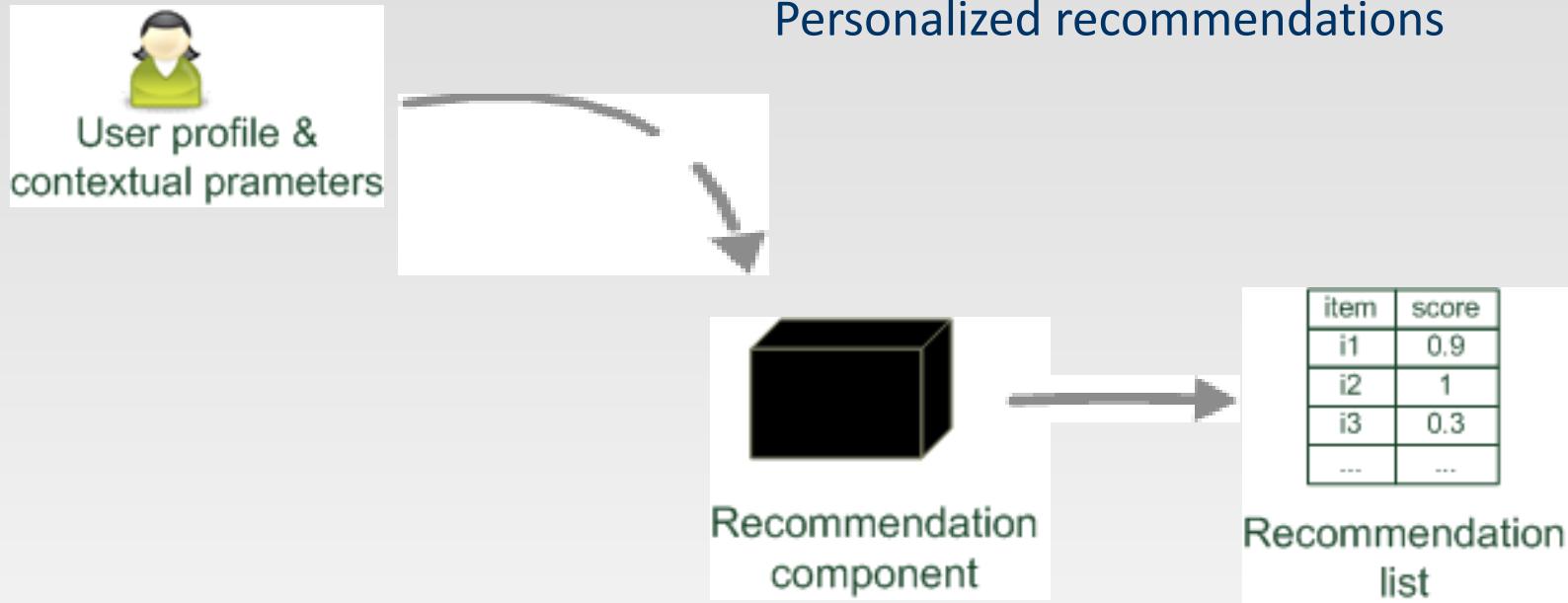
- ❖ Implicit
 - Learn ratings from user actions
 - ▶ E.g., purchase implies high rating

Paradigms of Recommender Systems

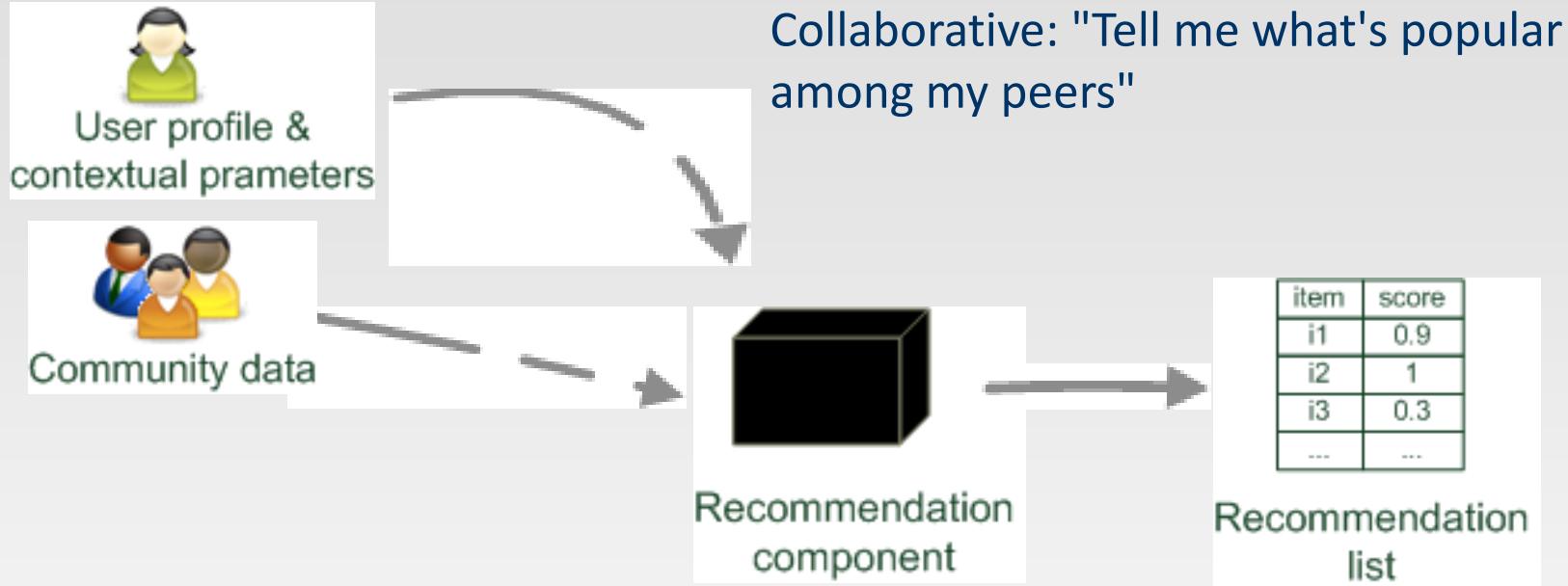
Recommender systems reduce information overload by estimating relevance



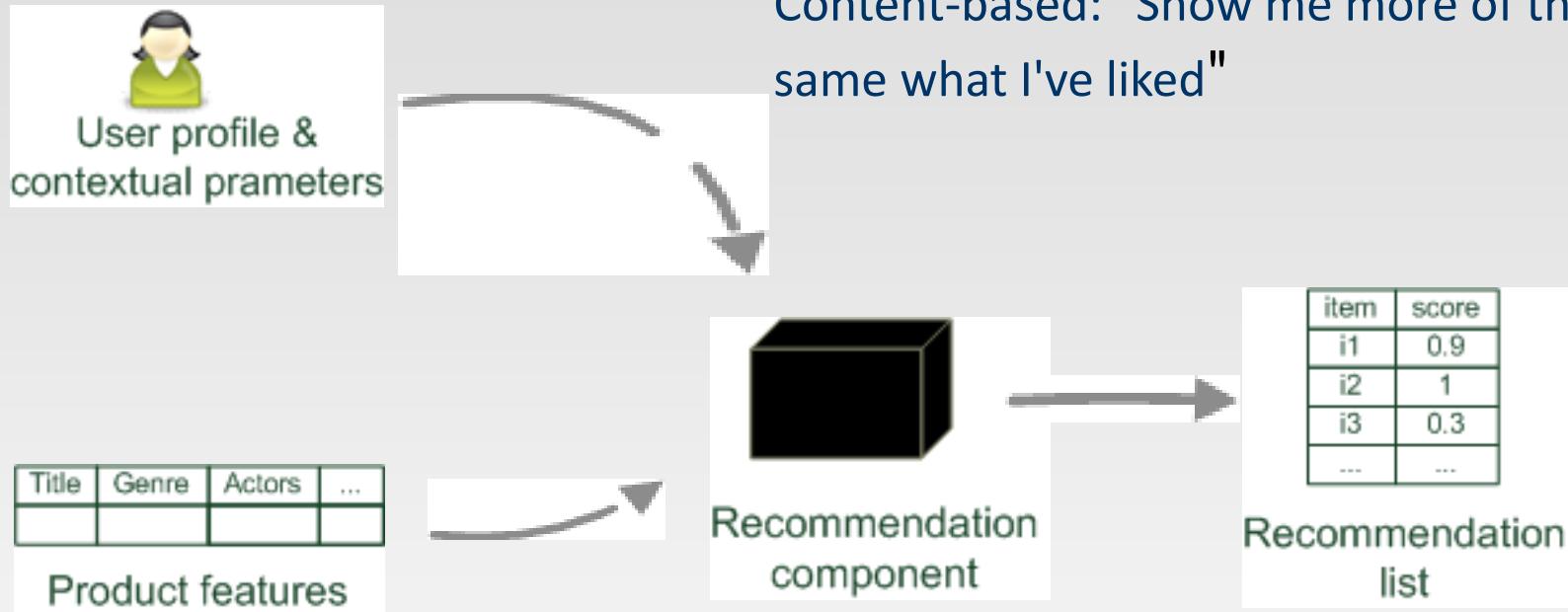
Paradigms of Recommender Systems



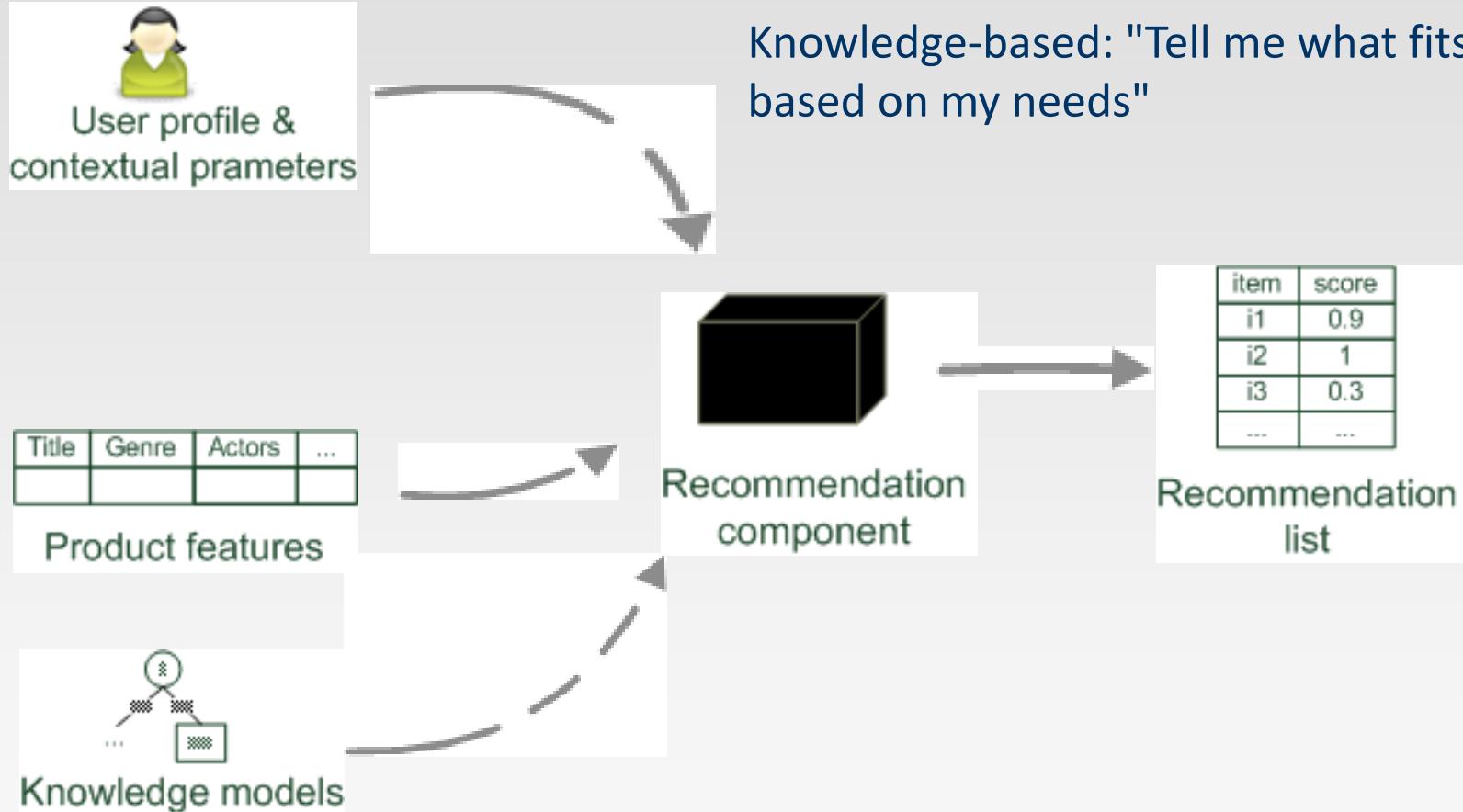
Paradigms of Recommender Systems



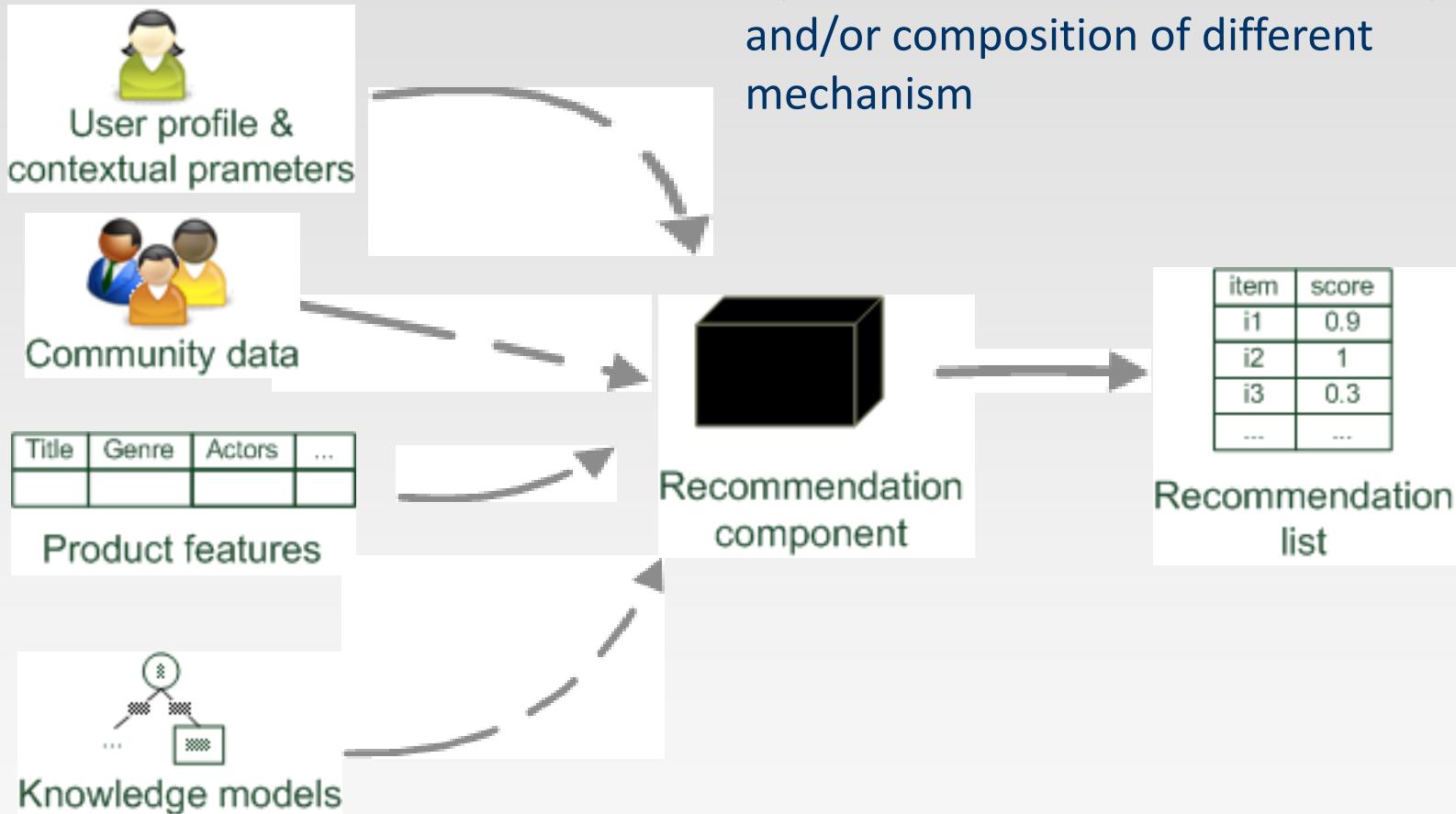
Paradigms of Recommender Systems



Paradigms of Recommender Systems



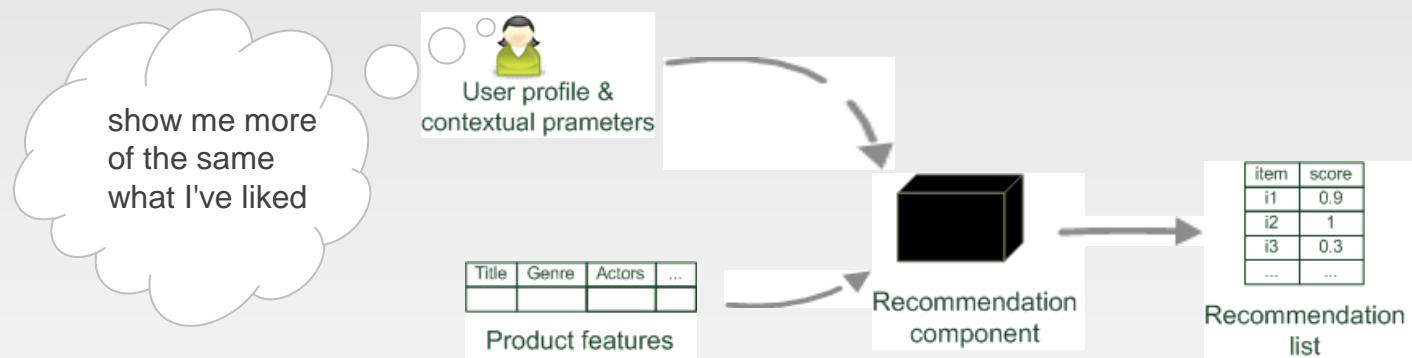
Paradigms of Recommender Systems



Paradigms of Recommender Systems

- ❖ **Content-based** – supervised machine learning used to induce a classifier to discriminate between interesting and uninteresting items for the user
- ❖ **Collaborative** – aggregation of consumers' preferences and recommendations to other users based on similarity in behavioral patterns
- ❖ **Knowledge-based** – knowledge about users and products used to reason what meets the user's requirements, using discrimination tree, decision support tools, case-based reasoning (CBR)

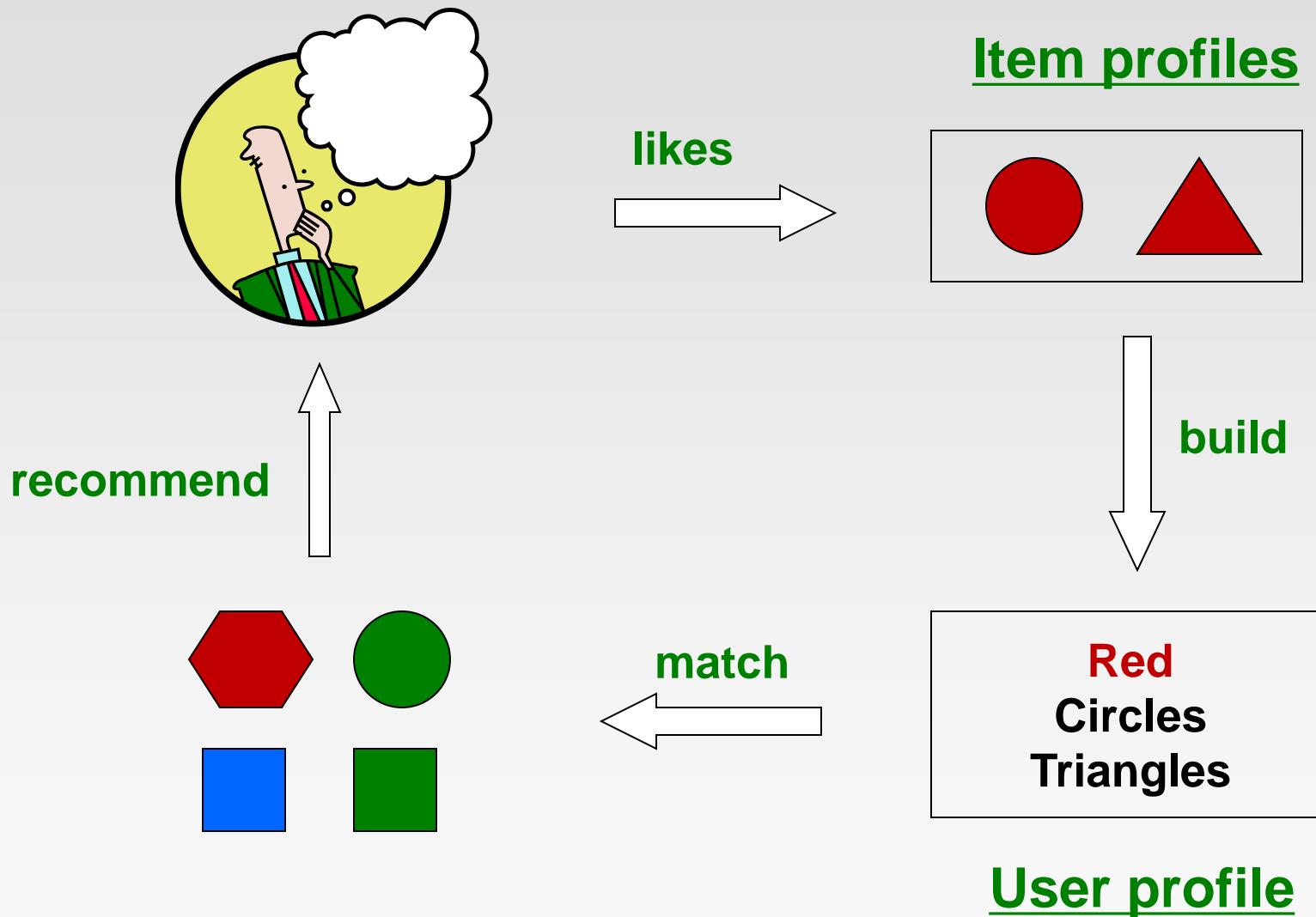
Part 2: Content-based Recommendation



Content-based Recommendations

- ❖ **Main idea:** Recommend items to customer x similar to previous items rated highly by x
- ❖ What do we need:
 - Some information about the available items such as the genre ("content")
 - Some sort of *user profile* describing what the user likes (the preferences)
- ❖ **Example:**
 - Movie recommendations:
 - ▶ Recommend movies with same actor(s), director, genre, ...
 - Websites, blogs, news:
 - ▶ Recommend other sites with “similar” content

Plan of Action



What is the “Content”?

- ❖ Most CB-recommendation techniques were applied to recommending text documents.
 - Like web pages or newsgroup messages for example.
- ❖ Content of items can also be represented as text documents.
 - With textual descriptions of their basic characteristics.
 - Structured: Each item is described by the same set of attributes



| Title | Genre | Author | Type | Price | Keywords |
|----------------------|-------------------|-------------------|-----------|-------|------------------------------------------------------------------|
| The Night of the Gun | Memoir | David Carr | Paperback | 29.90 | Press and journalism, drug addiction, personal memoirs, New York |
| The Lace Reader | Fiction, Mystery | Brunonia Barry | Hardcover | 49.90 | American contemporary fiction, detective, historical |
| Into the Fire | Romance, Suspense | Suzanne Brockmann | Hardcover | 45.90 | American fiction, murder, neo-Nazism |

- Unstructured: free-text description.

Item Profiles

- ❖ For each item, create an **item profile**
- ❖ Profile is a set (vector) of features
 - **Movies:** author, title, actor, director,...
 - **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**

Term-Frequency - Inverse Document Frequency (TF-IDF)

- ❖ Compute the overall importance of keywords

- Given a keyword i and a document j

$$TF-IDF(i,j) = TF(i,j) * IDF(i)$$

- ❖ Term frequency (TF)

- Let $freq(i,j)$ number of occurrences of keyword i in document j
 - Let $maxOthers(i,j)$ denote the highest number of occurrences of another keyword of j

$$TF(i,j) = \frac{freq(i,j)}{maxOthers(i,j)}$$

- ❖ Inverse Document Frequency (IDF)

- N: number of all recommendable documents
 - $n(i)$: number of documents in which keyword i appears

$$IDF(i) = \log \frac{N}{n(i)}$$

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

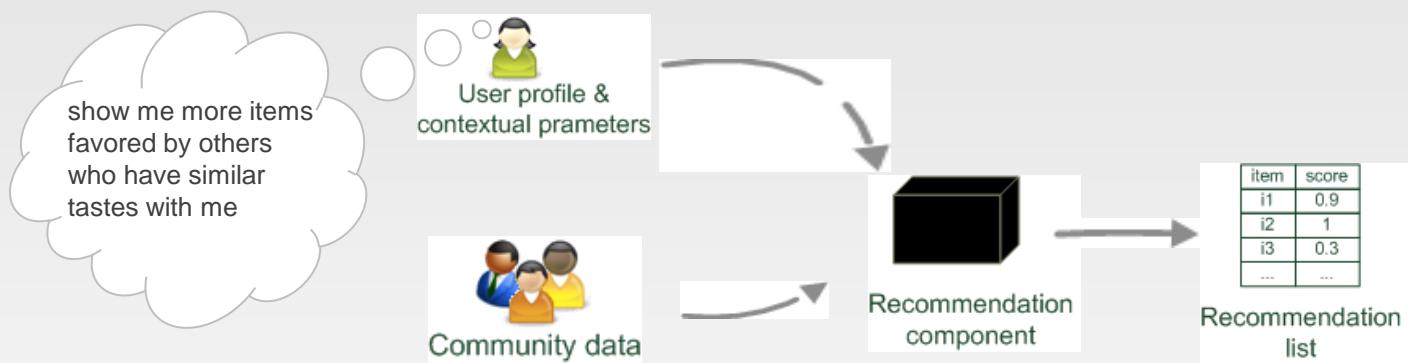
Pros: Content-based Approach

- ❖ +: No need for data on other users
- ❖ +: 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

Part 3: Collaborative Filtering



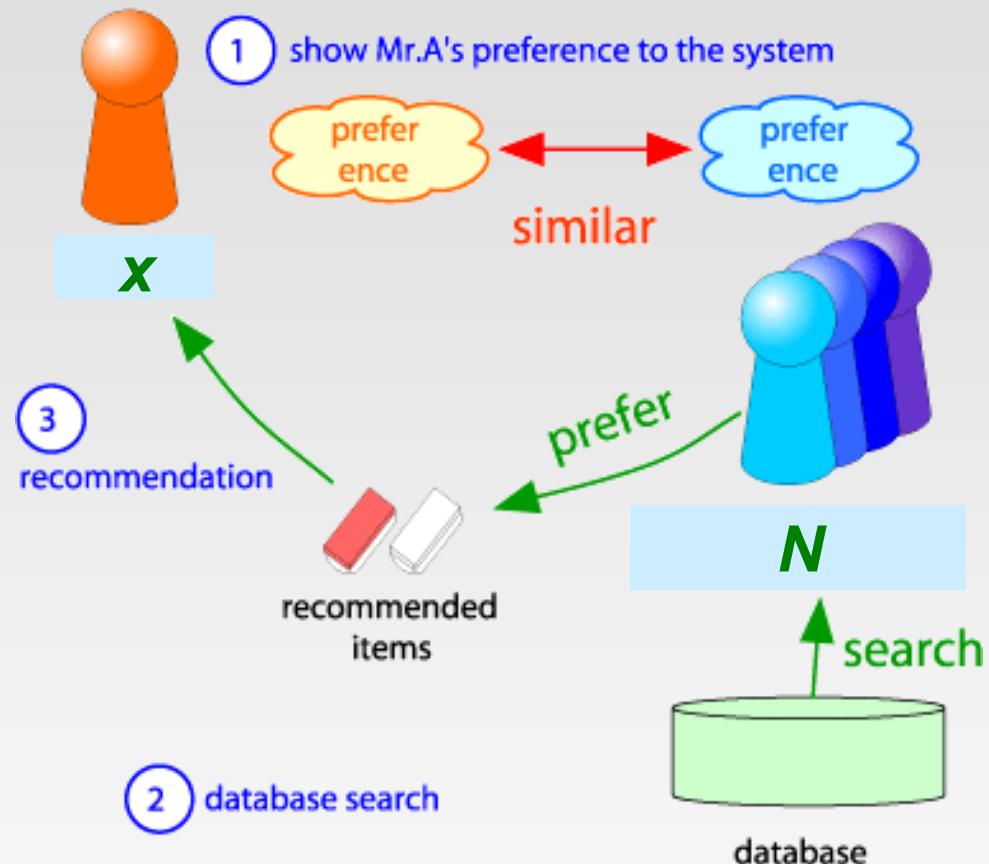
Collaborative Filtering (CF)

- ❖ The most prominent approach to generate recommendations
 - used by large, commercial e-commerce sites
 - well-understood, various algorithms and variations exist
 - applicable in many domains (book, movies, DVDs, ..)
- ❖ Approach
 - use the "wisdom of the crowd" to recommend items
- ❖ Basic assumption and idea
 - Users give ratings to catalog items (implicitly or explicitly)
 - Customers who had similar tastes in the past, will have similar tastes in the future



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



User-based Nearest-Neighbor Collaborative Filtering

❖ The basic technique

- Given an "active user" (Alice) and an item i not yet seen by Alice
 - ▶ find a set of users (peers/nearest neighbors) who liked the same items as Alice in the past **and** who have rated item i
 - ▶ use, e.g. the average of their ratings to predict, if Alice will like item i
 - ▶ do this for all items Alice has not seen and recommend the best-rated

❖ Basic assumption and idea

- If users had similar tastes in the past they will have similar tastes in the future
- User preferences remain stable and consistent over time

User-based Nearest-Neighbor Collaborative Filtering

❖ Example

- A database of ratings of the current user, Alice, and some other users is given:

| | Item1 | Item2 | Item3 | Item4 | Item5 |
|-------|-------|-------|-------|-------|-------|
| Alice | 5 | 3 | 4 | 4 | ? |
| User1 | 3 | 1 | 2 | 3 | 3 |
| User2 | 4 | 3 | 4 | 3 | 5 |
| User3 | 3 | 3 | 1 | 5 | 4 |
| User4 | 1 | 5 | 5 | 2 | 1 |

- Determine whether Alice will like or dislike *Item5*, which Alice has not yet rated or seen

User-based Nearest-Neighbor Collaborative Filtering

❖ Some first questions

- How do we measure similarity?
- How many neighbors should we consider?
- How do we generate a prediction from the neighbors' ratings?



| | Item1 | Item2 | Item3 | Item4 | Item5 |
|-------|-------|-------|-------|-------|-------|
| Alice | 5 | 3 | 4 | 4 | ? |
| User1 | 3 | 1 | 2 | 3 | 3 |
| User2 | 4 | 3 | 4 | 3 | 5 |
| User3 | 3 | 3 | 1 | 5 | 4 |
| User4 | 1 | 5 | 5 | 2 | 1 |

Finding “Similar” Users

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

用集合jaccard

$$r_x = [* , __, __, *, ***]$$

$$r_y = [* , __, **, **, __]$$

- Jaccard similarity measure** $\frac{||r_x \cap r_y||}{||r_x \cup r_y||}$ 但是会忽略实际
评分分数
 - Problem:** Ignores the value of the rating

r_x, r_y as sets:

$$r_x = \{1, 4, 5\}$$

$$r_y = \{1, 3, 4\}$$

- Cosine similarity measure**

$$\text{sim}(x, y) = \cos(r_x, r_y) = \frac{r_x \cdot r_y}{||r_x|| \cdot ||r_y||}$$

- Problem:** Treats missing ratings as “negative”

- Pearson correlation coefficient**

- S_{xy} = items rated by both users x and y

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

\bar{r}_x, \bar{r}_y ... avg.
rating of x, y

Similarity Metric

Cosine similarity:

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

| | 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: $1/5 < 2/4$ 直接错了

❖ Cosine similarity: $0.380 > 0.322$ 过于接近

- Considers missing ratings as “negative”
- 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 |

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

Notice cosine sim. is correlation when data is centered at 0

Similarity Metric (Cont')

- ❖ A popular similarity measure in user-based CF: **Pearson correlation**

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

- Possible similarity values between -1 and 1;

| | Item1 | Item2 | Item3 | Item4 | Item5 | |
|-------|-------|-------|-------|-------|-------|-------------|
| Alice | 5 | 3 | 4 | 4 | ? | |
| User1 | 3 | 1 | 2 | 3 | 3 | sim = 0,85 |
| User2 | 4 | 3 | 4 | 3 | 5 | sim = 0,00 |
| User3 | 3 | 3 | 1 | 5 | 4 | sim = 0,70 |
| User4 | 1 | 5 | 5 | 2 | 1 | sim = -0,79 |

Rating Predictions

From similarity metric to recommendations:

- ❖ Let r_x be the vector of user x 's ratings
- ❖ Let N be the set of k users most similar to x who have rated item i
- ❖ Prediction for item s of user x :
 - $r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi}$
 - $r_{xi} = \frac{\sum_{y \in N} s_{xy} \cdot r_{yi}}{\sum_{y \in N} s_{xy}}$
 - Other options?
- ❖ Many other tricks possible...

Shorthand:

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

CF: Common Practice

Before:

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

- ❖ Define **similarity** s_{ij} of users i and j
- ❖ Select k nearest neighbors $N(i; x)$
 - Users most similar to i , that have ratings on 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}

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

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

Memory-based and Model-based Approaches

- ❖ User-based CF is said to be "memory-based"
 - the rating matrix is directly used to find neighbors / make predictions
 - does not scale for most real-world scenarios
 - large e-commerce sites have tens of millions of customers and millions of items
- ❖ Model-based approaches
 - based on an offline pre-processing or "model-learning" phase
 - at run-time, only the learned model is used to make predictions
 - models are updated / re-trained periodically
 - large variety of techniques used
 - model-building and updating can be computationally expensive

Item-Item Collaborative Filtering

- ❖ So far: User-user collaborative filtering
- ❖ Another view: Item-item
 - Basic idea:
 - ▶ Use the similarity between items (and not users) to make predictions
 - For item i , find other similar items
 - 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 u on item j

$N(i; x)$... set items rated by x similar to i

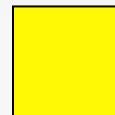
Item-Item Collaborative Filtering

- ❖ Example:
 - Look for items that are similar to Item5
 - Take Alice's ratings for these items to predict the rating for Item5

| | Item1 | Item2 | Item3 | Item4 | Item5 |
|-------|-------|-------|-------|-------|-------|
| Alice | 5 | 3 | 4 | 4 | ? |
| User1 | 3 | 1 | 2 | 3 | 3 |
| User2 | 4 | 3 | 4 | 3 | 5 |
| User3 | 3 | 3 | 1 | 5 | 4 |
| User4 | 1 | 5 | 5 | 2 | 1 |

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

| | 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 ($|N|=2$)

| | 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 | |



- estimate rating of movie 1 by user 5

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

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

$\text{sim}(1,m)$
 1.00
 -0.18
0.41
 -0.10
 -0.31
0.59

Neighbor selection:

Identify movies similar to movie 1, rated by user 5

Here 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

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

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

$\text{sim}(1,m)$

1.00

-0.18

0.41

-0.10

-0.31

0.59

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 | |

Predict by taking weighted average:

$$r_{1,5} = (0.41 \cdot 2 + 0.59 \cdot 3) / (0.41 + 0.59) = 2.6$$

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

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 item-item 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
- ❖ - 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
 - New items, Esoteric items
- ❖ - Popularity bias:
 - Cannot recommend items to someone with unique taste
 - Tends to recommend popular items

More on Ratings – Explicit Ratings

- ❖ Probably the most precise ratings
- ❖ Most commonly used (1 to 5, 1 to 7 Likert response scales)
- ❖ Main problems
 - Users not always willing to rate many items
 - ▶ number of available ratings could be too small → sparse rating matrices → poor recommendation quality
 - How to stimulate users to rate more items?

More on Ratings – Implicit Ratings

- ❖ Typically collected by the web shop or application in which the recommender system is embedded
- ❖ When a customer buys an item, for instance, many recommender systems interpret this behavior as a positive rating
- ❖ Clicks, page views, time spent on some page, demo downloads ...
- ❖ Implicit ratings can be collected constantly and do not require additional efforts from the side of the user
- ❖ Main problem
 - One cannot be sure whether the user behavior is correctly interpreted
 - For example, a user might not like all the books he or she has bought; the user also might have bought a book for someone else
- ❖ Implicit ratings can be used in addition to explicit ones; question of correctness of interpretation

Collaborative Filtering: Complexity

- ❖ Expensive step is finding k most similar customers: $\text{O}(|X|)$
- ❖ **Too expensive to do at runtime**
 - Could pre-compute
- ❖ Naïve pre-computation takes time $\text{O}(k \cdot |X|)$
 - X ... set of customers
- ❖ Ways of doing this:
 - Near-neighbor search in high dimensions (**LSH**)
 - Clustering
 - Dimensionality reduction
 -
- ❖ Supported by Hadoop: Apache Mahout
<https://mahout.apache.org/users/basics/algorithms.html>

What is a Good Recommendation in Practice?

- ❖ Total sales numbers
- ❖ Promotion of certain items
- ❖ ...
- ❖ Click-through-rates
- ❖ Interactivity on platform
- ❖ ...
- ❖ Customer return rates
- ❖ Customer satisfaction and loyalty



Evaluation

| movies | | | | | |
|--------|---|---|---|---|---|
| users | 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

Evaluating Predictions

- ❖ **Compare predictions with known ratings**

- Root-mean-square error (RMSE)

- ▶ $\sqrt{\sum_{xi} (r_{xi} - r_{xi}^*)^2}$ where r_{xi} is predicted, r_{xi}^* is the true rating of x on i

- Precision at top 10:

- ▶ % of those in top 10

- Rank Correlation:

- ▶ Spearman's *correlation* between system's and user's complete rankings

- ❖ **Another approach: 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

The Netflix Prize

❖ Training data

- 100 million ratings, 480,000 users, 17,770 movies
- 6 years of data: 2000-2005

❖ Test data

- Last few ratings of each user (2.8 million)
- **Evaluation criterion:** Root Mean Square Error (RMSE) =
$$\frac{1}{|R|} \sqrt{\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$$
- **Netflix's system RMSE: 0.9514**

❖ Competition

- 2,700+ teams
- **\$1 million** prize for 10% improvement on Netflix

The Netflix Utility Matrix R

Matrix R

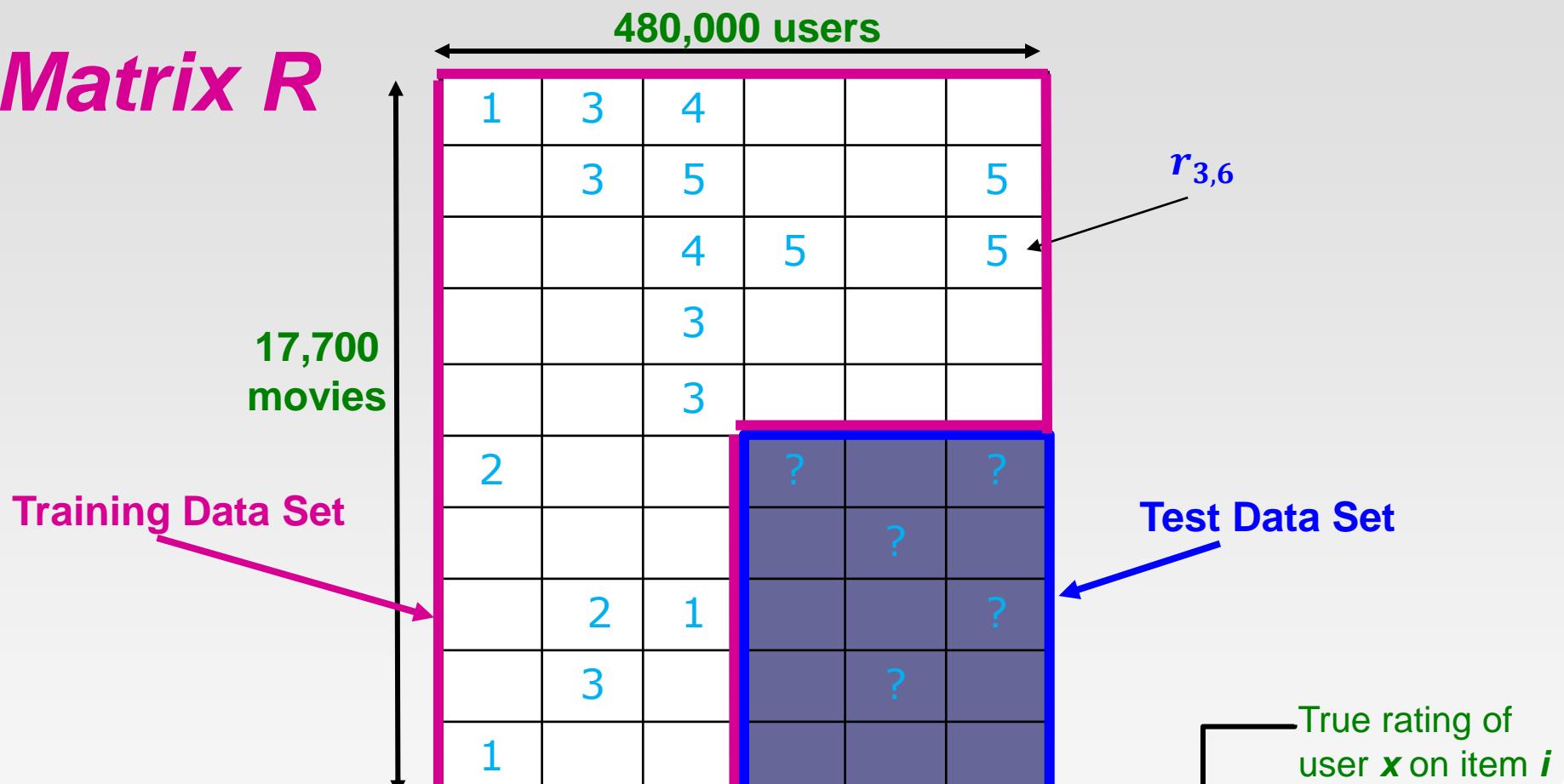
17,700
movies

480,000 users

| | | | | | |
|---|---|---|---|---|---|
| 1 | 3 | 4 | | | |
| | 3 | 5 | | | 5 |
| | | 4 | 5 | | 5 |
| | | | 3 | | |
| | | | 3 | | |
| 2 | | | 2 | | 2 |
| | | | | 5 | |
| | 2 | 1 | | | 1 |
| | | 3 | | 3 | |
| 1 | | | | | |

Utility Matrix R : Evaluation

Matrix R



$$\text{RMSE} = \frac{1}{|R|} \sqrt{\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$$

BellKor Recommender System

❖ The winner of the Netflix Challenge!

❖ Multi-scale modeling of the data:

Combine top level, “regional” modeling of the data, with a refined, local view:

➤ Global:

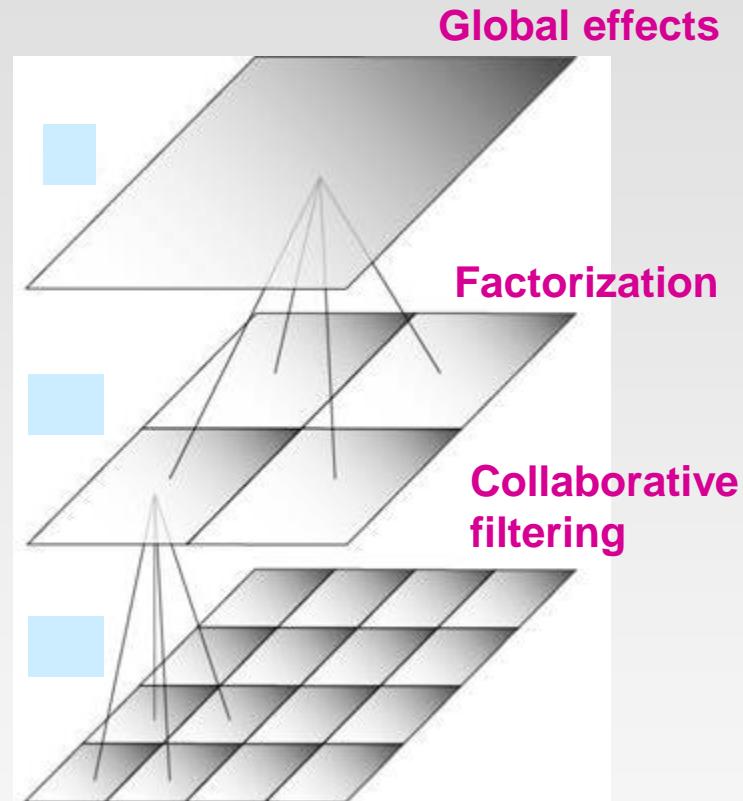
▶ Overall deviations of users/movies

➤ Factorization:

▶ Addressing “regional” effects

➤ Collaborative filtering:

▶ Extract local patterns



Modeling Local & Global Effects

❖ Global:

- Mean movie rating: **3.7 stars**
- *The Sixth Sense* is **0.5** stars above avg.
- Joe rates **0.2** stars below avg.

⇒ **Baseline estimation:**

Joe will rate *The Sixth Sense* 4 stars



❖ Local neighborhood (CF/NN):

- Joe didn't like related movie *Signs*
- ⇒ **Final estimate:**

Joe will rate *The Sixth Sense* 3.8 stars



Modeling Local & Global Effects

- ❖ In practice we get better estimates if we model deviations:

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

baseline estimate for r_{xi}

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

μ = overall mean rating

b_x = rating deviation of user x
= (avg. rating of user x) – μ

b_i = (avg. rating of movie i) – μ

Problems/Issues:

- 1) Similarity measures are “arbitrary”
- 2) Pairwise similarities neglect interdependencies among users
- 3) Taking a weighted average can be restricting

Solution: Instead of s_{ij} use w_{ij} that we estimate directly from data

Idea: Interpolation Weights w_{ij}

- ❖ Use a **weighted sum** rather than **weighted avg.**:

$$\widehat{r}_{xi} = b_{xi} + \sum_{j \in N(i;x)} w_{ij} (r_{xj} - b_{xj})$$

- ❖ **A few notes:**

- $N(i; x)$... set of movies rated by user x that are similar to movie i
- w_{ij} is the interpolation weight (some real number)
 - ▶ We allow: $\sum_{j \in N(i,x)} w_{ij} \neq 1$
- w_{ij} models interaction between pairs of movies (it does not depend on user x)

Idea: Interpolation Weights w_{ij}

❖ $\widehat{r}_{xi} = b_{xi} + \sum_{j \in N(i, x)} w_{ij} (r_{xj} - b_{xj})$

❖ How to set w_{ij} ?

- Remember, error metric is: $\frac{1}{|R|} \sqrt{\sum_{(i,x) \in R} (\widehat{r}_{xi} - r_{xi})^2}$ or equivalently
SSE: $\sum_{(i,x) \in R} (\widehat{r}_{xi} - r_{xi})^2$
- Find w_{ij} that minimize SSE on training data!
 - ▶ Models relationships between item i and its neighbors j
- w_{ij} can be learned/estimated based on x and all other users that rated i

Recommendations via Optimization

- ❖ **Goal:** Make good recommendations
 - Quantify goodness using **RMSE**:
Lower RMSE \Rightarrow better recommendations
 - Want to make good recommendations on items that user has not yet seen. **Can't really do this!**
 - **Let's build a system such that it works well on known (user, item) ratings**
And **hope** the system will also predict well the **unknown ratings**

Recommendations via Optimization

- ❖ Idea: Let's set values w such that they work well on known (user, item) ratings
- ❖ How to find such values w ?
- ❖ Idea: Define an objective function and solve the optimization problem
- ❖ Find w_{ij} that minimize SSE on training data!

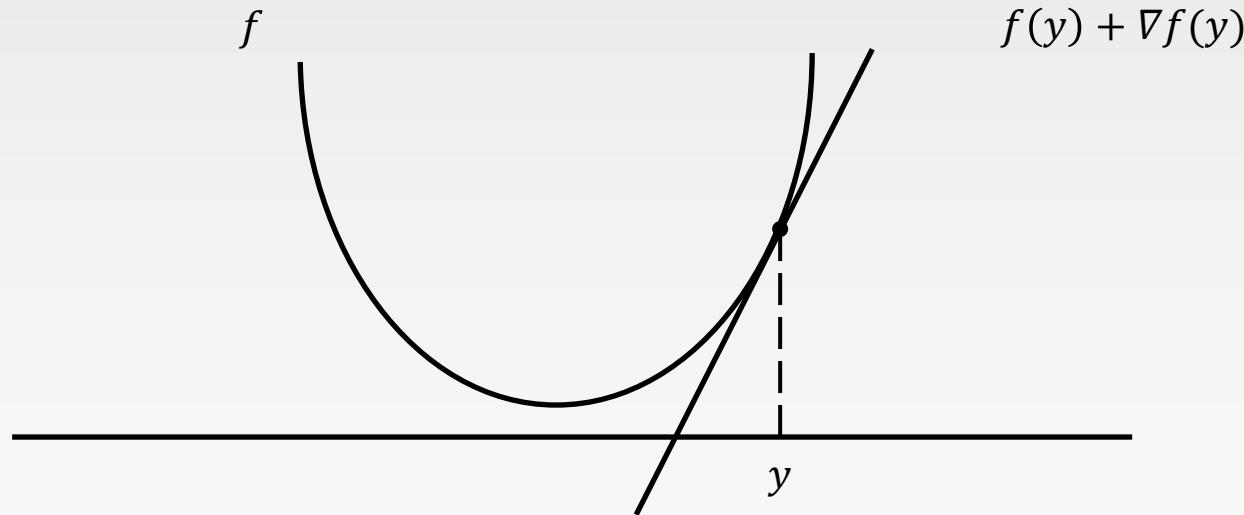
$$J(w) = \sum_{x,i} \left(\underbrace{\left[b_{xi} + \sum_{j \in N(i;x)} w_{ij}(r_{xj} - b_{xj}) \right]}_{\text{Predicted rating}} - r_{xi} \right)^2$$

True rating

- ❖ Think of w as a vector of numbers

Detour: Minimizing a function

- ❖ A simple way to minimize a function $f(x)$:
 - Compute the take a derivative ∇f
 - Start at some point y and evaluate $\nabla f(y)$
 - Make a step in the reverse direction of the gradient: $y = y - \nabla f(y)$
 - Repeat until converged



Interpolation Weights

- ❖ We have the optimization problem, now what?
- ❖ Gradient decent:

➤ Iterate until convergence: $w \leftarrow w - \eta \nabla_w J$

➤ where $\nabla_w J$ is the gradient (derivative evaluated on data):

$$\nabla_w J = \left[\frac{\partial J(w)}{\partial w_{ij}} \right]$$

$$= 2 \sum_{x,i} \left(\left[b_{xi} + \sum_{k \in N(i;x)} w_{ik} (r_{xk} - b_{xk}) \right] - r_{xi} \right) (r_{xj} - b_{xj})$$

for $j \in \{N(i; x), \forall i, \forall x\}$

else $\frac{\partial J(w)}{\partial w_{ij}} = 0$

➤ Note: We fix movie i , go over all r_{xi} , for every movie $j \in N(i; x)$, we compute $\frac{\partial J(w)}{\partial w_{ij}}$

while $|w_{new} - w_{old}| > \varepsilon$:

$$w_{old} = w_{new}$$

$$w_{new} = w_{old} - \eta \cdot \nabla w_{old}$$

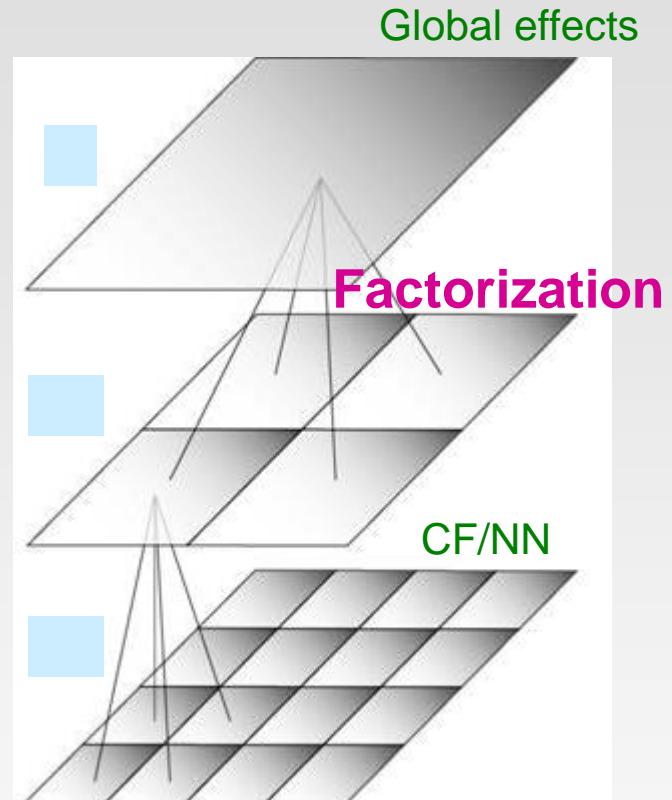
Interpolation Weights

❖ So far: $\widehat{r_{xi}} = b_{xi} + \sum_{j \in N(i;x)} w_{ij}(r_{xj} - b_{xj})$

- Weights w_{ij} derived based on their role; **no use of an arbitrary similarity measure** ($w_{ij} \neq s_{ij}$)
- Explicitly account for interrelationships among the neighboring movies

❖ Next: Latent factor model

- Extract “regional” correlations



More model-based approaches

- ❖ Plethora of different techniques proposed in the last years, e.g.,
 - Matrix factorization techniques, statistics
 - ▶ singular value decomposition, principal component analysis
 - Association rule mining
 - ▶ compare: shopping basket analysis
 - Probabilistic models
 - ▶ clustering models, Bayesian networks, probabilistic Latent Semantic Analysis
 - Various other machine learning approaches
- ❖ Costs of pre-processing
 - Usually not discussed
 - Incremental updates possible?

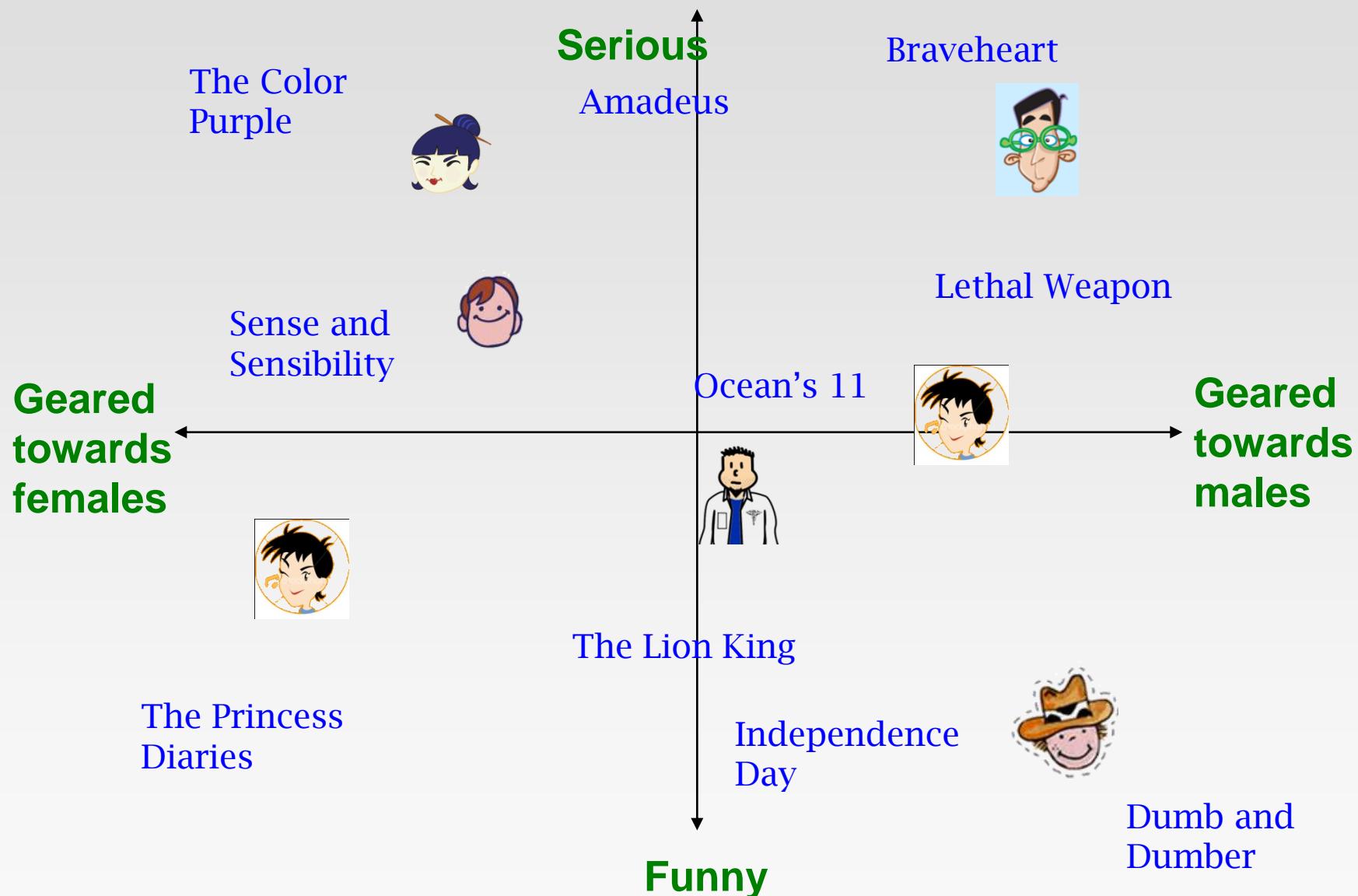
Matrix Factorization

- ❖ Informally, the SVD theorem (Golub and Kahan 1965) states that a given matrix M can be decomposed into a product of three matrices as follows

$$M = U \times \Sigma \times V^T$$

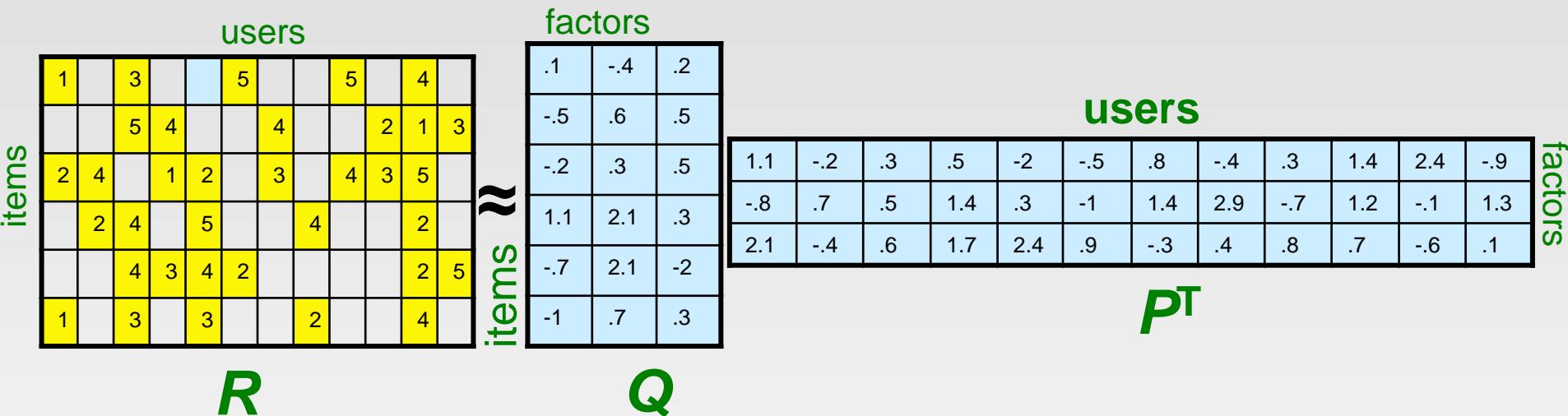
- where U and V are called *left* and *right singular vectors* and the values of the diagonal of Σ are called the *singular values*
- ❖ We can approximate the full matrix by observing only the most important features – those with the largest singular values
- ❖ In the example, we calculate U , V , and Σ (with the help of some linear algebra software) but retain only the two most important features by taking only the first two columns of U and V^T

Latent Factor Models



Matrix Factorization

- ❖ “SVD” on Netflix data: $\mathbf{R} \approx \mathbf{Q} \cdot \mathbf{P}^T$



- ❖ For now let's assume we can approximate the rating matrix \mathbf{R} as a product of “thin” $\mathbf{Q} \cdot \mathbf{P}^T$
 - \mathbf{R} has missing entries but let's ignore that for now!
 - ▶ Basically, we will want the reconstruction error to be small on known ratings and we don't care about the values on the missing ones

Ratings as Products of Factors

- ❖ How to estimate the missing rating of user x for item i ?

| users | | | | | | | | | |
|-------|---|---|---|-----|---|---|---|---|---|
| items | 1 | 3 | | 5 | | 5 | | 4 | |
| 1 | | 5 | 4 | 2.4 | 4 | | 2 | 1 | 3 |
| 2 | 4 | | 1 | 2 | | 3 | 4 | 3 | 5 |
| | 2 | 4 | | 5 | | 4 | | 2 | |
| | | 4 | 3 | 4 | 2 | | | 2 | 5 |
| 1 | 3 | 3 | | | 2 | | | 4 | |

≈

$$\hat{r}_{xi} = q_i \cdot p_x^T$$

$$= \sum_f q_{if} \cdot p_{xf}$$

q_i = row i of Q
 p_x = column x of P^T

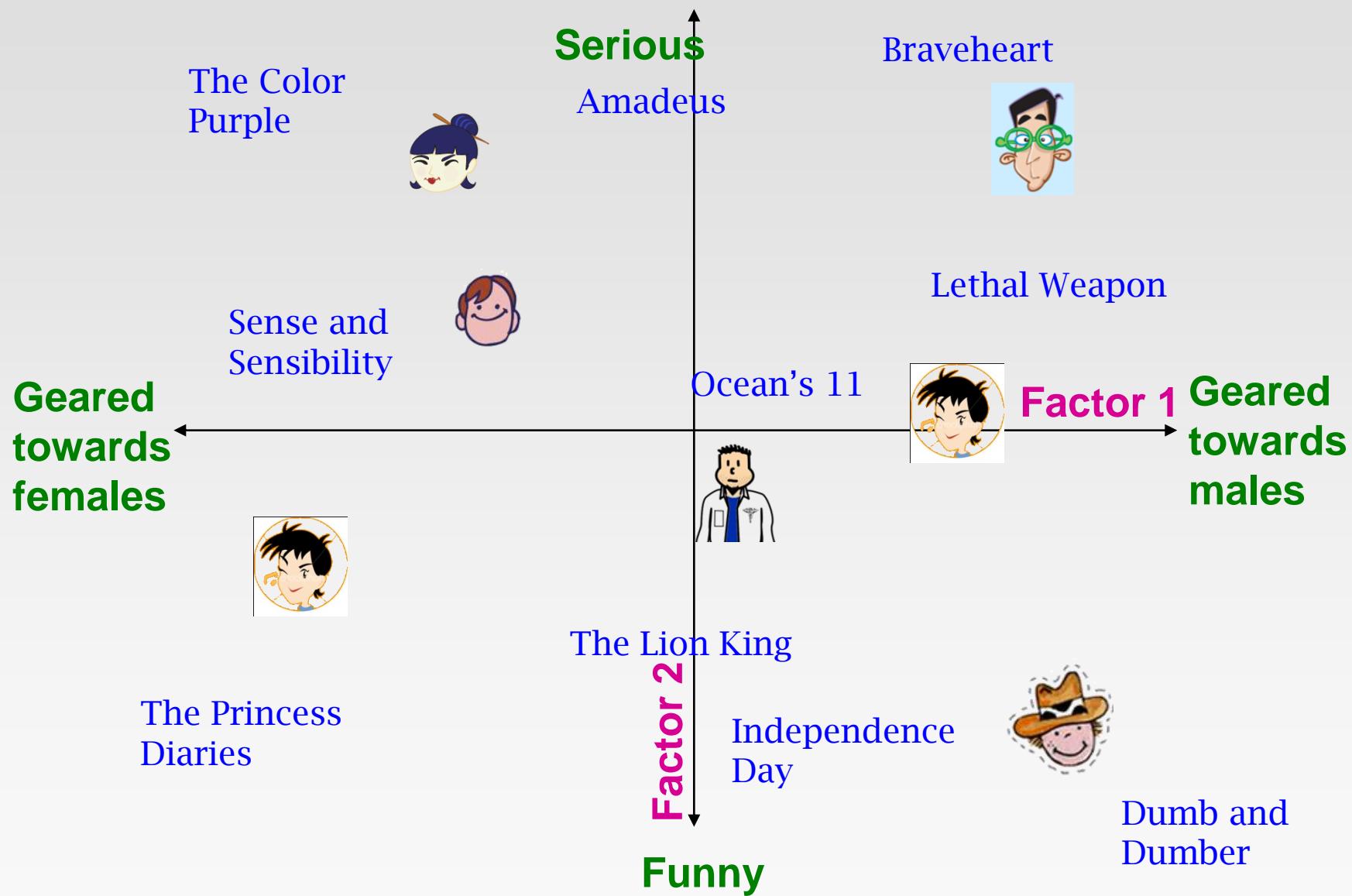
| items | .1 | -.4 | .2 |
|-----------|-----|-----|----|
| • factors | -.5 | .6 | .5 |
| | -.2 | .3 | .5 |
| | 1.1 | 2.1 | .3 |
| | -.7 | 2.1 | -2 |
| | -1 | .7 | .3 |

Q

| users | 1.1 | -.2 | .3 | .5 | -2 | -.5 | .8 | -.4 | .3 | 1.4 | 2.4 | -.9 |
|-----------|-----|-----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| • factors | -.8 | .7 | .5 | 1.4 | .3 | -1 | 1.4 | 2.9 | -.7 | 1.2 | -.1 | 1.3 |
| | 2.1 | -.4 | .6 | 1.7 | 2.4 | .9 | -.3 | .4 | .8 | .7 | -.6 | .1 |

P^T

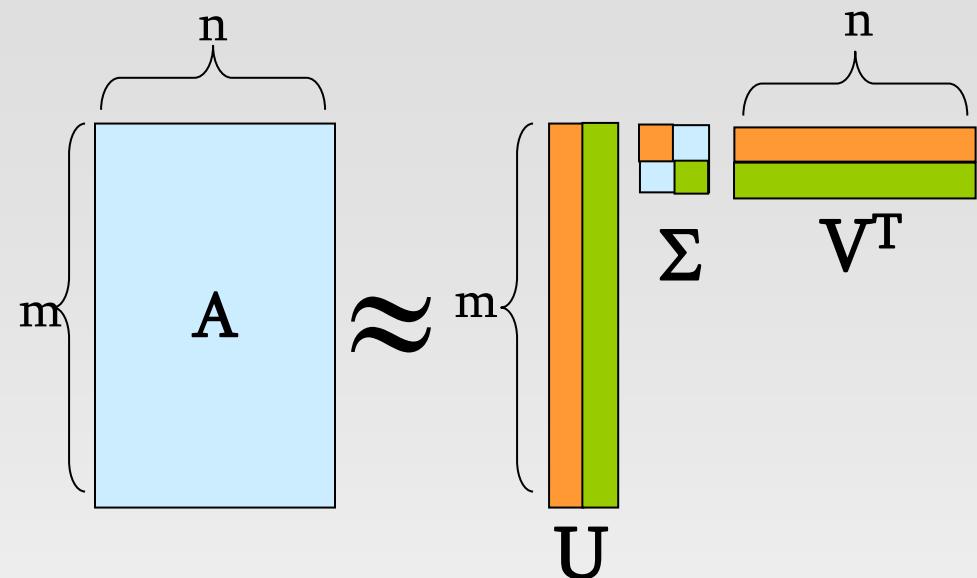
Latent Factor Models



Recap: SVD

❖ Remember SVD:

- A: Input data matrix
- U: Left singular vecs
- V: Right singular vecs
- Σ: Singular values



❖ So in our case:

“SVD” on Netflix data: $R \approx Q \cdot P^T$

$$A = R, \quad Q = U, \quad P^T = \Sigma \quad V^T$$

$$\hat{r}_{xi} = q_i \cdot p_x$$

SVD: More good stuff

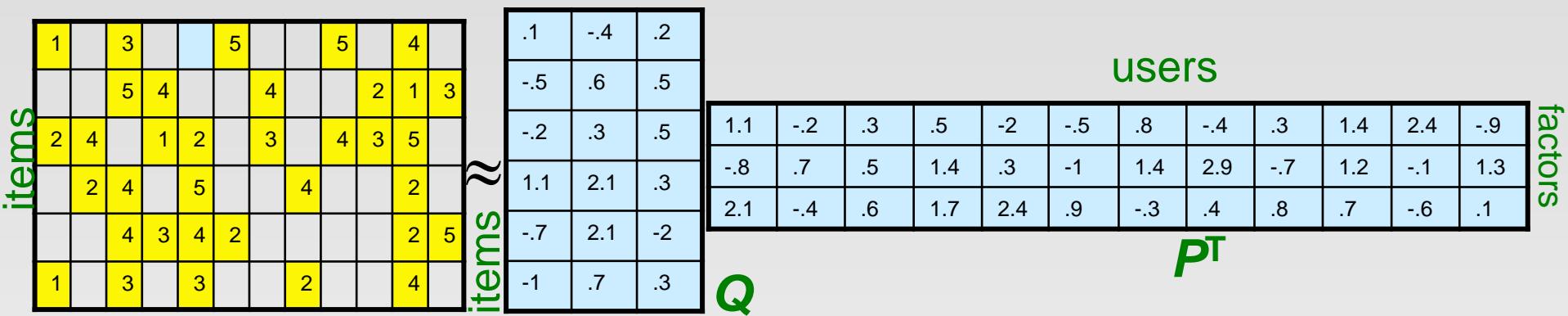
- ❖ We already know that SVD gives minimum reconstruction error (Sum of Squared Errors):

$$\min_{U,V,\Sigma} \sum_{ij \in A} (A_{ij} - [U\Sigma V^T]_{ij})^2$$

- ❖ Note two things:

- SSE and RMSE are monotonically related:
 - ▶ $RMSE = \frac{1}{c}\sqrt{SSE}$ Great news: SVD is minimizing RMSE
- Complication: The sum in SVD error term is over all entries (no-rating interpreted as zero-rating).
But our R has missing entries!

Matrix Factorization



- ❖ SVD isn't defined when entries are missing!
- ❖ Use specialized methods to find P, Q

➤ $\min_{P,Q} \sum_{(i,x) \in R} (r_{xi} - q_i \cdot p_x^T)^2$

➤ Note:

- We don't require cols of P, Q to be orthogonal/unit length
- P, Q map users/movies to a latent space

$$\hat{r}_{xi} = q_i \cdot p_x^T$$

Modeling Biases and Interactions

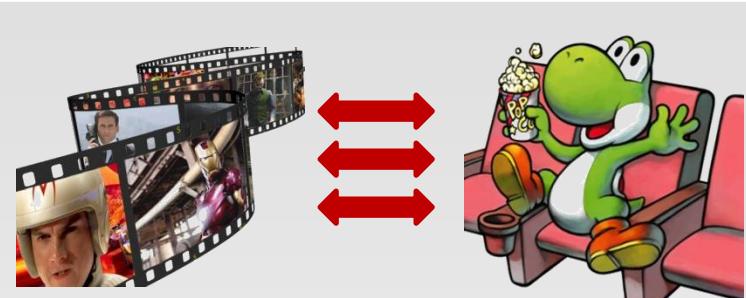
user bias



movie bias



user-movie interaction



Baseline predictor

- Separates users and movies
- Benefits from insights into user's behavior
- Among the main practical contributions of the competition

User-Movie interaction

- Characterizes the matching between users and movies
- Attracts most research in the field
- Benefits from algorithmic and mathematical innovations

μ = overall mean rating

b_x = bias of user x

b_i = bias of movie i

Baseline Predictor

- ❖ We have expectations on the rating by user x of movie i , even without estimating x 's attitude towards movies like i



+



- Rating scale of user x
- Values of other ratings user gave recently (day-specific mood, anchoring, multi-user accounts)

- (Recent) popularity of movie i
- Selection bias; related to number of ratings user gave on the same day (“frequency”)

Putting It All Together

$$r_{xi} = \mu + b_x + b_i + q_i \cdot p_x$$

Overall mean rating Bias for user x Bias for movie i User-Movie interaction

❖ Example:

- Mean rating: $\mu = 3.7$
- You are a critical reviewer: your ratings are 1 star lower than the mean: $b_x = -1$
- Star Wars gets a mean rating of 0.5 higher than average movie: $b_i = +0.5$
- Predicted rating for you on Star Wars:
 $= 3.7 - 1 + 0.5 = 3.2$

Fitting the New Model

- ❖ **Solve:**

$$\min_{Q,P} \sum_{(x,i) \in R} (r_{xi} - (\mu + b_x + b_i + q_i p_x))^2$$

goodness of fit

$$+ \left(\lambda_1 \sum_i \|q_i\|^2 + \lambda_2 \sum_x \|p_x\|^2 + \lambda_3 \sum_x \|b_x\|^2 + \lambda_4 \sum_i \|b_i\|^2 \right)$$

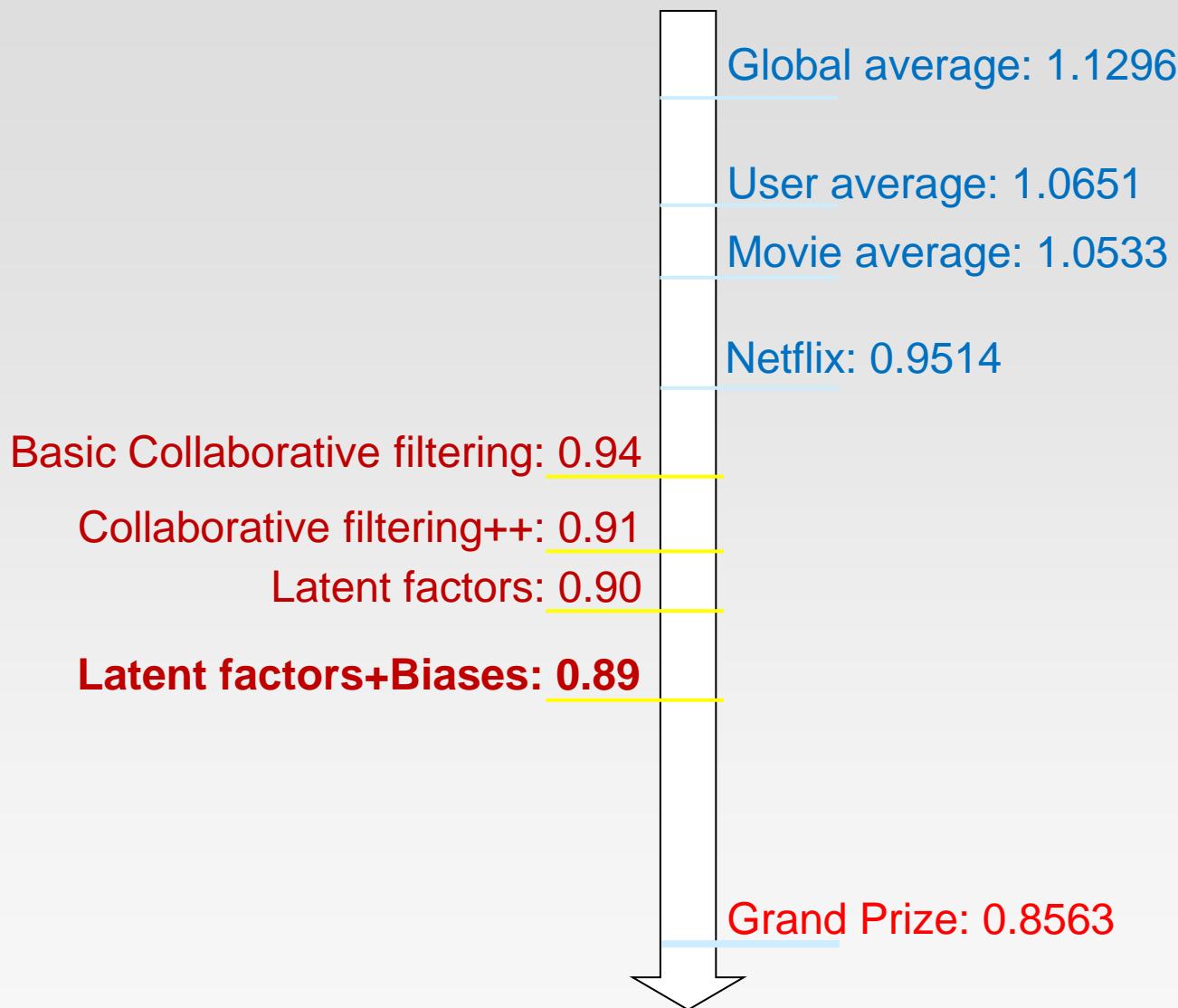
regularization

λ is selected via grid-search
on a validation set

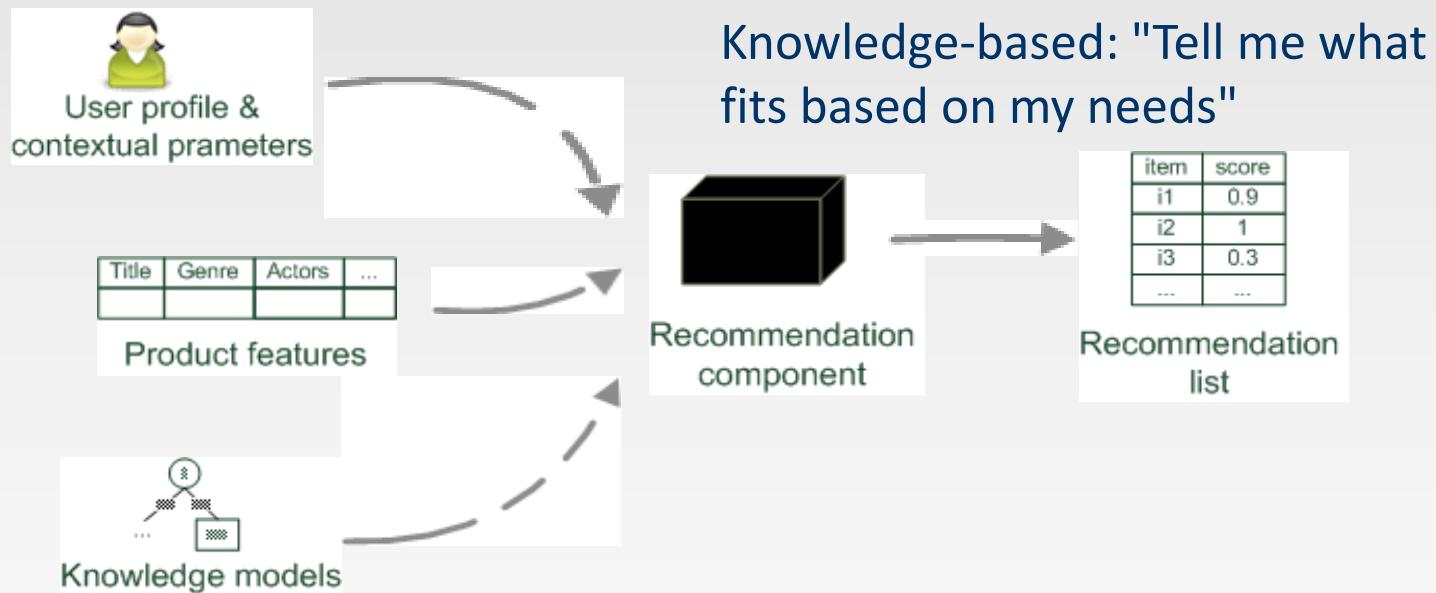
- ❖ **Stochastic gradient decent to find parameters**

- **Note:** Both biases b_x , b_i as well as interactions q_i , p_x are treated as parameters (we estimate them)

Performance of Various Methods



Part 4: Knowledge-Based Recommendation



Why do we need knowledge-based recommendation?

- ❖ Products with low number of available ratings



- ❖ Time span plays an important role
 - Five-year-old ratings for computers
 - User lifestyle or family situation changes
- ❖ Customers want to define their requirements explicitly
 - “The color of the car should be black”

Knowledge-based Recommendation

- ❖ Constraint-based
 - based on explicitly defined set of recommendation rules
 - fulfill recommendation rules
- ❖ Case-based
 - based on different types of similarity measures
 - retrieve items that are similar to specified requirements
- ❖ Both approaches are similar in their **conversational** recommendation process
 - users specify the requirements
 - systems try to identify solutions
 - if no solution can be found, users change requirements

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

References

- ❖ Chapter 9 of Mining of Massive Datasets.
- ❖ Tutorial: Recommender Systems. IJCAI 2013.

End of Chapter 9