

Information Retrieval

CS 547/DS 547

Worcester Polytechnic Institute
Department of Computer Science
Instructor: Prof. Kyumin Lee

Upcoming Schedule

- March 1: Midterm Exam
- March 3: due date of HW3
- March 17: due date of project proposal
- March 21: due date of proposal presentation slides

Midterm

- The exam will be held at 6pm next Wednesday in class.
- The exam is closed book.
- You may prepare and use one standard 8.5" by 11" piece of paper with any notes you think appropriate or significant.
- You may use a calculator if it make you feel comfortable. But no other electronic devices are allowed (e.g., cell phone, tablet and computer).

Previous Class...

Boolean Retrieval
Model

Previous Class...

Boolean Retrieval
Model

Inverted index

Previous Class...

Query Optimization

Previous Class...

Query Optimization

Preprocessing
Documents
→Tokenization,
Normalization,
Stemming, Stop words

Previous Class...

Query Optimization

Preprocessing
Documents
→Tokenization,
Normalization,
Stemming, Stop words

Skip pointers &
Positional index

Previous Class...

Wild-card queries
→ Permuterm Index

Previous Class...

TF and IDF

Previous Class...

TF and IDF

tf-idf weighting

$$w_{t,d} = (1 + \log \text{tf}_{t,d}) \cdot \log \frac{N}{\text{df}_t}$$

Previous Class...

Cosine Similarity

Previous Class...

Cosine Similarity

Computing scores in a
complete search system

Previous Class...

Statistical Language
Models

Crawler

Previous Class...

Statistical Language
Models

Crawler

Web APIs

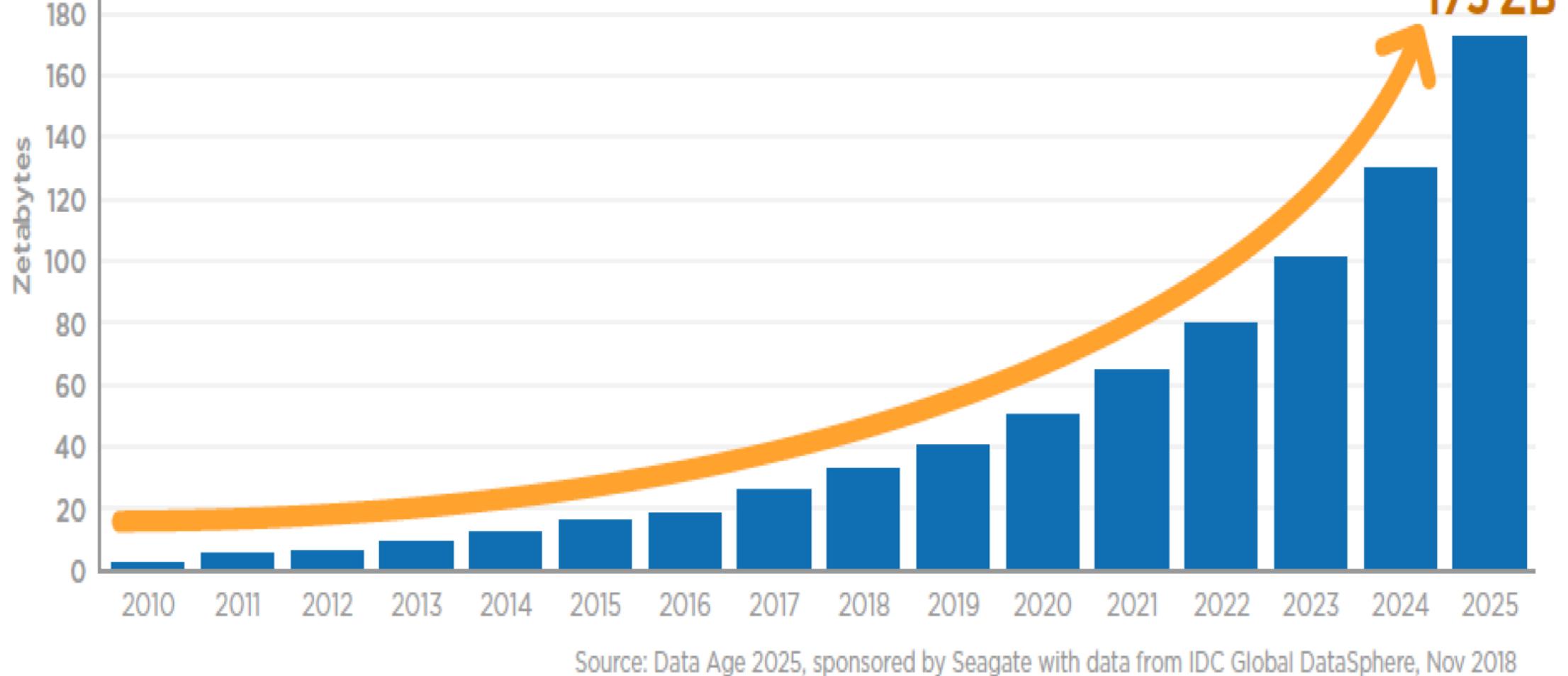
Previous Class...

PageRank

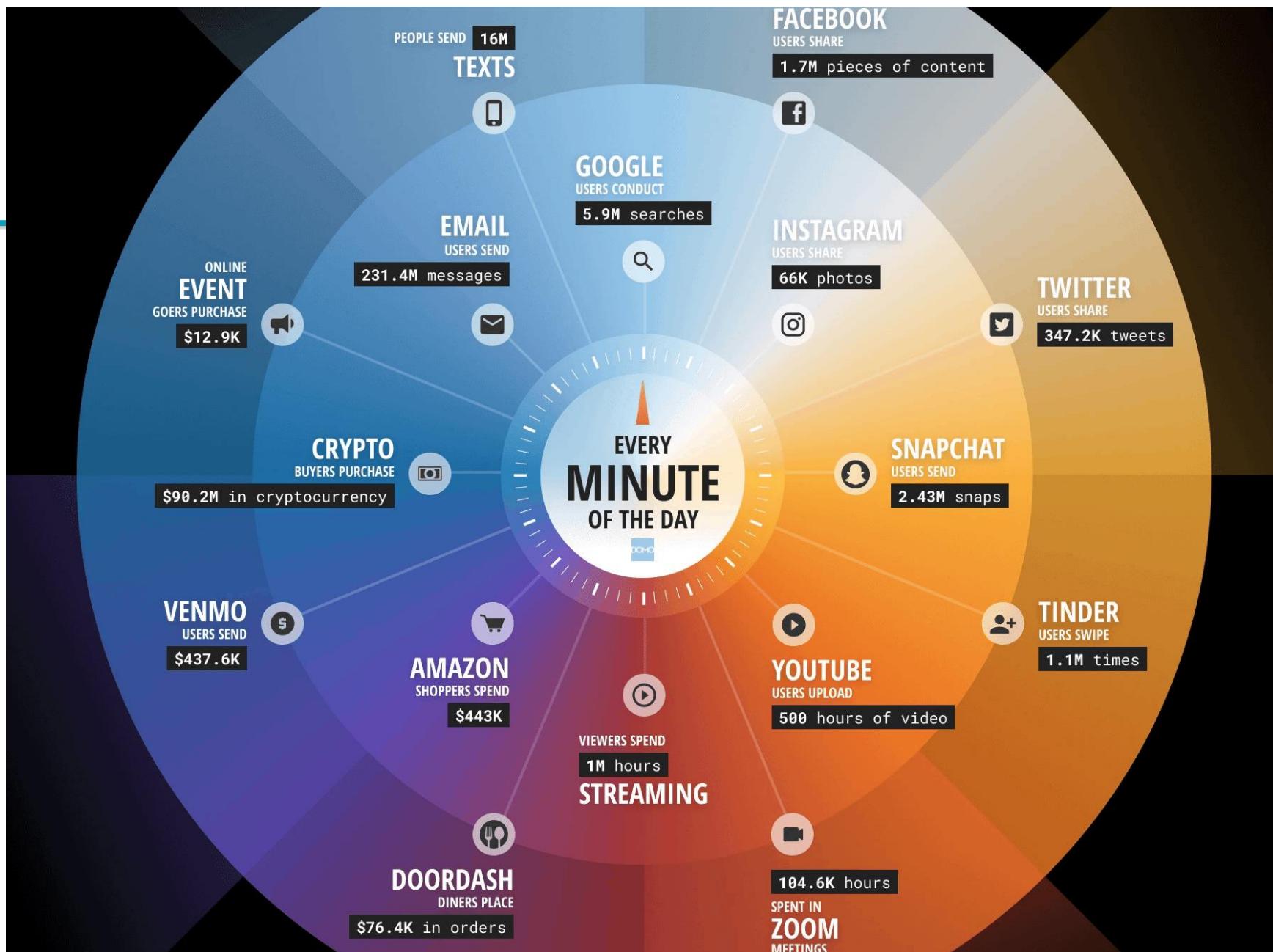
Evaluation

Recommenders

Annual Size of the Global Datasphere



In 2022



How Many Products Does Amazon Sell?



Amazon itself sells over 12 million products. If you take into account all products sold on the Amazon marketplace by third-party sellers, that number rises to more than 353 million products.

How Many Orders Does Amazon Get A Day?

Amazon ships approximately **1.6 million packages a day**.

That works out to more than **66 thousand orders per hour**, and **18.5 orders per second**.

Amazon Sales Statistics: How Much Amazon Makes in a Day

In 2019, Amazon made \$141.25 billion in retail product sales. This comes out to an average of \$385 million each day.



About Spotify

Spotify transformed music listening forever when it launched in 2008. Discover, manage and share over 100 million tracks and 5 million podcasts titles, for free, or upgrade to Spotify Premium to access exclusive features for music including improved sound quality and an on-demand, offline, and ad-free music listening experience. Today, Spotify is the world's most popular audio streaming subscription service with 489 million users, including 205 million subscribers in more than 180 markets.

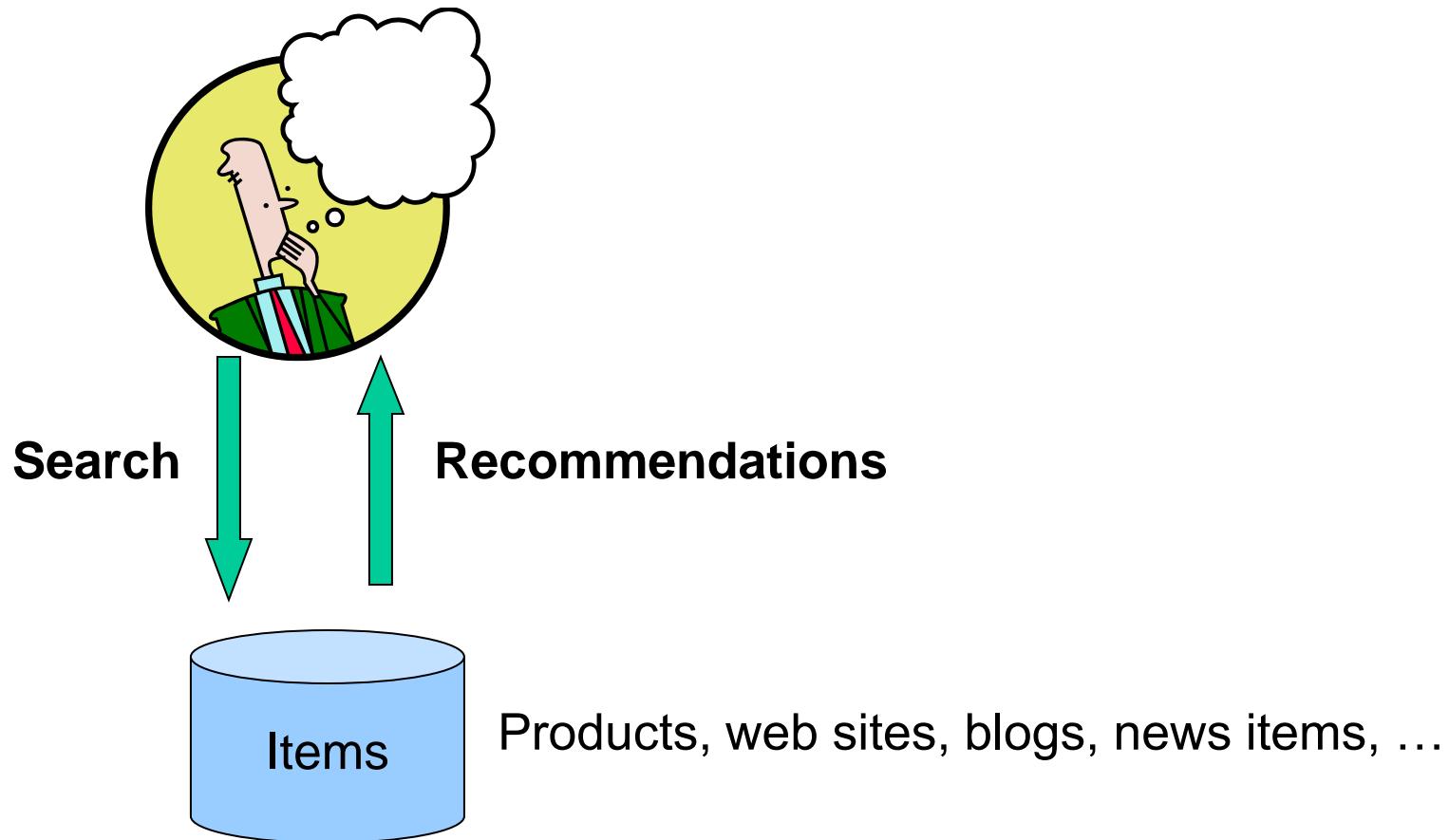


YouTube Usage Statistics

YouTube users consume 1 billion hours' worth of video every day. This figure translates to approximately 5 billion videos every day.

- The average time spent on YouTube by a user is 11.24 minutes.
- It is estimated that YouTube has around 5 billion videos.
- 500 hours of video is uploaded to YouTube every minute.

Recommendations



Amazon.com: Online Shopping for Electronics, Apparel, Computers, Books, DVDs & more

www.amazon.com

amazon Prime

James's Amazon.com Today's Deals Gift Cards Sell Help

Shop by Department Search All Go Hello, James Your Account Your Prime Cart Wish List

Back-to-School Deals Sponsored by Lysol > Shop now

Instant Video Digital Music Store Cloud Drive Kindle Appstore for Android Digital Games & Software Audible Audiobooks

eTextbooks for your iPad

Save up to 80% with eTextbooks and carry less

> Learn more

Top Picks: Women's Denim Really Great Toys Off to College

Building on a Family Legacy

Kathryn Parish's son added Internet sales to build on his family's six-decade tradition of delighting kids and parents with quality toys.

> "Mom got a big kick out of it"

> Shop Really Great Toys

One of thousands of small businesses thriving because of Amazon customers.

Related to Items You've Viewed

You viewed	Customers who viewed this also viewed
Introduction to Information Retrieval Hinrich Schütze, Christopher D. Manning, Prabhakar Raghavan Hardcover ★★★★★ (21) \$69.00 \$57.42	Taming Text: How to Find, Organize... Grant S. Ingersoll, Thomas S. Morton, ... Paperback ★★★★★ (9) \$44.99 \$31.65
Foundations of Statistical Natural... Christopher D. Manning and Hinrich Schütze	Speech and Language Processing, 2nd... Daniel Jurafsky, James H. Martin Hardcover ★★★★★ (18) \$177.29 \$141.37
Information Retrieval: Implementing... Stefan Buetzner, Charles L. A. Ongaro, ... Hardcover ★★★★★ (4) \$62.99 \$40.00	Search Engines: Information Retrieval... Bruce Croft, Donald Metzler, Trevor... Hardcover ★★★★★ (9) \$133.00 \$111.84

> View or edit your browsing history

kindle fire HDX From \$229 > Shop now

GOOGLE HANGOUTS BUILT-IN chromebook Buy Now

School Lists > Suggested supplies by grade

Prep Your Fuel for School Nerf N-Strike Elite Strong Arm

Speed and mobility are yours with the quick draws and fast firing of the Strongarm blaster. [Read more](#) \$42.99 \$9.44

Best Sellers Sports & Outdoors : Golf Clubs

www.netflix.com/WiHome

Reader

NETFLIX

Browse Taste Profile **KIDS** DVDs

Titles, People, Genres

James

Recently Watched

Popular on Netflix

Top Picks for James

Home

All

Music

Boston Dynamics

Live

Gaming

Culinary arts

Tools

Background music

Game shows

Driving

Basketball

Algorithms

Classical Music

History

>

Shorts

Subscriptions

Library

History

Your videos

Watch later

Liked videos

Subscriptions

AI Coffee Break ...

Alexander Amini

Browse channels

Explore

Trending

Shopping

Music

Movies & TV

**Do You Love Me?**

Boston Dynamics

38M views • 2 years ago

**Why we all need subtitles now**

Vox

8.8M views • 1 month ago

**How Singapore Airlines Makes 50,000 In-Flight Meals A Day | Big Business |...**

Insider Business

1.5M views • 3 days ago

**Glitterbomb Trap Catches Phone Scammer (who gets arrested)**

Mark Rober

70M views • 1 year ago

**Smith: "Call of Duty" more available after Nintendo deal**

Quest Means Business

2K views • 20 hours ago

**The 4 Reasons Why You're Poor**

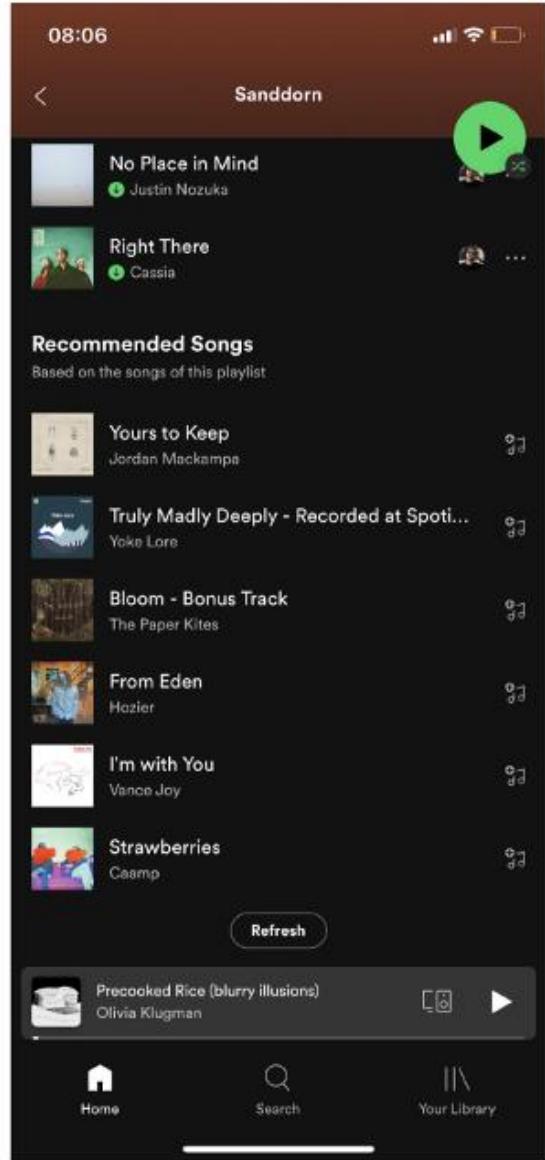
Iman Gadzhi

838K views • 5 months ago

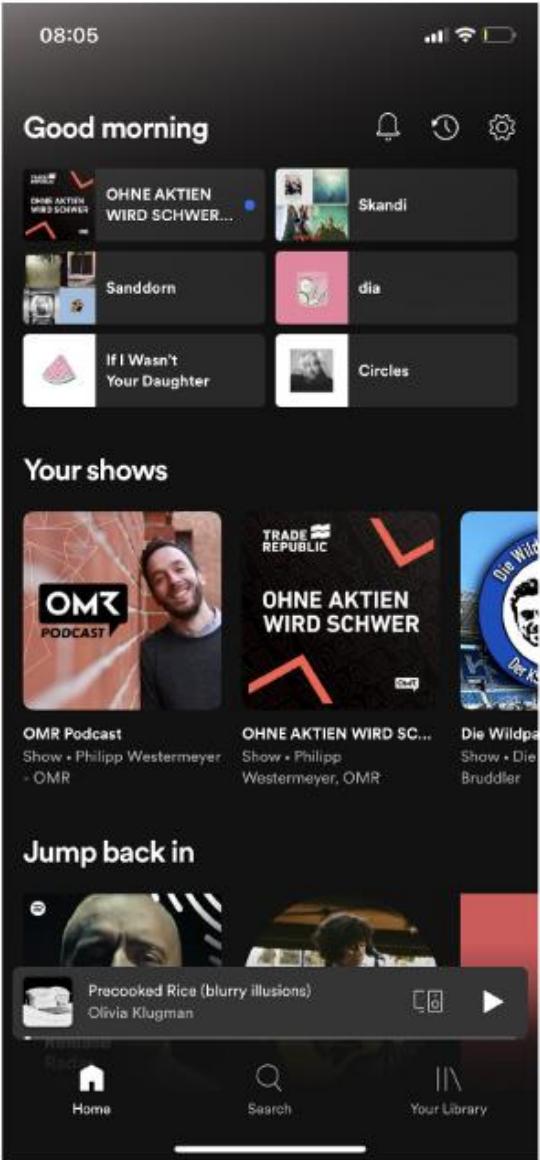
**NBA bloopers but they keep getting more embarrassing**

CoshReport

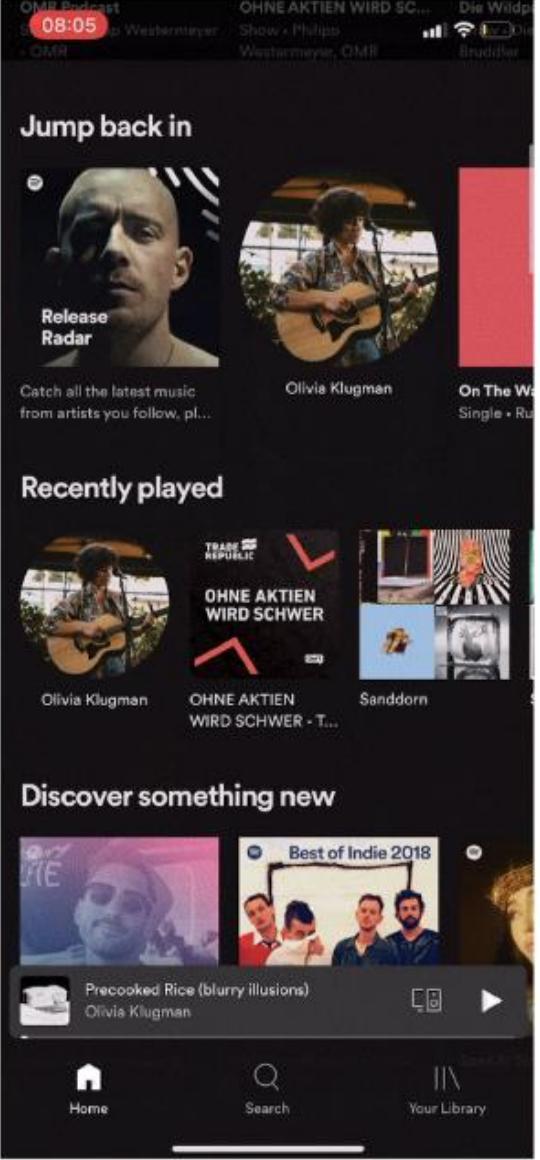
2.9M views • 10 months ago



Recommended Songs below



Home Screen with Shelves



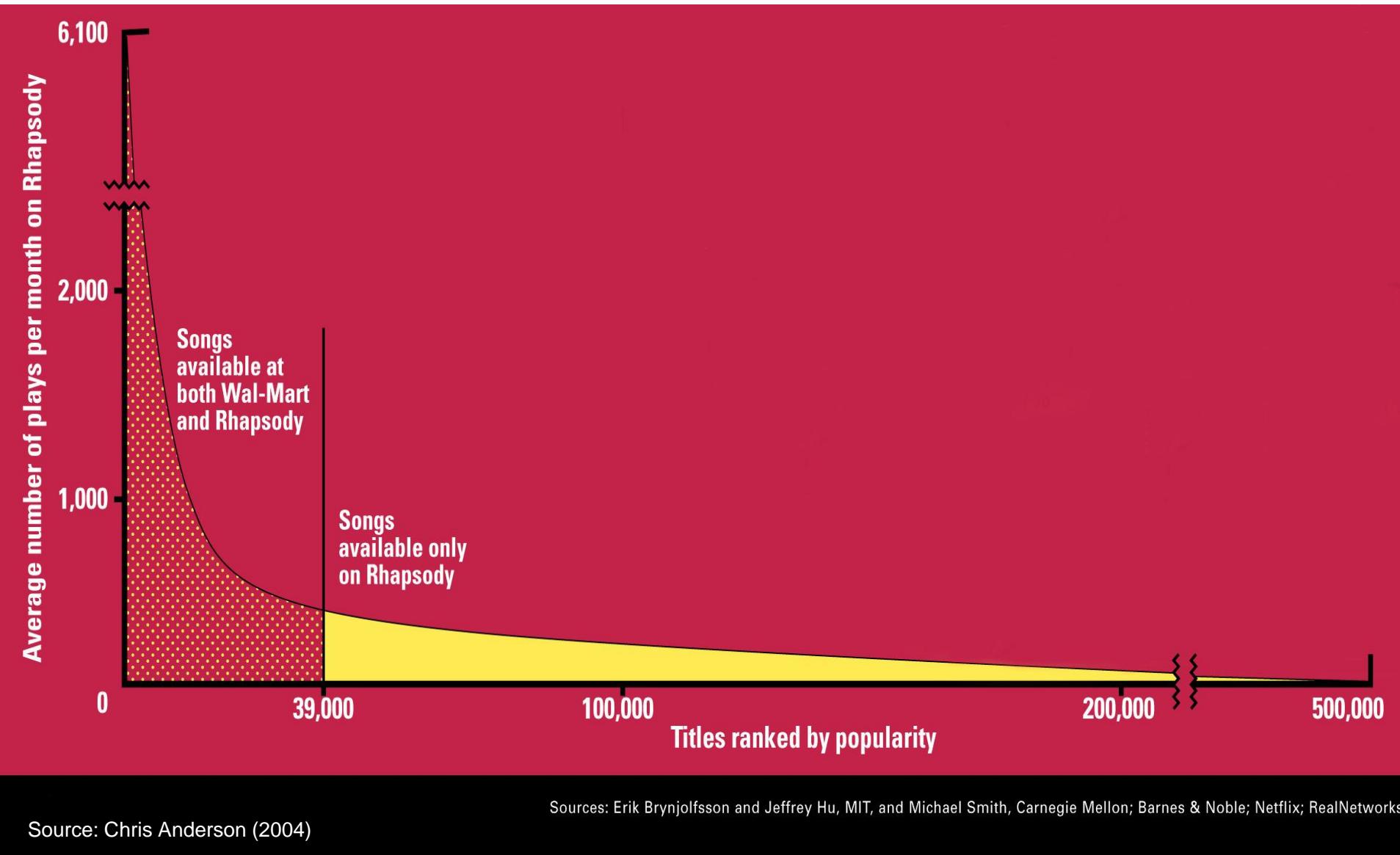
Jump back in, Recently played

Other examples of recommenders?

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>

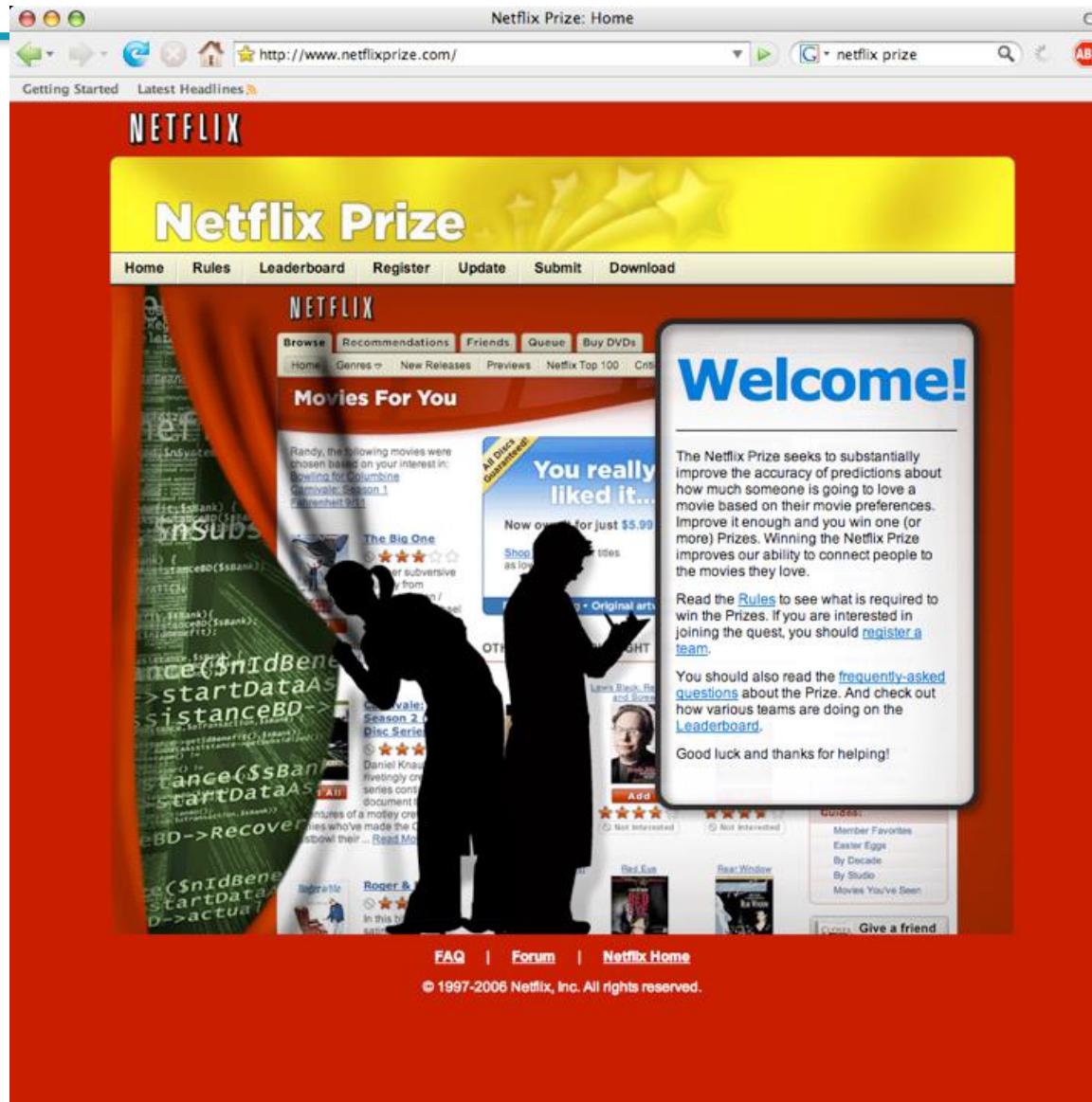
The Long Tail



Recommendation Types

- Editorial and hand curated
 - List of favorites
 - Lists of “essential” items
- Simple aggregates
 - Top 10, Most Popular, Recent Uploads
- Tailored to individual users
 - Amazon, Netflix, ...

\$\$\$



Formal Model

- X = set of Customers (Users)
- 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

	The Lego Movie	The Fault in Our Stars	Guardians of the Galaxy	Star Wars
Alice	1		0.2	
Bob		0.5		0.3
Carol	0.2		1	
David				0.4

Key Problems

- Gathering “known” ratings for matrix
- Extrapolate unknown ratings from the known ones
 - Mainly interested in high unknown ratings
 - Don’t care about finding what you **don’t like**, but rather what you like
- Evaluating extrapolation methods
 - How do we know if we’ve done a good job?

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, clicks, listen to a song
 - What about low ratings?
 - Called negative sampling
 - Random sampling, popularity-based sampling, and so on.

(2) Extrapolating Utilities

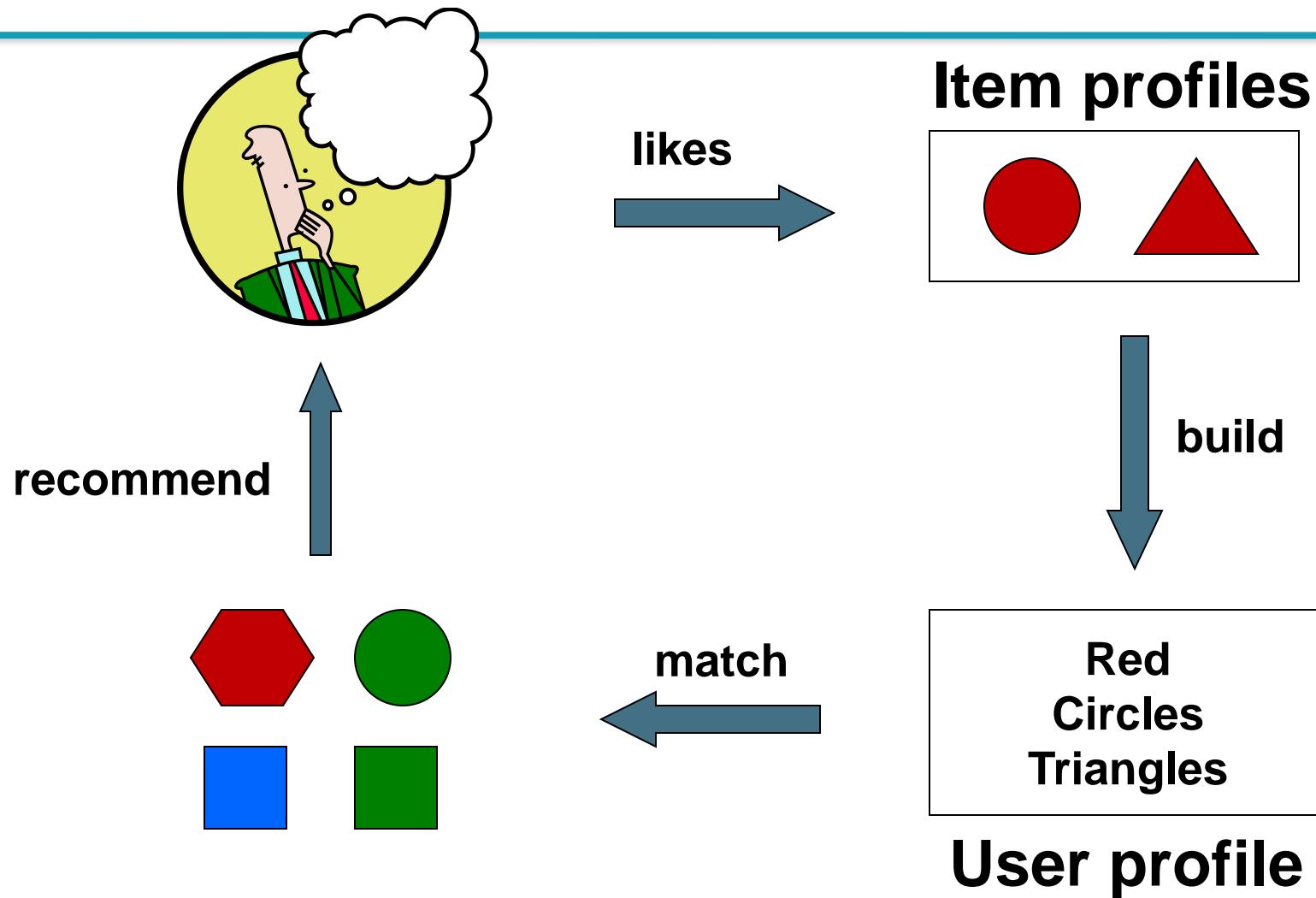
- Key problem: Utility matrix U is sparse
 - Most people have not rated most items
 - Cold start:
 - New items have no ratings
 - New users have no history
- Three main approaches
 - Content-based
 - Collaborative
 - Latent factor based
- These days, people focus on neural network/deep learning based approaches to capture more complex relation between users and items

Content-based Recommender Systems

Content-based Recommendations

- **Main idea:** recommend items to customer x similar to previous items rated highly by x
- Movie recommendations
 - recommend movies with same actor(s), director, genre, ...
- 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 of features (vectors!)
 - Movies: author, title, actor, director,...
 - Text: Set of “important” words in document
- How to pick important features?
 - Usual heuristic is TF-IDF

Sidenote: TF-IDF

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

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

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

← Whichever doc/item has the max # of frequency of term i

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/item 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 \mathbf{x} and item profile \mathbf{i} , estimate
 - $u(\mathbf{x}, \mathbf{i}) = \cos(\mathbf{x}, \mathbf{i}) = \mathbf{x} \cdot \mathbf{i} / (\|\mathbf{x}\| \|\mathbf{i}\|)$

Advantages of 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
- Can provide explanations of recommended items by listing content-features that caused an item to be recommended

Limitations of 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

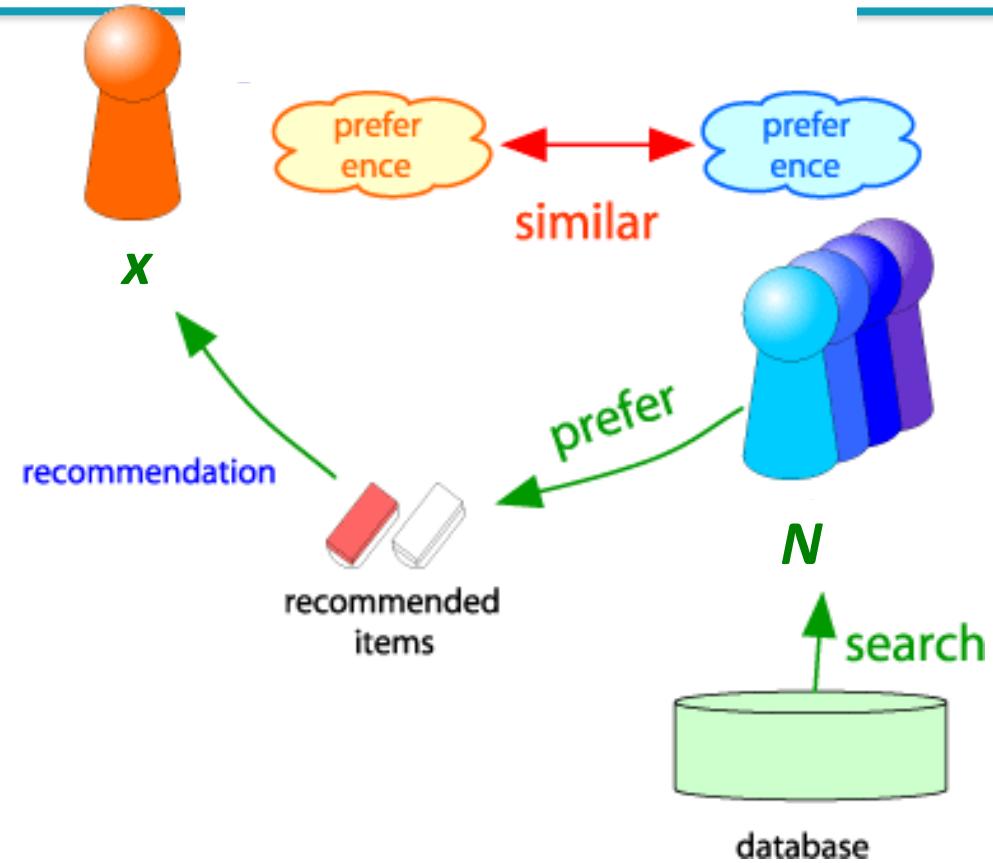
Collaborative Recommendations

Collaborative Recommendations

- User-based recommendation
- Item-based recommendation

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



Finding “Similar” Users

$$\begin{aligned} \mathbf{r}_x &= [* , _, _, *, *, ***] \\ \mathbf{r}_y &= [* , _, **, **, _, _] \end{aligned}$$

- Let \mathbf{r}_x be the vector of user x 's ratings
- Jaccard similarity measure**

- Problem:** Ignores the value of the rating

- Cosine similarity measure**

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

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

- Pearson correlation coefficient**

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

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

$\mathbf{r}_x, \mathbf{r}_y$ as sets:
 $\mathbf{r}_x = \{1, 4, 5\}$
 $\mathbf{r}_y = \{1, 3, 4\}$

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

\bar{r}_x, \bar{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: $\text{sim}(A,B) > \text{sim}(A,C)$
- Jaccard: $1/5 < 2/4$
- Cosine: $0.380 > 0.322$
 - Considers missing ratings as “negative”
 - Solution: subtract the (row) mean

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

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

- $\text{sim}(A,B)$ vs $\text{sim}(A,C)$
 - $0.092 > -0.559$

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

Rating Predictions

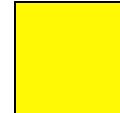
From similarity metric to recommendations:

- 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 i of user x :**
 - $r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi}$
 - $r_{xi} = \frac{\sum_{y \in N} s_{xy} \cdot r_{yi}}{\sum_{y \in N} s_{xy}}$

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

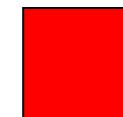
User-User CF ($|N|=2$)

	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

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

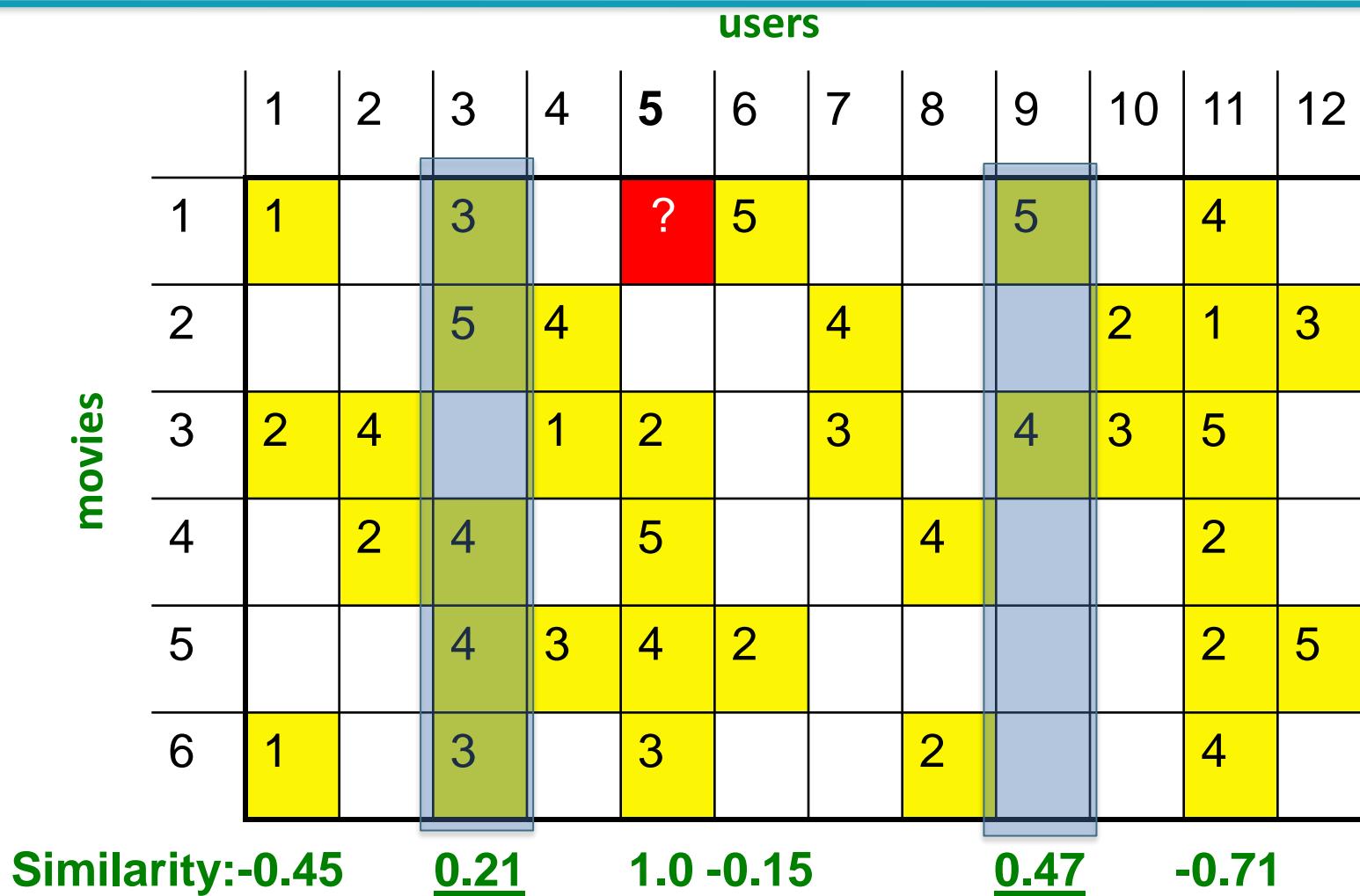
	users											
	1	2	3	4	5	6	7	8	9	10	11	12
1	1		3		?	5			5		4	
2			5	4			4			2	1	3
3	2	4		1	2		3		4	3	5	
4		2	4		5			4			2	
5			4	3	4	2					2	5
6	1		3		3			2			4	



- estimate rating of movie 1 by user 5

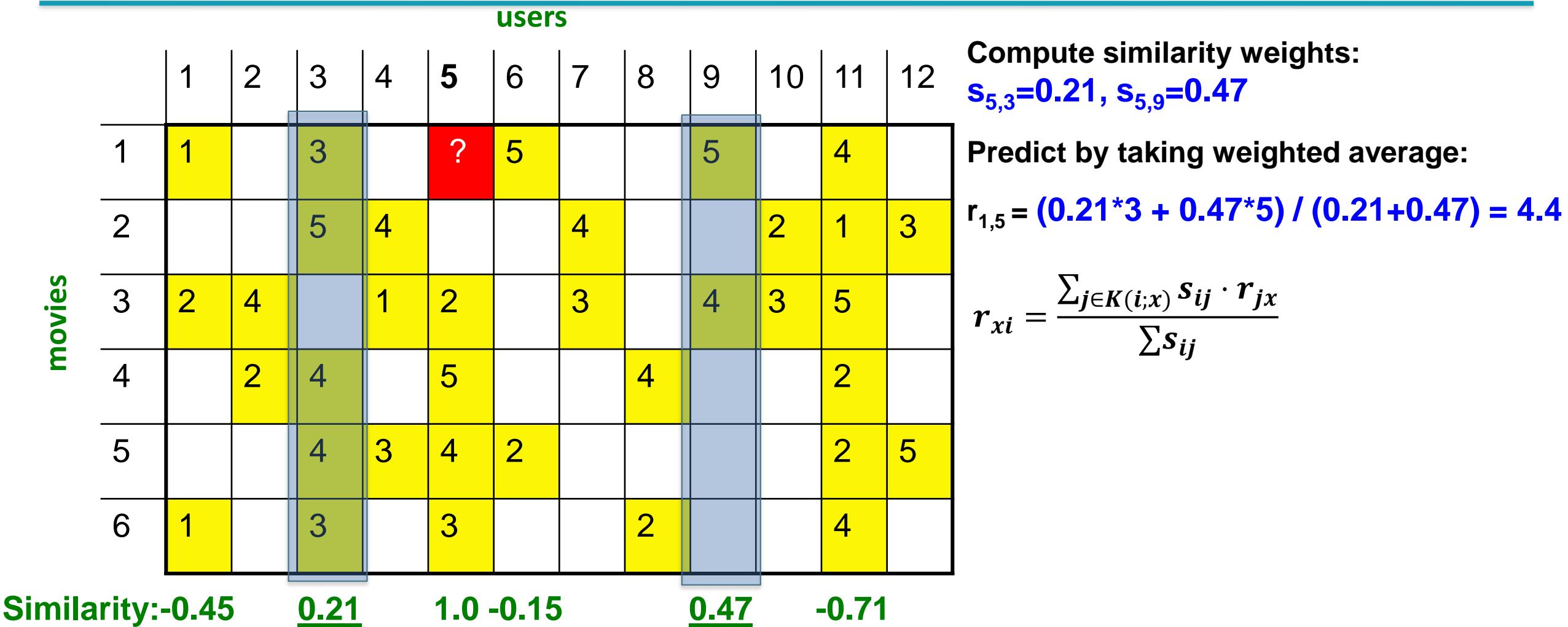
Neighbor selection:
Identify users similar to
user 5, and rated item 1

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



Neighbor selection:
Identify users similar to user 5, and rated item 1

User-based CF ($|K|=2$)



User-based CF ($|K|=2$)

	users											
	1	2	3	4	5	6	7	8	9	10	11	12
1	1		3		4.4	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	

Similarity:-0.45 0.21 1.0 -0.15 0.47 -0.71

Compute similarity weights:
 $s_{5,3}=0.21, s_{5,9}=0.47$

Predict by taking weighted average:

$$r_{1,5} = (0.21 * 3 + 0.47 * 5) / (0.21 + 0.47) = 4.4$$

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

Issue with the user-based CF

- So far: user-based collaborative filtering
- Another view is **item-based CF**.

Item-based CF

- The item-based approach works by comparing items based on their pattern of ratings across users. The similarity of items i and j is computed as follows:

$$sim(i, j) = \frac{\sum_{\mathbf{u} \in U} (r_{\mathbf{u}, i} - \bar{r}_{\mathbf{u}})(r_{\mathbf{u}, j} - \bar{r}_{\mathbf{u}})}{\sqrt{\sum_{\mathbf{u} \in U} (r_{\mathbf{u}, i} - \bar{r}_{\mathbf{u}})^2} \sqrt{\sum_{\mathbf{u} \in U} (r_{\mathbf{u}, j} - \bar{r}_{\mathbf{u}})^2}}$$

Recommendation phase

- After computing the similarity between items we select a set of k most similar items to the target item and generate a predicted value of user x 's rating

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

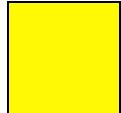
s_{ij} ... similarity of items i and j

r_{xj} ...rating of user x on item j

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

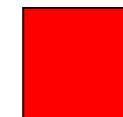
Item-Item CF ($|N|=2$)

	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

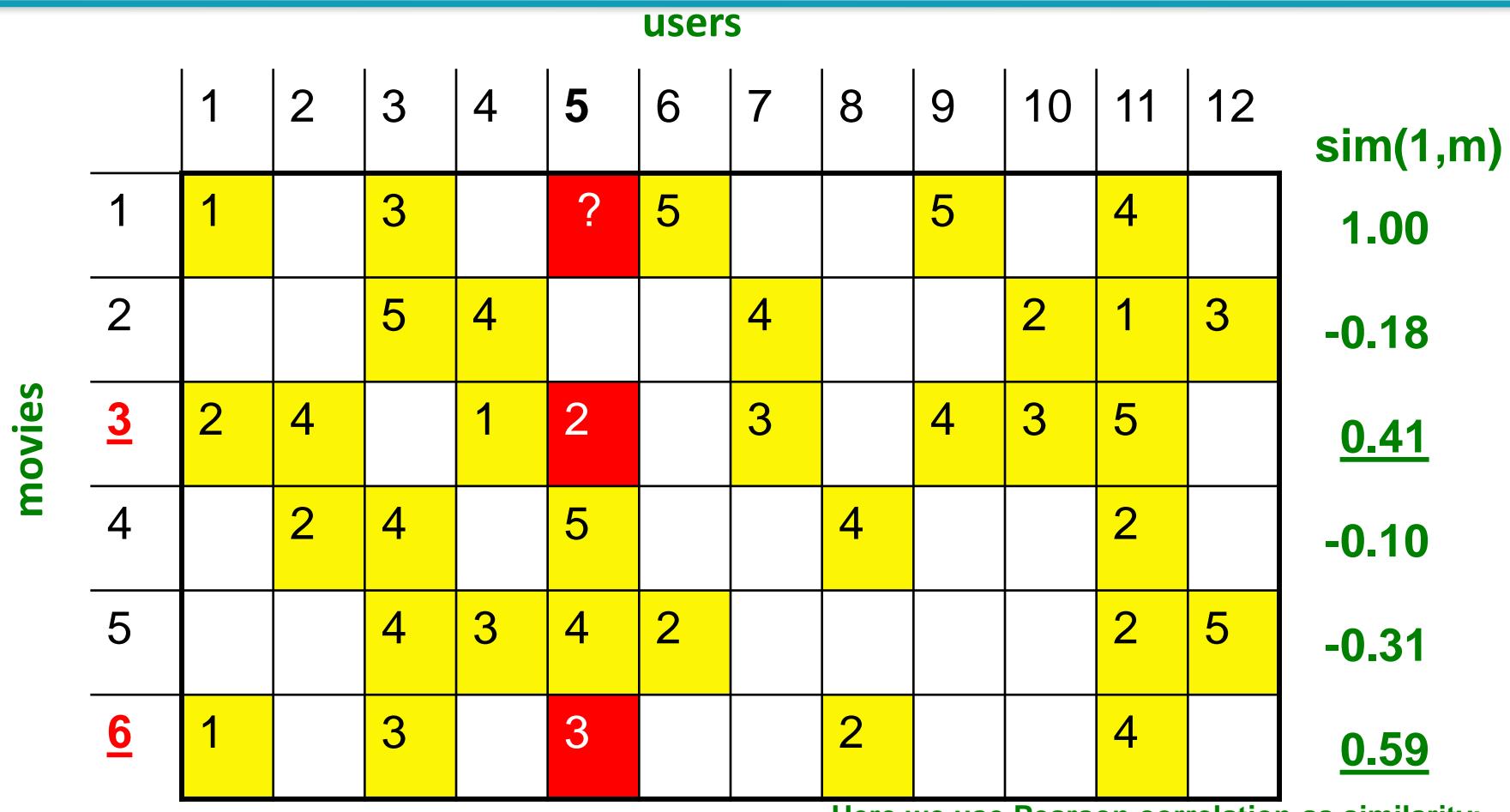
Item-Item CF ($|N|=2$)

	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 ($|N|=2$)



Neighbor selection:

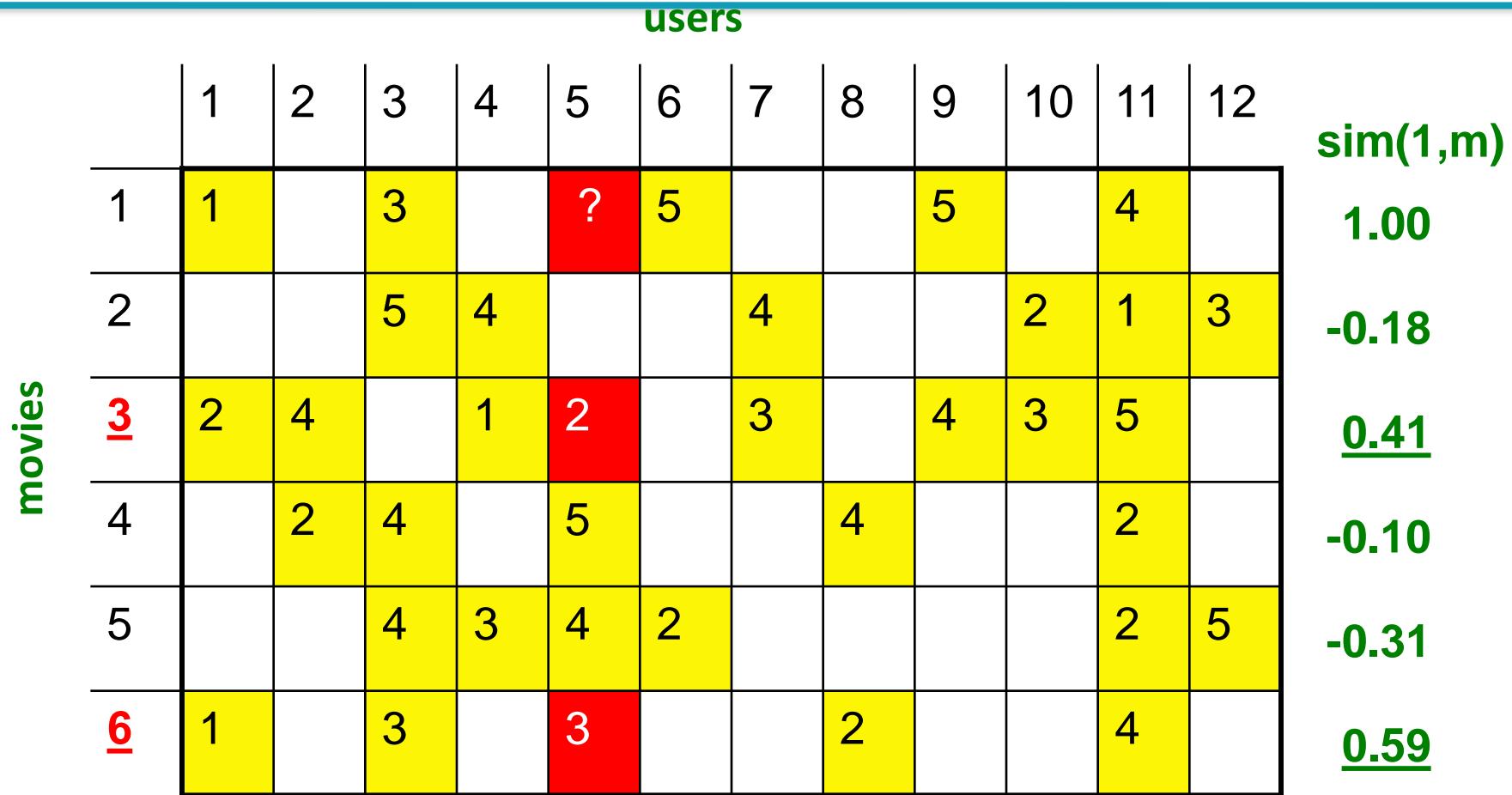
Identify movies similar to
movie 1, rated by user 5

Here we use Pearson correlation as similarity:

- 1) Subtract mean rating m_i from each movie i

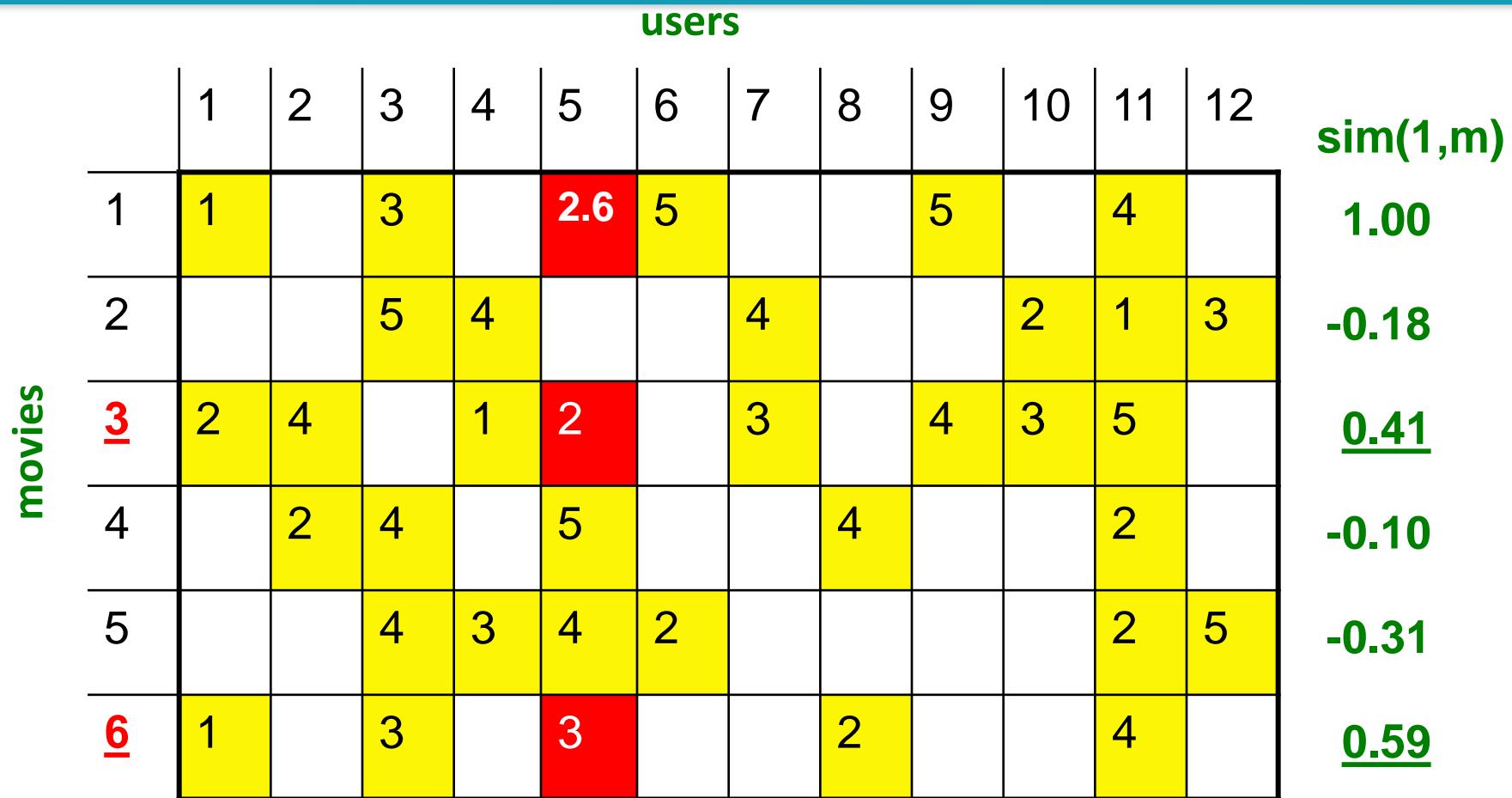
$$m_1 = (1+3+5+5+4)/5 = 3.6$$
row 1: [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]
- 2) Compute cosine similarities between rows

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



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

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



Predict by taking weighted average:

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

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

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

Hybrid Methods

Hybrid Methods

- Add content-based methods to collaborative filtering
 - item profiles for new item problem
 - demographics to deal with new user problem
- Implement two separate recommenders and combine predictions
 - Perhaps using a linear model
 - E.g., global baseline + collaborative filtering

Global baseline estimate

- Estimate Joe's rating for the movie **The Sixth Sense**
 - No feature selection needed
 - Problem: Joe has not rated any movie similar to **The Sixth Sense**
- **Global baseline estimate**
 - Mean movie rating: **3.7 stars**
 - **The Six Sense** is **0.5 stars** above avg
 - Joe rates **0.2 stars** below avg
 - Baseline estimate: **$3.7 + 0.5 - 0.2 = 4$ stars**

Combining Global Baseline with CF

- **Global Baseline estimate:**
 - Joe will give *The Sixth Sense* 4 stars
- **Local neighborhood (CF/NN):**
 - Joe didn't like related movie *Signs*
 - Rated it 1 star below his average rating
- **Final estimate**
 - Joe will rate *The Sixth Sense* $4 - 1 = 3$ stars

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 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}

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

- 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

Evaluation

Evaluation

movies

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

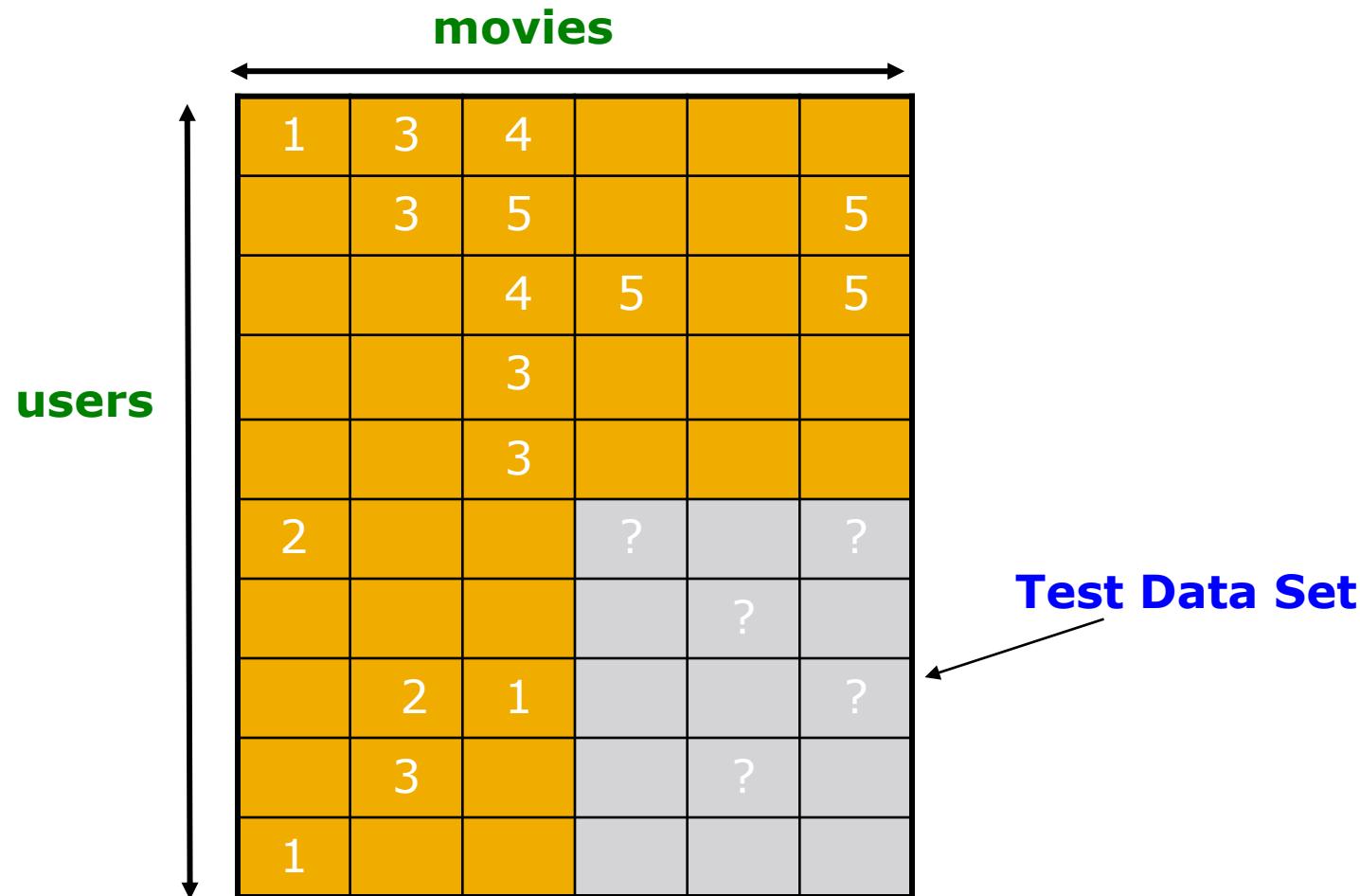
users



Evaluation

		movies					
		1	3	4			
			3	5			5
				4	5		5
				3			
				3			
users		2			?		?
					?		
			2	1			?
			3			?	
1							

Test Data Set



Evaluation Predictions

- Root-mean-square error (RMSE)

$\sqrt{\frac{1}{|R|} \sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$ where \hat{r}_{xi} is predicted, r_{xi} is the true rating of x on i

- Problems with Error Measures

- Narrowly focus on accuracy sometimes misses the point

- Order of predictions

- In practice, we care only to predict high ratings:

- RMSE might penalize a method that does well for high ratings and badly for others

- Alternative: Precision@ k , Hit Rate@ k (i.e., recall of positive interactions), and NDCG@ k

- E.g., $k=10$

Three ways running Jupyter Notebook

- 1. Install Jupyter Notebook
 - <https://jupyter.org/install>
- 2. install Anaconda (if you installed it before, Jupyter Notebook was already installed together)
 - <https://www.anaconda.com/products/individual>
 - Beyond all of the normal (non-data centric) packages that Python comes with, Anacoda comes with even more!
- 3. Use a cloud computing (e.g., Google Colab)
 - <https://colab.research.google.com/notebooks/intro.ipynb#recent=true>

Google Colab

- **Google Colaboratory** is a free online cloud-based **Jupyter notebook** environment that allows us to train our machine learning and deep learning models on CPUs, GPUs, and TPUs.

<https://www.youtube.com/watch?v=inN8seMm7UI>



Recommendation Demo

- recommenderDemo.ipynb

Latent Factor Models

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) =

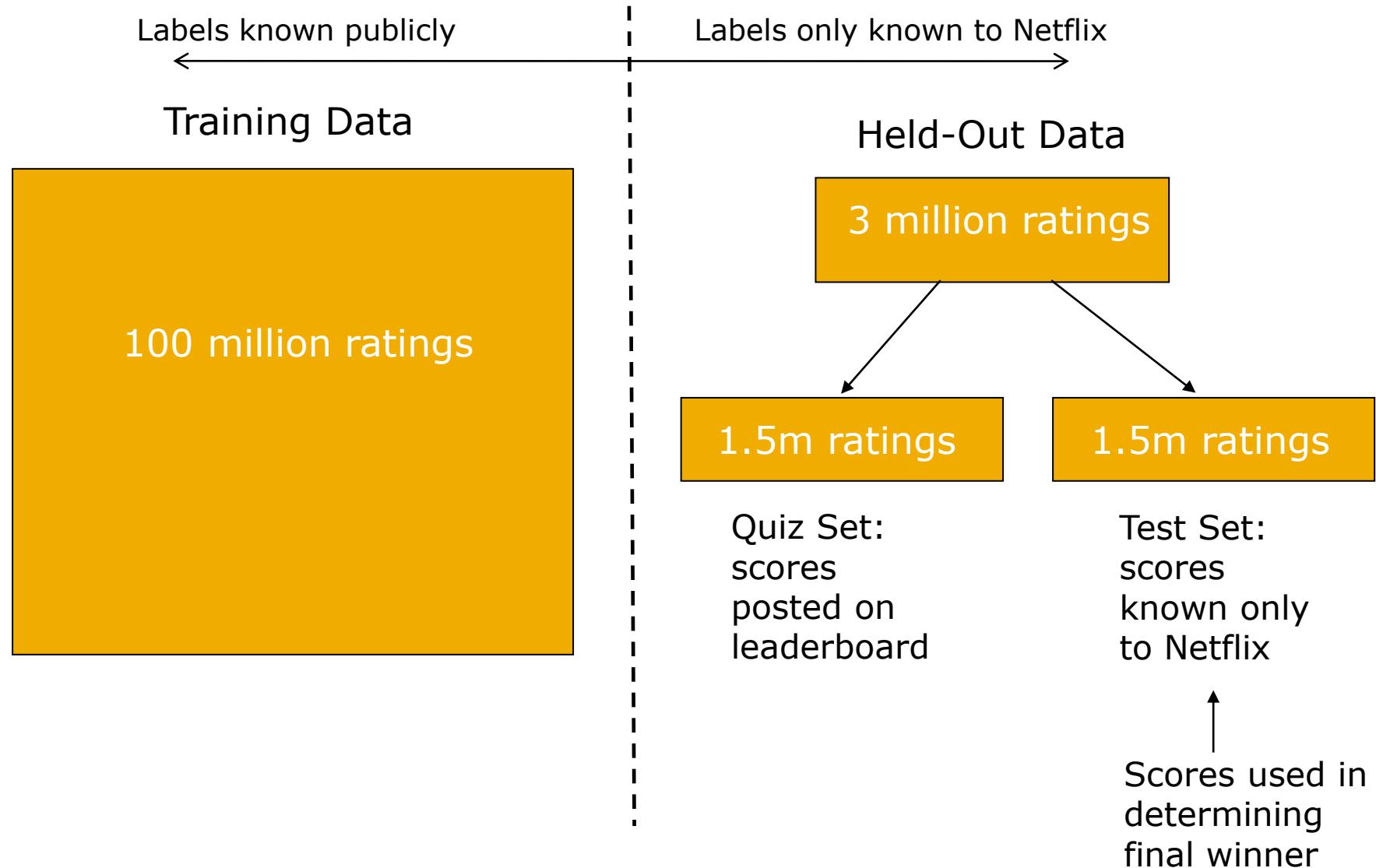
$$\sqrt{\frac{1}{|R|} \sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$$

- Netflix's system RMSE: 0.9514

■ Competition

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

Competition Structure



The Netflix Utility Matrix R

Matrix R

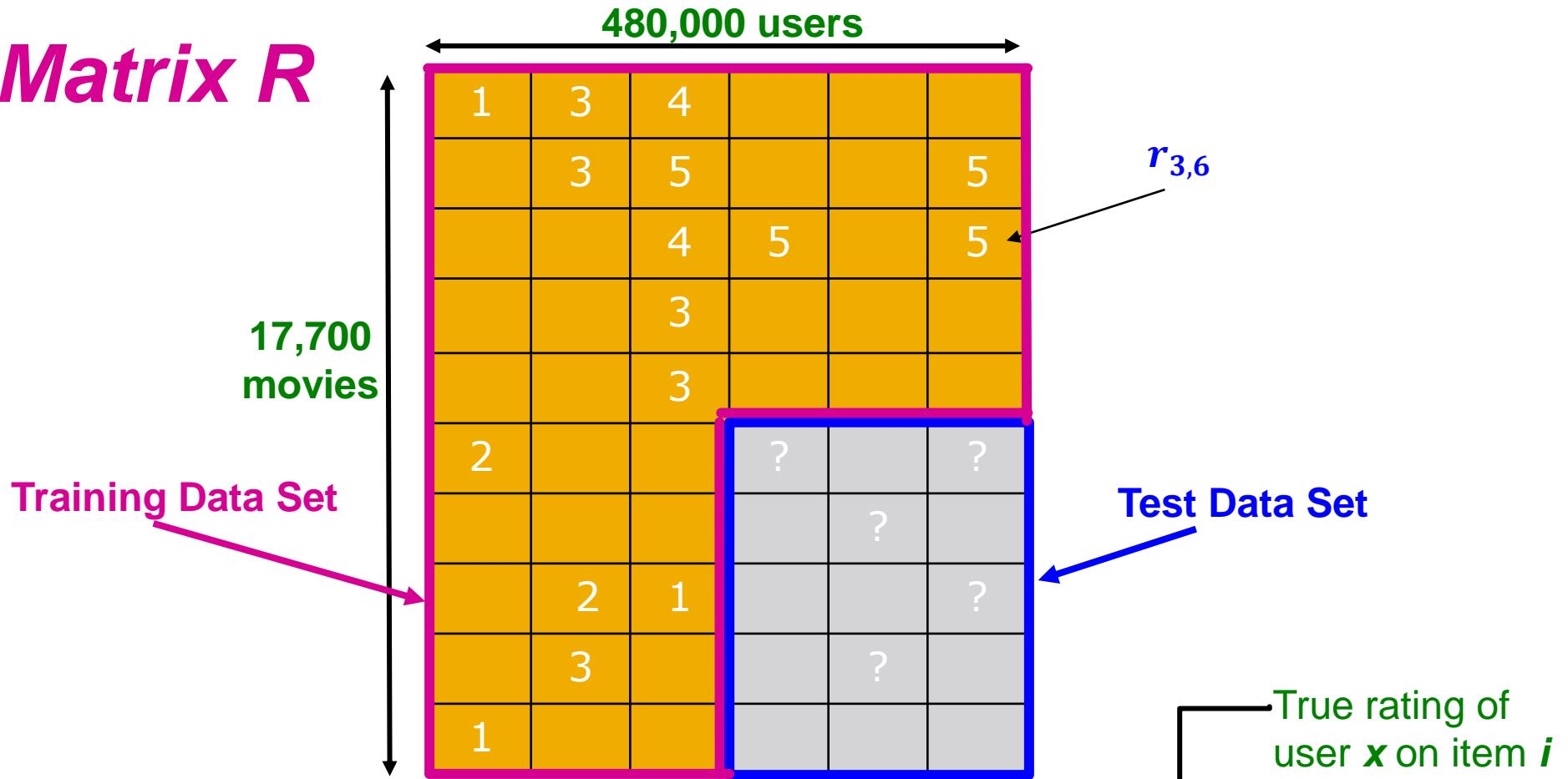
480,000 users

17,700 movies

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

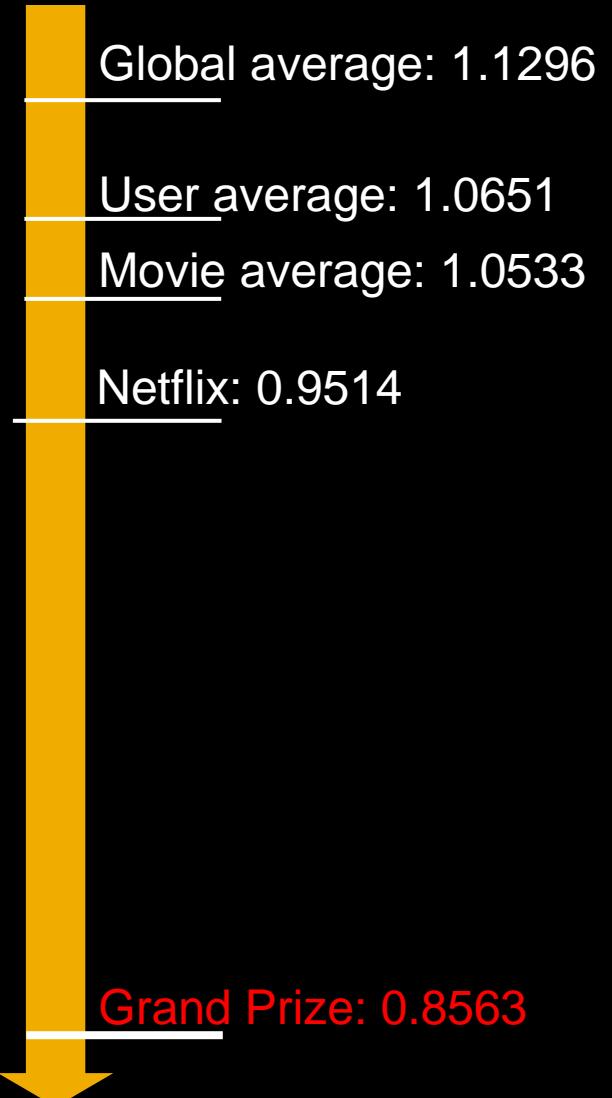
Utility Matrix R : Evaluation

Matrix R



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

Performance of Various Methods



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:

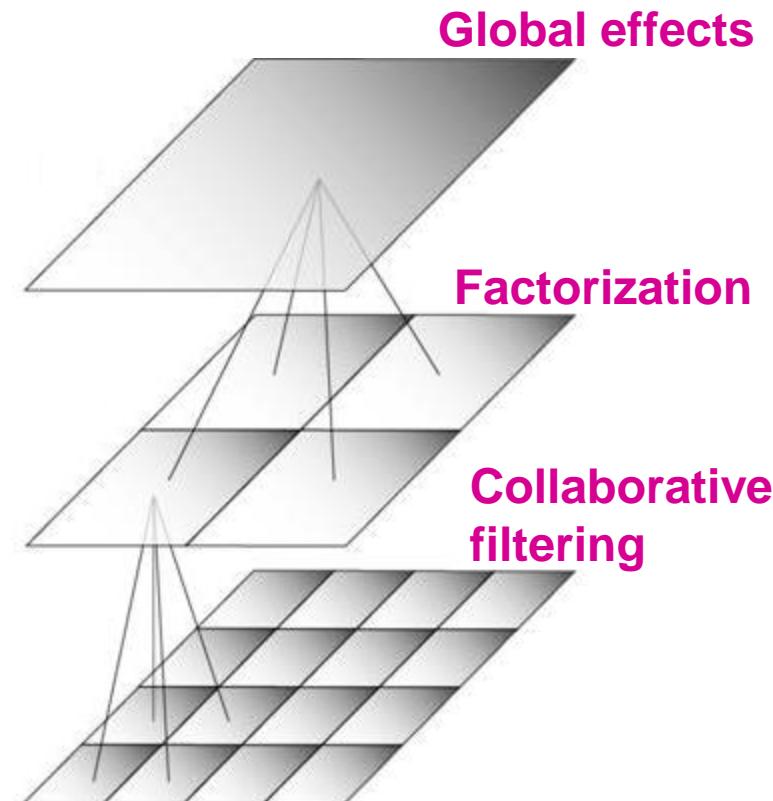
- Overall deviations of users/movies

- Factorization:

- Addressing “regional” effects

- Collaborative filtering:

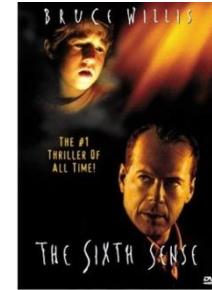
- Extract local patterns



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*
- Rated it **1 star** below his average rating



■ Final estimate

- Joe will rate ***The Sixth Sense* 4 – 1 = 3 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:

- Similarity measures are “arbitrary”
- Pairwise similarities neglect interdependencies among users
- 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)
 - Note, 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}

- $\hat{r}_{xi} = b_{xi} + \sum_{j \in N(i,x)} w_{ij}(r_{xj} - b_{xj})$
- How to set w_{ij} ?
 - Remember, error metric is:
$$\sqrt{\frac{1}{|R|} \sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2}$$
 or equivalently Sum of Squared Error (SSE): $\sum_{(i,x) \in R} (\hat{r}_{xi} - r_{xi})^2$
 - 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**

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

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!

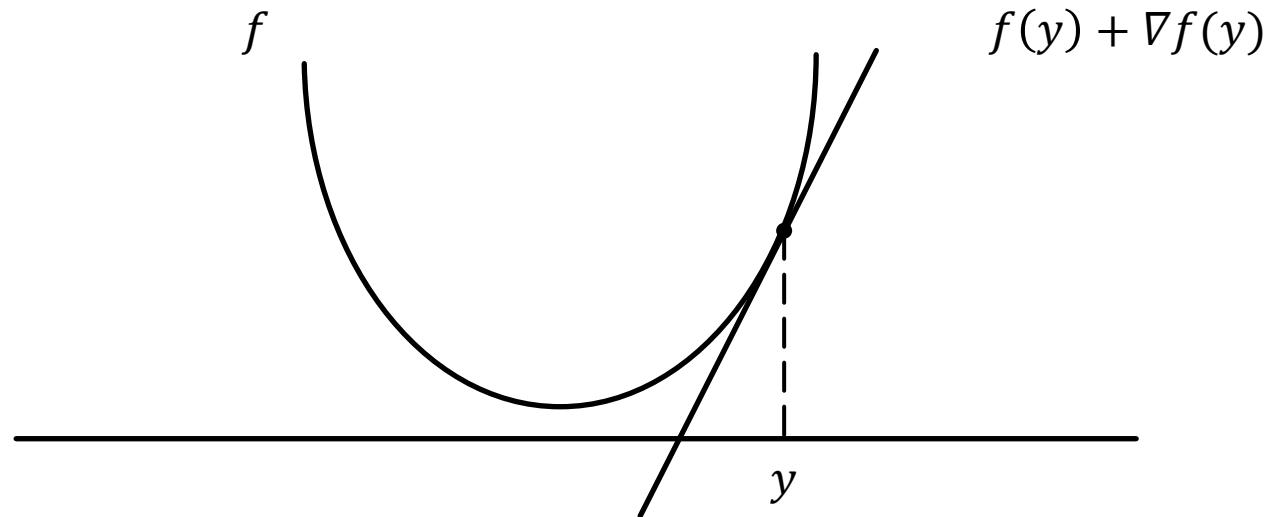
$$J(w) = \sum_{x,i \in R} \left(\underbrace{\left[b_{xi} + \sum_{j \in N(i;x)} w_{ij}(r_{xj} - b_{xj}) \right]}_{\text{Predicted rating}} - r_{xi} \right)^2$$

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 derivative $\nabla f(x)$
 - 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



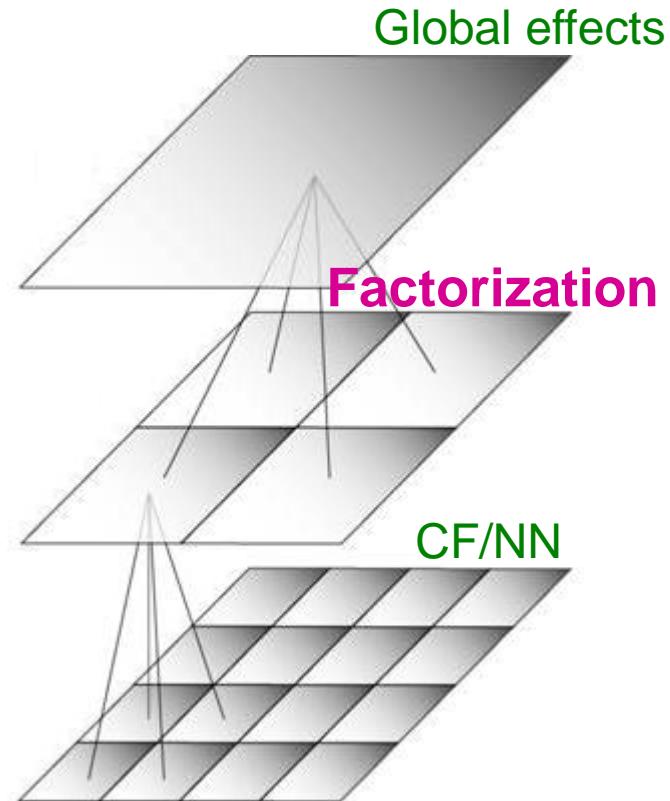
Interpolation Weights

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

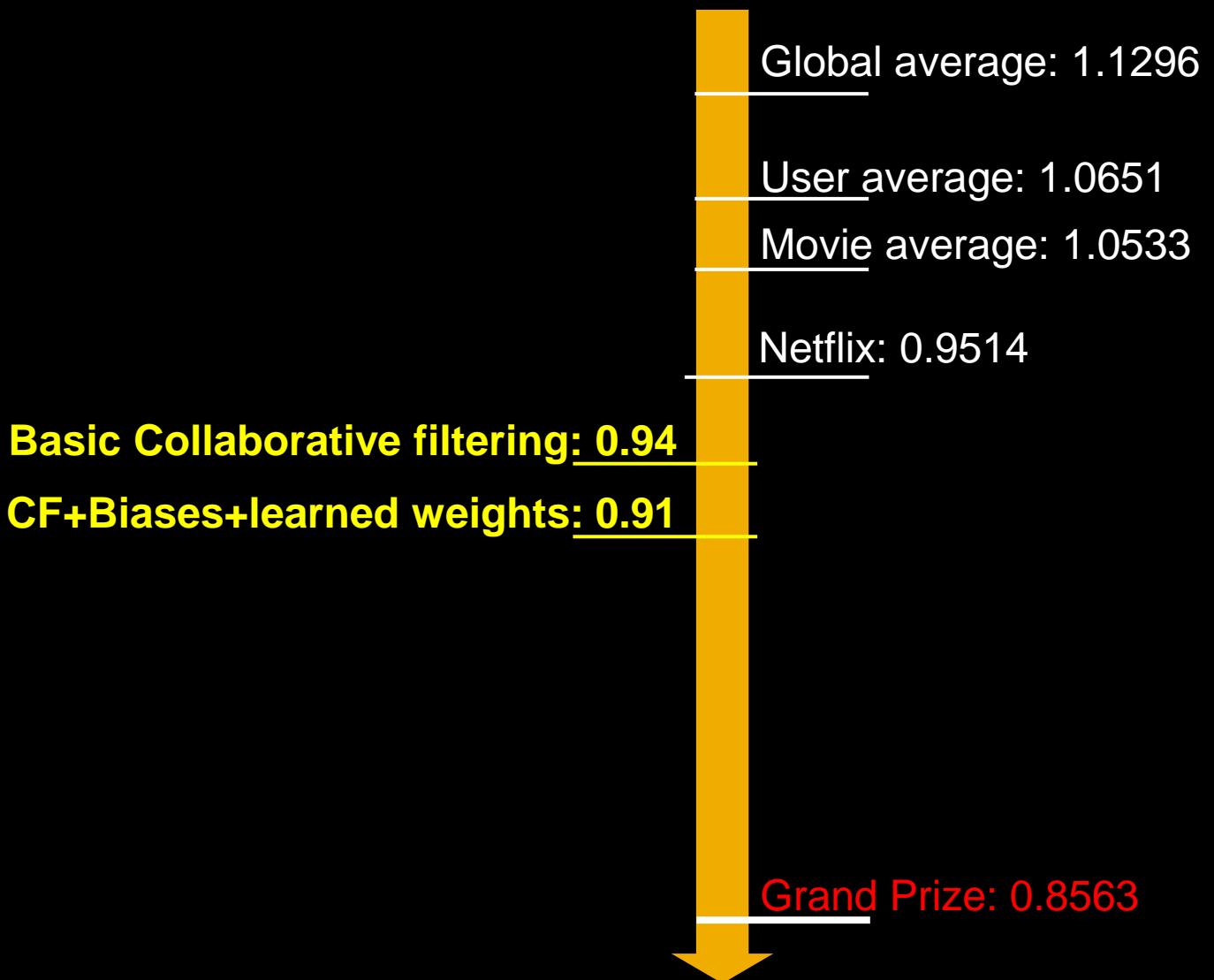
- Weights w_{ij} derived based on their role; **no use of an arbitrary similarity measure** ($w_{ij} \neq s_{ij}$)
- Explicitly account for interrelationships among the neighboring movies

- **Next: Latent factor model**

- Extract “regional” correlations



Performance of Various Methods



Latent Factor Models

- Latent Factor Model on Netflix data: $\mathbf{R} \approx \mathbf{Q} \cdot \mathbf{P}^T$

$$\begin{matrix} & \text{users} \\ \text{items} & \begin{matrix} 1 & 3 & 5 & 5 & 4 & 2 & 1 & 3 \\ 5 & 4 & 1 & 2 & 3 & 4 & 3 & 5 \\ 2 & 4 & 1 & 2 & 3 & 4 & 3 & 5 \\ 2 & 4 & 5 & 4 & 4 & 2 & 2 & 5 \\ 4 & 3 & 4 & 2 & 2 & 4 & 2 & 5 \\ 1 & 3 & 3 & 2 & 2 & 4 \end{matrix} \end{matrix} \approx \begin{matrix} & \text{factors} \\ \text{items} & \begin{matrix} .1 & -.4 & .2 \\ -.5 & .6 & .5 \\ -.2 & .3 & .5 \\ 1.1 & 2.1 & .3 \\ -.7 & 2.1 & -2 \\ -1 & .7 & .3 \end{matrix} \end{matrix}$$

\mathbf{R} \mathbf{Q}

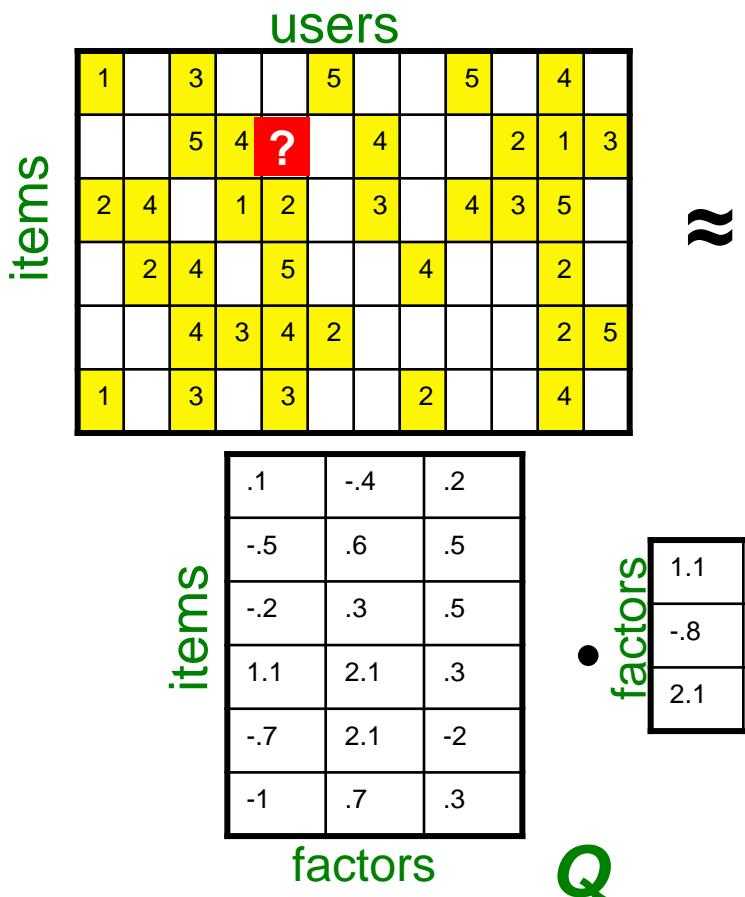
users factors

\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 ?



$$\begin{aligned}\hat{r}_{xi} &= q_i \cdot p_x \\ &= \sum_f q_{if} \cdot p_{xf}\end{aligned}$$

q_i = row i of Q

p_x = column x of P^T

Ratings as Products of Factors

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

users

1	3		5		5		4		
		5	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		

items

≈

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

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

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

Q

users

.1	-.4	.2									
-.5	.6	.5									
-.2	.3	.5									
1.1	2.1	.3									
-.7	2.1	-2									
-1	.7	.3									

items

• factors

P^T

Ratings as Products of Factors

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

users

1	3	5	5	5	4
	5	4	2.4	4	
2	4	1	2	3	4
	2	4	5	4	2
	4	3	4	2	
1	3	3		2	4

items

\approx

Q

$$\hat{r}_{xi} = q_i \cdot p_x$$
$$= \sum_f q_{if} \cdot p_{xf}$$

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

users

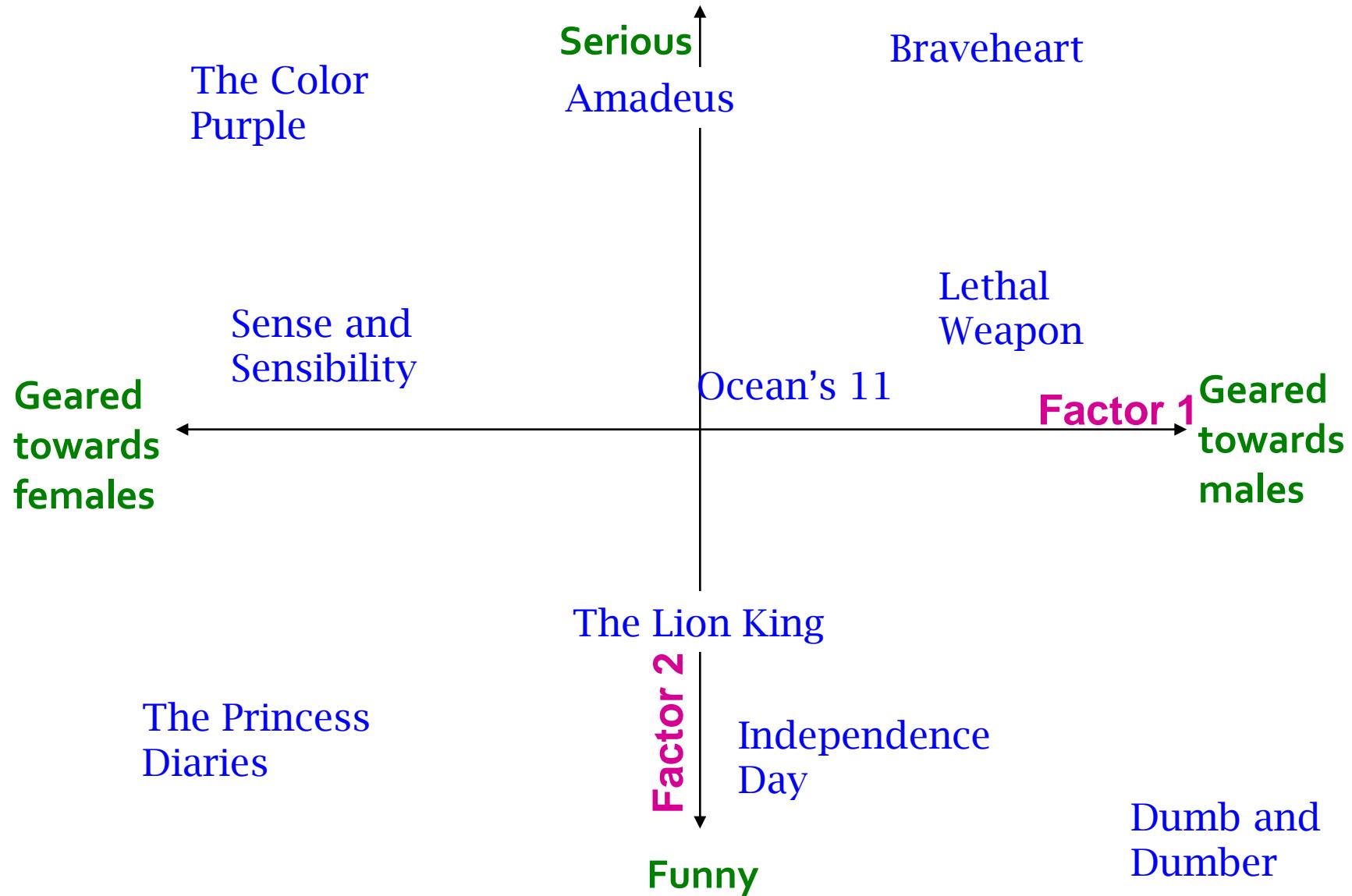
P^T

.1	-.4	.2									
-.5	.6	.5									
-.2	.3	.5									
1.1	2.1	.3									
-.7	2.1	-2									
-1	.7	.3									

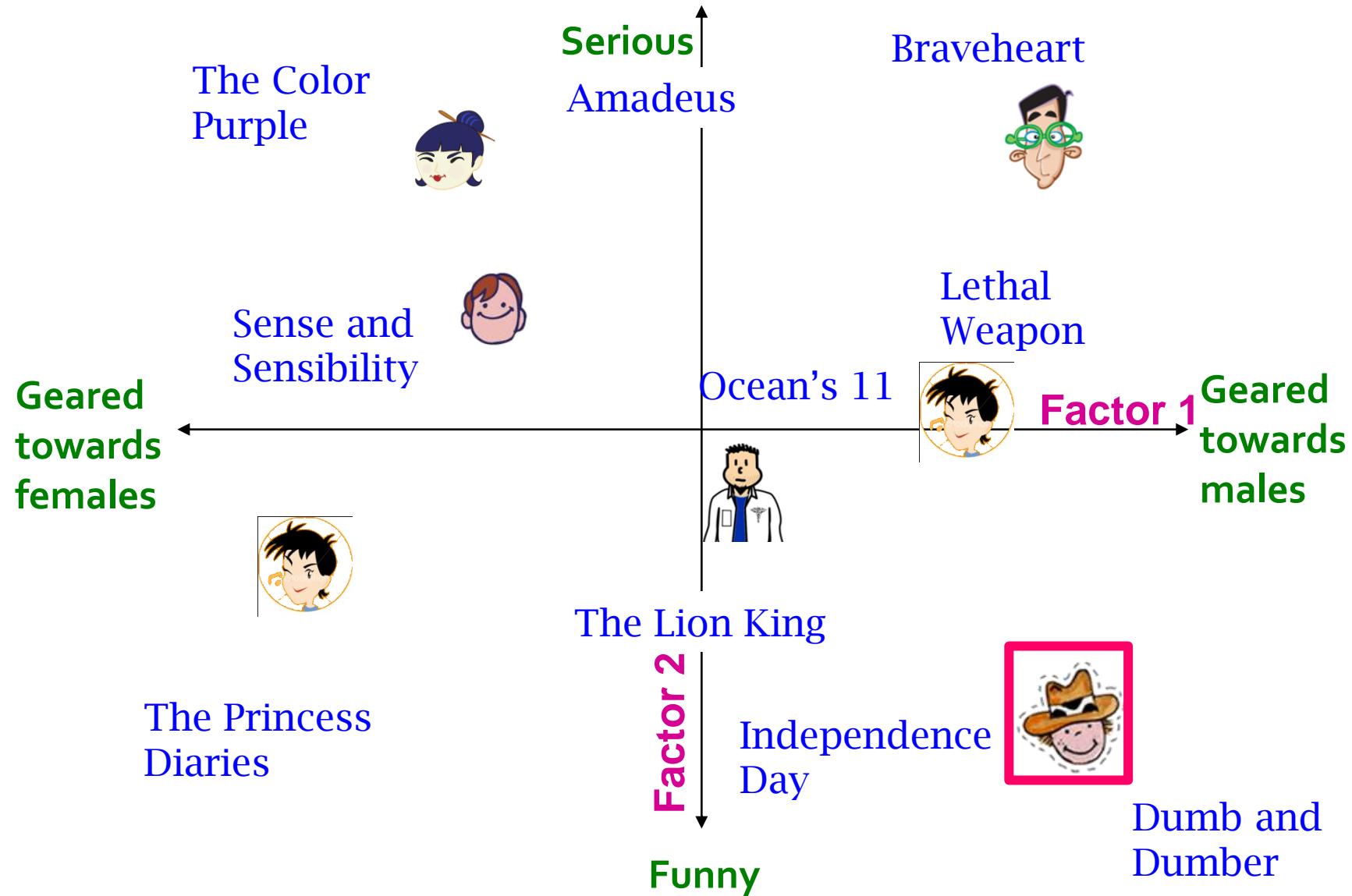
items

f factors

Latent Factor Models



Latent Factor Models

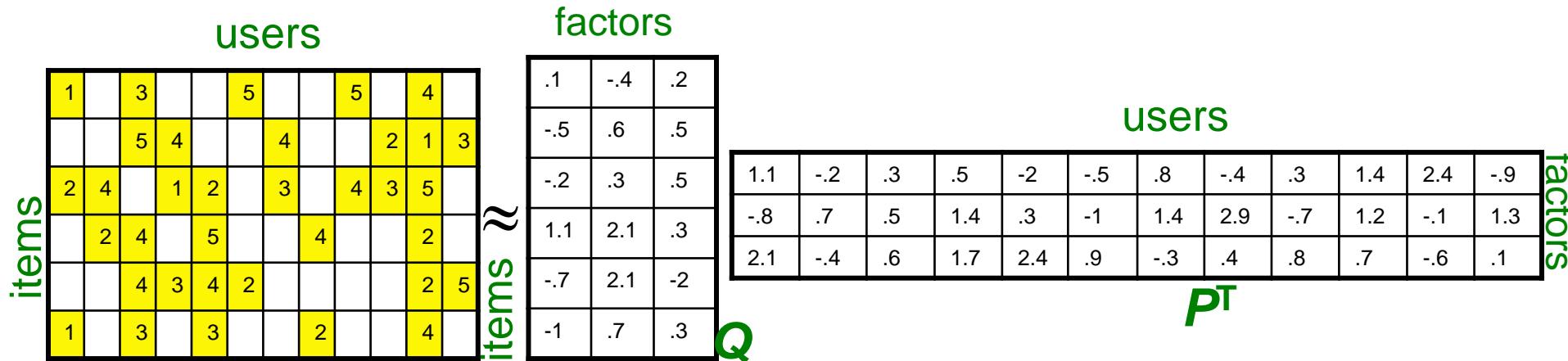


Finding the Latent Factors

Latent Factor Models

- Our goal is to find P and Q such that:

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



Back to Our Problem

- Want to minimize sum of the squared errors (SSE) for unseen test data
- Idea: Minimize SSE on training data
 - Want large k (# of factors) to capture all the signals
 - But, SSE on test data begins to rise for $k > 2$
- This is a classical example of **overfitting**:
 - With too much freedom (too many free parameters) the model starts fitting noise
 - That is it fits too well the training data and thus **not generalizing** well to unseen test data

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

Dealing with Missing Entries

- To solve overfitting we introduce **regularization**:

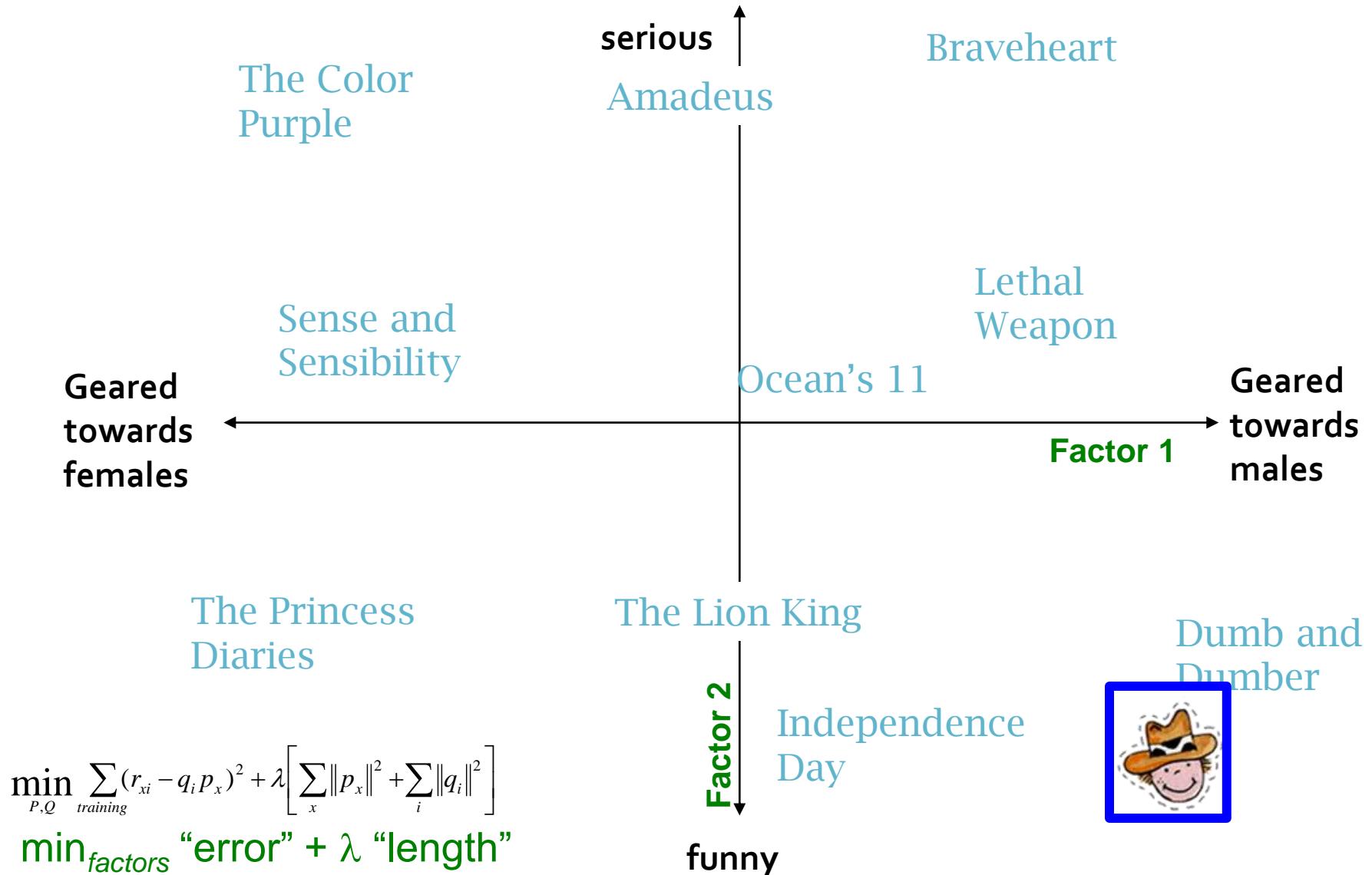
- Allow rich model where there is sufficient data
- Shrink aggressively where data is scarce

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

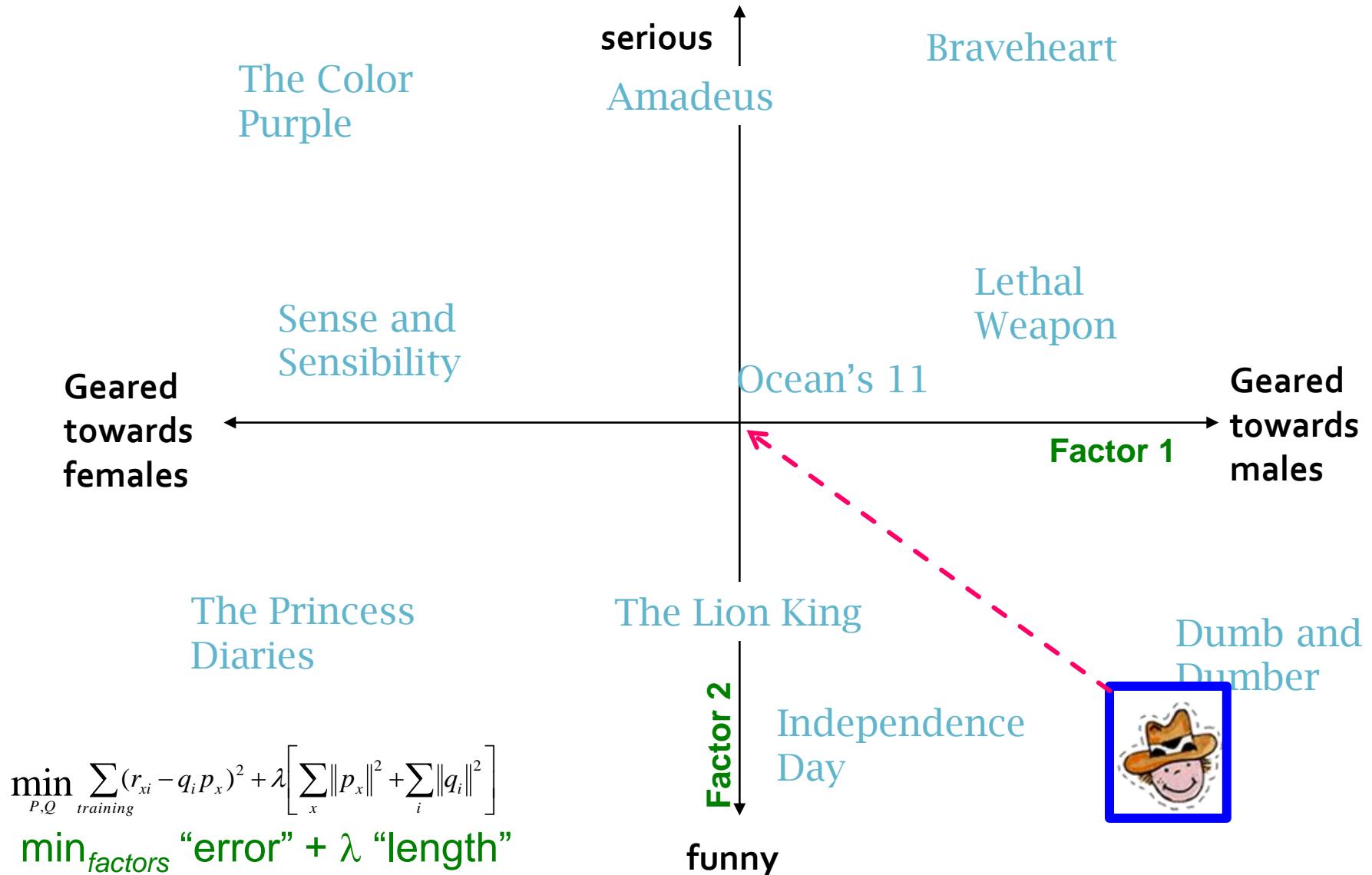
$$\min_{P,Q} \underbrace{\sum_{\text{training}} (r_{xi} - q_i p_x)^2}_{\text{"error"}} + \left[\underbrace{\lambda_1 \sum_x \|p_x\|^2}_{\text{"length"}, \text{ shrinkage}} + \underbrace{\lambda_2 \sum_i \|q_i\|^2}_{\text{"length"}, \text{ shrinkage}} \right]$$

$\lambda_1, \lambda_2 \dots$ user set regularization parameters

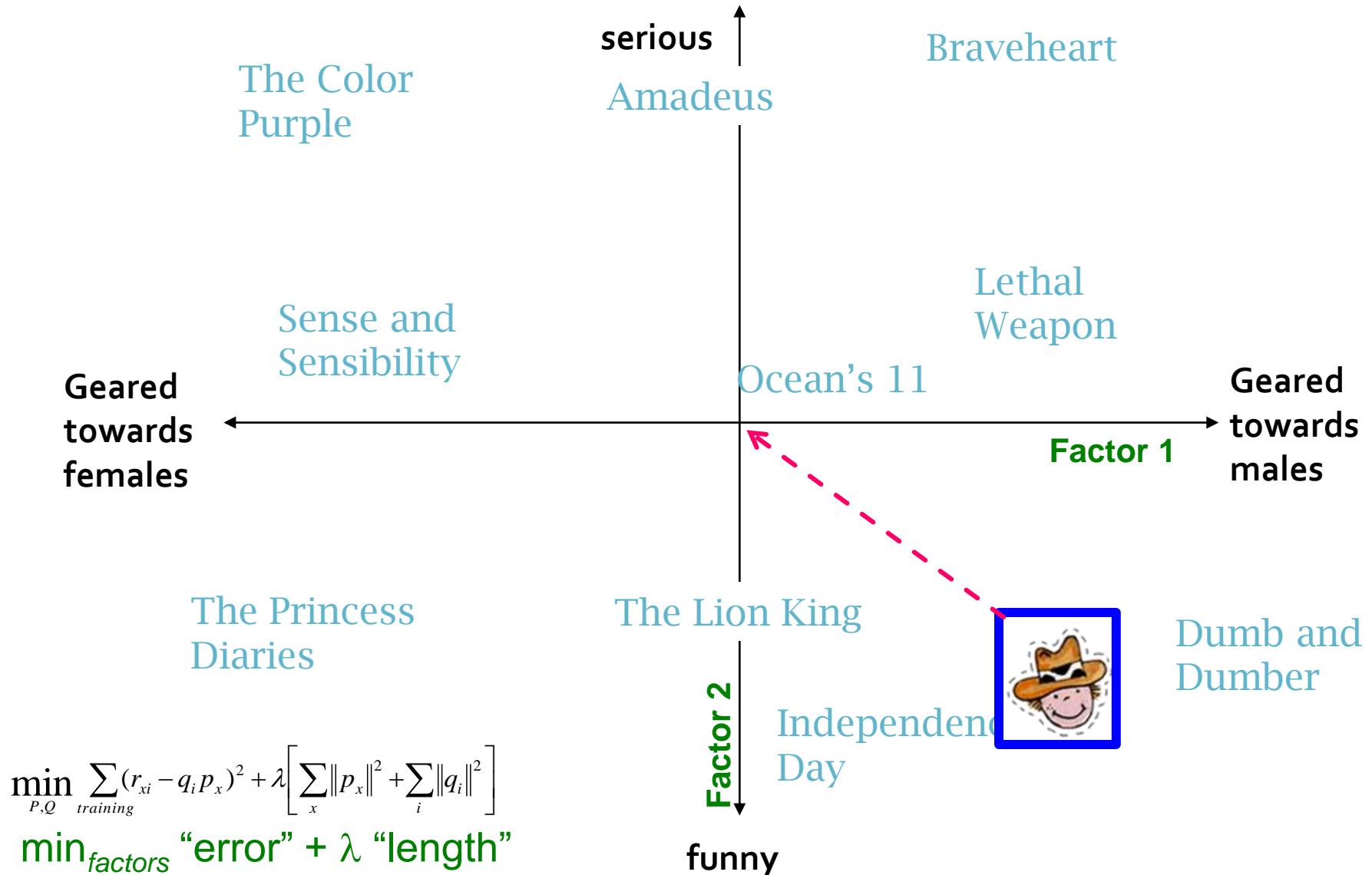
The Effect of Regularization



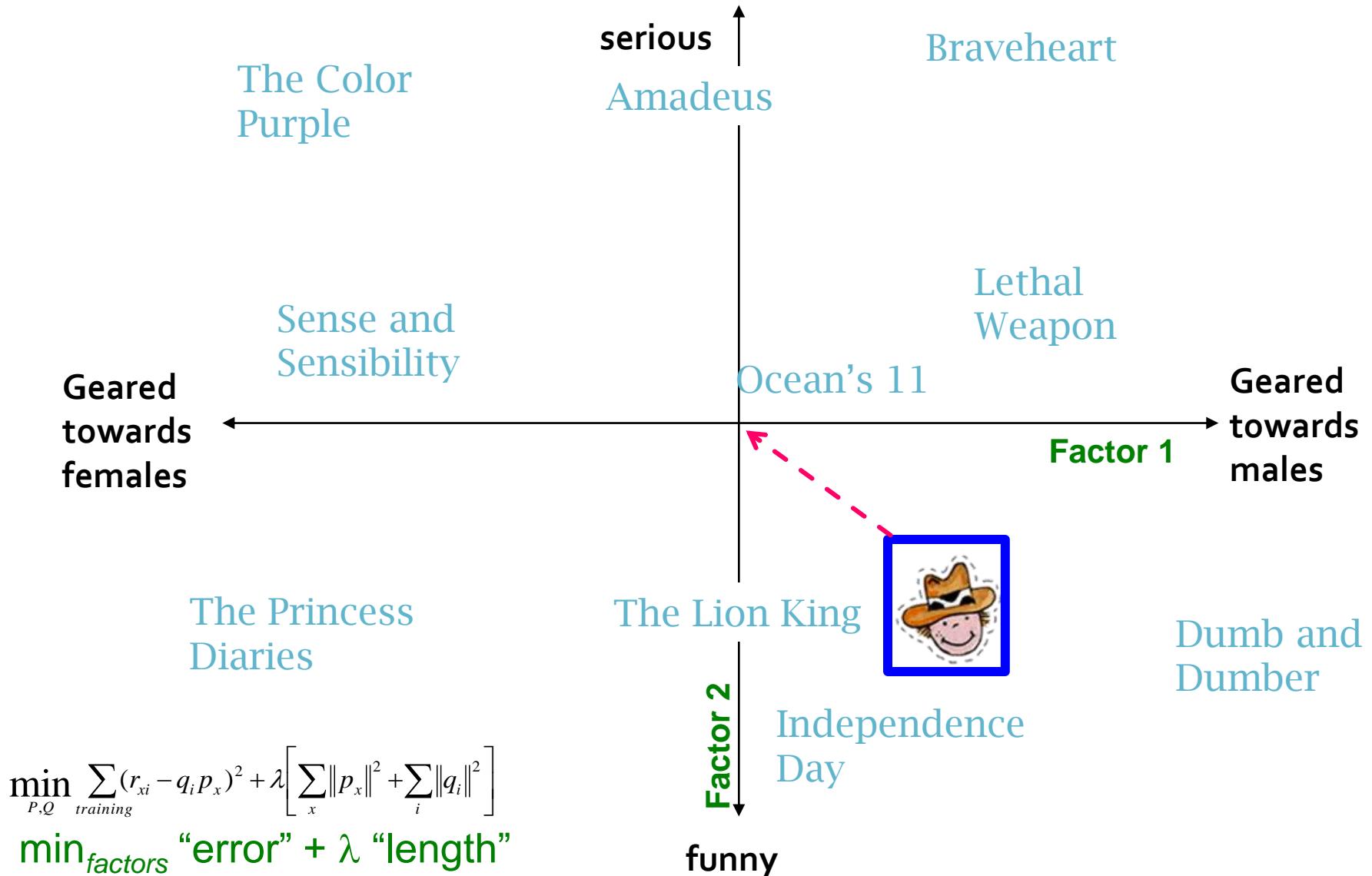
The Effect of Regularization



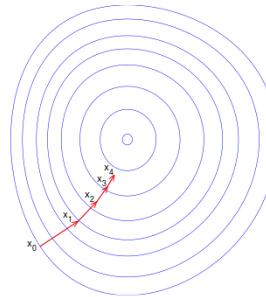
The Effect of Regularization



The Effect of Regularization



Stochastic Gradient Descent



- Want to find matrices P and Q :

$$\min_{P,Q} \sum_{\text{training}} (r_{xi} - q_i p_x)^2 + \left[\lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]$$

- Gradient descent:

- Initialize P and Q (using SVD, pretend missing ratings are 0)

- Do gradient descent:

- $P \leftarrow P - \eta \cdot \nabla P$

- $Q \leftarrow Q - \eta \cdot \nabla Q$

- where ∇Q is gradient/derivative of matrix Q :

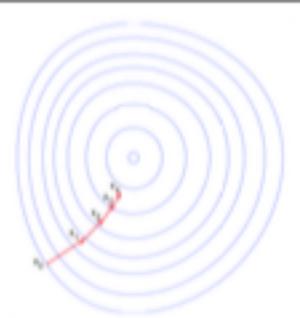
$$\nabla Q = [\nabla q_{if}] \text{ and } \nabla q_{if} = \sum_{x,i} -2(r_{xi} - q_i p_x) p_{xf} + 2\lambda_2 q_{if}$$

- Here q_{if} is entry f of row q_i of matrix Q

How to compute gradient of a matrix?
Compute gradient of every element independently!

- Observation: Computing gradients is slow!

Stochastic Gradient Descent



■ Gradient Descent (GD) vs. Stochastic GD

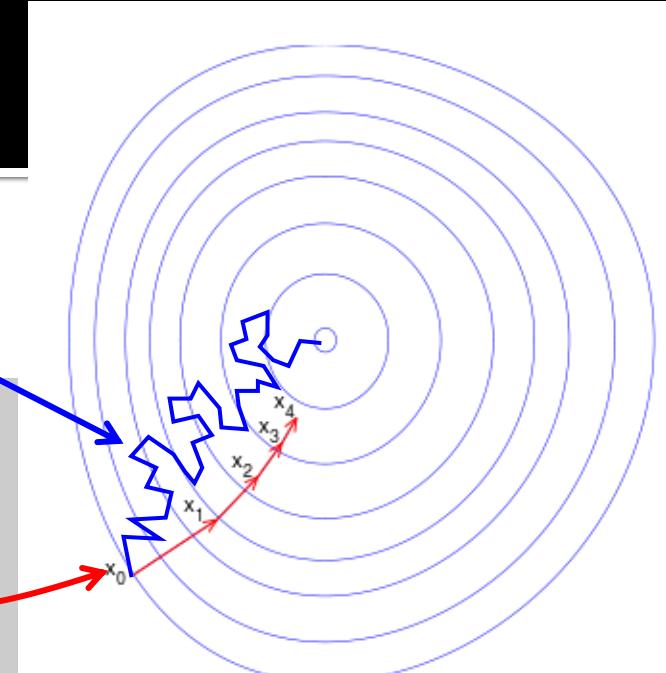
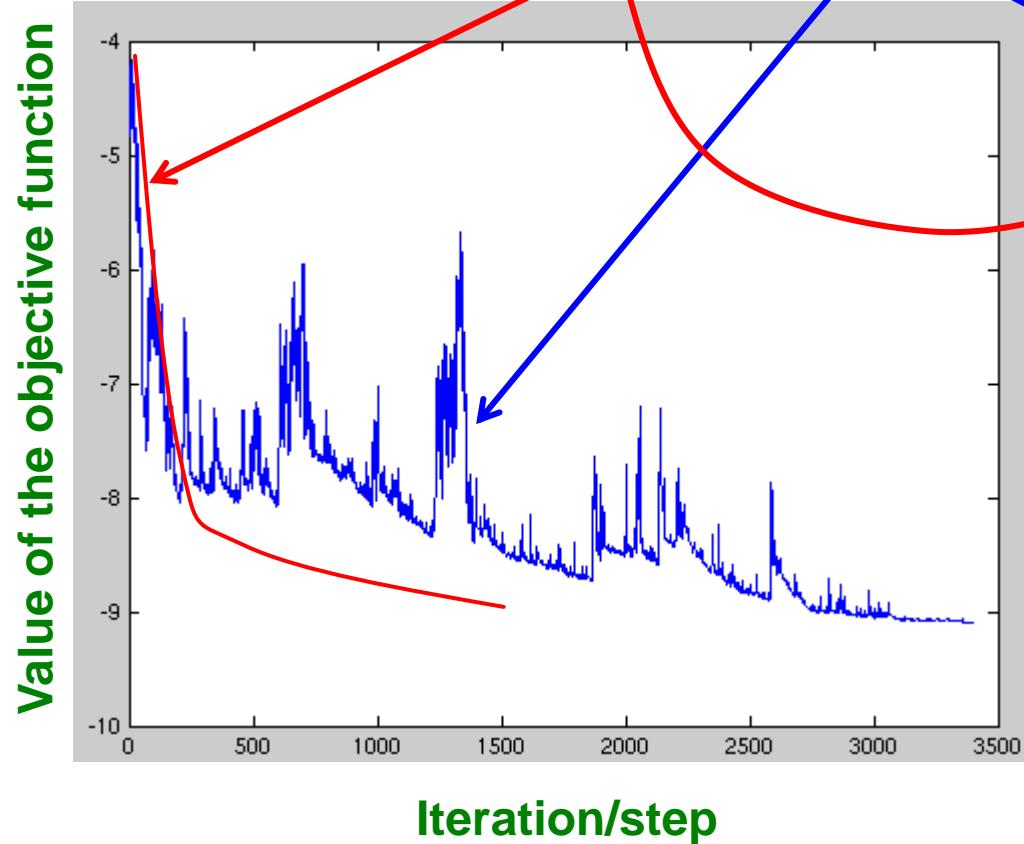
- **Observation:** $\nabla Q = [\nabla q_{if}]$ where

$$\nabla q_{if} = \sum_{x,i} -2(r_{xi} - q_{if} p_{xf}) p_{xf} + 2\lambda q_{if} = \sum_{x,i} \nabla Q(r_{xi})$$

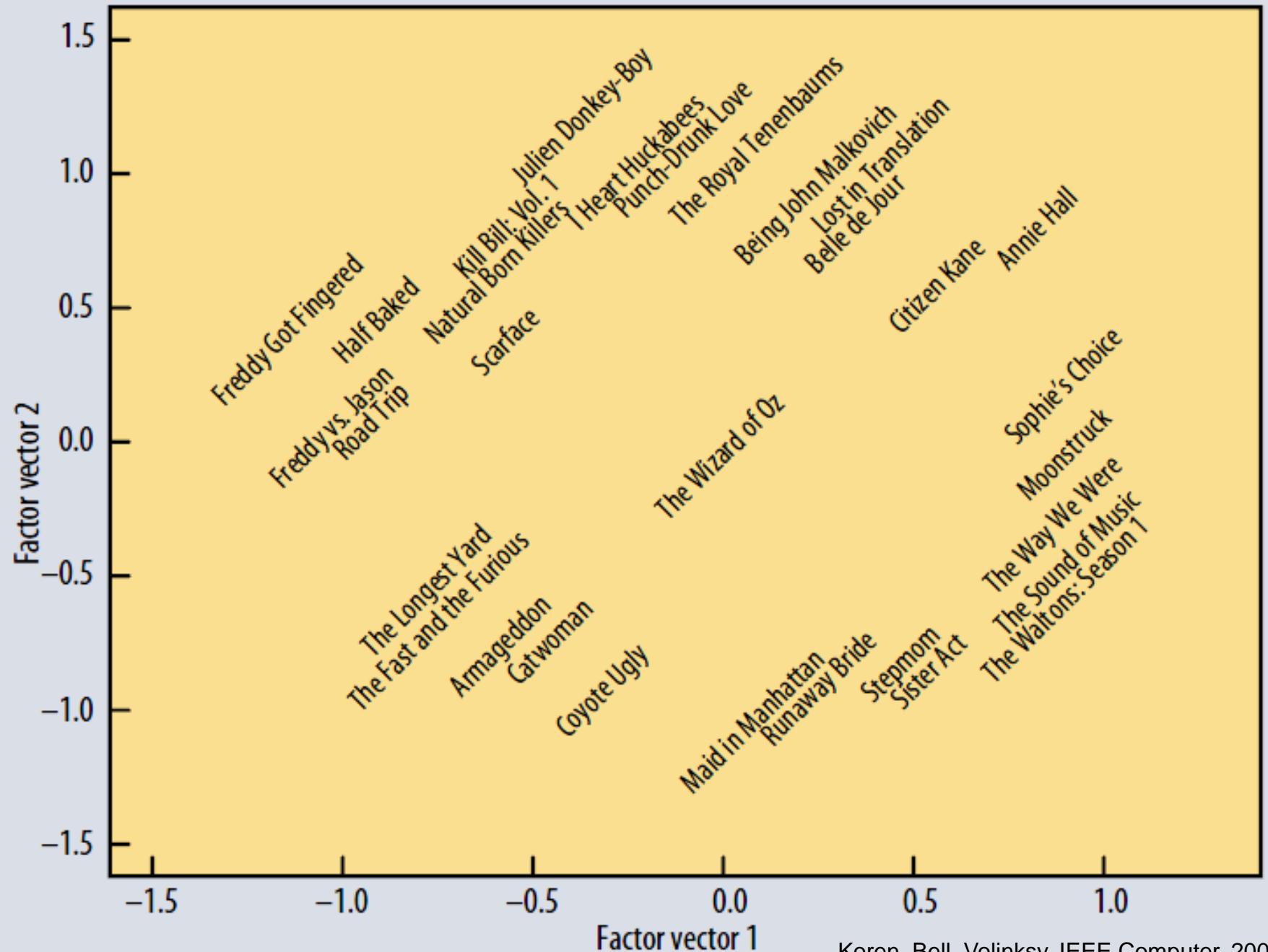
- Here q_{if} is entry f of row q_i of matrix Q
- $Q = Q - \eta \nabla Q = Q - \eta [\sum_{x,i} \nabla Q(r_{xi})]$
- **Idea:** Instead of evaluating gradient over all ratings evaluate it for each individual rating and make a step
- **GD:** $Q \leftarrow Q - \eta [\sum_{r_{xi}} \nabla Q(r_{xi})]$
- **SGD:** $Q \leftarrow Q - \mu \nabla Q(r_{xi})$
- **Faster convergence!**
 - Need more steps but each step is computed much faster

SGD vs. GD

Convergence of **GD** vs. **SGD**

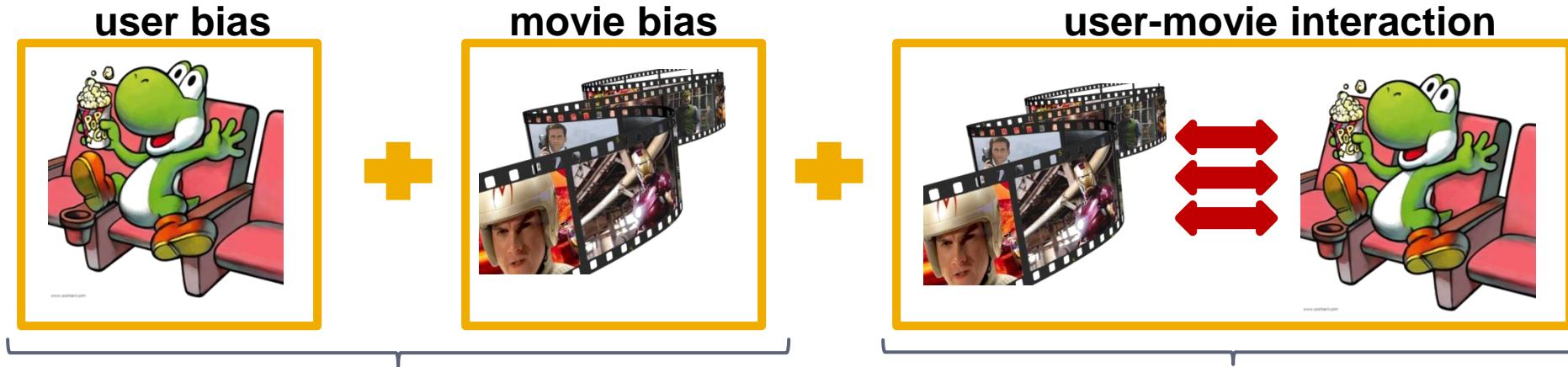


GD improves the value of the objective function at every step.
SGD improves the value but in a “noisy” way.
GD takes fewer steps to converge but each step takes much longer to compute.
In practice, **SGD** is much faster!



Extending Latent Factor Model to Include Biases

Modeling Biases and Interactions



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

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

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

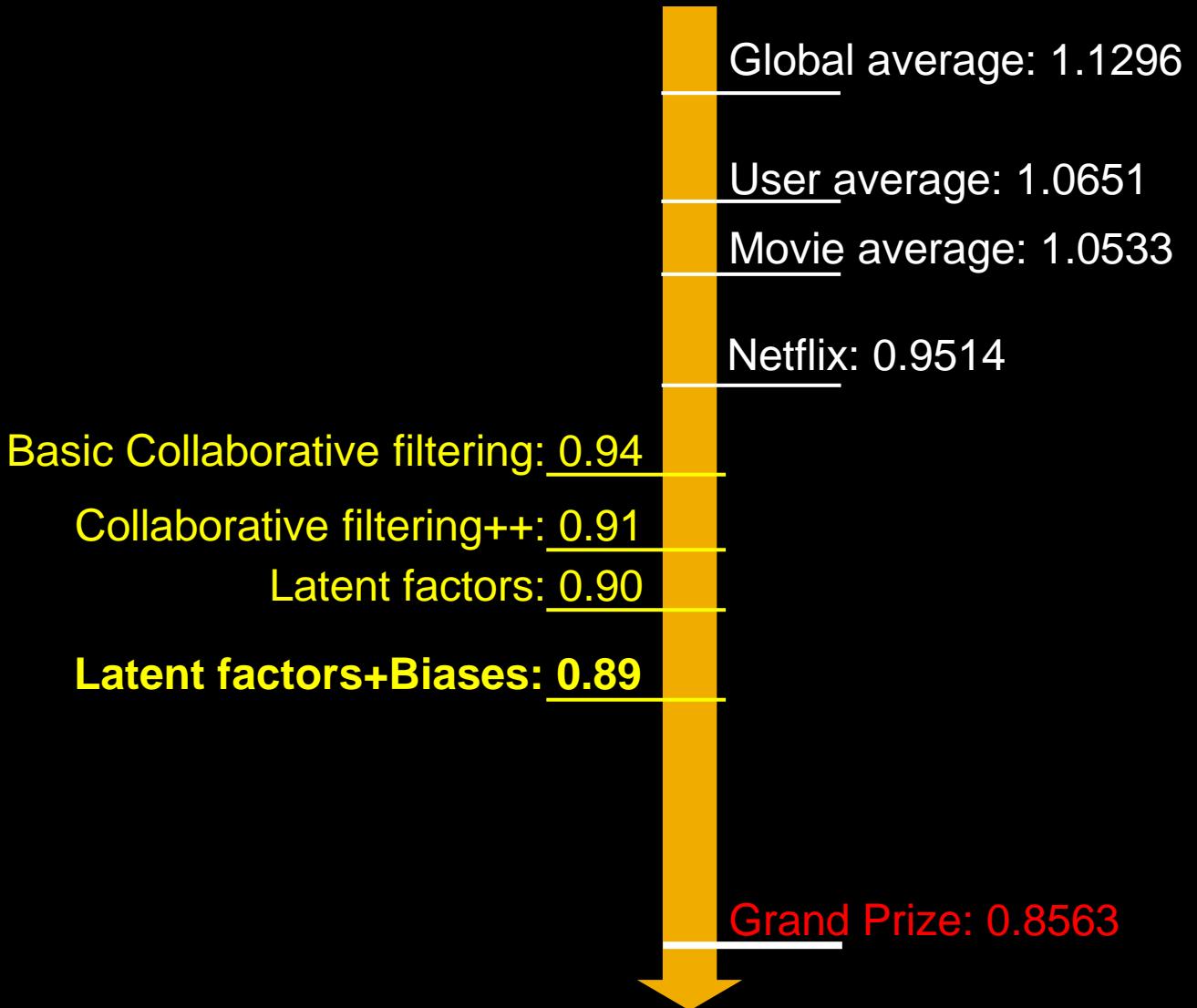
regularization

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 (and we learn them)

Performance of Various Methods

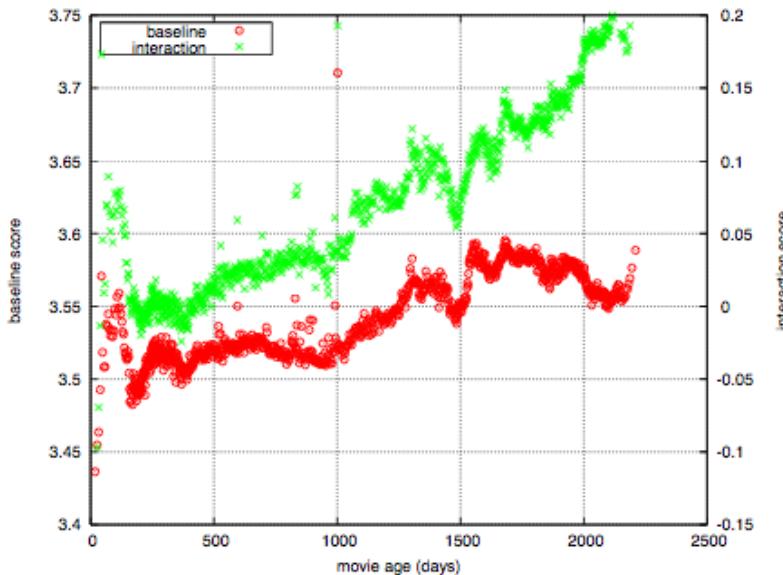
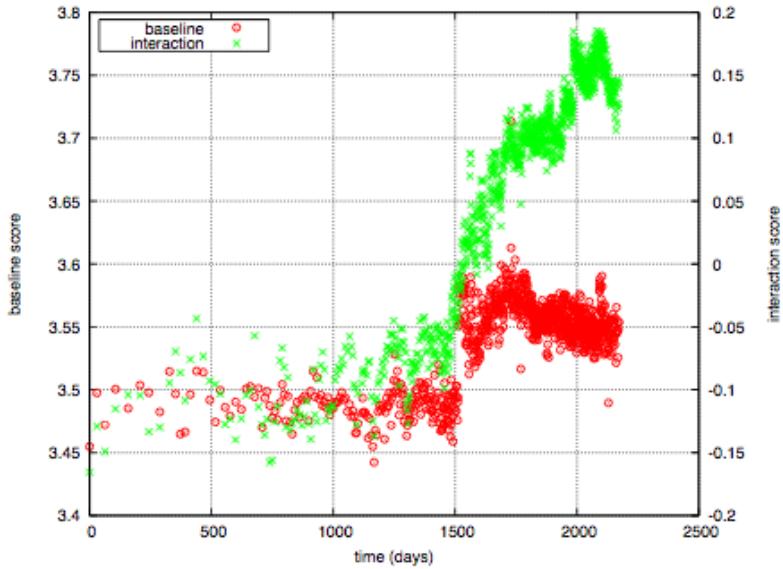


The Netflix Challenge: 2006-09

Temporal Biases Of Users

- **Sudden rise in the average movie rating (early 2004)**
 - Improvements in Netflix
 - GUI improvements
 - Meaning of rating changed
- **Movie age**
 - Users prefer new movies without any reasons
 - Older movies are just inherently better than newer ones

Y. Koren, Collaborative filtering with temporal dynamics, KDD '09



Temporal Biases & Factors

- Original model:

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

- Add time dependence to biases:

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

- Make parameters b_x and b_i to depend on time
 - (1) Parameterize time-dependence by linear trends
 - (2) Each bin corresponds to 10 consecutive weeks

$$b_i(t) = b_i + b_{i,\text{Bin}(t)}$$

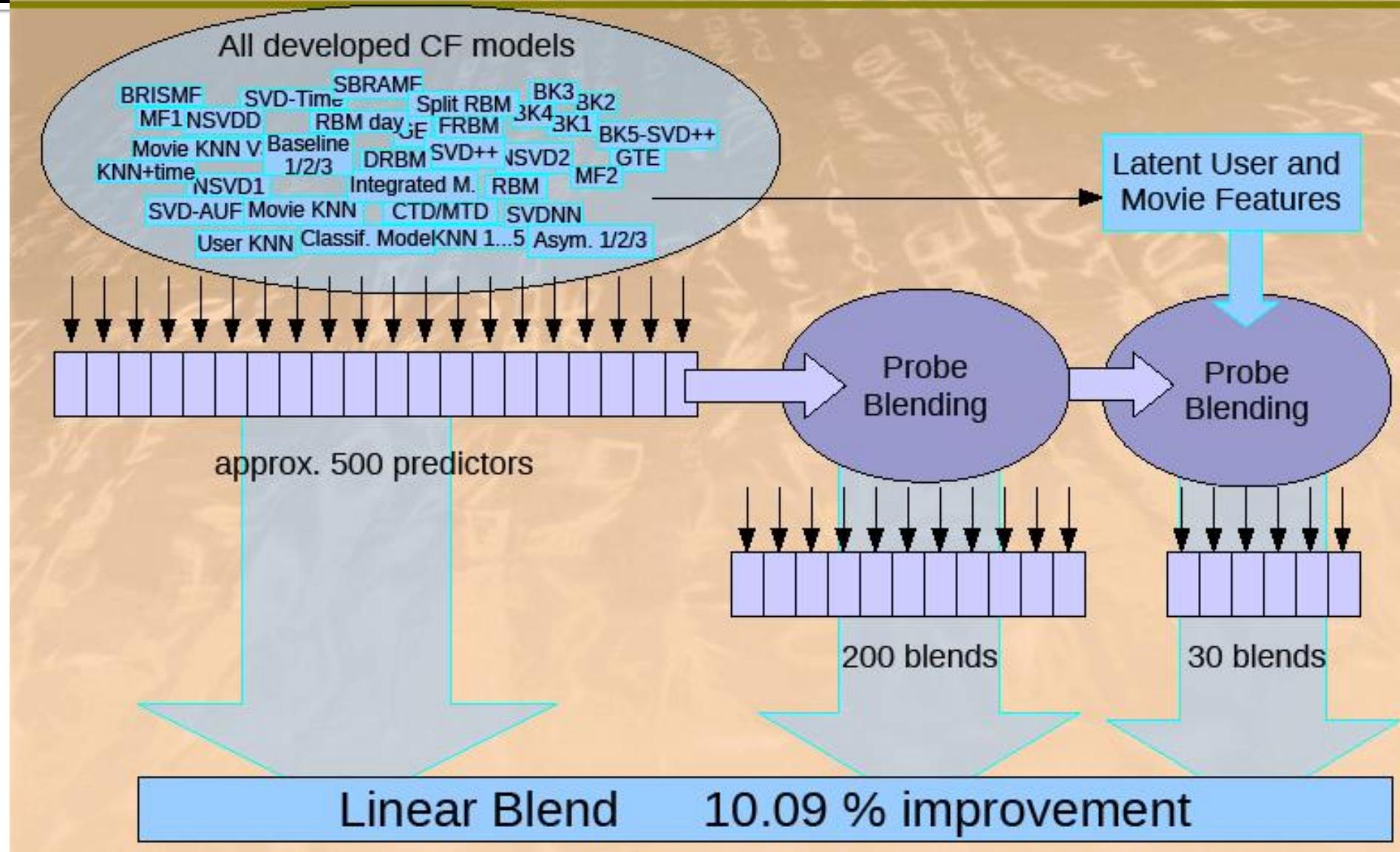
- Add temporal dependence to factors

- $p_x(t)$... user preference vector on day t

Performance of Various Methods



The big picture Solution of BellKor's Pragmatic Chaos



Standing on June 26th 2009

NETFLIX

Netflix Prize

Home Rules Leaderboard Register Update Submit Download

Leaderboard

Display top 20 leaders.

Rank	Team Name	Best Score	% Improvement	Last Submit Time
1	BellKor's Pragmatic Chaos	0.8558	10.05	2009-06-26 18:42:37
Grand Prize - RMSE <= 0.8563				
2	PragmaticTheory	0.8582	9.80	2009-06-25 22:15:51
3	BellKor in BigChaos	0.8590	9.71	2009-05-13 08:14:09
4	Grand Prize Team	0.8593	9.68	2009-06-12 08:20:24
5	Dace	0.8604	9.56	2009-04-22 05:57:03
6	BigChaos	0.8613	9.47	2009-06-23 23:06:52
Progress Prize 2008 - RMSE = 0.8616 - Winning Team: BellKor in BigChaos				
7	BellKor	0.8620	9.40	2009-06-24 07:16:02
8	Gravity	0.8634	9.25	2009-04-22 18:31:32
9	Opera Solutions	0.8638	9.21	2009-06-26 23:18:13
10	BruceDengDiaoCiYiYou	0.8638	9.21	2009-06-27 00:55:55
11	pengpengzhou	0.8638	9.21	2009-06-27 01:06:43
12	xvector	0.8639	9.20	2009-06-26 13:49:04
13	xiangliang	0.8639	9.20	2009-06-26 07:47:34

June 26th submission triggers 30-day “last call”

Netflix Prize

COMPLETED[Home](#) | [Rules](#) | [Leaderboard](#) | [Update](#) | [Download](#)

Leaderboard

Showing Test Score. [Click here to show quiz score](#)Display top leaders.

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8562	9.90	2009-07-10 21:44:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries!	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos

1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8562	9.90	2009-07-10 21:44:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries!	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07
12	BellKor	0.8624	9.46	2009-07-26 17:19:11

Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos

13	xiangliang	0.8642	9.27	2009-07-15 14:53:22
14	Gravity	0.8643	9.26	2009-04-22 18:31:32
15	Ces	0.8651	9.18	2009-06-21 19:24:53
16	Invisible Ideas	0.8653	9.15	2009-07-15 15:53:04
17	Just a guy in a garage	0.8662	9.06	2009-05-24 10:02:54
18	J Dennis Su	0.8666	9.02	2009-03-07 17:16:17
19	Craig Carmichael	0.8666	9.02	2009-07-25 16:00:54
20	acmehill	0.8668	9.00	2009-03-21 16:20:50

Progress Prize 2007 - RMSE = 0.8723 - Winning Team: KorBell

Million \$ Awarded Sept 21st 2009



Acknowledgments

- Some slides and plots borrowed from
Yehuda Koren, Robert Bell and Padhraic Smyth, Jure Leskovec
- **Further reading:**
 - Y. Koren, Collaborative filtering with temporal dynamics, KDD '09
 - [Matrix Factorization Techniques for Recommender Systems](#)
 - [How the Netflix Prize was won](#)