

# Information Storage and Retrieval

CSCE 670  
Texas A&M University  
Department of Computer Science & Engineering  
Instructor: Prof. James Caverlee

**Recommenders**  
**6/8 March 2018**

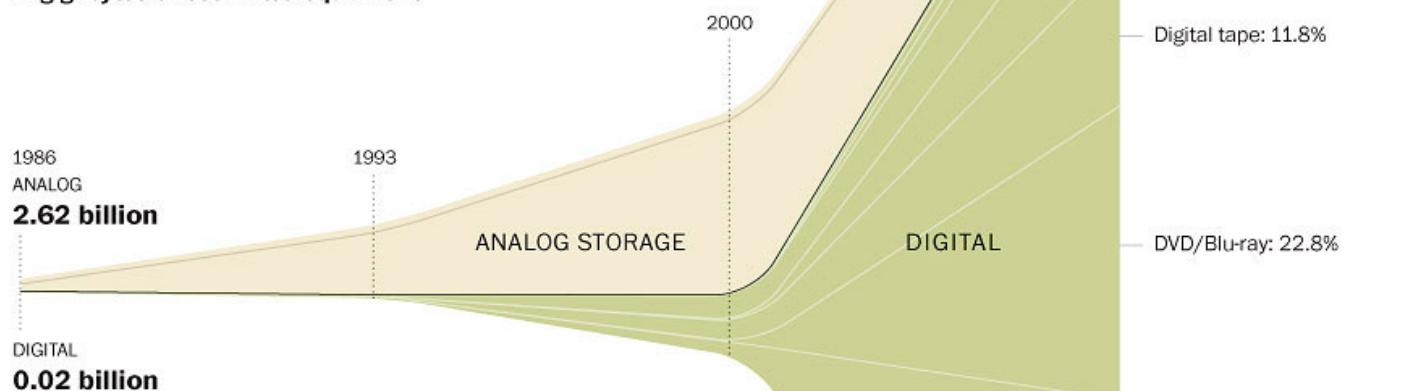
# Today

- Recommendation systems: Intro
- Collaborative Filtering
  - Finding nearest neighbors
  - How to aggregate ratings?
- Content-based Recommenders
- Netflix Prize (start)

### THE WORLD'S CAPACITY TO STORE INFORMATION

This chart shows the world's growth in storage capacity for both analog data (books, newspapers, videotapes, etc.) and digital (CDs, DVDs, computer hard drives, smartphone drives, etc.)

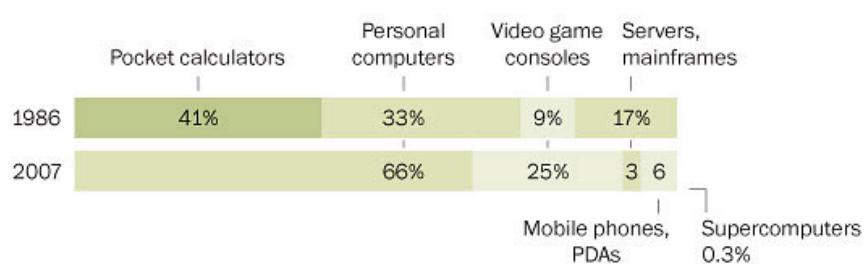
**In gigabytes or estimated equivalent**



### COMPUTING POWER

In 1986, pocket calculators accounted for much of the world's data-processing power.

**Percentage of available processing power by device:**



THE  
MOBILE WEB  
RECEIVES

217

NEW USERS.

WORDPRESS  
USERS PUBLISH

347 NEW  
BLOG  
POSTS.

571

NEW WEBSITES  
ARE CREATED.

FOURSQUARE USERS  
PERFORM

2,083

CHECK-INS.

FLICKR

YOUTUBE  
USERS UPLOAD

48

HOURS  
OF NEW VIDEO.



EMAIL  
USERS  
SEND  
204,166,667  
MESSAGES.



GOOGLE  
RECEIVES  
OVER  
2,000,000  
SEARCH QUERIES.



FACEBOOK  
USERS

SHARE  
684,478

PIECES OF CONTENT.

CONSUMERS  
SPEND  
\$272,070

ON WEB SHOPPING.



TWITTER USERS  
SEND OVER

100,000  
TWEETS.

ADDIE

EVERY  
MINUTE  
*of the*  
DAY

# How Many Products Does Amazon Sell?

by Paul Grey on 15 December 2013 in E-Commerce

Amazon.com is the self-styled "Greatest Store on Earth." It's been said that Amazon aims one day to sell everything to everyone.



Exactly how much choice do you have when shopping with Amazon?

Today Amazon sells over 200 million products in the USA, which are categorised into 35 departments. There are almost 5 million items in the Clothing department, almost 20 million in Sports & Outdoors, and over 4 million Office Products. There 7 million items in the Amazon Jewelry department, 24 million in Electronics, 1.4 million products in the Beauty department, 570 thousand Baby products, and 600 thousand Grocery items.

That's in the USA. This table lists my estimates of the number of products offered on the main Amazon websites around the world.

Amazon.com	USA	232 million
Amazon.co.uk	UK	132 million
Amazon.de	Germany	118 million

## **How many tracks does it have?**

**Spotify: 30 million+**

**Apple Music: 30 million+**

The two services have roughly the same number of tracks; some 30 million at their disposal.

Apple Music users can listen to their own personal iTunes library through the app as well as the Apple Music streaming library. Spotify allows you to listen to non-Spotify music too, but adding local files has to be done through a desktop Mac or PC. For example, you can't listen to your iTunes library through Spotify itself.

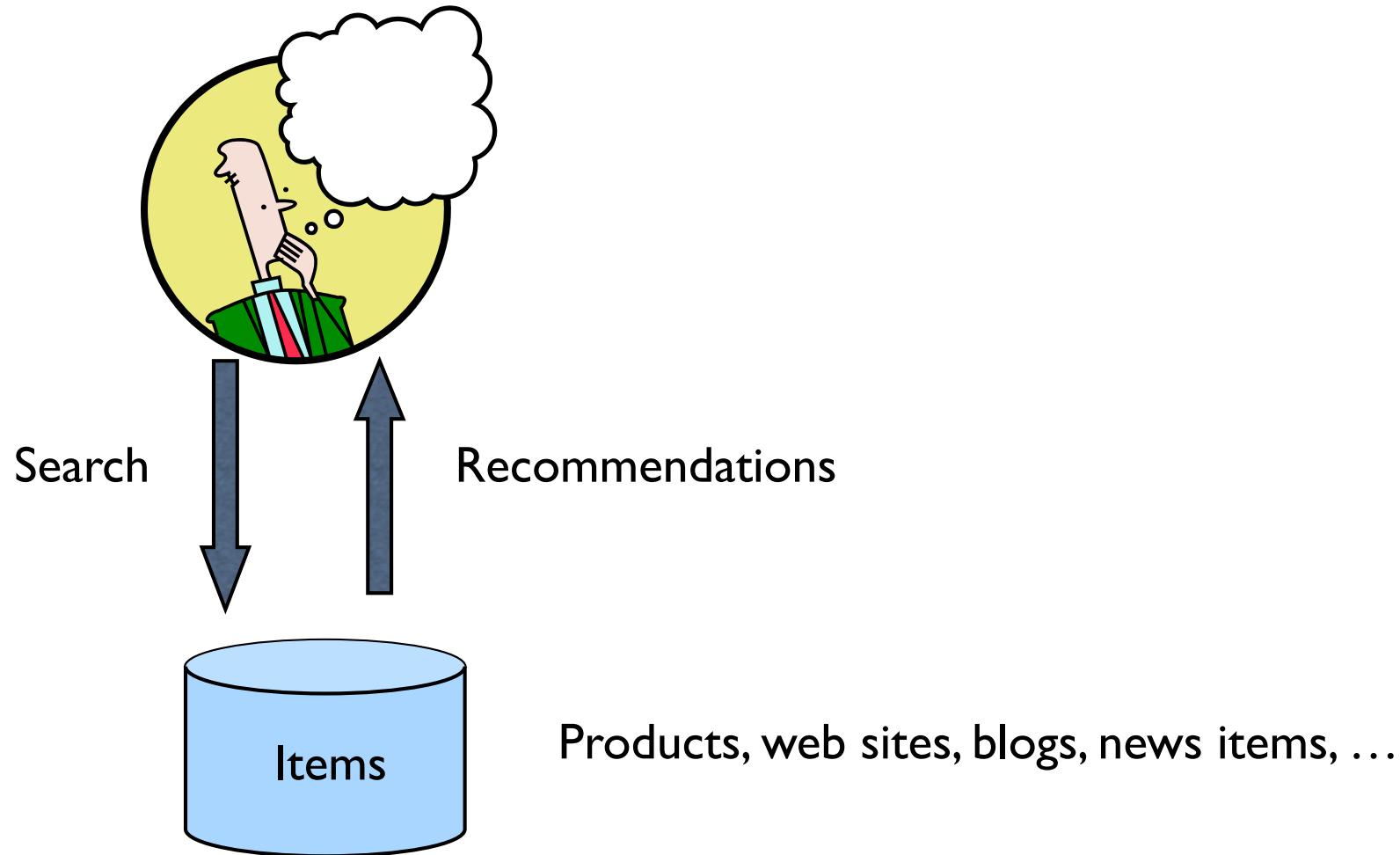


## Statistics

### Viewership

- More than 1 billion unique users visit YouTube each month
- Over 6 billion hours of video are watched each month on YouTube—that's almost an hour for every person on Earth
- 100 hours of video are uploaded to YouTube every minute
- 80% of YouTube traffic comes from outside the US
- YouTube is localized in 61 countries and across 61 languages
- According to Nielsen, YouTube reaches more US adults ages 18-34 than any cable network
- Millions of subscriptions happen each day. The number of people subscribing daily is up more than 3x since last year. The number of daily subscriptions is up more than 4x since last year

# Recommendations



# Why recommendation?

The goal of recommender systems is...

- To help people discover new content

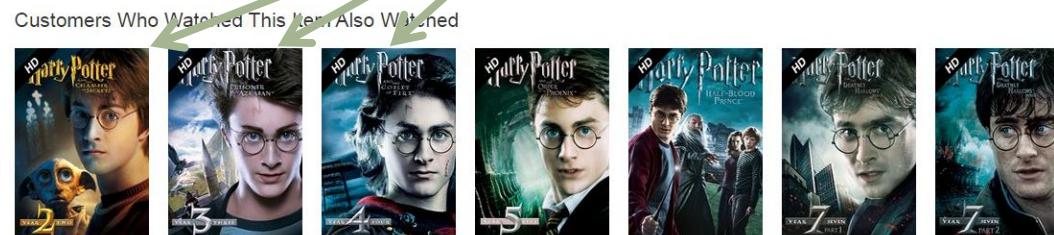
Recommendations for You in Amazon Instant Video [See more](#)



# Why recommendation?

The goal of recommender systems is...

- To help us find the content we were already looking for



# Why recommendation?

The goal of recommender systems is...

- To discover which things go together



Calvin Klein Men's Relaxed Straight Leg Jean In Cove

★★★★★ 5+ 20 customer reviews

Price: \$48.16 - \$69.99 & FREE Returns. Details

Size:

Select ▾ Sizing info | Fit: As expected (55%) ▾

Color: Cove

- 98% Cotton/2% Elastane
- Imported
- Button closure
- Machine Wash
- Relaxed straight-leg jean in light-tone denim featuring whiskering and five-pocket styling
- Zip fly with button
- 10.25-inch front rise, 19-inch knee, 17.5-inch leg opening

Frequently Bought Together



Calvin Klein Jeans  
\$57.94 - \$69.50



Calvin Klein Jeans  
\$49.92



Calvin Klein Jeans  
\$50.67 - \$69.99



Levi's  
\$23.99 - \$68.00

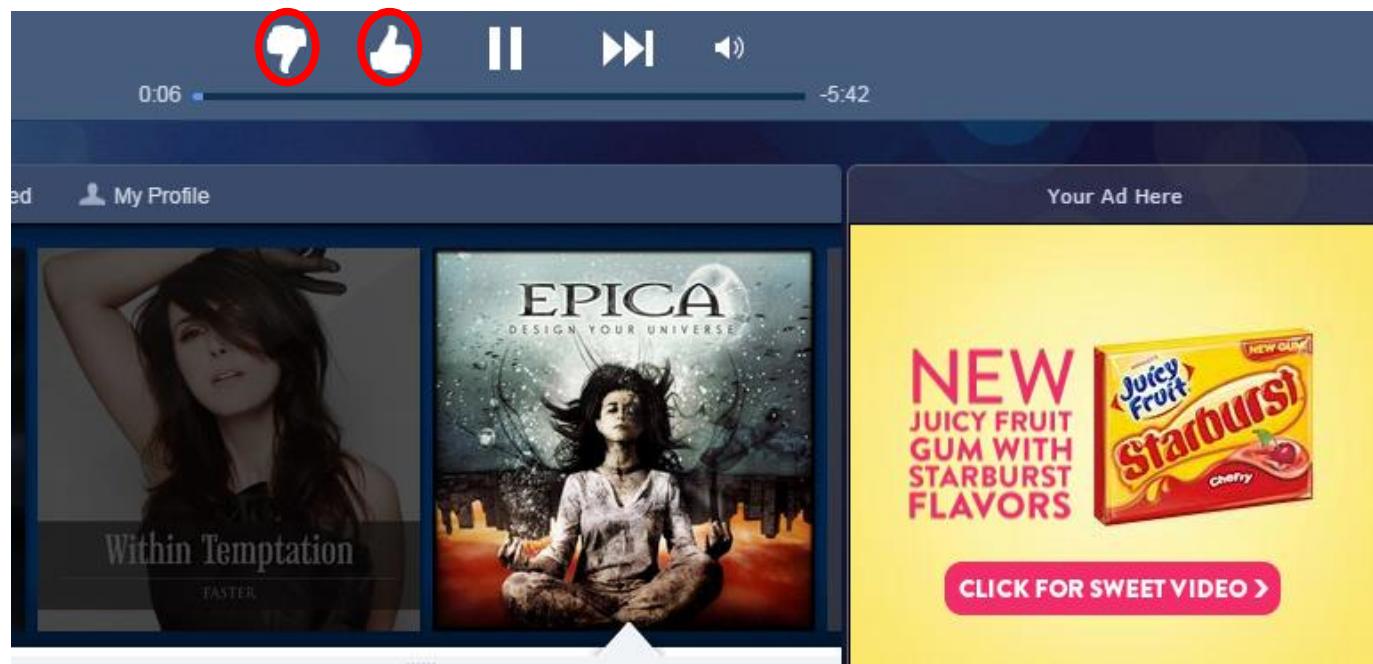
Customers Who Bought This Item Also Bought



# Why recommendation?

The goal of recommender systems is...

- To personalize user experiences in response to user feedback



# Why recommendation?

The goal of recommender systems is...

- To identify things that we **like**

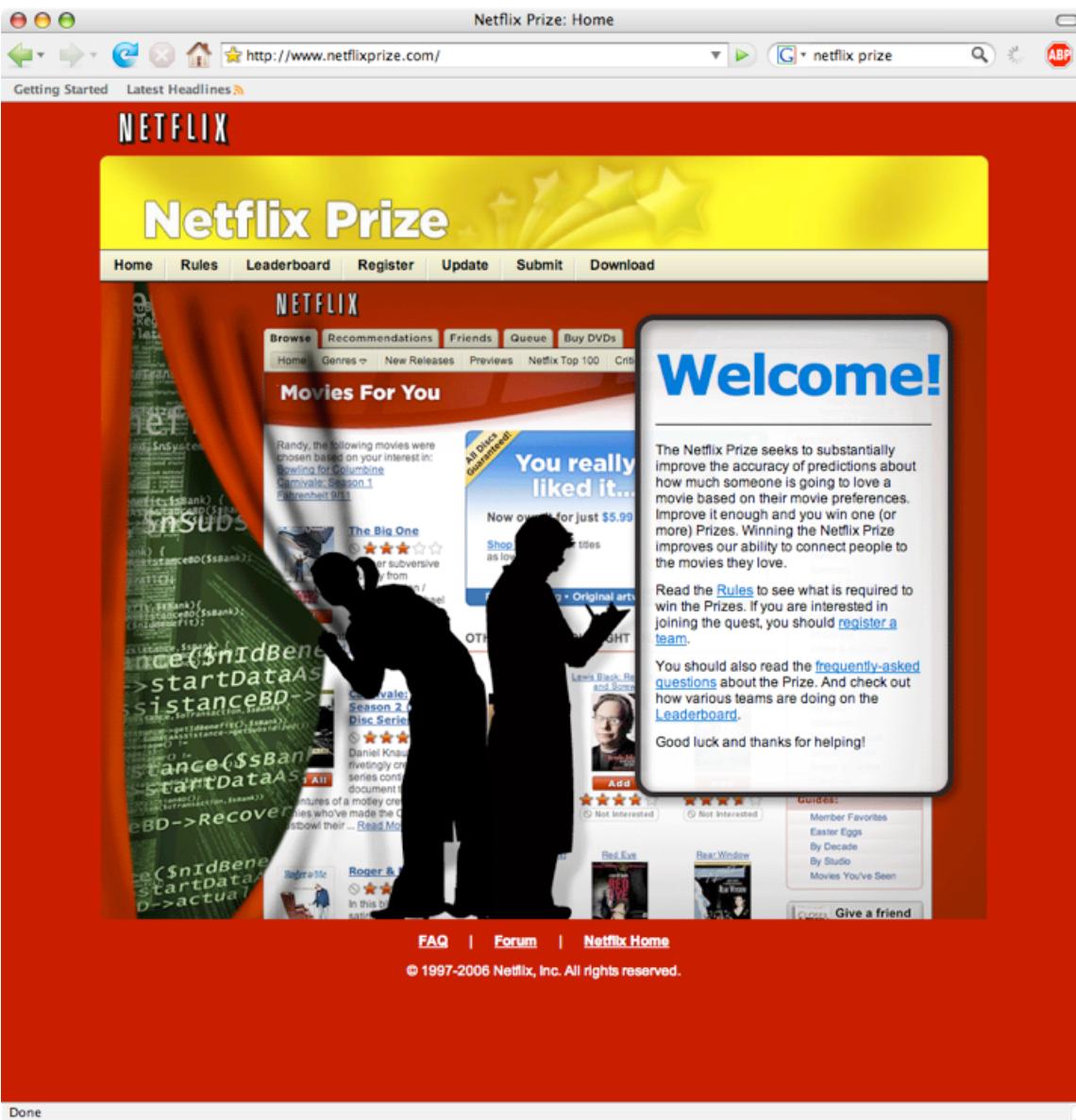


# Examples of recommenders?

# 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
- $S$  = 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

	Spotlight	The Revenant	Mad Max: Fury Road	Room
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 known ratings
  - Mainly interested in high unknown ratings
  - Don’t care about finding what you **don’t like**, but rather what you do 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
  - What about low ratings?

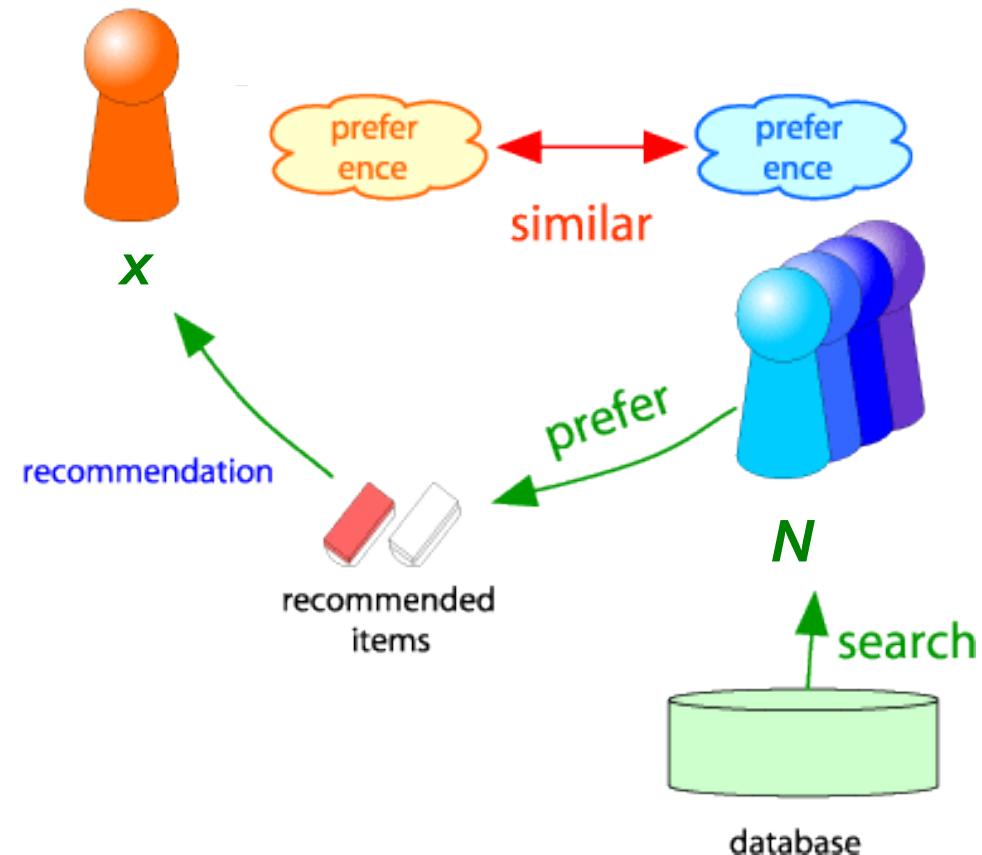
# Extrapolating Utilities

- Key problem: 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
  - Collaborative
  - Content-based
  - Latent factor models

# Collaborative Recommendations

# 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

- Let  $r_x$  be the vector of user  $x$ 's ratings
- Jaccard!
  - But it ignores the value of the rating
- Cosine similarity measure
  - $\text{sim}(x,y) = \cos(r_x, r_y)$
  - But it treats missing ratings as “negative”
- Pearson correlation coefficient
  - $S_{xy} = \text{items rated by both users } x \text{ and } y$

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

# Jaccard

$\text{Jaccard}(A, B) =$

$\text{Jaccard}(U_i, U_j) =$

→ Maximum of 1 if the two users purchased **exactly the same** set of items

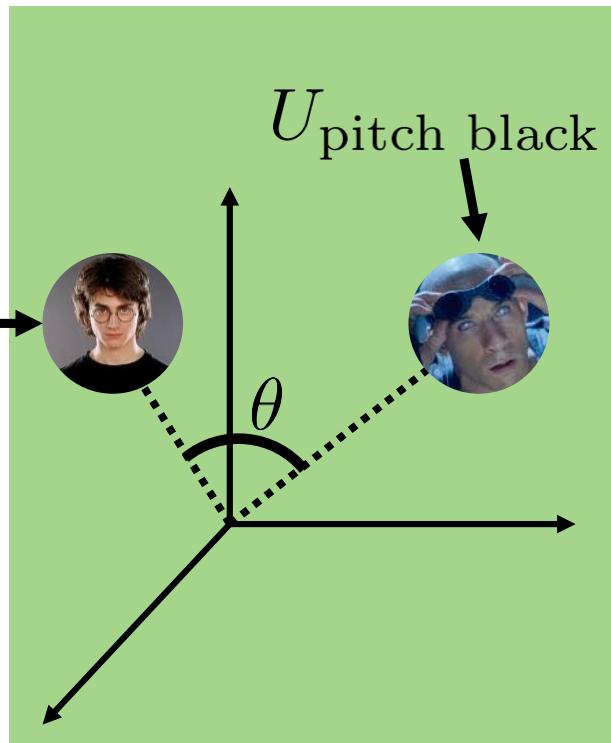
(or if two items were purchased by the same set of users)

→ Minimum of 0 if the two users purchased **completely disjoint** sets of items

(or if the two items were purchased by completely disjoint sets of users)

# Cosine

$U_{\text{harry potter}}$   
(vector representation of  
users who purchased  
harry potter)



$$\cos(\theta) = 1$$

(theta = 0)  $\rightarrow$   $A$  and  $B$  point in  
exactly the same direction

$$\cos(\theta) = -1$$

(theta = 180)  $\rightarrow$   $A$  and  $B$  point  
in opposite directions (won't  
actually happen for 0/1 vectors)

$$\cos(\theta) = 0$$

(theta = 90)  $\rightarrow$   $A$  and  $B$  are  
orthogonal

# Cosine

## Why cosine?

- Unlike Jaccard, works for arbitrary vectors
- E.g. what if we have **opinions** in addition to purchases?

$$R = \begin{pmatrix} 1 & 0 & \cdots & 1 \\ 0 & 0 & & 1 \\ \vdots & & \ddots & \vdots \\ 1 & 0 & \cdots & 1 \end{pmatrix} \xrightarrow{\hspace{1cm}} \begin{pmatrix} -1 & 0 & \cdots & 1 \\ 0 & 0 & & -1 \\ \vdots & & \ddots & \vdots \\ 1 & 0 & \cdots & -1 \end{pmatrix}$$

bought and **liked**

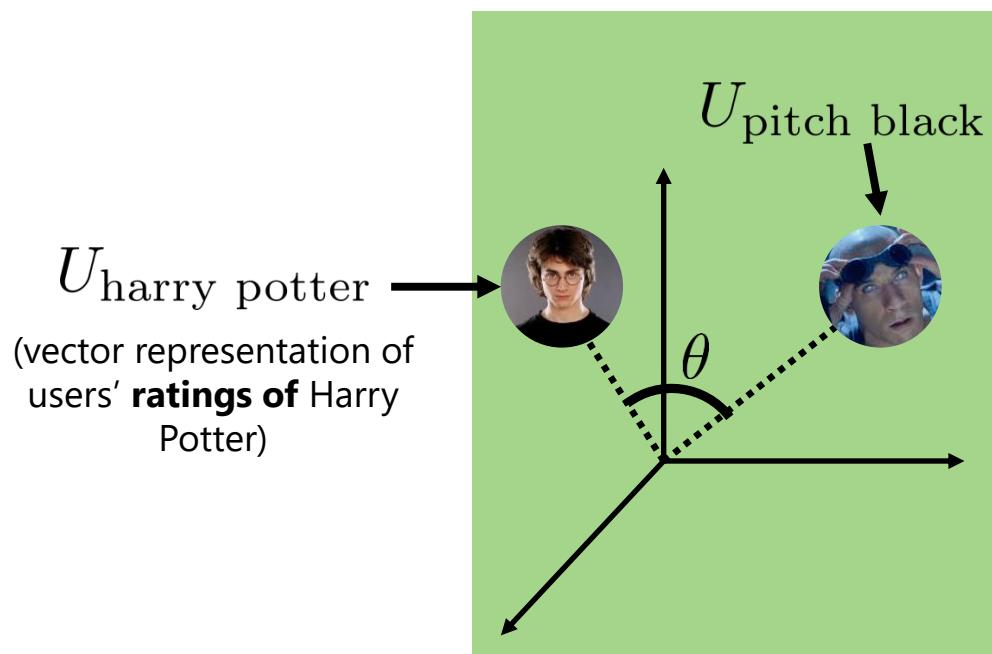
didn't buy

bought and **hated**

The diagram illustrates the transformation of a binary matrix  $R$  into a signed matrix. The original matrix  $R$  is a binary matrix where rows represent users and columns represent items. The transformed matrix has the same structure but with signed values: 1 for 'bought and liked', 0 for 'didn't buy', and -1 for 'bought and hated'. This transformation allows cosine similarity to measure the angle between user profiles, even when they have different purchase patterns.

# Cosine (with ratings)

E.g. our previous example, now with “thumbs-up/thumbs-down” ratings



$$\cos(\theta) = 1$$

(theta = 0) → Rated by the same users, and they all agree

$$\cos(\theta) = -1$$

(theta = 180) → Rated by the same users, but they **completely disagree** about it

$$\cos(\theta) = 0$$

(theta = 90) → Rated by different sets of users

# Pearson correlation

What if we have numerical ratings  
(rather than just thumbs-up/down)?

$$R = \begin{pmatrix} -1 & 0 & \dots & 1 \\ 0 & 0 & & -1 \\ \vdots & & \ddots & \vdots \\ 1 & 0 & \dots & -1 \end{pmatrix} \xrightarrow{\hspace{1cm}} \begin{pmatrix} 4 & 0 & \dots & 2 \\ 0 & 0 & & 3 \\ \vdots & & \ddots & \vdots \\ 5 & 0 & \dots & 1 \end{pmatrix}$$

bought and **liked**  
didn't buy  
bought and **hated**

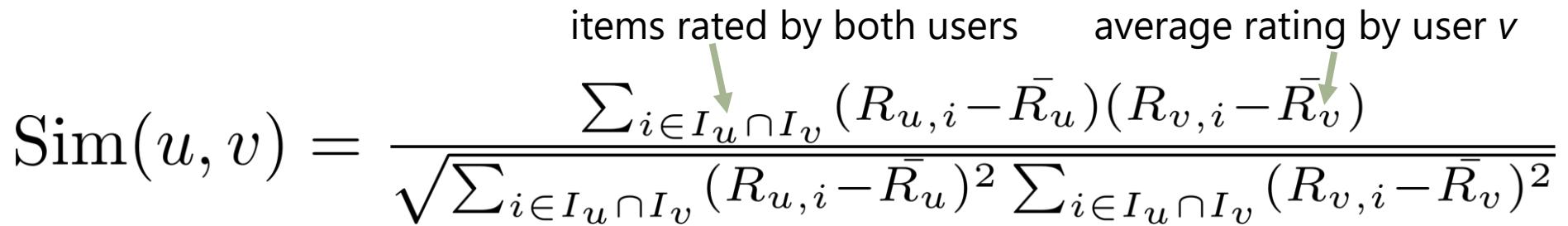
# Pearson correlation

What if we have numerical ratings  
(rather than just thumbs-up/down)?

- We wouldn't want 1-star ratings to be parallel to 5-star ratings
  - So we can subtract the average – values are then **negative** for below-average ratings and **positive** for above-average ratings

$$\text{Sim}(u, v) = \frac{\sum_{i \in I_u \cap I_v} (R_{u,i} - \bar{R}_u)(R_{v,i} - \bar{R}_v)}{\sqrt{\sum_{i \in I_u \cap I_v} (R_{u,i} - \bar{R}_u)^2} \sqrt{\sum_{i \in I_u \cap I_v} (R_{v,i} - \bar{R}_v)^2}}$$

items rated by both users      average rating by user  $v$



# Example

	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:  $0.386 > 0.322$ 
  - 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

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

- So far, we can find a rating, but how do we actually generate recommendations?

# How to aggregate ratings?

$$r_{c,s} = \text{aggr } r_{c',s},$$

$c' \in \hat{C}$

“similar”  
users to c

# How to aggregate ratings?

$$(a) \ r_{c,s} = \frac{1}{N} \sum_{c' \in \hat{C}} r_{c',s},$$

$$(b) \ r_{c,s} = k \sum_{c' \in \hat{C}} sim(c, c') \times r_{c',s},$$

$$(c) \ r_{c,s} = \bar{r}_c + k \sum_{c' \in \hat{C}} sim(c, c') \times (r_{c',s} - \bar{r}_{c'}),$$

Many other tricks possible

# Item-Item Collaborative Filtering

- So far: User-user collaborative filtering
- Another view
  - For item  $s$ , find other similar items
  - Estimate rating for item based on ratings for similar items
  - Can use same similarity metrics and prediction functions as in user-user model
- In practice, it has been observed that item-item often works better than user-user

# Pros/cons of collaborative filtering

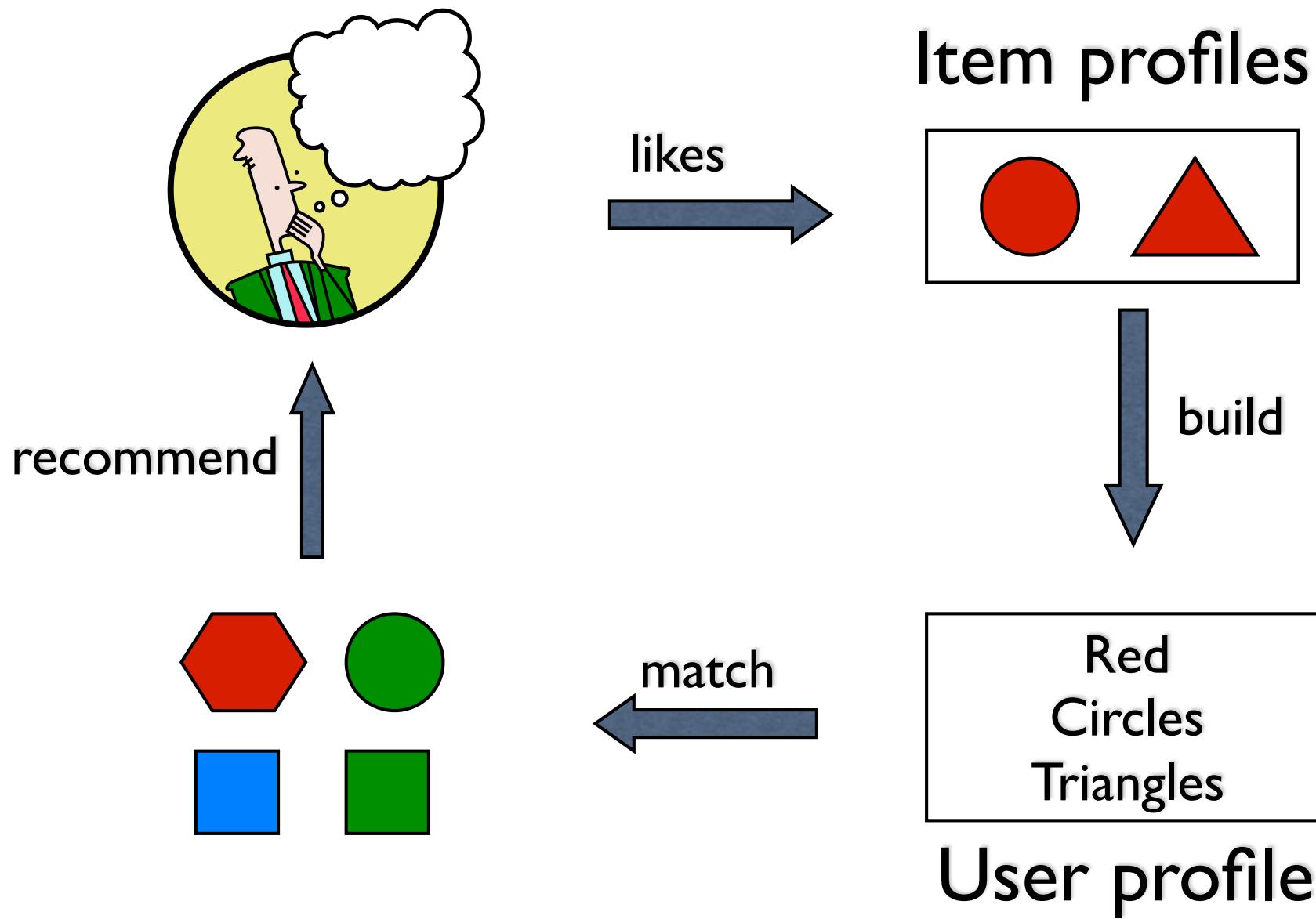
- Works for any kind of item
  - No feature selection needed
- Cold start:
  - Need enough users in the system to find a match
- Sparsity:
  - The user/ratings matrix is sparse
  - Hard to find users that have rated the same items
- First rater:
  - Cannot recommend an item that has not been previously rated
  - New items, esoteric items
- Popularity bias:
  - Cannot recommend items to someone with unique taste
  - Tends to recommend popular items

# Content-based recommendations

# 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 words?
  - Usual heuristic is TF.IDF

# 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 Recs?

- 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 and unpopular items
  - No first-rater problem
- Can provide explanations of recommended items by listing content-features that caused item to be recommended

# Limitations of content-based approach

- Finding the appropriate features
  - e.g., images, movies, music
- Recommendations for new users
  - How to build a profile?
- Overspecialization
  - Never recommends items outside user's content profile
  - People might have multiple interests
  - Unable to exploit quality judgments of other users

**Hybrid: Content +  
Collaborative**

# Hybrid Methods

- Implement two separate recommenders and combine predictions
- Add content-based methods to collaborative filtering
  - item profiles for new item problem
  - demographics to deal with new user problem

Recommendation Approach	Recommendation Technique	
	Heuristic-based	Model-based
Content-based	<p>Commonly used techniques:</p> <ul style="list-style-type: none"> <li>• TF-IDF (information retrieval)</li> <li>• Clustering</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Lang 1995</li> <li>• Balabanovic &amp; Shoham 1997</li> <li>• Pazzani &amp; Billsus 1997</li> </ul>	<p>Commonly used techniques:</p> <ul style="list-style-type: none"> <li>• Bayesian classifiers</li> <li>• Clustering</li> <li>• Decision trees</li> <li>• Artificial neural networks</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Pazzani &amp; Billsus 1997</li> <li>• Mooney et al. 1998</li> <li>• Mooney &amp; Roy 1999</li> <li>• Billsus &amp; Pazzani 1999, 2000</li> <li>• Zhang et al. 2002</li> </ul>
Collaborative	<p>Commonly used techniques:</p> <ul style="list-style-type: none"> <li>• Nearest neighbor (cosine, correlation)</li> <li>• Clustering</li> <li>• Graph theory</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Resnick et al. 1994</li> <li>• Hill et al. 1995</li> <li>• Shardanand &amp; Maes 1995</li> <li>• Breese et al. 1998</li> <li>• Nakamura &amp; Abe 1998</li> <li>• Aggarwal et al. 1999</li> <li>• Delgado &amp; Ishii 1999</li> <li>• Pennock &amp; Horwitz 1999</li> <li>• Sarwar et al. 2001</li> </ul>	<p>Commonly used techniques:</p> <ul style="list-style-type: none"> <li>• Bayesian networks</li> <li>• Clustering</li> <li>• Artificial neural networks</li> <li>• Linear regression</li> <li>• Probabilistic models</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Billsus &amp; Pazzani 1998</li> <li>• Breese et al. 1998</li> <li>• Ungar &amp; Foster 1998</li> <li>• Chien &amp; George 1999</li> <li>• Getoor &amp; Sahami 1999</li> <li>• Pennock &amp; Horwitz 1999</li> <li>• Goldberg et al. 2001</li> <li>• Kumar et al. 2001</li> <li>• Pavlov &amp; Pennock 2002</li> <li>• Shani et al. 2002</li> <li>• Yu et al. 2002, 2004</li> <li>• Hofmann 2003, 2004</li> <li>• Marlin 2003</li> <li>• Si &amp; Jin 2003</li> </ul>
Hybrid	<p>Combining content-based and collaborative components using:</p> <ul style="list-style-type: none"> <li>• Linear combination of predicted ratings</li> <li>• Various voting schemes</li> <li>• Incorporating one component as a part of the heuristic for the other</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Balabanovic &amp; Shoham 1997</li> <li>• Claypool et al. 1999</li> <li>• Good et al. 1999</li> <li>• Pazzani 1999</li> <li>• Billsus &amp; Pazzani 2000</li> <li>• Tran &amp; Cohen 2000</li> <li>• Melville et al. 2002</li> </ul>	<p>Combining content-based and collaborative components by:</p> <ul style="list-style-type: none"> <li>• Incorporating one component as a part of the model for the other</li> <li>• Building one unifying model</li> </ul> <p>Representative research examples:</p> <ul style="list-style-type: none"> <li>• Basu et al. 1998</li> <li>• Condliff et al. 1999</li> <li>• Soboroff &amp; Nicholas 1999</li> <li>• Ansari et al. 2000</li> <li>• Popescul et al. 2001</li> <li>• Schein et al. 2002</li> </ul>

# Netflix Prize

# The Netflix Prize

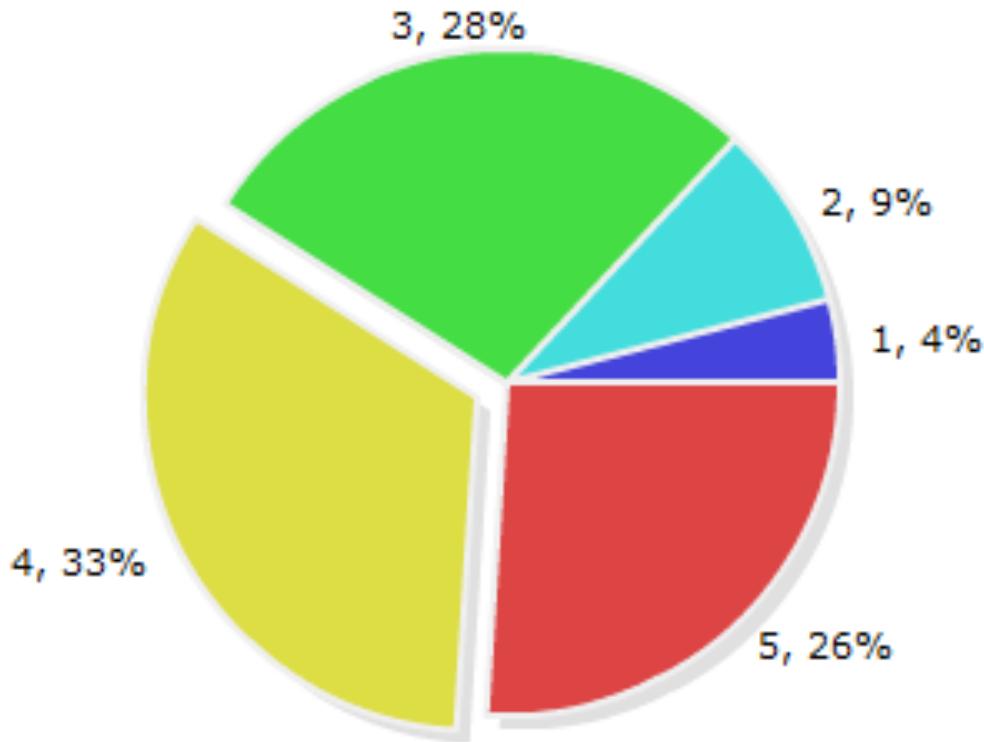
- Training data
  - 100 million ratings
  - 480,000 users
  - 17,770 movies
  - 6 years of data: 2000-2005
- Test data
  - Last few ratings of each user (2.8 million)
  - Evaluation criterion: root mean squared error (RMSE)
  - Netflix Cinematch system RMSE: 0.9514
- Competition
  - 2700+ teams
  - \$1 million grand prize for 10% improvement over Netflix

# RMSE

- » Compare predictions with known ratings
- » My system predicted you would rate
  - » **The Shawshank Redemption as 4.3 stars**
    - » In reality, you gave it **5 stars**
  - » **The Matrix with 3.9 stars**
    - » In reality, you gave it **4 stars**
- »  $\text{RMSE} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$

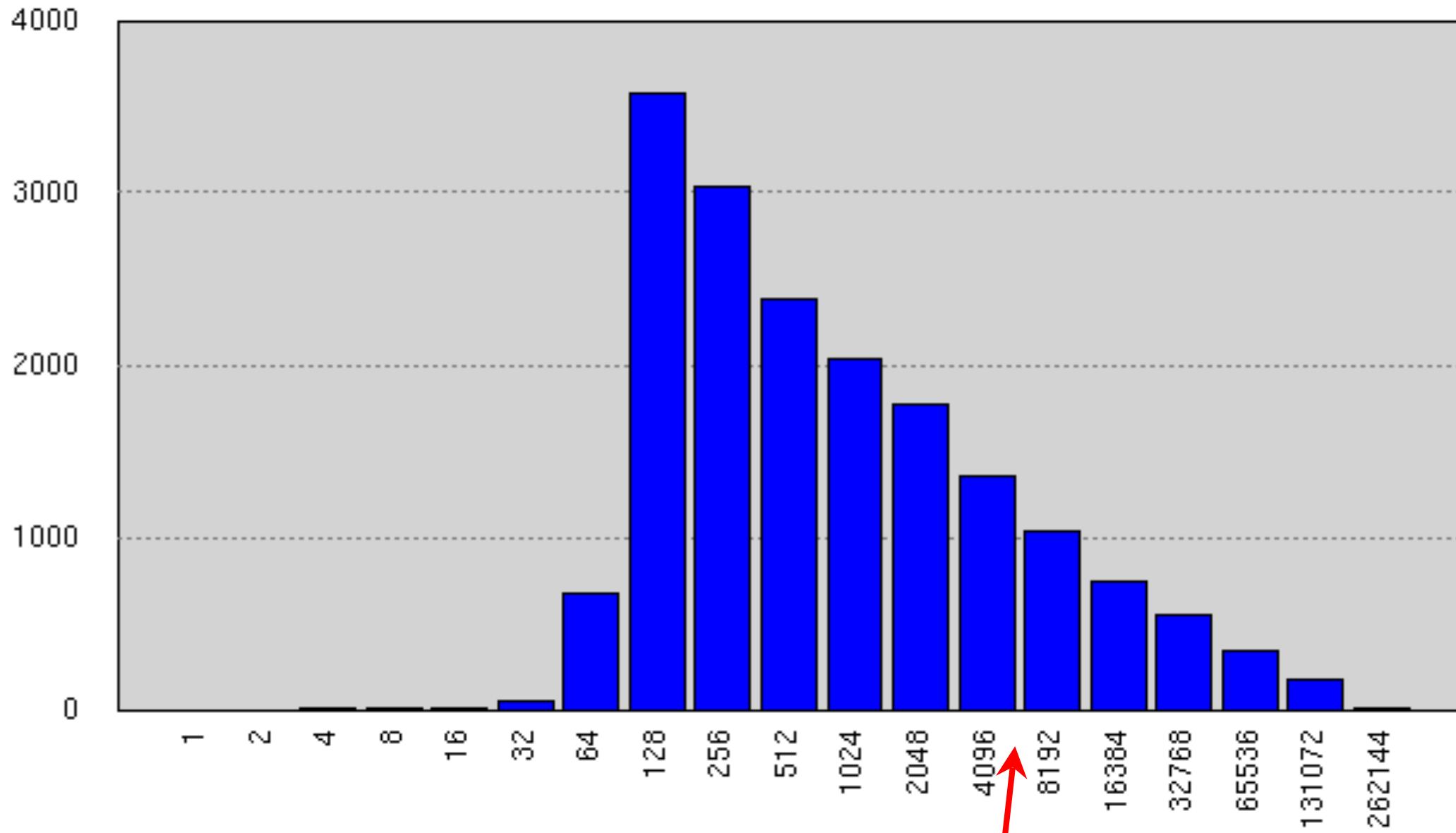
$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

# Overall rating distribution



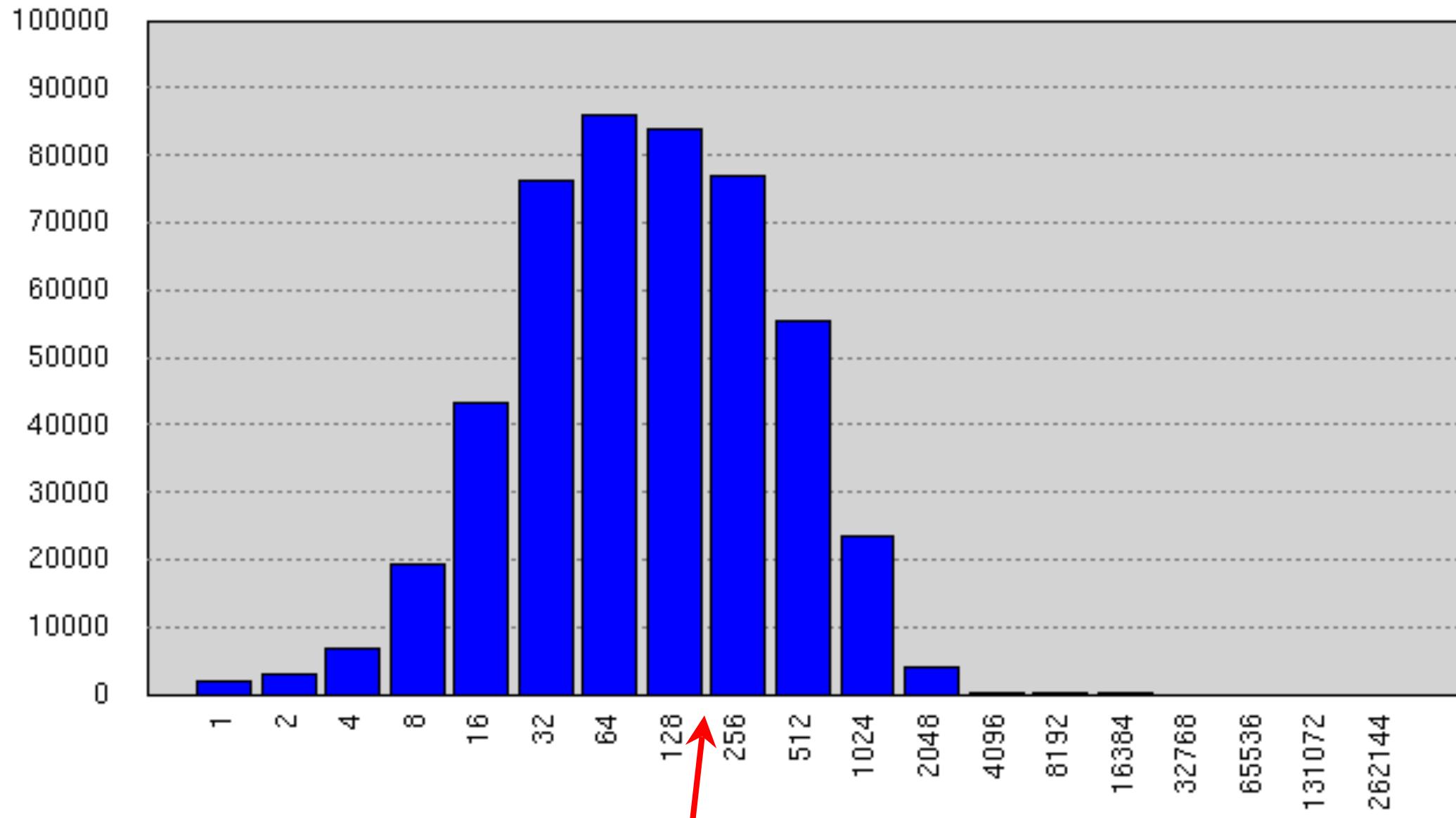
- Third of ratings are 4s
- Average rating is 3.68

# #ratings per movie



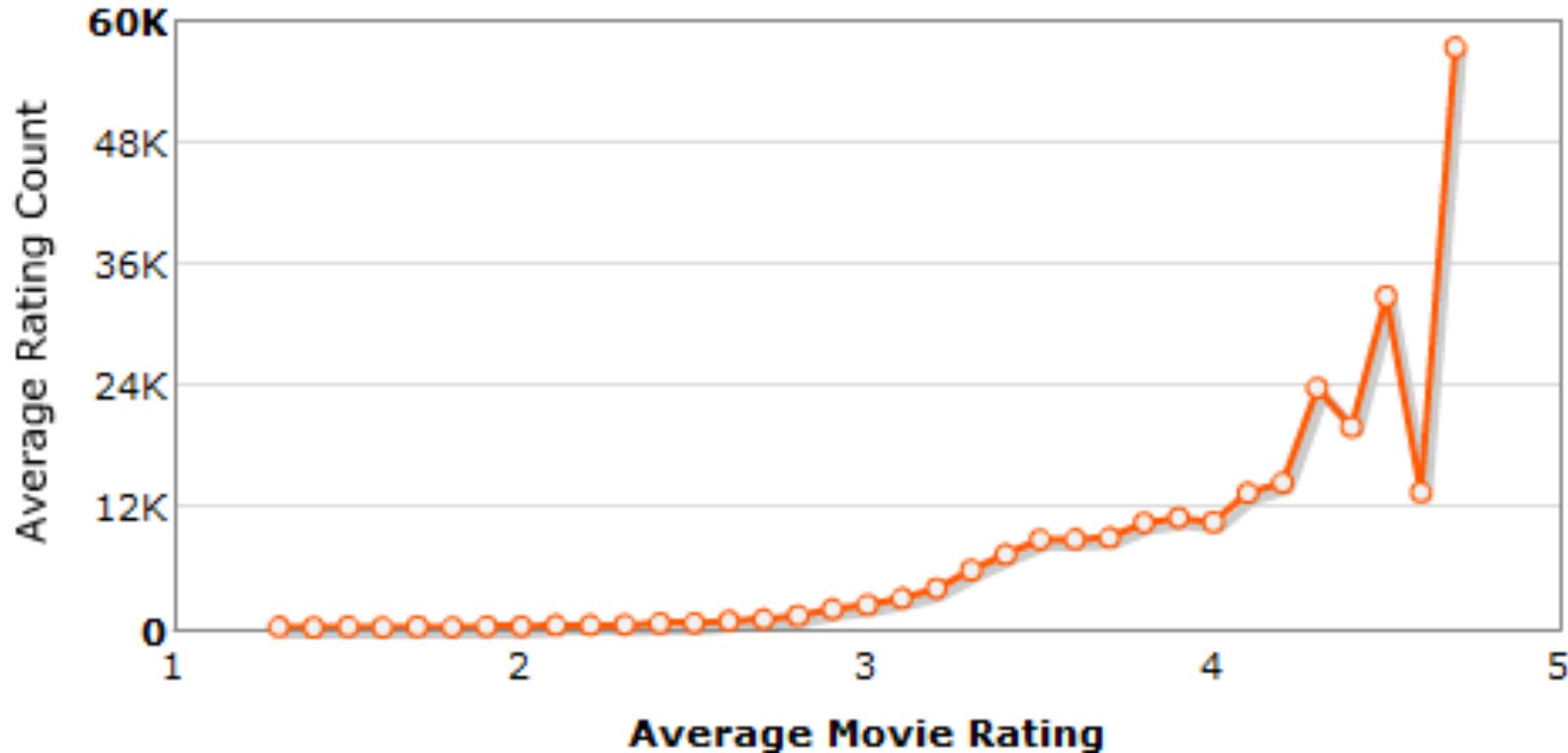
- Avg #ratings/movie: **5627**

# #ratings per user



- Avg #ratings/user: 208

# Average movie rating by movie count



- More ratings to better movies

# Most loved movies

Title	Avg rating	Count
The Shawshank Redemption	4.593	137812
Lord of the Rings: The Return of the King	4.545	133597
The Green Mile	4.306	180883
Lord of the Rings: The Two Towers	4.460	150676
Finding Nemo	4.415	139050
Raiders of the Lost Ark	4.504	117456
Forrest Gump	4.299	180736
Lord of the Rings: The Fellowship of the ring	4.433	147932
The Sixth Sense	4.325	149199
Indiana Jones and the Last Crusade	4.333	144027

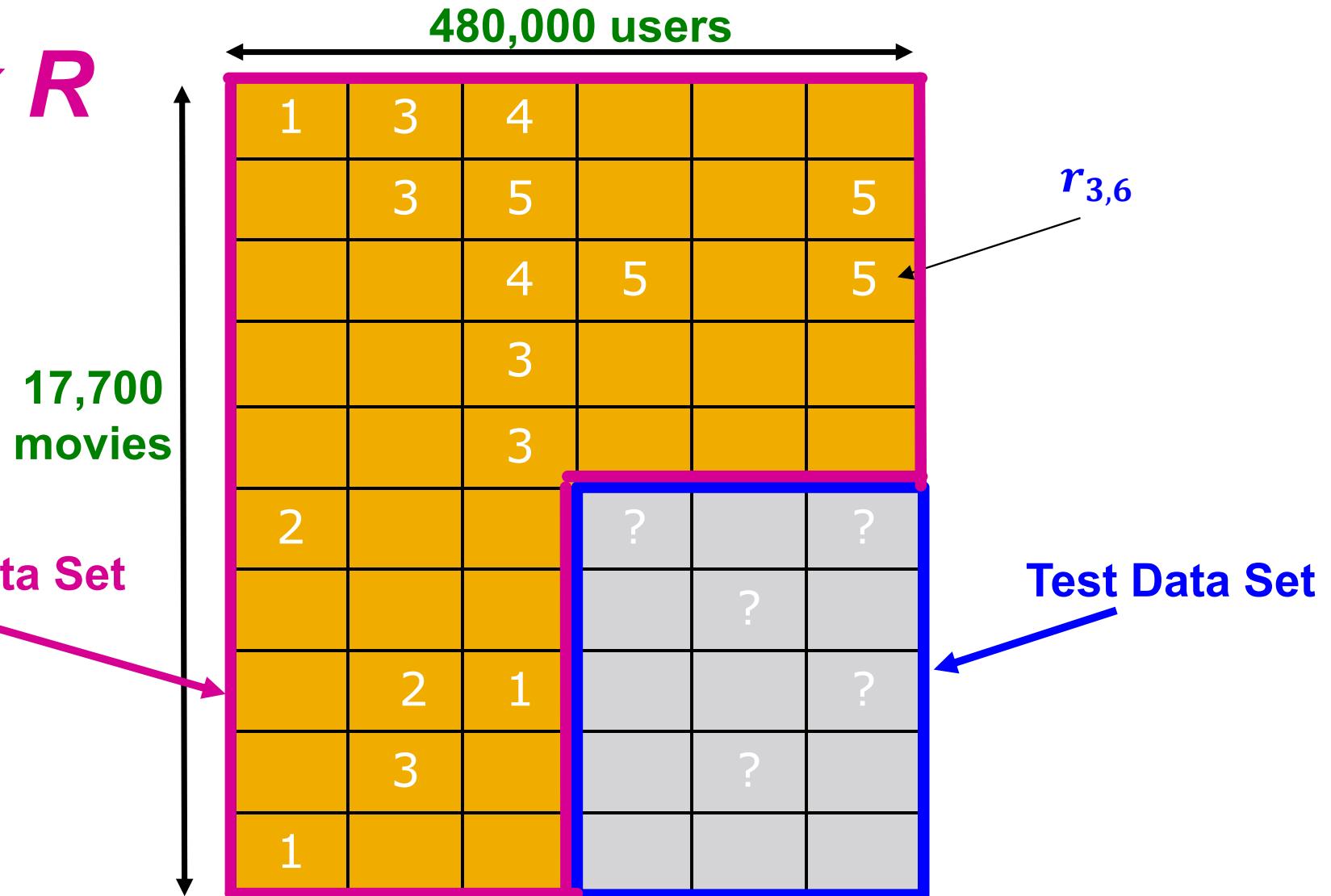
# The Netflix Utility Matrix R

**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**





# Netflix Prize

[Home](#) [Rules](#) [Leaderboard](#) [Register](#) [Update](#) [Submit](#) [Download](#)

## Leaderboard

Display top  leaders.

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

June 26<sup>th</sup> submission triggers 30-day “last call”

## ■ Ensemble team formed

- Group of other teams on leaderboard forms a new team
- Relies on combining their models
- Quickly also get a qualifying score over 10%

## ■ BellKor

- Continue to get small improvements in their scores
- Realize that they are in direct competition with Ensemble

## ■ Strategy

- Both teams carefully monitoring the leaderboard
- Only sure way to check for improvement is to submit a set of predictions
  - This alerts the other team of your latest score

- **Submissions limited to 1 a day**
  - Only 1 final submission could be made in the last 24h
- **24 hours before deadline...**
  - **BellKor** team member in Austria notices (by chance) that **Ensemble** posts a score that is slightly better than BellKor's
- **Frantic last 24 hours for both teams**
  - Much computer time on final optimization
  - Carefully calibrated to end about an hour before deadline
- **Final submissions**
  - **BellKor** submits a little early (on purpose), 40 mins before deadline
  - **Ensemble** submits their final entry 20 mins later
  - ....and everyone waits....

# 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
------	-----------	-----------------	---------------	------------------

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

1	<a href="#">BellKor's Pragmatic Chaos</a>	0.8567	10.06	2009-07-26 18:18:28
2	<a href="#">The Ensemble</a>	0.8567	10.06	2009-07-26 18:38:22
3	<a href="#">Grand Prize Team</a>	0.8562	9.88	2009-07-10 21:24:00
4	<a href="#">Opera Solutions and Vandelay United</a>	0.8588	9.84	2009-07-10 01:12:31
5	<a href="#">Vandelay Industries !</a>	0.8591	9.81	2009-07-10 00:32:20
6	<a href="#">PragmaticTheory</a>	0.8594	9.77	2009-06-24 12:06:56
7	<a href="#">BellKor in BigChaos</a>	0.8601	9.70	2009-05-13 08:14:09
8	<a href="#">Dace</a>	0.8612	9.59	2009-07-24 17:18:43
9	<a href="#">Feeds2</a>	0.8622	9.48	2009-07-12 13:11:51
10	<a href="#">BigChaos</a>	0.8623	9.47	2009-04-07 12:33:59
11	<a href="#">Opera Solutions</a>	0.8623	9.47	2009-07-24 00:34:07
12	<a href="#">BellKor</a>	0.8624	9.46	2009-07-26 17:19:11

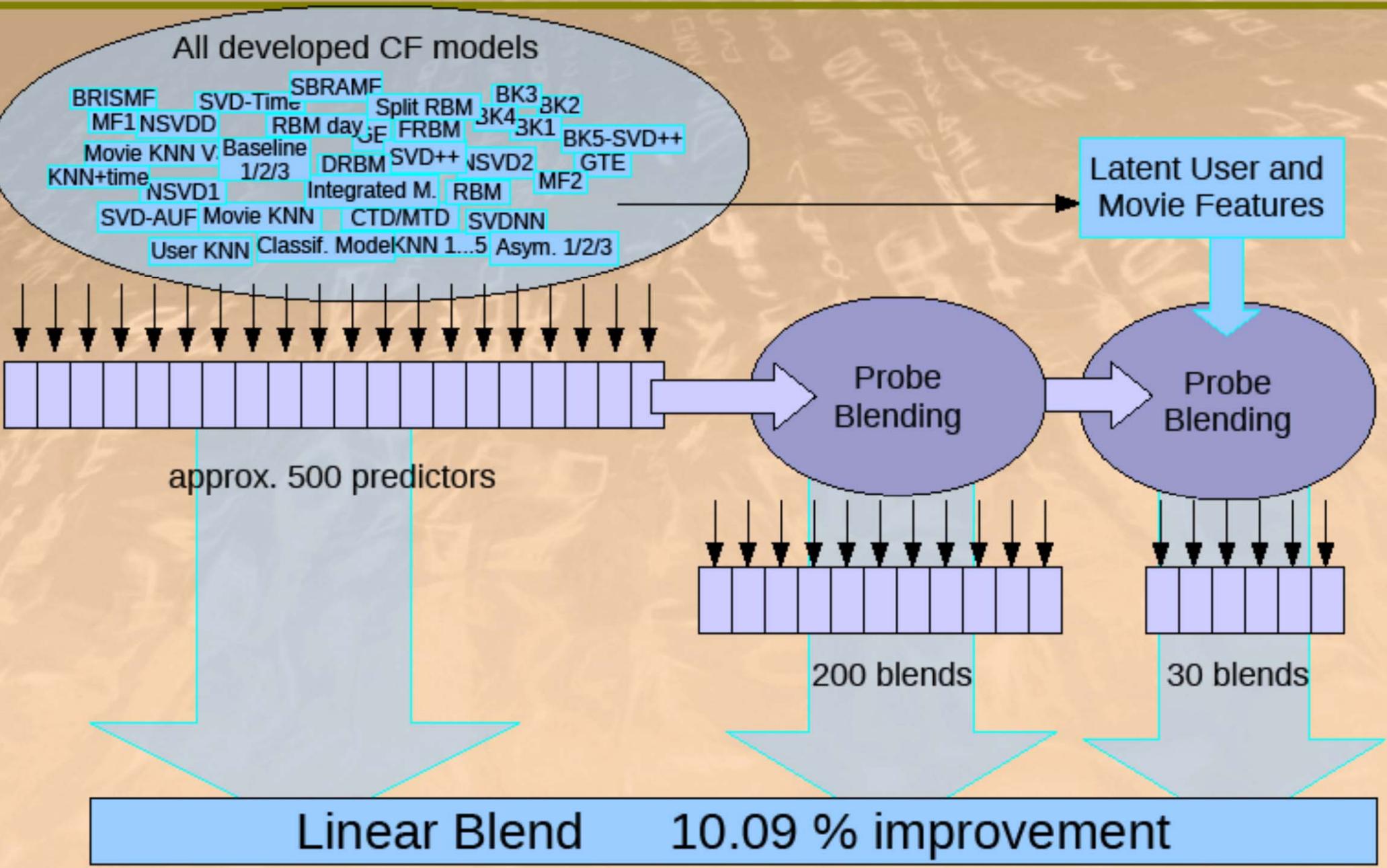
**Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos**

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



## The big picture

# Solution of BellKor's Pragmatic Chaos



CASEY JOHNSTON, ARS TECHNICA BUSINESS 04.16.12 08:20 AM

## SHARE

SHARE  
13

TWEET



COMMENT



EMAIL

# NETFLIX NEVER USED ITS \$1 MILLION ALGORITHM DUE TO ENGINEERING COSTS

Rank	Team Name	Best Score	Improvement	Last Submit Time
No Grand Prize candidates yet				
1	PragmaticTheo	0.8884	9.78	2009-09-15 01:04:47
2	Believe In BigChase	0.8890	9.71	2009-05-13 08:14:06
3	GrandPrizeTeam	0.8893	9.68	2009-05-32 08:20:24
4	Coda	0.8894	9.68	2009-04-22 08:57:03
5	BigChase	0.8813	9.47	2009-09-15 18:03:55
Previous Prize, 2008 - RMSE = 0.9416 - Winning Team: Believer In BigChase				
6	Believer	0.8820	9.48	2009-09-17 13:41:48
7	Groth	0.8834	9.25	2009-04-22 18:31:32
8	Omega Solutions	0.8840	9.18	2009-09-09 02:24:53
9	Winton	0.8840	9.18	2009-09-17 12:47:27
10	SmartDeepCachingYou	0.8841	9.18	2009-09-02 17:09:21
11	Coda	0.8842	9.17	2009-09-12:23:04:25
12	mapa2	0.8842	9.17	2009-09-01 02:35:00
13	Xiangliang	0.8842	9.17	2009-09-13:11:35:35
14	Fredzie	0.8847	9.11	2009-09-19 22:21:18
15	JeffLackInAPlane	0.8850	9.68	2009-05-24 18:02:54
16	Team ESP	0.8853	9.05	2009-09-05:25:25
17	getachewabu	0.8854	9.04	2009-09-05 18:10:03
18	NewHedgeTeam	0.8857	9.01	2009-05-13 07:30:22
19	J.Donna.Bu	0.8859	9.00	2009-03-11 08:41:54
20	VAMPIRE INDUSTRIES	0.8858	9.00	2009-09-11 08:43:14

Netflix awarded a \$1 million prize to a developer team in 2009 for an algorithm that increased the accuracy of the company's recommendation engine by 10 percent. But it doesn't use the million-dollar code, and has no plans to implement it in the future, Netflix announced on its blog

# What the Failed \$1M Netflix Prize Says About Business Advice



Ryan Holiday, CONTRIBUTOR

An inside look at the hidden side of marketing, PR and strategy. [FULL BIO ▾](#)

Opinions expressed by Forbes Contributors are their own.

## Recommended by Forbes

MOST POPULAR	TRENDING ON LINKEDIN	Forbes BrandVoice
The 10 Most Dangerous U.S. Cities	Then Silence - - How Should I Follow Up?	A Degree Really Required For Entry?



Netflix (Photo credit: Wikipedia)

If you've picked up a business book in the last four years, you undoubtedly heard it: the endless championing of [Netflix's \\$1 Million Prize](#).

It made for great "the future of business" fodder. A technology company crowdsources its R&D by ponying up \$1,000,000 to anyone who can improve their recommendation engine. [Netflix](#), stickler for user-experience as it is, was willing to pay all this money even if it meant making their algorithm just 10% better. Amazing! If you were looking for an anecdote to cheerlead Web 2.0 to, this was almost too good to be true.

DIRECTV  
NOW

firetv

**GET AN AMAZON FIRE TV ON US**

WHEN YOU PREPAY 2 MONTHS OF DIRECTV NOW

NO ANNUAL CONTRACT

**GET IT NOW**

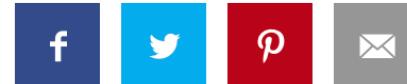
# THE NETFLIX PRIZE: HOW A \$1 MILLION CONTEST CHANGED BINGE-WATCHING FOREVER



By DAN JACKSON

Published On 07/07/2017

@danielvjackson



```
STRUCT GROUP_INFO INIT_GROUPS = { .USAGE = ATOMIC_INIT(2) };

STRUCT GROUP_INFO *GROUPS_ALLOC(INT GIDSETSIZE){

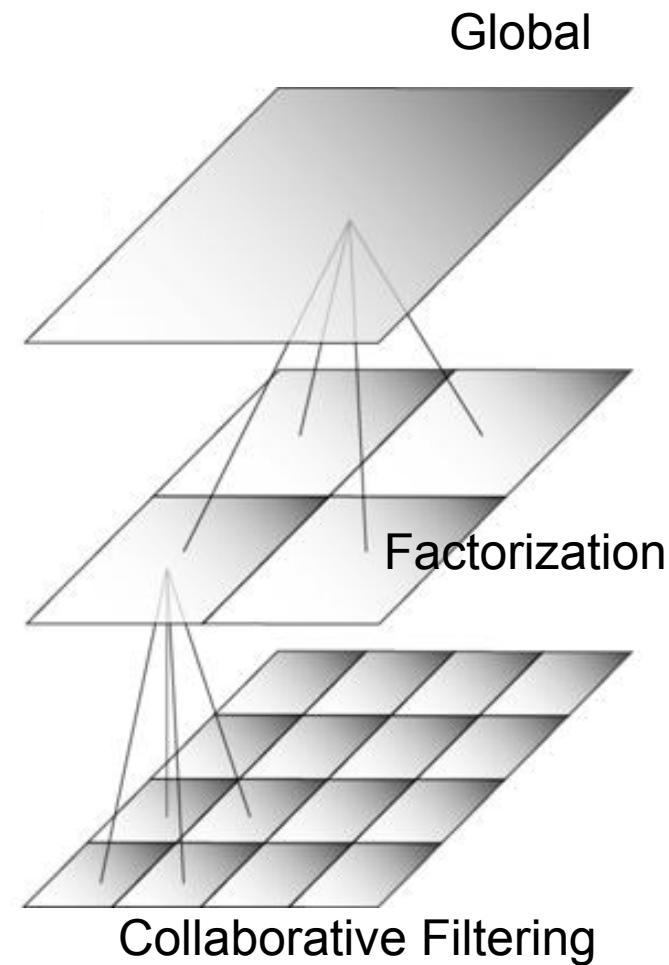
    STRUCT GROUP_INFO *GROUP_INFO;
    INT NBLOCKS;
    INT I;

    NBLOCKS = (GIDSETSIZE + NGROUPS_PER_BLOCK - 1) / NGROUPS_PER_BLOCK;
    /* MAKE SURE WE ALWAYS ALLOCATE AT LEAST ONE INDIRECT BLOCK POINTER */
    NBLOCKS = NBLOCKS ? : 1;
    GROUP_INFO = KMALLOC(SIZEOF(*GROUP_INFO) + NBLOCKS*SIZEOF(GID_T *), GFP_USER);
    IF (!GROUP_INFO)
```



# BellKor Recommender System

- The winner of the Netflix Challenge
- Multi-scale modeling of the data
  - **Global:**
    - Overall deviations of users/movies
  - **Factorization:**
    - Addressing “regional” effects
  - **Collaborative Filtering:**
    - Extract local patterns



# Modeling Local and Global Effects

- **Global:**

- Mean movie rating: **3.7 stars**
- *The Sixth Sense* is **0.5** stars above average
- Joe rates **0.2** stars below average
- => Baseline estimate: Joe will rate *The Sixth Sense* **4 stars**
- **Local Neighborhood (Collaborative Filtering)**
  - Joe didn't like related movie *Signs*
  - =>Final estimate: Joe will rate *The Sixth Sense* **3.8 stars**

# Recap: Collaborative Filtering

- Earliest and most popular **collaborative filtering method**
- Derive unknown ratings from those of “similar” movies (item-item variant)
- Define **similarity measure  $s_{ij}$**  of items  $i$  and  $j$
- Select  $k$ -nearest neighbors, compute the rating
  - **$N(i; x)$ :** items most similar to  $i$  that were rated by  $x$

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

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

# Modeling Local and 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  $r_{xi}$

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

$\mu$  = overall mean rating

$b_x$  = rating deviation of user  $x$

= (avg. rating of user  $x$ ) –  $\mu$

$b_i$  = (avg. rating of movie  $i$ ) –  $\mu$

## 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