

Information Storage and Retrieval

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

Recommenders
21 March 2017

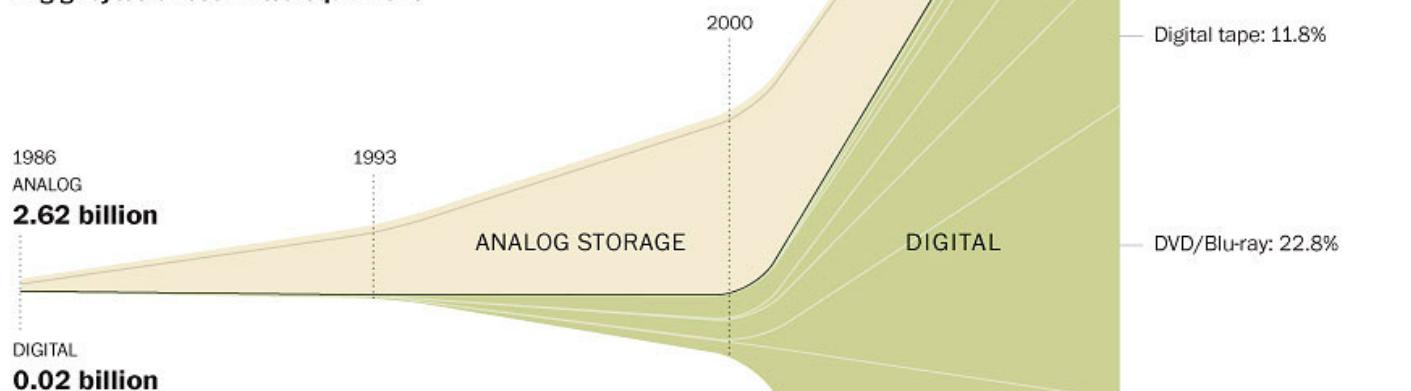
Today

- Recommendation systems: Intro
- Collaborative Filtering
 - Finding nearest neighbors
 - How to aggregate ratings?

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.)

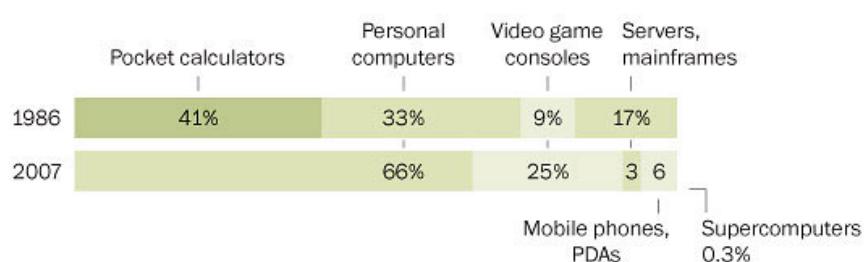
In gigabytes or estimated equivalent



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

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

Apple: iTunes Now Has 20M Songs; Over 16B Downloads

Posted Oct 4, 2011 by Leena Rao (@leenarao)

0 [Share](#) 6 [Share](#) 0 [Tweet](#) 147 ▾

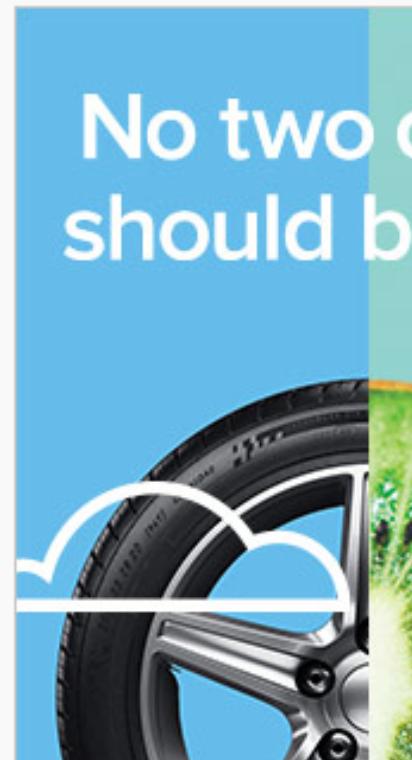


At today's Apple iPhone event, newly appointed CEO Tim Cook [revealed some staggering numbers](#) on Apple's music downloads and songs in iTunes. iTunes now offers 20 million songs, which is up from 200,000 when iTunes was first launched eight years ago.

iTunes is the number one music store in the world with over 16 billion songs downloaded. It looks like Apple has seen about a billion downloads in the past four months, as the company revealed [15 billion](#) song downloads in June. At the time, Apple revealed that it 225 million credit card accounts listed with iTunes.

Interestingly, streaming service Spotify has [around 15 million songs](#) available, which isn't that much less.

ADVERTISEMENT



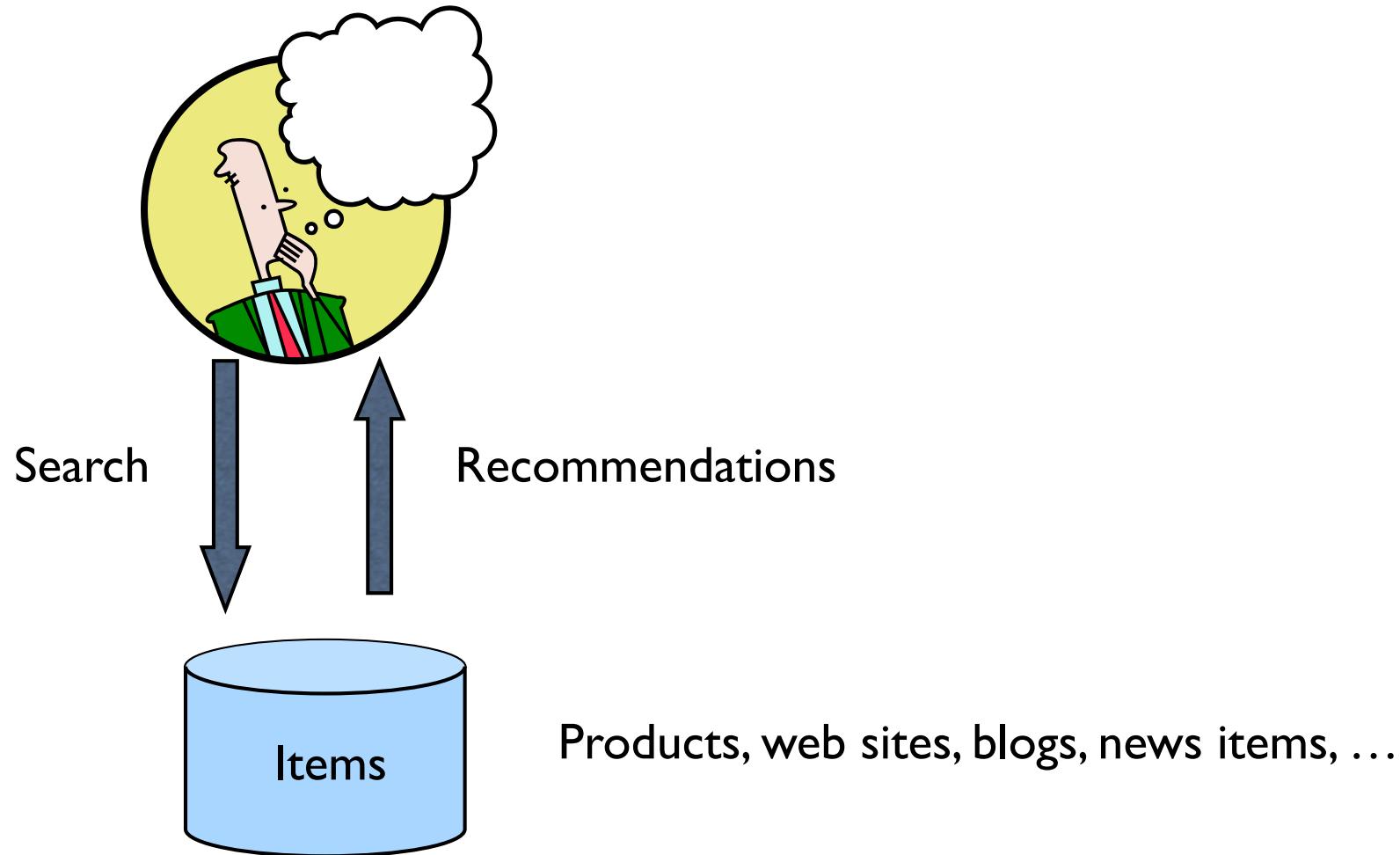


Statistics

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

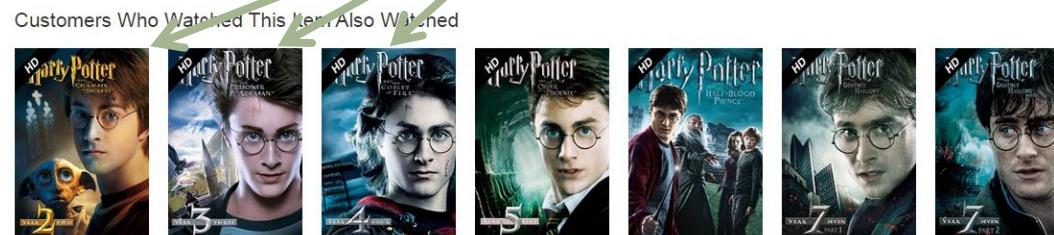
Recommendations for You in Amazon Instant Video [See more](#)



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 stars - 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



Item	Price Range
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

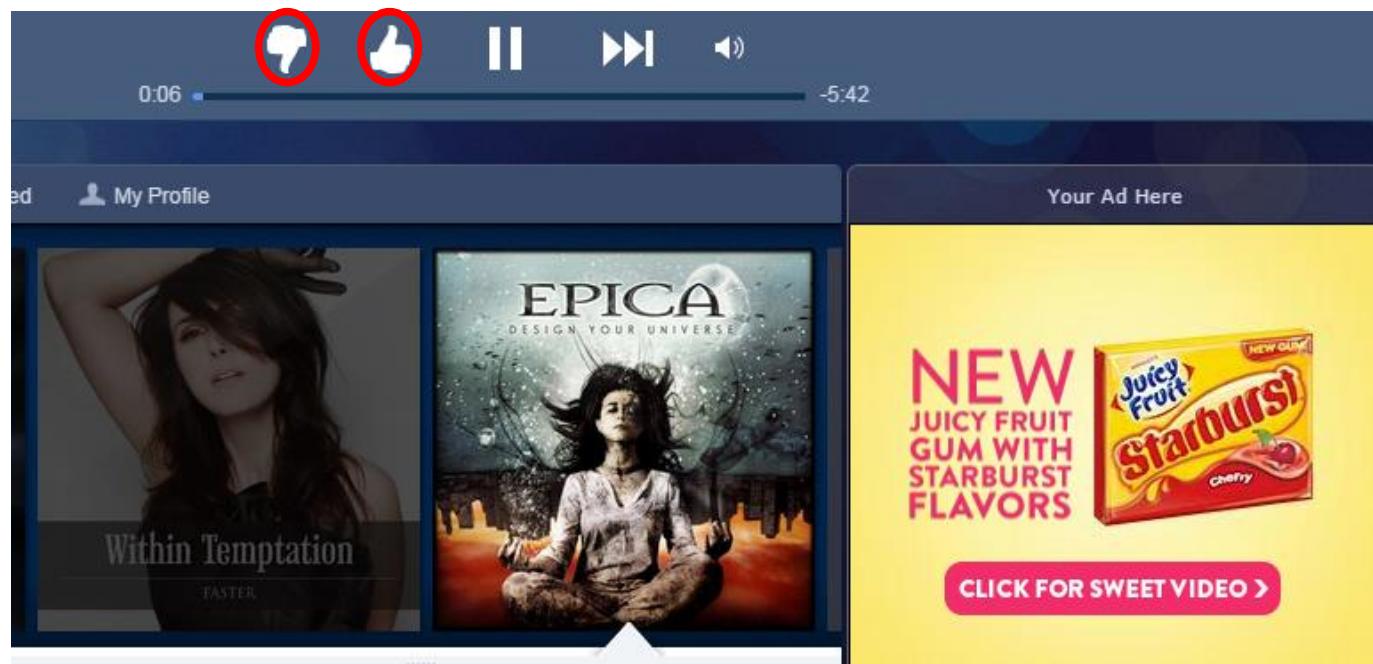
Customers Who Bought This Item Also Bought



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?

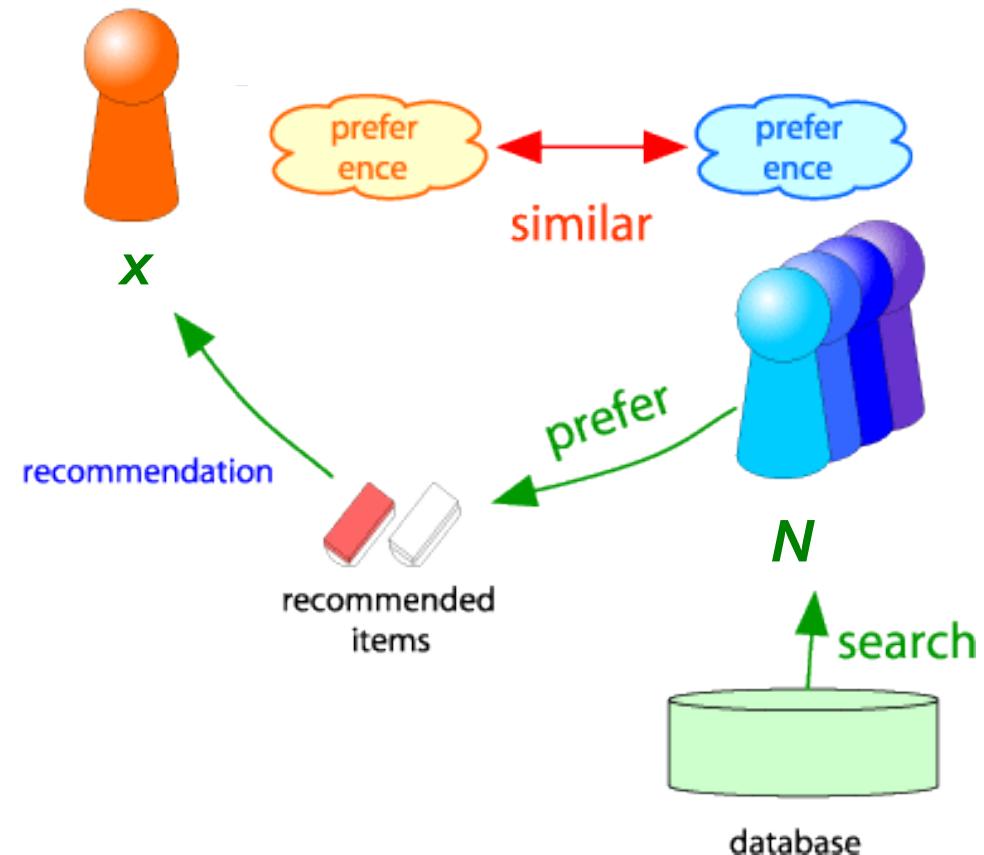
Extrapolating Utilities

- 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

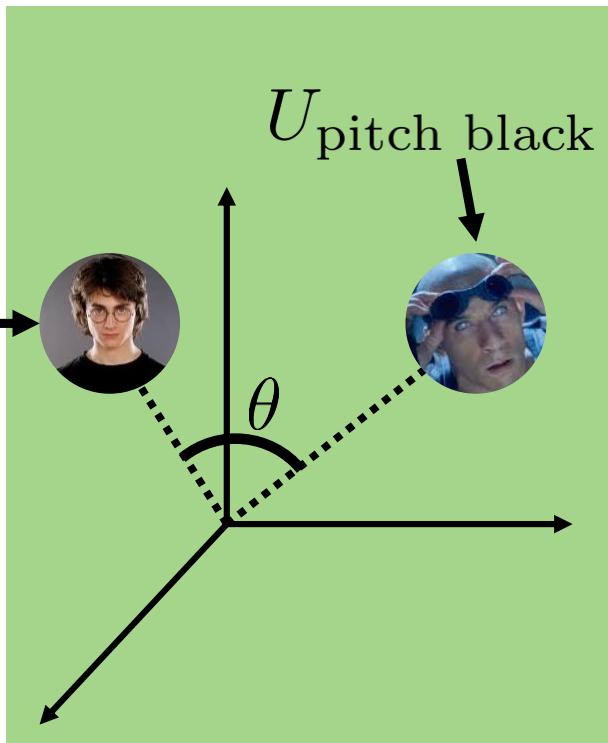
(or if two items were purchased by the same set of users)

→ 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**

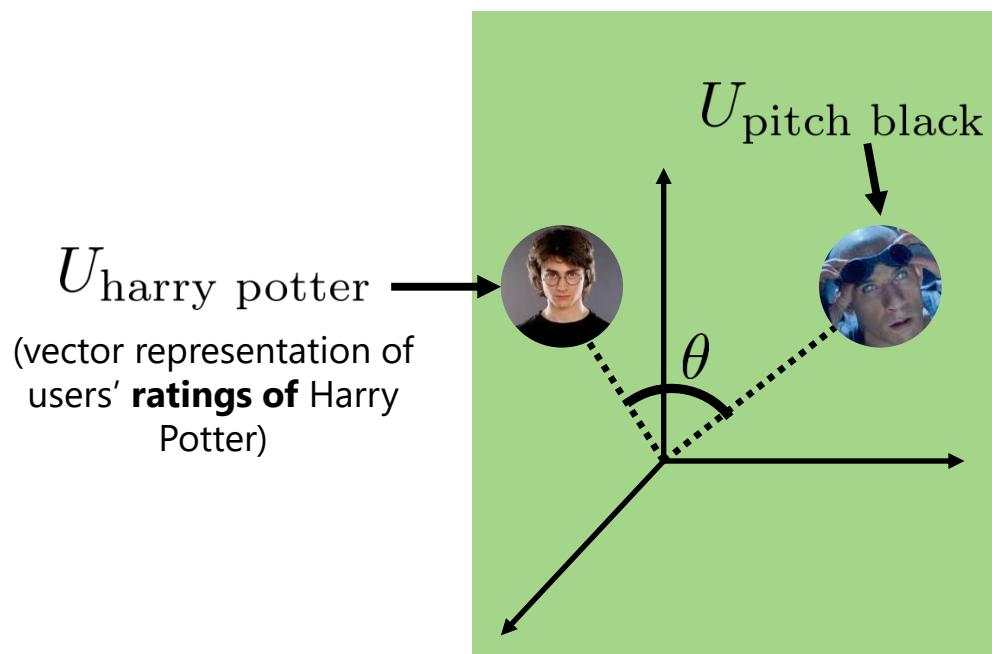
didn't buy

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 histories.

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**

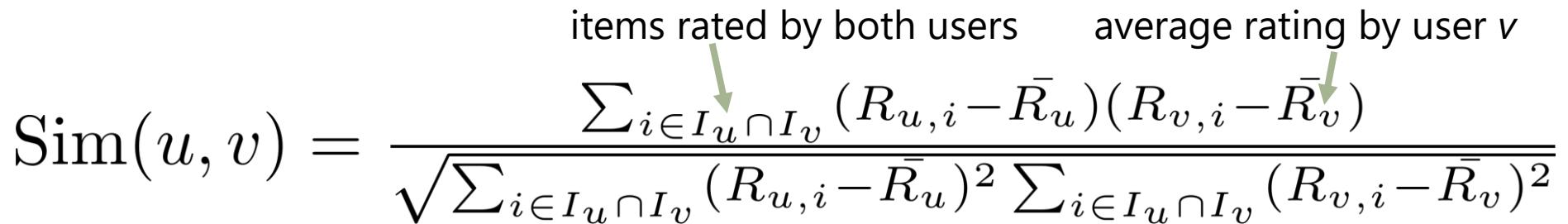
Pearson correlation

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