

Recommendation Systems

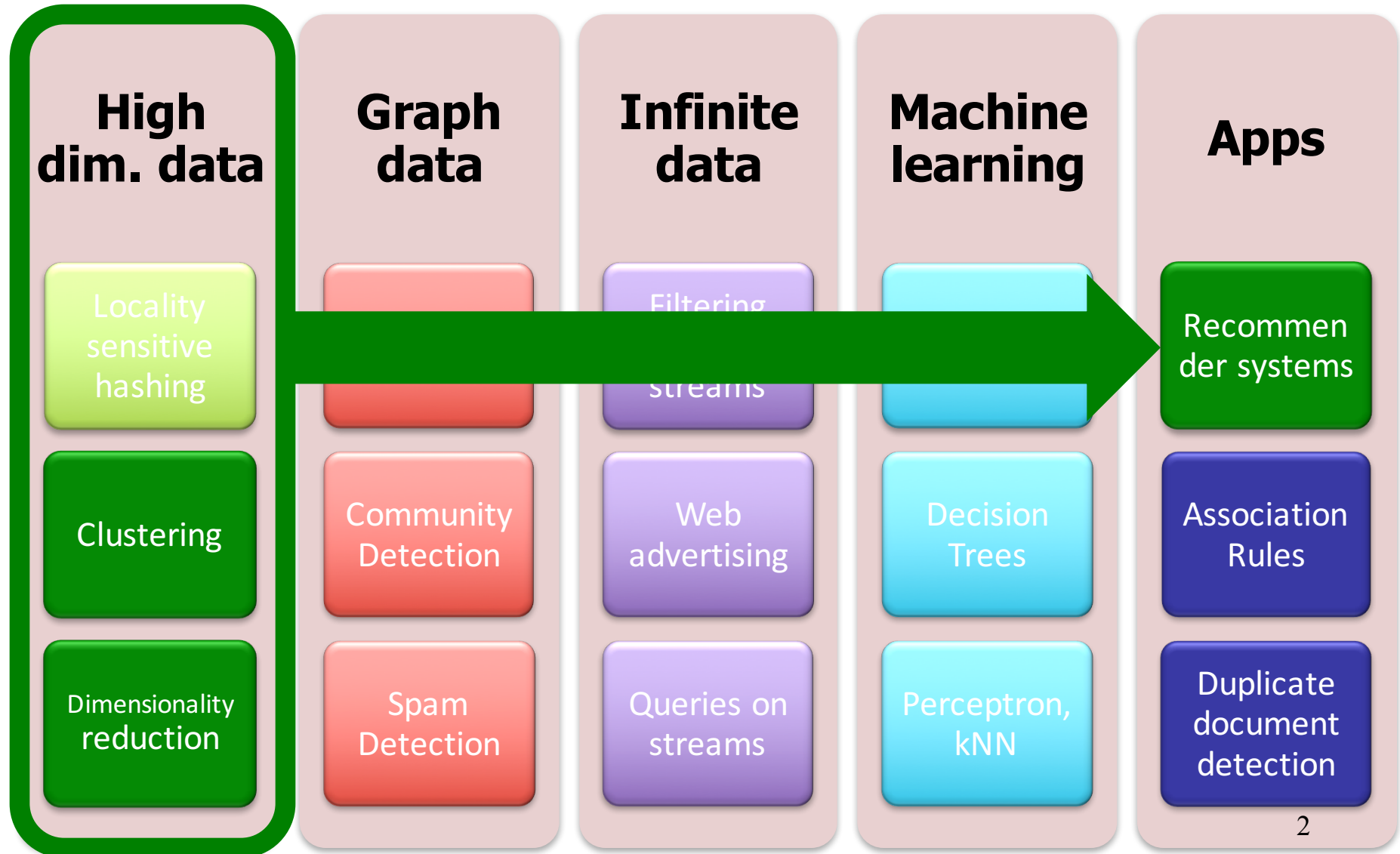
Content-Based Recommendations

Collaborative Filtering

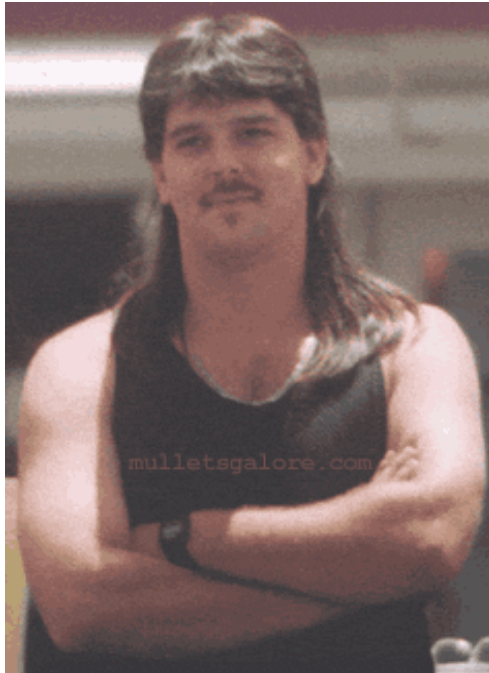
Hybrid Systems

Anna Farzindar, Ph.D.

High Dimensional Data



Example: Recommender Systems



◆ Customer X

- Buys Metallica CD
- Buys Megadeth CD



◆ Customer Y

- Does search on Metallica
- Recommender system suggests Megadeth from data collected about customer X

Recommendations



Examples:

amazon.com.



movielens
helping you find the *right* movies

last.fm™
the social music revolution

Google
News



Motivation: The Long Tail

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

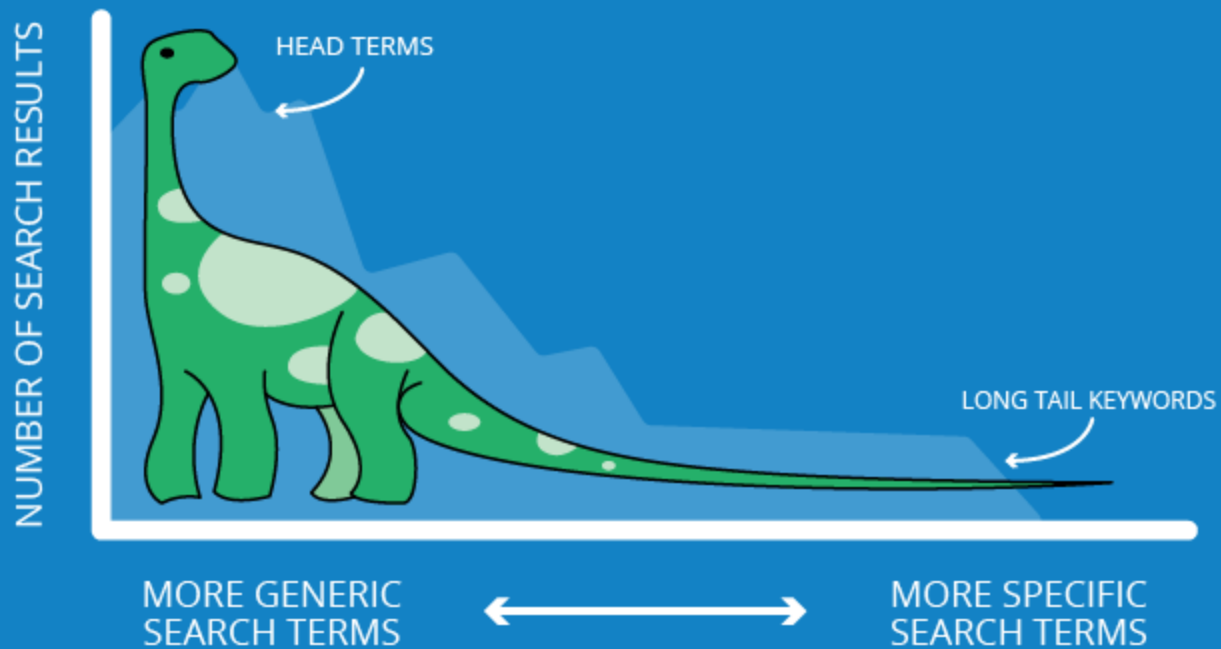
◆ More choice necessitates better filters

➤ Recommendation engines

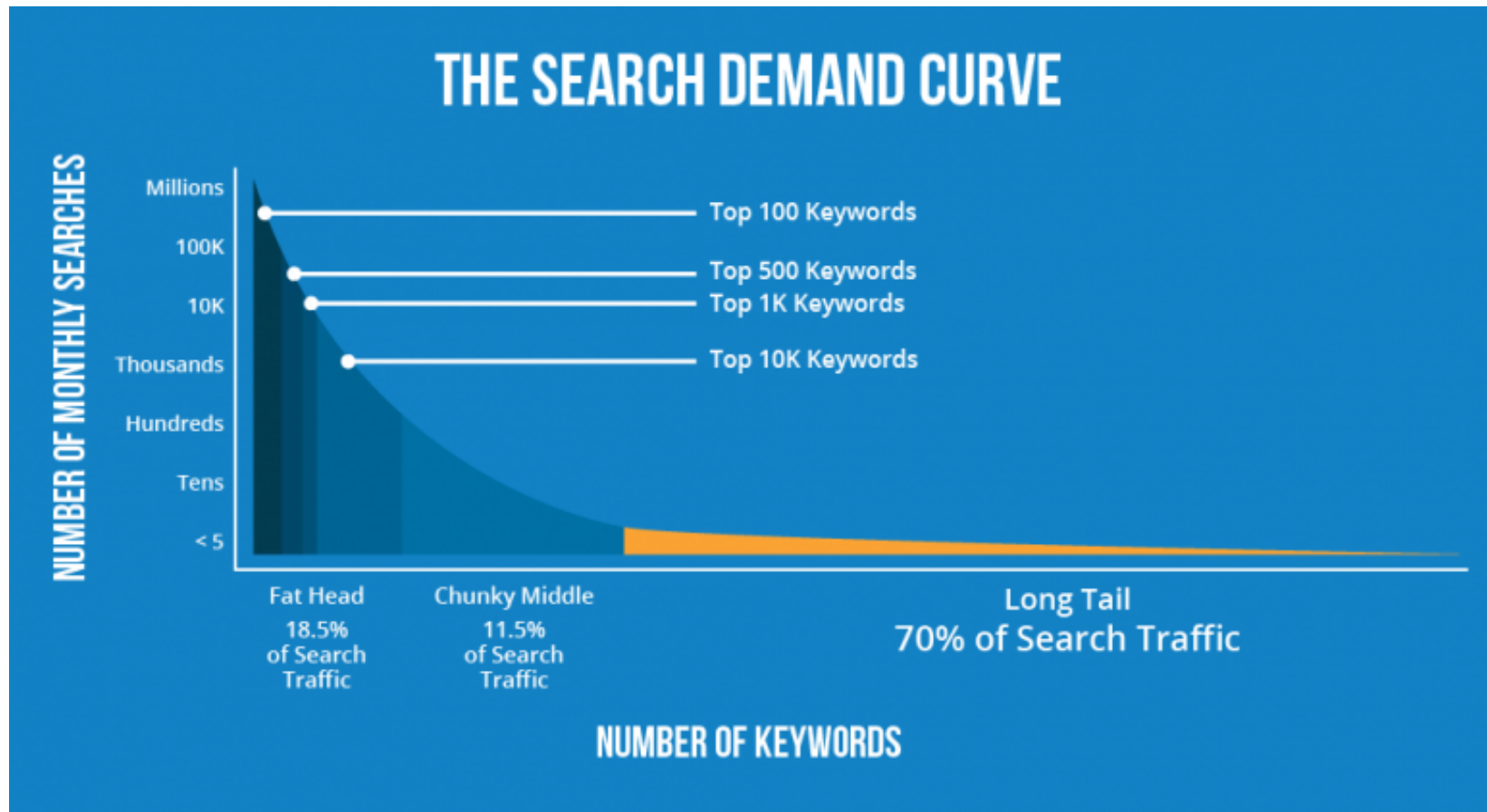
➤ How **Into Thin Air** made **Touching the Void** a bestseller: <http://www.wired.com/wired/archive/12.10/tail.html>

The Long Tail

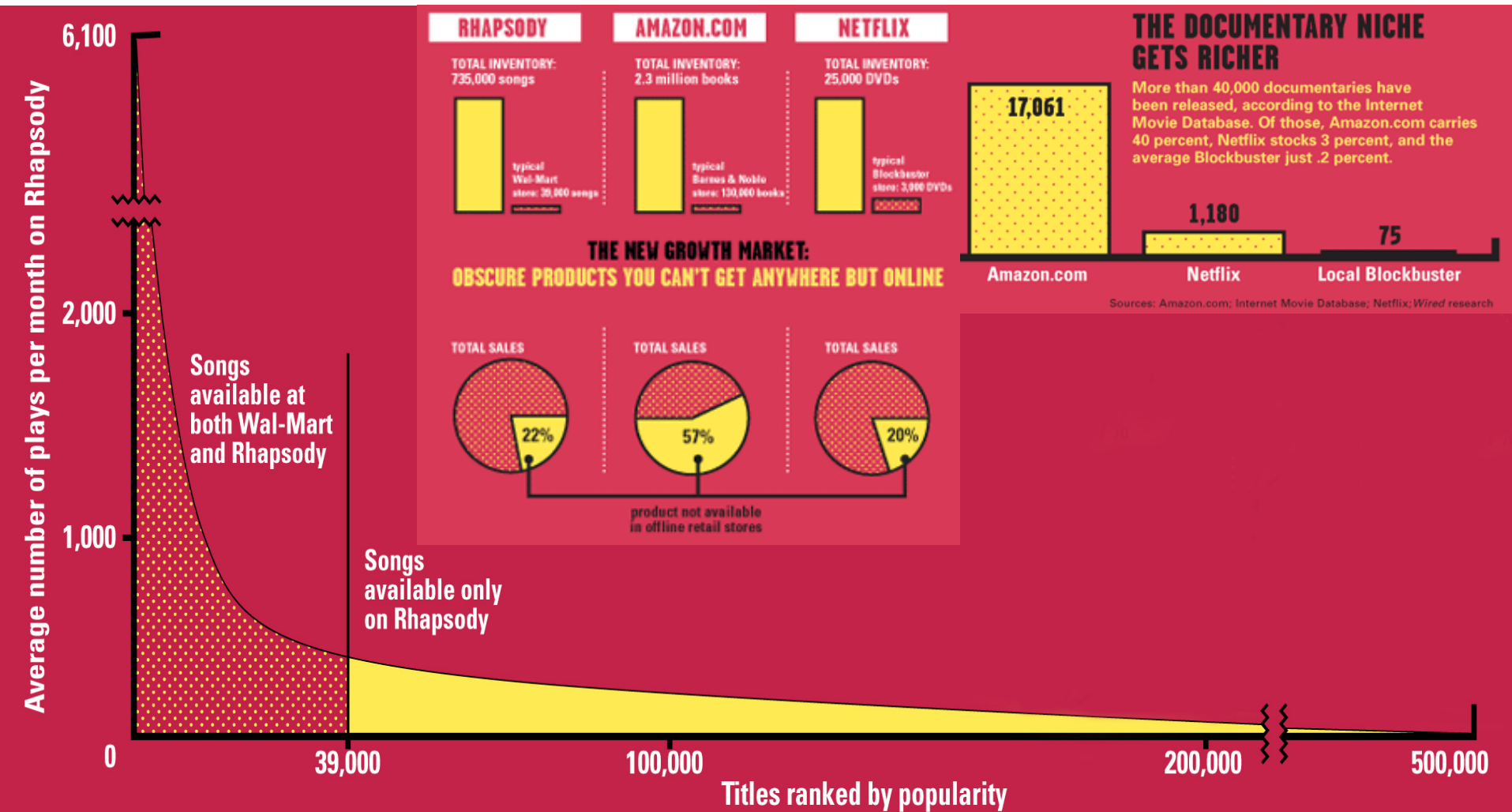
LONG TAIL KEYWORDS



The Long Tail



The Long Tail



Change in thinking compared with online stores

- ◆ "What **percentage** of the top 10,000 titles in any online media store (Netflix, iTunes, Amazon, or any other) will **rent or sell** at least once a month?"
- ◆ Most people guess 20 percent
 - 80-20 rule, also known as Pareto's principle (1896)
 - Only 20 percent of major studio films, TV shows, books, etc. will be hits
- ◆ The right answer: 99 percent
 - Demand for nearly every one of those top 10,000 titles.

Counterintuitive to old way of thinking

- ◆ The 20 percent rule in the entertainment industry is about ***hits***, not sales of any sort
 - **Hit-driven mindset**: think that if something isn't a **hit**, it won't make money
 - Makes sense with **scarce shelf space** in a retail store
 - iTunes, Amazon, and Netflix: discovered that "**misses**" **usually make money, too**
 - And because **there are so many more of them**, that money can add up quickly to a huge new market
- ◆ Industry has a poor sense of what people want
 - Turns out that people like a wide range of things when they are easily available.

Rules of thumb

- ◆ **Rule 1: Make everything available**
 - Embrace under-served markets, niches (e.g., obscure video genres)
 - What matters is not where customers are, or how many of them are seeking a particular title, but that some number of them exist, anywhere
 - As a result, almost anything is worth offering on the chance it will find a buyer.

Rules of thumb (cont.)

◆ Rule 2: Lower costs

- Have taken away the unnecessary **costs of the retail channel**: manufacturing, distribution, and retail overheads
- Leaves the costs of finding, making, and marketing content
- Ensure that the people on the creative and business side still make money

◆ Rule 3: Help users find new content

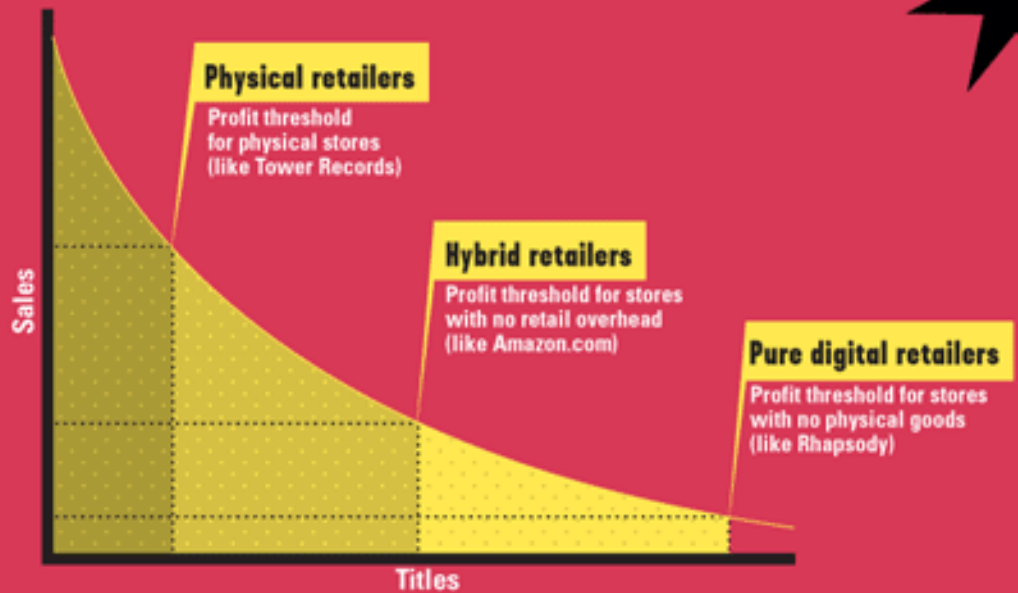
- Edited recommendations
- Content-based recommendations
- Collaborative filtering: uses browsing and purchasing patterns of users to guide those who follow them
- "Customers who bought this also bought ..."
- **Use recommendations to drive demand down the Long Tail.**

Physical vs. Online

THE BIT PLAYER ADVANTAGE

Beyond bricks and mortar there are two main retail models – one that gets halfway down the Long Tail and another that goes all the way. The first is the familiar hybrid model of Amazon and Netflix, companies that sell physical goods online. Digital catalogs allow them to offer unlimited selection along with search, reviews, and recommendations, while the cost savings of massive warehouses and no walk-in customers greatly expands the number of products they can sell profitably.

Pushing this even further are pure digital services, such as iTunes, which offer the additional savings of delivering their digital goods online at virtually no marginal cost. Since an extra database entry and a few megabytes of storage on a server cost effectively nothing, these retailers have no economic reason not to carry *everything* available.



“IF YOU LIKE BRITNEY, YOU’LL LOVE ...”

Just as lower prices can entice consumers down the Long Tail, recommendation engines drive them to obscure content they might not find otherwise.



Source: Amazon.com

<http://www.wired.com/wired/archive/12.10/tail.html>

Types of Recommendations

Types of Recommendations

◆ 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**
- ◆ Users have preferences for certain items
- ◆ **Want to extract preferences from data**
- ◆ **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

- For each user-item pair, value represents degree of preference of that user for that item (e.g., rating)
- Matrix is sparse (most entries unknown)

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

Predictions

- ◆ Goal of a Recommendation System is to **predict the blanks in the utility matrix**
 - Would Alice like Pirates?
 - Would David like Avatar?
- ◆ Not necessary to predict every blank entry
- ◆ **Want to discover some entries in each row that are likely to be high**
- ◆ **Usually recommend a few items the user should value highly**
 - **Most likely to generate additional revenue.**

Key Problems

◆ (1) Gathering “known” ratings for matrix

- How to collect the data in the utility matrix,

◆ (2) Extrapolate unknown ratings from the known ones

- Mainly interested in **high unknown ratings**
 - We are not interested in knowing what you don't like but what you like
 - To generate revenue,

◆ (3) Evaluating extrapolation methods

- How to measure success/performance of recommendation methods.

(1) Gathering Ratings

◆ Explicit

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

◆ Implicit

- Learn ratings from user actions
 - E.g., purchase implies high rating
- What about low ratings?

(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

◆ **How to extrapolate missing entries?**

Three Approaches to Recommendation Systems

◆ 1) Content-based

- Use characteristics of an item
- Recommend items that have similar content to items user liked in the past
- Or items that match pre-defined attributes of the user

◆ 2) Collaborative filtering

- Build a model from a user's past behavior (items previously purchased or rated) and similar decisions made by other users
- Use the model to predict items that the user may like
- Collaborative: suggestions made to a user utilize information across the entire user base

◆ 3) Hybrid approaches.

Content-based Recommender Systems

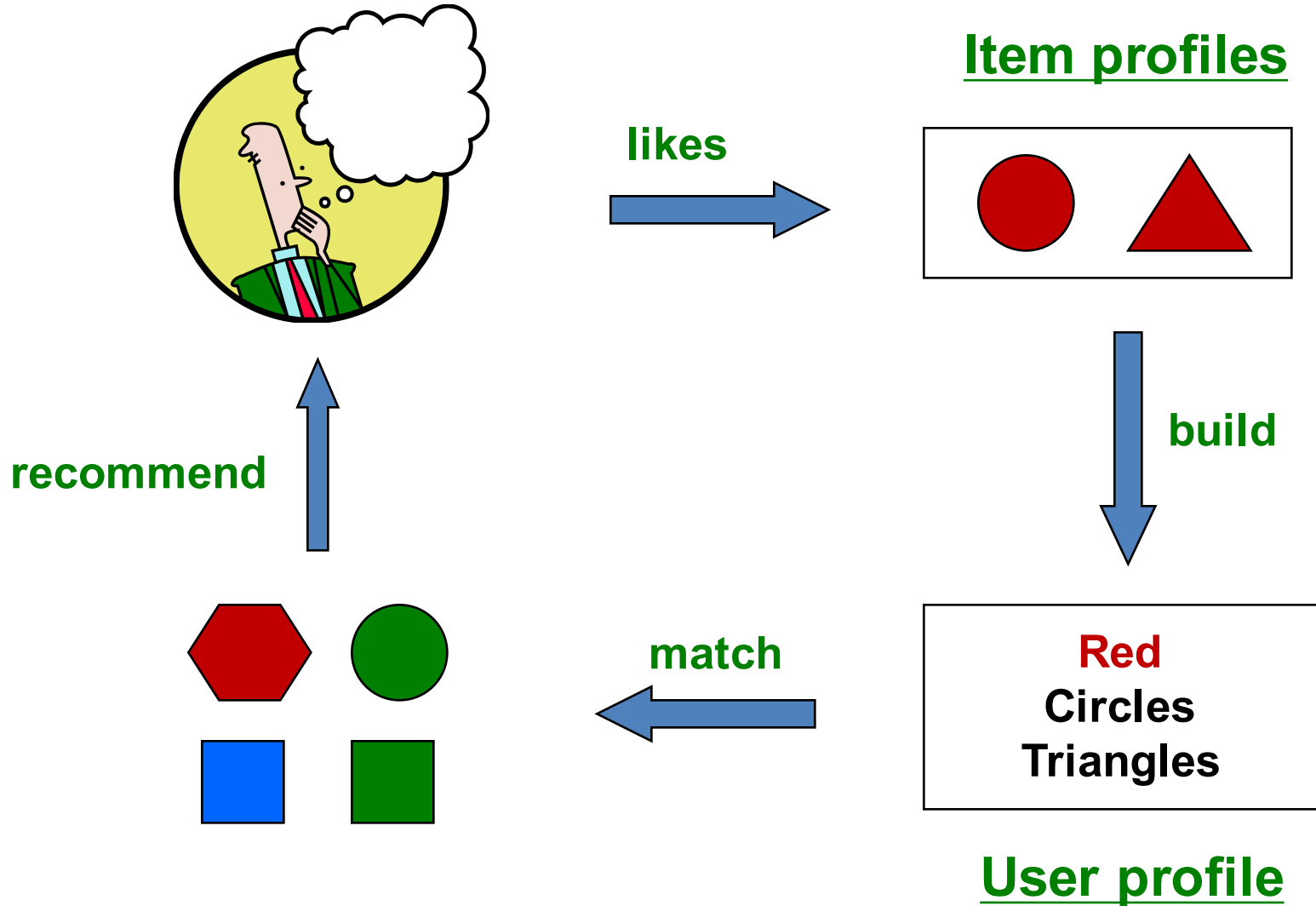
Content-based Recommendations

- ◆ **Main idea: Recommend items to customer x that are similar to previous items rated highly by x**
 - *Requires characterizing the content of items in some way*

Examples:

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

Plan of Action



General Strategy for Content-Based Recommendations

◆ Construct item profiles

- Explicit features in a database, discovering features in documents, Tags
- **Create vectors representing items**
 - Boolean vectors indicate occurrence of high TF.IDF word
 - Numerical vectors might contain ratings

◆ Construct user profiles

- **Create vectors with same components that describe user's preferences**

◆ Recommend items to users based on content

- **Calculate cosine distance between item and user vectors**
- Classification algorithms.

ITEM PROFILES

Item Profiles

- ◆ For each item, create an **item profile**

- ◆ **Profile is a set (vector) of features**

- **Movies:** screenwriter, title, actor, director,...
- **Text:** Set of “important” words in document

- ◆ **Example 9.2**

- Features of movies are set of actors and average rating
- Each movie has 5 actors, two in both movies
- Average ratings 3 and 4 (with unknown scaling factor alpha)
- **Must scale non-boolean components so they are not dominant or irrelevant**

0	1	1	0	1	1	0	1	3α
1	1	0	1	0	1	1	0	4α

Cosine Distance (Section 3.5.4)

- ◆ **Cosine distance is used in spaces with dimensions, including Euclidean spaces**

- Where points are vectors with integer or Boolean components

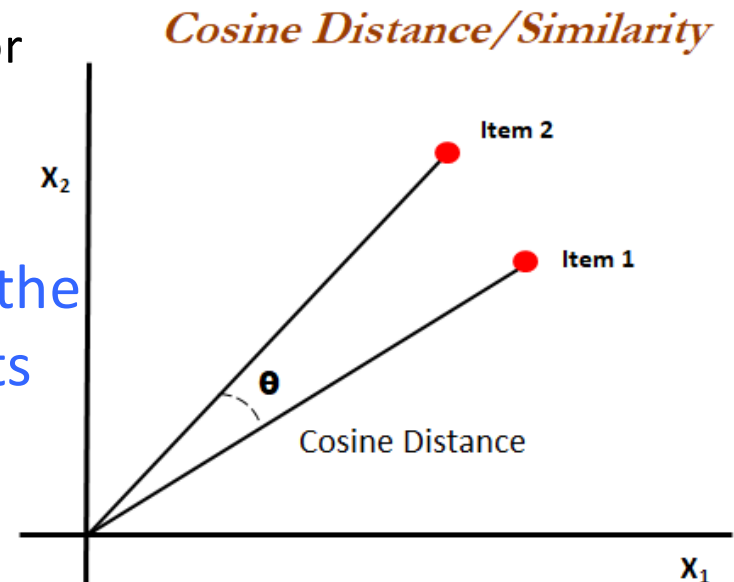
- ◆ Points thought of as directions

- ◆ **Cosine distance between 2 points is the angle that the vectors to those points make**

- Range from 0 to 180 degrees

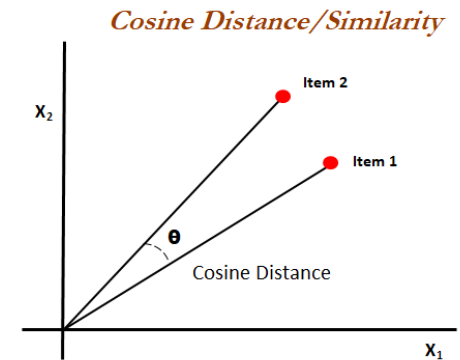
- ◆ **First compute cosine of the angle between vectors x and y**

- ◆ Then apply **arc-cosine** function to translate to **0-180 degrees**: the cosine distance.



Cosine Similarity

(Used to Calculate Cosine Distance)



$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

◆ Dot product $x \cdot y$ divided by Euclidean distance of x and y from origin (Section 3.5.2)

➤ Dot product of vectors:

$$[x_1, x_2, \dots, x_n] \cdot [y_1, y_2, \dots, y_n] \text{ is } \sum_{i=1}^n x_i y_i.$$

➤ Euclidean distance of two vectors x, y

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Example 3.1.3: Cosine Distance

- ◆ $x = [1, 2, -1]$ $y = [2, 1, 1]$
- ◆ **Dot product** $x \cdot y = 1 \times 2 + 2 \times 1 + (-1) \times 1 = 3$
- ◆ **Euclidean distance** of two vectors x, y

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- ◆ **Euclidean distance of a vector from the origin: vector for origin is all zeros**
- ◆ Distance of x from origin is: $\sqrt{1^2 + 2^2 + (-1)^2} = \sqrt{6}$.
- ◆ Distance of y from origin also $\sqrt{6}$
- ◆ Cosine of angle from x to y is $3/(\sqrt{6}\sqrt{6})$ or $1/2$.
- ◆ **Arccosine** of $1/2$ is 60 degrees
- ◆ That is the **cosine distance between x and y .**

Back to Item Profiles Example (9.2)

- ◆ Features of movies are set of actors and average rating (scaled)

$$\begin{array}{cccccccccc} 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 3\alpha \\ 1 & 1 & 0 & 1 & 0 & 1 & 1 & 0 & 4\alpha \end{array}$$

- ◆ Dot product of these two item vectors is:

$$[x_1, x_2, \dots, x_n] \cdot [y_1, y_2, \dots, y_n] \text{ is } \sum_{i=1}^n x_i y_i.$$

- Or $2 + 12\alpha^2$

- ◆ Distance from origin using

$$d([x_1, x_2, \dots, x_n], [y_1, y_2, \dots, y_n]) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- ◆ Lengths of vectors (distance from origin) are:

$$\sqrt{5 + 9\alpha^2} \text{ and } \sqrt{5 + 16\alpha^2}.$$

- ◆ Cosine of angle between vectors is $\frac{2 + 12\alpha^2}{\sqrt{25 + 125\alpha^2 + 144\alpha^4}}$

- ◆ Scaling factor alpha affects how similar items are.

ITEM PROFILES BASED ON TEXTUAL CONTENT

Item Profiles based on Textual Content

- ◆ Much research on content-based recommendations focuses on **textual content**
 - Recommend items (web pages, books, movies) based on associated textual content
 - Descriptions, user reviews
- ◆ Can treat this as an **Information Retrieval task (IR)**
- ◆ **How to identify whether two documents are about similar things?**
- ◆ **How to pick important features of documents?**
- ◆ Want to **identify the significant words** in documents.

TF.IDF: Measure of Word Importance

- ◆ Classification of documents as being about similar things starts with finding significant words in those documents
- ◆ **Not most frequent words**
 - (the, and, a, ...) – called “stop words”
- ◆ **Not just rare words either**
- ◆ **Want concentration of useful words in just a few documents**
- ◆ Usual heuristic from text mining is **TF-IDF:**
(Term frequency * Inverse Doc Frequency)
- ◆ **Words with highest TF.IDF score are often the terms that best characterize the topic of a document**
- ◆ When constructing an item profile for Recommender system:
 - Term ... Feature
 - Document ... Item.

TF-IDF: From Section 1.3.1

(Term frequency * Inverse Document Frequency)

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

Term Frequency: $TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}}$

- ◆ Term frequency of term i in document j is **normalized**
- ◆ **Divide by maximum occurrences of any term in document j**
- ◆ Most frequent term has $TF=1$

n_i = number of docs that mention term i

N = total number of docs

Inverse Document Frequency: $IDF_i = \log_2(N/n_i)$

TF-IDF score: $w_{ij} = TF_{ij} \times IDF_i$

Item profile for a document = set of words with highest TF-IDF scores, together with their scores.

TF.IDF Example (Example 1.3)

$$TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}}$$
$$IDF_i = \log_2(N/n_i)$$

- ◆ Repository of $2^{20} = 1,048,576$ documents
- ◆ Suppose word w appears in $2^{10} = 1024$ documents
- ◆ **Inverse document frequency:**

$$IDF_w = \log_2(2^{20}/2^{10}) = \log_2(2^{10}) = 10 \text{ (logarithm scaled)}$$

- ◆ Consider document j in which **word w appears 20 times**
 - This is the maximum number of times any word appears in document j (after eliminating stop words)
 - So $TF_{wj} = 20/20 = 1$
 - **TF.IDF score for w in document j is 10**
- ◆ Consider document k where word **w appears once**
 - Maximum number of occurrences of any word in k is 20
 - So $TF_{wk} = 1/20$
 - **TF.IDF score for w in document k is $\frac{1}{2}$.**

Recommender Systems: Make Recommendations Based on Features of Documents

- ◆ **Want to suggest articles, pages, blogs a user might want to see**
- ◆ **Hard to classify items by topic**
- ◆ In practice, **try to identify words that characterize the topic of a document**
- ◆ **Eliminate stop words:** several hundred most common words
- ◆ **For remaining words, calculate the TF.IDF score for each word in the document**
- ◆ **The words with the highest TF.IDF scores characterize the document.**

Represent documents by a set of words

- ◆ Take as features of the document the n words with highest TF.IDF scores
 - Could pick same n for all documents
 - Or let n be fixed percentage of words in the document
 - Could also make all words with TF.IDF scores above a given threshold are part of feature set
- ◆ Documents then represented by set of words
- ◆ Expect these words to express subjects or main ideas of documents
- ◆ Then can measure the similarity of two documents using:
 - Cosine distance between the sets, treated as vectors
 - Jaccard distance (Ch. 3) between sets of word.

Document Similarity using Cosine Distance

- ◆ First compute cosine of the angle between vectors A and B

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

- ◆ Then apply arc-cosine function to translate to 0-180 degrees: the cosine distance.

Document Similarity using Cosine Distance

- ◆ Think of set of **high-TF.IDF words as a vector**, with one component for each possible word
- ◆ **Vector has 1 if word is in the set** for that document and 0 if not
- ◆ Between two documents, only a finite number of words among their two sets
- ◆ **Almost all components are 0**; do not affect dot product
- ◆ **Dot products** are size of **intersection of the two sets** of words
- ◆ **Lengths of vectors** are square roots of number of words in each set
- ◆ Cosine of angle between vectors: dot product divided by product of vector lengths.

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

Document Similarity: Two Kinds

◆ Chapter 3: find (nearly) identical documents

- **Lexical similarity:** Documents are similar if they contain large fraction of identical sequences of characters
- **Shingling:** convert documents to sets
- **Minhashing:** convert large sets to short signatures, preserving similarity
- **Locality-sensitive hashing (LSH):** Focus on parts of signatures likely to be from similar documents to **identify candidate pairs**

◆ For Recommendation Systems (Chapter 9):

- Interested in **occurrences of many important words in both documents** (even if little lexical similarity between documents)

Similar methodology:

- **High TF.IDF words form a vector**, with a component for each possible word set to 1 or 0 (*analogous to sets of shingles*)
- Based on a **distance measure** (Jaccard or cosine distance).

Another Option to Describe Item Content: Obtaining Item Profile Features from Tagging Systems

- ◆ Useful for content-based recommendations for **images**
- ◆ E.g., Flickr photosharing, Facebook, Instagram, Snapchat
- ◆ Users **enter words or phrases that describe items**
- ◆ **GPS information/geofilters**: e.g., automatically add location information when a photo is uploaded
- ◆ **Can use tags as a recommendation system**
 - E.g., if user retrieves or bookmarks pages with certain tags, **recommend other pages with same tags**
- ◆ Only works if users create tags or allow automatic geotagging.

USER PROFILES

General Strategy for Content-Based Recommendations

◆ Construct item profiles

- Source we discussed:
 - Explicit features in a database
 - Discovering features in documents
 - Tags,
- Create vectors representing items
 - Boolean vectors indicate occurrence of high TF.IDF word
 - Numerical vectors might contain ratings,

◆ Construct user profiles

- Create vectors with same components that describe user's preferences,

◆ Recommend items to users based on content

- Calculate cosine distance between item and user vectors
- Classification algorithms.

Example 1: Boolean Utility Matrix

- Items are movies, only feature is “Actor”
 - Item profile: vector with 0 or 1 for each Actor
- Suppose user x has watched 5 movies
 - 2 movies featuring actor A
 - 3 movies featuring actor B
- User profile = mean of item profiles
 - Feature A's weight = $2/5 = 0.4$
 - Feature B's weight = $3/5 = 0.6$

Example 2: Star Ratings

- Same example, 1-5 star ratings
 - Actor A's movies rated 3 and 5
 - Actor B's movies rated 1, 2 and 4
- Useful step: Normalize ratings by subtracting user's mean rating (3)
 - Actor A's normalized ratings = 0, +2
 - Profile weight = $(0 + 2)/2 = 1$
 - Actor B's normalized ratings = -2, -1, +1
 - Profile weight = $-2/3$

User Profiles (Examples 9.3 and 9.4)

◆ Construct user profiles

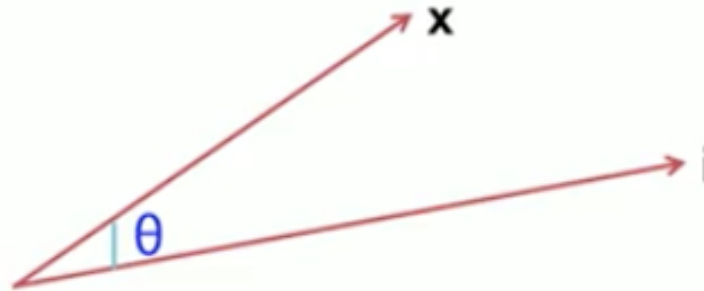
- Create **vectors** with same components that **describe user's preferences**
- Best estimate regarding which items a user likes is **some aggregation of the profiles of those items**

◆ User profile possibilities:

- **Boolean utility matrix:** average the components of vectors representing item profiles for the items in which utility matrix has a 1 for that user
 - E.g., 20% of movies that user U likes have actor A (has a 1)
 - User profile for U will have 0.2 in component for actor A
- **Non-boolean utility matrix: (e.g., ratings)** weight the vectors representing profiles of items by utility (rating) value
 - U gives **average rating of 3**; rates three movies with actor A: 3, 4, 5
 - **Normalize** by subtracting user's average rating: new ratings 0, 1, 2
 - Then user profile component for actor A will have value of 1
 - **Negative weights for below-average ratings, positive for above-avg.**

Making Predictions

- User profile \mathbf{x} , Item profile \mathbf{i}
- Estimate $U(\mathbf{x}, \mathbf{i}) = \cos(\theta) = (\mathbf{x} \cdot \mathbf{i}) / (|\mathbf{x}| |\mathbf{i}|)$



Technically, the cosine distance is actually the angle θ
And the cosine similarity is the angle $180-\theta$

For convenience, we use $\cos(\theta)$ as our similarity measure
and call it the “cosine similarity” in this context.

Content-Based Recommendations

◆ Prediction heuristic

- Given user profile x and item profile i
- Estimate degree to which a user would prefer an item by computing cosine distance between x and i vectors,

◆ Classification Algorithms

- Use machine learning techniques
- Regard given data as a training set
- For each user, build a classifier that predicts the rating of all items
- Ratings on a scale of 1 to k can be directly mapped to k classes
- Many different classifiers
 - Naïve Bayes classifier
 - K-nearest neighbor
 - Decision trees
 - Neural networks.

Example 9.5

- ◆ Building User Profile in Previous Example (9.4):
 - U gives average rating of 3; rates three movies with actor A: 3, 4, 5
 - Normalize by subtracting user's average rating: new ratings 0, 1, 2
 - Negative weights for below-average ratings, positive for above-average,
- ◆ Movie with many actors the user likes:
 - Cosine of angle will be large positive fraction
 - After applying the arccosine, will have a small cosine distance between vectors (angle close to 0 degrees),
- ◆ Movie with mix of actors the user likes/doesn't like:
Cosine of angle will be around 0 (angle close to 90 degrees)
- ◆ Movie with many actors user doesn't like: Cosine will be a large negative fraction => large cosine distance between vectors (angle close to 180 degrees).

Recommending items based on cosine distance

- ◆ Estimate degree to which a user would prefer an item by **computing cosine distance** between x and i vectors
- ◆ Scale components with values that are not boolean (e.g., ratings)
- ◆ Use Random hyperplanes (RH)* and Locality Sensitive Hashing (LSH) techniques to place item profiles (i vectors) in buckets
- ◆ **For a given user (x vector), apply RH and LSH techniques: identify in which bucket we look for items that might have a small cosine distance from user.**
- ◆ *Locality Sensitive Hashes (LSH) for Cosine Similarity are generated using randomly selected hyperplanes. The random hyperplanes slice up the space, and the more hyperplanes you have the more tightly the LSH of v constrains v .

DECISION TREES

Example 9.6: Decision Trees

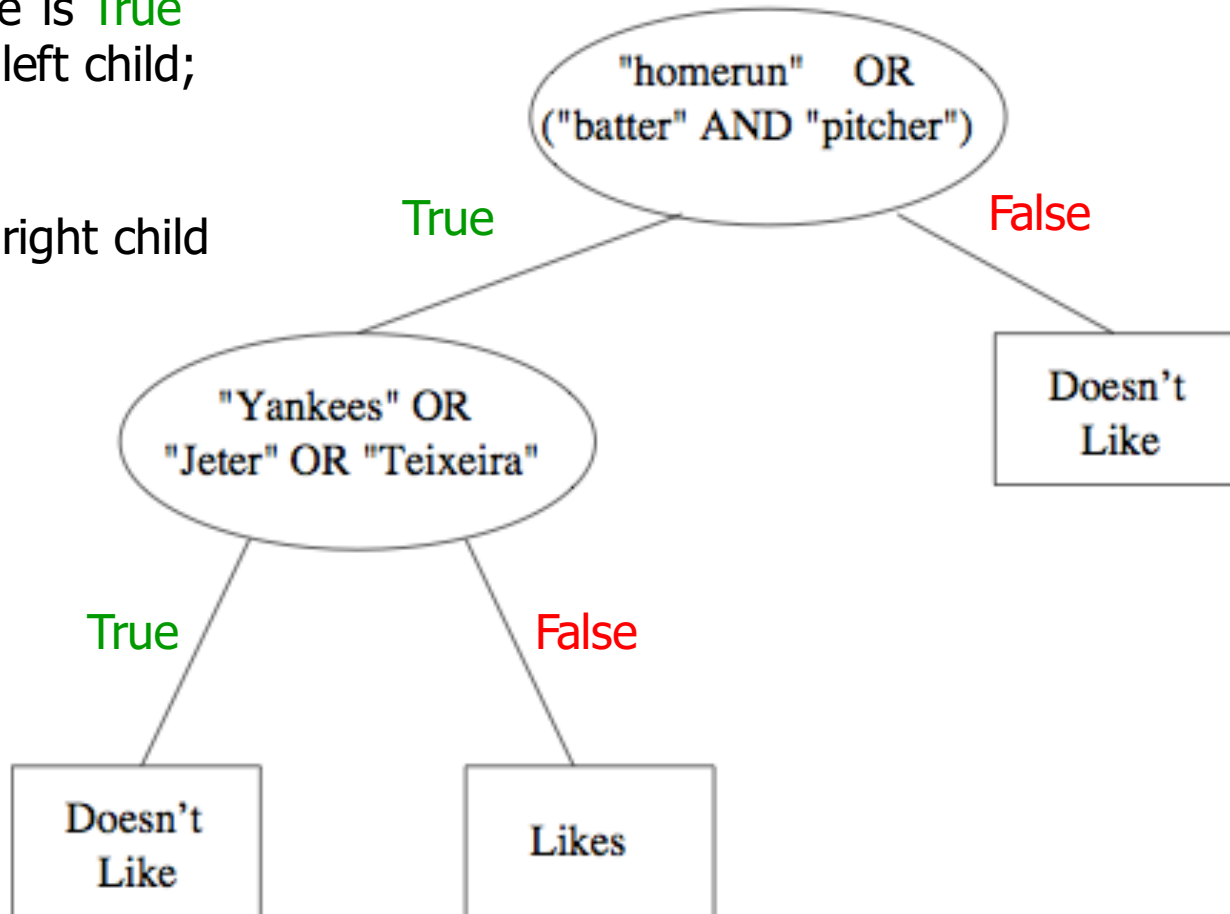
- ◆ Items are **news articles**
- ◆ **Features are high TF.IDF words (keywords) in these documents**
- ◆ Suppose user U likes articles about baseball **except** articles about **New York Yankees**
- ◆ Row of utility matrix for U has a 1 if U has read the article, blank otherwise
- ◆ **Construct decision tree:**
 - select a predicate for each interior node
- ◆ **Classify an item:**
 - Start at root and apply predicate to the item
 - If predicate is true, go to left child; if false, go to right child
 - Repeat until a leaf is reached: leaf classified liked or not liked.

Example 9.6 (cont.)

Condition to check:

-If predicate is **True**
then go to left child;

-if **False**,
then go to right child



Classifiers

- ◆ **Classifiers of all types take a long time to construct**
 - E.g., for decision trees: need one tree per user
- ◆ Constructing a tree requires looking at all item profiles
- ◆ Have to consider many different predicates
- ◆ Could involve complex combinations of features
- ◆ Typically applied only to small problem sizes.

SUMMARY OF CONTENT-BASED APPROACH

Pros: Content-based Approach

◆ +: No need for data on other users

- No cold-start (for item) or sparsity problems (i.e., new items can receive recommendations)

◆ +: Able to recommend to users with unique tastes

◆ +: Able to recommend new & unpopular items

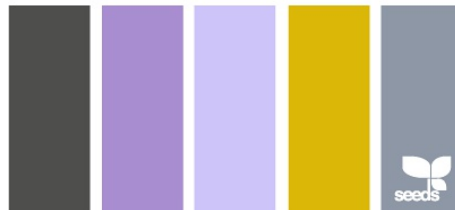
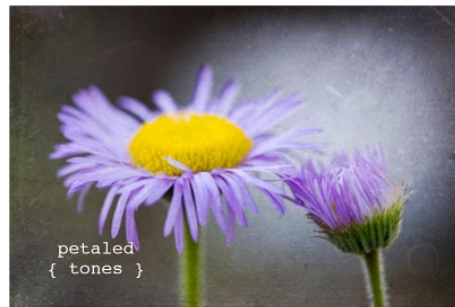
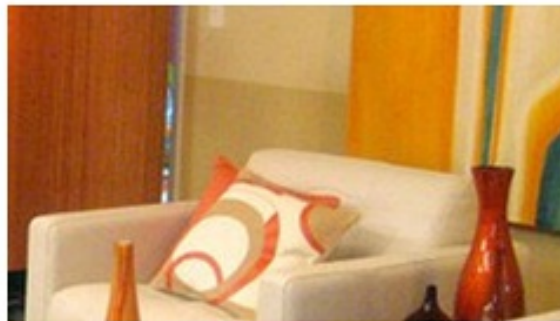
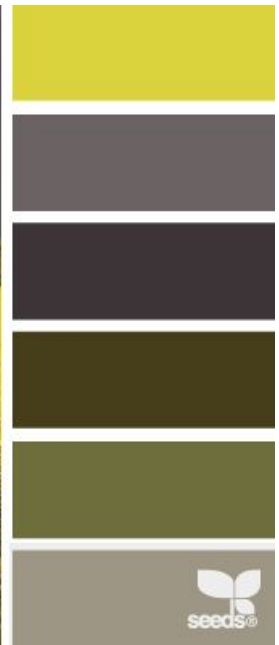
- No first-rater problem (i.e., new products never have been rated, therefore they cannot be recommended)

◆ +: 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 (don't use ratings!).**



Dive Into Color

February, 2010

22

Recommendation Systems

Collaborative Filtering

Part 2

**Harnessing quality judgments of
other users**

Three Approaches to Recommendation Systems

◆ 1) Content-based

- Use characteristics of an item
- Recommend items that
 - have similar content to items user liked in the past
 - match pre-defined attributes of the user

◆ 2) Collaborative filtering

- Build a model from
 - a user's past behavior (e.g., items previously purchased or rated) and
 - similar decisions made by other users
- Use the model to predict items that the user may like
- Collaborative: suggestions made to a user utilizing information across the entire user base

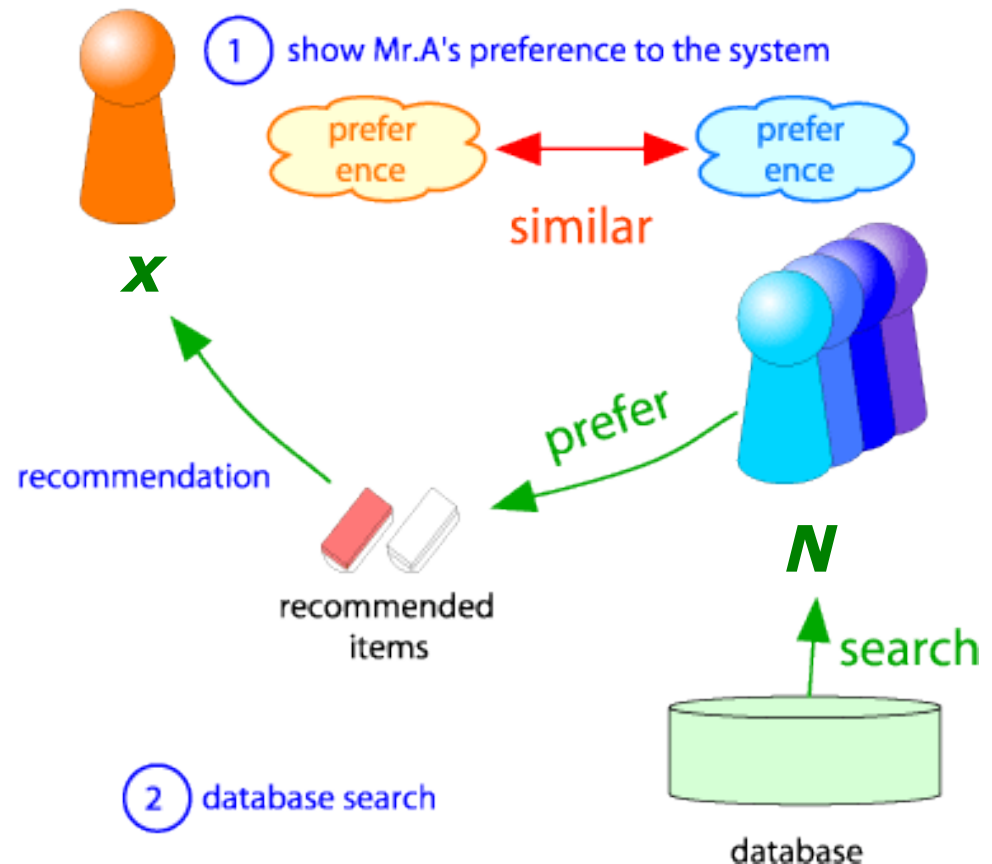
◆ 3) Hybrid approaches.

Collaborative Filtering Example

◆ User-based 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



Collaborative Filtering: Overview

- ◆ CF works by **collecting user feedback**: **ratings for items**
 - Exploit **similarities** in **rating behavior** among **users** in determining **recommendations**
- ◆ **Two classes of CF algorithms:**
 1. **Neighborhood-based or Memory-based approaches**
 - User-based CF
 - Item-based CF,
 2. **Model-based approaches**
 - Estimate parameters of statistical models for user ratings
 - Latent factor and matrix factorization models.

NEIGHBORHOOD-BASED OR MEMORY-BASED COLLABORATIVE FILTERING

USER-BASED CF

Neighborhood-based / Memory-based Collaborative Filtering

- ◆ Active user: the user we want to make predictions for
- ◆ **User-based CF:** A subset of other users is chosen based on their similarity to the active user
- ◆ A weighted combination of their ratings is used to make predictions for the active user
- ◆ **Steps:**
 1. Assign a weight to all users w.r.t. **similarity with the active user**
 2. Select **k** users that have the **highest similarity** with active user (the neighborhood)
 3. Compute a prediction from **a weighted combination of the selected neighbors' ratings.**

Similarity between users: by what measure?

- ◆ Weight $w_{x,y}$ is measure of similarity between user x and active user y

Let r_x be the vector of user x 's ratings

Possible similarity metrics:

- ◆ Jaccard similarity

- (Intersection / union) of two sets
- Doesn't use non-boolean values: e.g., ratings

- ◆ Cosine similarity

- Treat ratings as vectors to points
- Cosine similarity between points

r_x, r_y as sets:

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

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

r_x, r_y as points:

$$r_x = \{1, 0, 0, 1, 3\}$$

$$r_y = \{1, 0, 2, 2, 0\}$$

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$

Similarity between users:
by what measure?

JACCARD SIMILARITY

Jaccard Similarity and Distance of Sets

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

$$Sim (C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$

- ◆ *Jaccard distance* = 1 – Jaccard Similarity

How well do these similarity metrics work?

	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: $\text{sim}(A, B) > \text{sim}(A, C)$ (why?)

- For A and B: One movie rated in common with similar ratings
- For A and C: Two movies rated in common but with dissimilar ratings

◆ Jaccard similarity (Example 9.7)

- Ignores values (rates) in matrix
- Only look at which items are rated in matrix

◆ Intersection / union: $\text{sim}(A, B) = 1/5$, $\text{sim}(A, C) = 2/4$

◆ $1/5 < 2/4$ indicates $\text{sim}(A, C) > \text{sim}(A, B)$

◆ **Not a good similarity metric for ratings data!**

Similarity between users:
by what measure?

COSINE SIMILARITY

Same matrix with cosine similarity (Example 9.8)

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

- ◆ Treat blanks as 0 value: **questionable**, since it treats no rating as more similar to disliking a movie than liking it
- ◆ Cosine similarity of A and B is:

$$\frac{4 \times 5}{\sqrt{4^2 + 5^2 + 1^2} \sqrt{5^2 + 5^2 + 4^2}} = 0.380$$

- ◆ Cosine similarity of A and C is: 0.322
- ◆ $0.380 > 0.322$: **Indicates A,B slightly more similar than A,C**
- ◆ **in this case, cosine similarity better than Jacquard similarity.**

Normalizing ratings (Example 9.9)

◆ To deal better with non-rated items

→ ◆ **Subtract the average rating of that user from each rating**

◆ low ratings -> negative numbers; high ratings -> positive numbers

	HP1	HP2	HP3	TW	SW1	SW2	SW3	Avg. Rating
A	4			5	1			10/3
B	5	5	4					14/3
C				2	4	5		11/3
D		3					3	6/2

	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

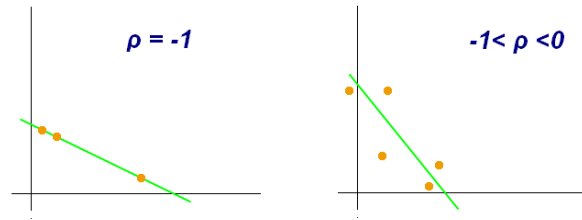
Cosine sim A,B vs. A,C: $0.092 > -0.559$ (85 degrees > 124 degrees)

A, C are much further apart than A, B, but neither is close.

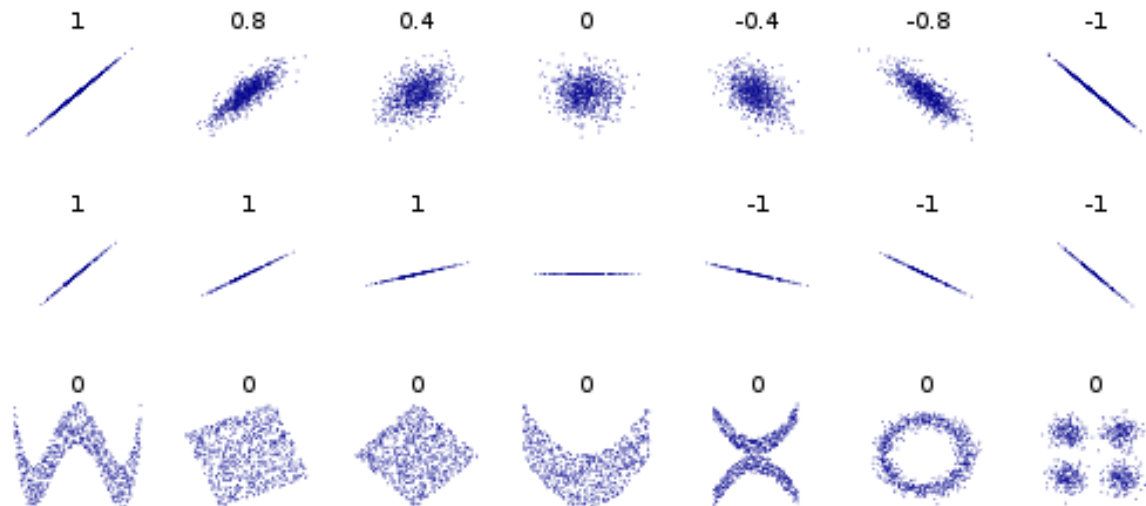
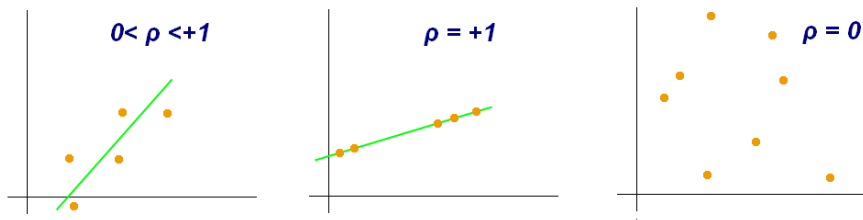
Similarity between users:
by what measure?

PEARSON CORRELATION

Pearson Correlation Examples



Different values of
correlation coefficient (ρ)



Most commonly used measure of similarity:

Pearson Correlation Coefficient

- ◆ Pearson correlation measures extent to which two variables **linearly relate**
- ◆ For users u, v : Pearson correlation is

$$w_{u,v} = \frac{\sum_{i \in I} (r_{u,i} - \bar{r}_u) (r_{v,i} - \bar{r}_v)}{\sqrt{\sum_{i \in I} (r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{i \in I} (r_{v,i} - \bar{r}_v)^2}}$$

where the $i \in I$ summations are over the items that both the users u and v have rated and \bar{r}_u is the average rating of the co-rated items of the u th user.

➤ And $r_{u,i}$ is rating of item i by user u

- ◆ **Note: When calculating these similarities, look only at the co-rated items.**

Example: Pearson Correlation Coefficient

$$w_{u,v} = \frac{\sum_{i \in I} (r_{u,i} - \bar{r}_u) (r_{v,i} - \bar{r}_v)}{\sqrt{\sum_{i \in I} (r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{i \in I} (r_{v,i} - \bar{r}_v)^2}}$$

	I_1	I_2	I_3	I_4
U_1	4	?	5	5
U_2	4	2	1	
U_3	3		2	4
U_4	4	4		
U_5	2	1	3	5

- ◆ Want correlation of user 1, user 5: $w_{1,5}$ Set I of co-rated movies = $\{I_1, I_3, I_4\}$
- ◆ For user 1: **average rating on co-rated movies** is **14/3**; For user 5: **10/3**
- ◆ Number: $(4 - 14/3)(2 - 10/3) + (5 - 14/3)(3 - 10/3) + (5 - 14/3)(5 - 10/3) = 12/9 = 1.3333$
- ◆ Denominator: $\text{sqrt}((4 - 14/3)^2 + (5 - 14/3)^2 + (5 - 14/3)^2) * \text{sqrt}((2 - 10/3)^2 + (3 - 10/3)^2 + (5 - 10/3)^2) = 1.76383$
- ◆ **Pearson correlation $w_{1,5} = 1.3333 / 1.76383 = 0.756$**

Making User-based CF Predictions with Pearson: Weighted Sum of Other Users' Ratings

- ◆ **Weighted average of their ratings is used to generate predictions**
- ◆ **To make a prediction for an active user a on an item i :**

$$P_{a,i} = \bar{r}_a + \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_u) \cdot w_{a,u}}{\sum_{u \in U} |w_{a,u}|}$$

where \bar{r}_a and \bar{r}_u are the average ratings for the user a and user u on all other rated items, and $w_{a,u}$ is the weight between the user a and user u . The summations are over all the users $u \in U$ who have rated the item i .

- ◆ **Note: When making predictions, calculate average of ALL co-rated items for users a and u**
- ◆ **Summation is over all users who rated item i .**

Continued Example: User-Based CF Prediction with Pearson Correlation Coefficient

	I_1	I_2	I_3	I_4
U_1	4	?	5	5
U_2	4	2	1	
U_3	3		2	4
U_4	4	4		
U_5	2	1	3	5

Average rating for each user
(not I_2)

U_1 14/3 = 4.67

U_2 5/2=2.5

U_4 4/1=4

U_5 10/3=3.33

- ◆ Want to predict rating for user U_1 on item I_2 : users U_2 , U_4 and U_5 rated I_2
- ◆ Similarity of U_1 to these users: $w_{1,5} = 0.756$ (see page 21), $w_{1,4} = 0$, $w_{1,2} = -1$

Predict rating for user U_1 on item I_2 \longrightarrow

$$P_{1,2} = \bar{r}_1 + \frac{\sum_u (r_{u,2} - \bar{r}_u) \cdot w_{1,u}}{\sum_u |w_{1,u}|}$$

$$= \bar{r}_1 + \frac{(r_{2,2} - \bar{r}_2)w_{1,2} + (r_{4,2} - \bar{r}_4)w_{1,4} + (r_{5,2} - \bar{r}_5)w_{1,5}}{|w_{1,2}| + |w_{1,4}| + |w_{1,5}|}$$

$$= 4.67 + \frac{(2 - 2.5)(-1) + (4 - 4)0 + (1 - 3.33)0.756}{1 + 0 + 0.756}$$

$$= 3.95.$$

Neighborhood-Based algorithms

- ◆ In neighborhood-based CF algorithms, **a subset of nearest neighbors** of the active user are **chosen based on their similarity with active user**
- ◆ Use these for predictions rather than all users who have rated the item.

Three Approaches to Recommendation Systems

◆ 1) Content-based

- Use characteristics of an item
- Recommend items that
 - have similar content to items user liked in the past
 - match pre-defined attributes of the user

◆ 2) Collaborative filtering

- Build a model from
 - a user's past behavior (e.g., items previously purchased or rated) and
 - similar decisions made by other users
- Use the model to predict items that the user may like
- Collaborative: suggestions made to a user utilizing information across the entire user base
- User-Based CF
- **Item-based CF**

◆ 3) Hybrid approaches.

RECOMMENDATION SYSTEMS

COLLABORATIVE FILTERING

ITEM-BASED CF

Item-based Collaborative Filtering

- ◆ Neighborhood-based CF algorithms **do not scale well** when applied to millions of users & items
 - Due to computational complexity of search for similar users (possible solutions?)
- ◆ **Item-to-item collaborative filtering**
 - Rather than matching similar users
 - **Match user's rated items to similar items**
- ◆ In practice, often **leads to faster online systems and better recommendations**
- ◆ **Similarities between pairs of items i and j are computed off-line**
- ◆ **Predict rating of user a on item i with a simple weighted average.**

Vector Cosine Similarity for two items

- ◆ Vector $A = \{x_1, y_1\}$
- ◆ Vector $B = \{x_2, y_2\}$
- ◆ Vector cosine similarity between A and B is:

$$w_{A,B} = \cos(\vec{A}, \vec{B}) = \frac{\vec{A} \bullet \vec{B}}{\|\vec{A}\| * \|\vec{B}\|} = \frac{x_1 x_2 + y_1 y_2}{\sqrt{x_1^2 + y_1^2} \sqrt{x_2^2 + y_2^2}}$$

Pearson Correlation between items i, j (Cont'd)

	1	2	...	i		j	...	$m-1$	m
1				R		?			
2				R		R			
⋮									
l				R		R			
⋮									
$n-1$?		R			
n				R		R			

FIGURE 2: item-based similarity ($w_{i,j}$) calculation based on the co-rated items i and j from users 2, l and n .

Source: Su, X., & Khoshgoftaar, T. M. (2009). A survey of collaborative filtering techniques. *Advances in artificial intelligence*, 2009, 4.

Recall: Pearson Correlation Coefficient

- ◆ Pearson correlation measures extent to which two variables linearly relate
- ◆ For users u, v : Pearson correlation is

$$w_{u,v} = \frac{\sum_{i \in I} (r_{u,i} - \bar{r}_u) (r_{v,i} - \bar{r}_v)}{\sqrt{\sum_{i \in I} (r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{i \in I} (r_{v,i} - \bar{r}_v)^2}}$$

where the $i \in I$ summations are over the items that both the users u and v have rated and \bar{r}_u is the average rating of the co-rated items of the u th user.

➤ And $r_{u,i}$ is rating of item i by user u

- ◆ Note: When calculating these similarities, look only at the co-rated items.

Pearson Correlation between items i, j

For the item-based algorithm, denote the set of users $u \in U$ who rated both items i and j , then the *Pearson Correlation* will be

$$w_{i,j} = \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_i) (r_{u,j} - \bar{r}_j)}{\sqrt{\sum_{u \in U} (r_{u,i} - \bar{r}_i)^2} \sqrt{\sum_{u \in U} (r_{u,j} - \bar{r}_j)^2}}, \quad (2)$$

where $r_{u,i}$ is the rating of user u on item i , \bar{r}_i is the average rating of the i th item by those users, see Figure 2 [40].

- ◆ **Note: Sum over set of users U who rated both items i, j**
- ◆ $r_{u,i}$ is rating of user u on item i
- ◆ \bar{r}_i is average rating of i^{th} item by those users.

Make Item-Based Predictions Using a Simple Weighted Average

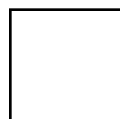
- ◆ Predict rating for user u on item i
- ◆ $w_{i,n}$ is weight between items i and n
- ◆ $r_{u,n}$ is rating for user u on item n
- ◆ Summation over **neighborhood set N of items** rated by u **that are most similar to i**

$$P_{u,i} = \frac{\sum_{n \in N} r_{u,n} w_{i,n}}{\sum_{n \in N} |w_{i,n}|}$$

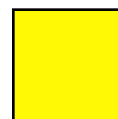
where the summations are over all other rated items $n \in N$ for user u , $w_{i,n}$ is the weight between items i and n , $r_{u,n}$ is the rating for user u on item n .

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	



- 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
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	



- estimate rating of movie 1 by user 5

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

First: what is similarity between items?

		users												Similarity (made up for example):	
		1	2	3	4	5	6	7	8	9	10	11	12	$w_{1,j}$	
movies	1	1		3		?	5			5		4		1.00	
	2			5	4			4			2	1	3	-0.18	
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>	
	4		2	4		5			4			2		-0.10	
	5			4	3	4	2					2	5	-0.31	
	<u>6</u>	1		3		3			2			4		<u>0.59</u>	

Neighbor selection: Identify movies most similar to movie 1 and rated by user 5

Neighborhood size is 2: pick movies 3 and 6

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

movies	users												Similarity (made up): $w_{1,j}$
	1	2	3	4	5	6	7	8	9	10	11	12	
	1		3		?	5			5		4		1.00
	2		5	4			4			2	1	3	-0.18
	<u>3</u>	2	4		1	2	3		4	3	5		<u>0.41</u>
	4		2	4		5		4			2		-0.10
	5			4	3	4	2				2	5	-0.31
	<u>6</u>	1		3		3		2			4		<u>0.59</u>

Similarity weights:
 $w_{1,3}=0.41$, $w_{1,6}=0.59$

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												Similarity: $w_{1,j}$
		1	2	3	4	5	6	7	8	9	10	11	12	
movies	1	1		3		2.6	5			5		4		1.00
	2			5	4			4			2	1	3	-0.18
	<u>3</u>	2	4		1	2		3		4	3	5		<u>0.41</u>
	4		2	4		5			4			2		-0.10
	5			4	3	4	2					2	5	-0.31
	<u>6</u>	1		3		3			2			4		<u>0.59</u>

Predict by taking weighted average:

$$P_{1.5} = (0.41 \cdot 2 + 0.59 \cdot 3) / (0.41 + 0.59) = 2.6$$

$$P_{u,i} = \frac{\sum_{n \in N} r_{u,n} w_{i,n}}{\sum_{n \in N} |w_{i,n}|}$$

Example: Item-to-Item Collaborative Filtering with Pearson Similarity

	I1	I2	I3	I4
U1	2	1		3
U2	3	?	5	2
U3		4	2	3
U4	5	3	1	

$$w_{i,j} = \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_i) (r_{u,j} - \bar{r}_j)}{\sqrt{\sum_{u \in U} (r_{u,i} - \bar{r}_i)^2} \sqrt{\sum_{u \in U} (r_{u,j} - \bar{r}_j)^2}}$$

- ◆ What does set U represent?
 - Set of users U who rated both items i, j
- ◆ What are the members of the set U for $w_{1,2}$?
 - U1 and U4 rated items I1 and I2
- ◆ When calculating average ratings for item i, which ratings do we use?
All ratings or just for co-rated items?
 - We showed examples of co-rated items
 - We can use all ratings as well: based on all ratings: $\text{avg}(r_1) = 10/3$, $\text{avg}(r_2)$ ⁹⁶ = 8/3

Example: Item-to-Item Collaborative Filtering with Pearson Similarity

	I1	I2	I3	I4
U1	2	1		3
U2	3	?	5	2
U3		4	2	3
U4	5	3	1	

$$w_{i,j} = \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_i)(r_{u,j} - \bar{r}_j)}{\sqrt{\sum_{u \in U} (r_{u,i} - \bar{r}_i)^2} \sqrt{\sum_{u \in U} (r_{u,j} - \bar{r}_j)^2}}$$

- ◆ Based on all ratings: average ratings for items I1, I2: $r_1 = 10/3$, $r_2 = 8/3$
- ◆ Pearson correlation for $w_{1,2}$: *Similarity between items 1 and 2*
- ◆ Set U includes U1 and U4
- ◆ $w_{1,2} = ((r_{U1,I1} - 10/3)(r_{U1,I2} - 8/3) + (r_{U4,I1} - 10/3)(r_{U4,I2} - 8/3)) /$
 $(\text{sqrt}((r_{U1,I1} - 10/3)^2 + (r_{U4,I1} - 10/3)^2) * \text{sqrt}((r_{U1,I2} - 8/3)^2 + (r_{U4,I2} - 8/3)^2))$
- ◆ $w_{1,2} = \frac{(2-10/3)(1-8/3) + (5-10/3)(3-8/3)}{\text{sqrt}((2-10/3)^2 + (5-10/3)^2) * \text{sqrt}((1-8/3)^2 + (3-8/3)^2)}$

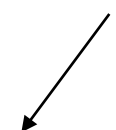
$$= 2.778 / 3.628 = 0.765$$

Example: Item-to-Item Collaborative Filtering with Pearson Similarity

	I1	I2	I3	I4
U1	2	1		3
U2	3	?	5	2
U3		4	2	3
U4	5	3	1	

$$w_{i,j} = \frac{\sum_{u \in U} (r_{u,i} - \bar{r}_i) (r_{u,j} - \bar{r}_j)}{\sqrt{\sum_{u \in U} (r_{u,i} - \bar{r}_i)^2} \sqrt{\sum_{u \in U} (r_{u,j} - \bar{r}_j)^2}}$$

co-rated items only



- ◆ Based on co-ratings: average ratings for items I1, I2: $r_1 = 7/2$, $r_2 = 4/2$
- ◆ Pearson correlation for $w_{1,2}$: *Similarity between items 1 and 2*
- ◆ Set U includes U1 and U4
- ◆ $w_{1,2} = ((r_{U1,I1} - 7/2)(r_{U1,I2} - 4/2) + (r_{U4,I1} - 7/2)(r_{U4,I2} - 4/2)) /$
 $(\text{sqrt}((r_{U1,I1} - 7/2)^2 + (r_{U4,I1} - 7/2)^2) * \text{sqrt}((r_{U1,I2} - 4/2)^2 + (r_{U4,I2} - 4/2)^2))$
- ◆ $w_{1,2} = \frac{(2-7/2)(1-4/2) + (5-7/2)(3-4/2)}{\text{sqrt}((2-7/2)^2 + (5-7/2)^2) * \text{sqrt}((1-4/2)^2 + (3-4/2)^2)}$

$= 3/3 = 1$

Item-Based Prediction

	I1	I2	I3	I4
U1	2	1		3
U2	3	?	5	2
U3		4	2	3
U4	5	3	1	

$$P_{u,i} = \frac{\sum_{n \in N} r_{u,n} w_{i,n}}{\sum_{n \in N} |w_{i,n}|}$$

- ◆ If we have the following item similarities: $w_{2,1} = 0.5$, $w_{2,3} = 0.2$, $w_{2,4} = 0.3$
- ◆ Which items are in the neighborhood N for item 2 if $|N| = 2$?
 - Items I1 and I4
- ◆ **Predict the rating of user U2 on I2:** user = 2, item = 2, $|N| = \text{items } 1,4$

$$P_{2,2} = \frac{(r_{2,1} * w_{2,1}) + (r_{2,4} * w_{2,4})}{0.5 + 0.3} = \frac{3*0.5 + 2*0.3}{0.8} = 2.625$$

Summary

- ◆ What is RecSys
- ◆ The Long Tail
- ◆ Three Approaches to Recommendation Systems

1) Content-based

- Use characteristics of an item
- Recommend items that have similar content to items user liked in the past
- Or items that match pre-defined attributes of the user

2) Collaborative filtering

- Build a model from a user's past behavior (items previously purchased or rated) and similar decisions made by other users
- Use the model to predict items that the user may like
- Collaborative: suggestions made to a user utilize information across the entire user base

3) Hybrid approaches.