

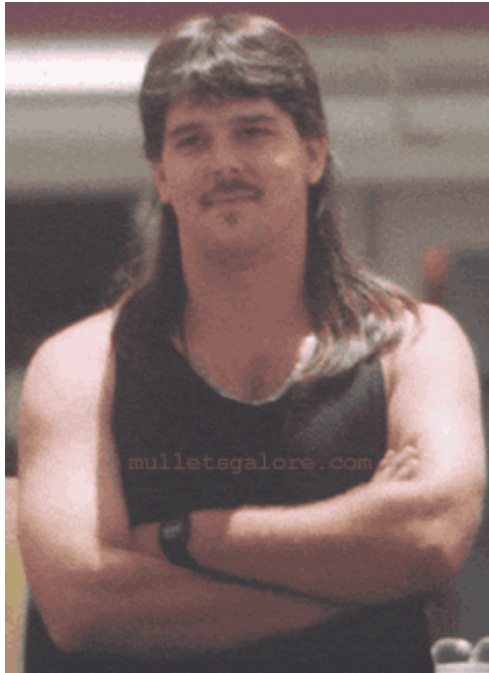
# Recommender System

---

Book source: Mining of Massive Datasets (<http://www.mmds.org>)  
by  
Jure Leskovec, Anand Rajaraman, Jeff Ullman



# Example: Recommender Systems



- **Customer X**

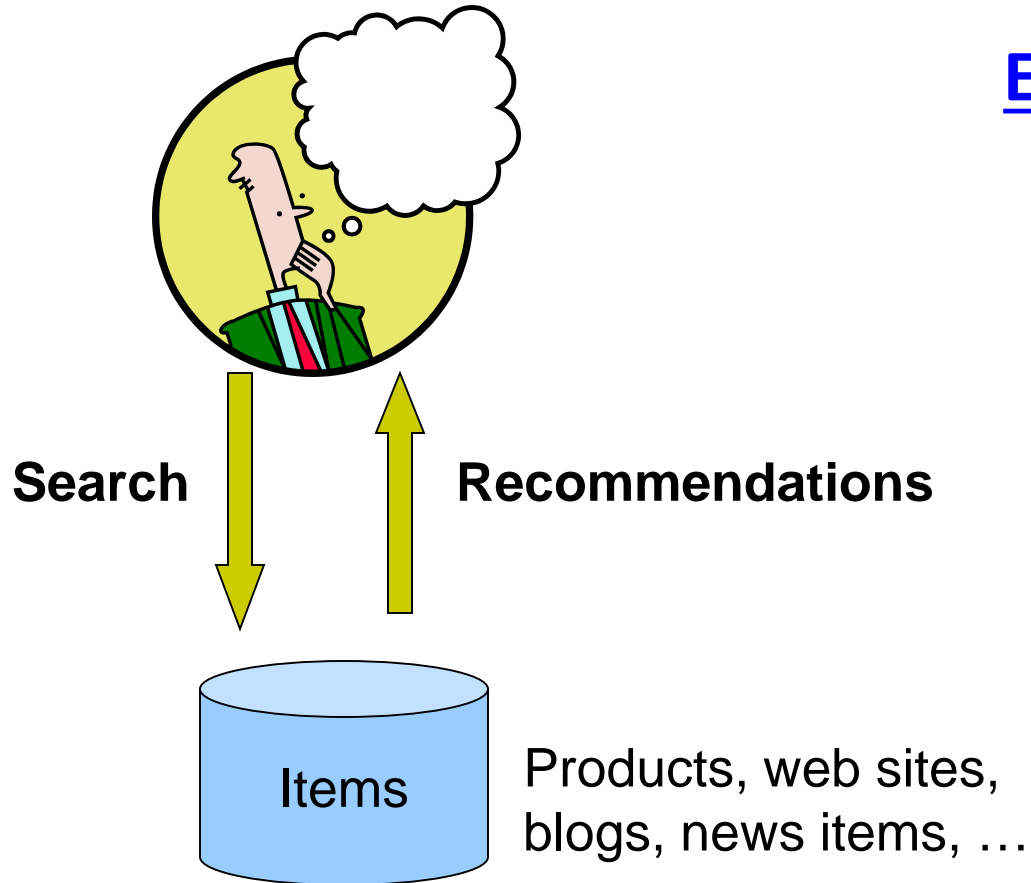
- Buys Metallica CD
- Buys Megadeth CD



- **Customer Y**

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

# Recommendations



## Examples:

amazon.com



**movie lens**  
helping you find the *right* movies



# Types of Recommendation Systems



- Content based
  - Based on content similarity
- Collaborative filtering

# Content-based Recommendations



- **Main idea:** Recommend items to customer  $x$  similar to previous items rated highly by  $x$

## *Example:*

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



# Item Profiles

- For each item, create an **item profile**
- **Profile is a set (vector) of features**
  - **Movies:** author, title, actor, director,...
  - **Text:** Set of “important” words in document



# User Profiles and Prediction

- User profile possibilities:
  - Weighted average of rated item profiles
  - Variation: weight by difference from average rating for item
  - ...
- Prediction heuristic:
  - Given user profile  $x$  and item profile  $i$ , estimate
$$u(x, i) = \cos(x, i) = \frac{x \cdot i}{||x|| \cdot ||i||}$$



# Content-based : Pros and Cons

## ● Pros

- No need for data on other users
- Able to recommend to users with unique tastes
- +: Able to recommend new & unpopular items
- +: Able to provide explanations

## ● Cons

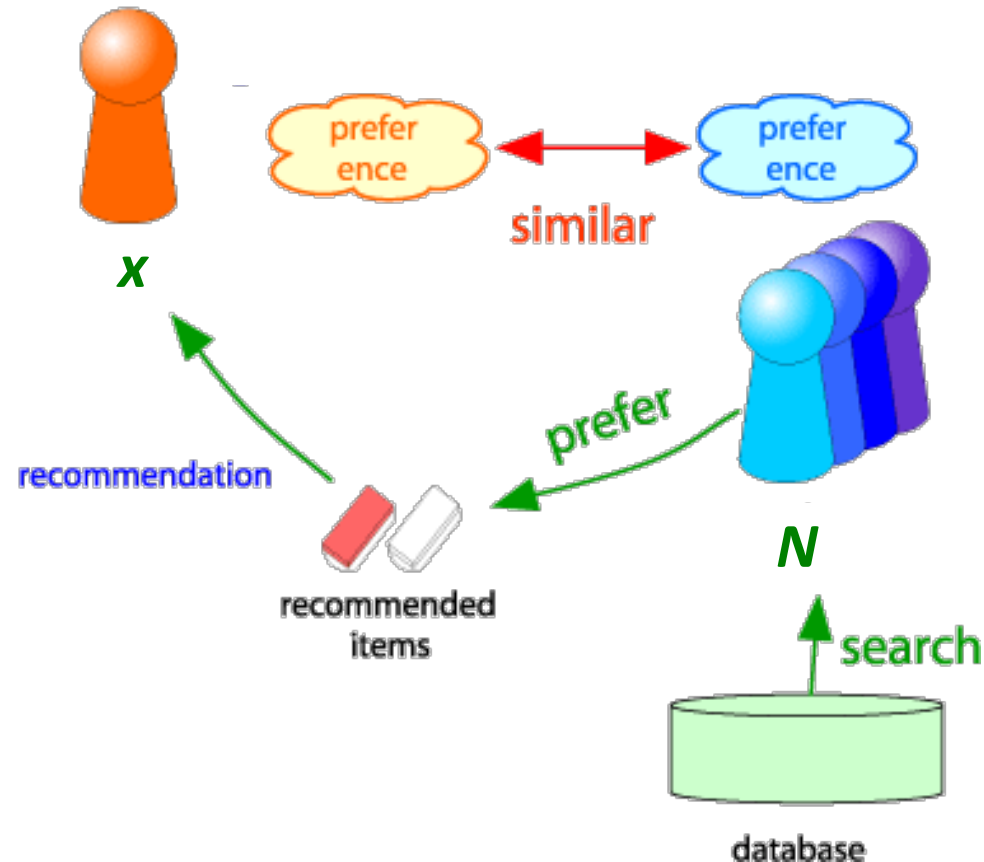
- Finding the appropriate features is hard
- Recommendations for new users
- How to build a user profile?
- Overspecialization
  - Never recommends items outside user's content profile
  - Unable to exploit quality judgments of other users





# 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 similarity measure**
  - **Problem:** Ignores the value of the rating
- **Cosine similarity measure**
  - $\text{sim}(x, y) = \cos(r_x, r_y) = \frac{r_x \cdot r_y}{||r_x|| \cdot ||r_y||}$
  - **Problem:** Treats missing ratings as “negative”
- **Pearson correlation coefficient**
  - $S_{xy}$  = items rated by both users  $x$  and  $y$

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

# 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 similarity:**  $1/5 < 2/4$
- **Cosine similarity:**  $0.386 > 0.322$ 
  - Considers missing ratings as “negative”

- **Solution: subtract the (row) mean**

	HP1	HP2	HP3	TW	SW1	SW2	SW3
A	2/3			5/3	-7/3		
B	1/3	1/3	-2/3				
C				-5/3	1/3	4/3	
D		0					0

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

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



# Rating Predictions

## From similarity metric to recommendations:

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

Shorthand:

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