

# Recommender Systems

Ashwin Malshé

“There is an extensive class of Web applications that involve predicting user responses to options. Such a facility is called a recommendation system.”

– *Jure Leskovec, Stanford University,  
Anand Rajaraman, Milliway Labs,  
Jeffrey D. Ullman, Stanford University*

# Netflix

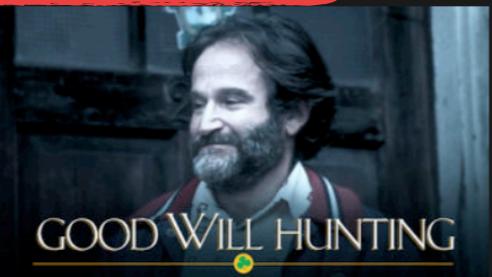
Because you watched Stranger Things



Trending Now



Because you watched Forrest Gump



# Amazon

Inspired by your shopping trends



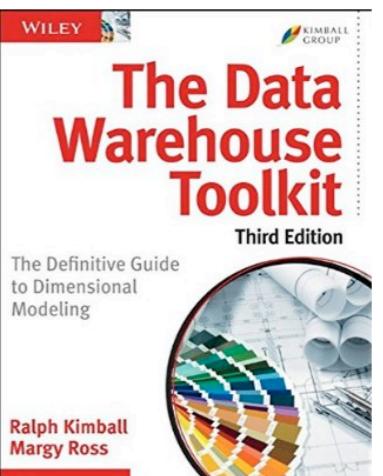
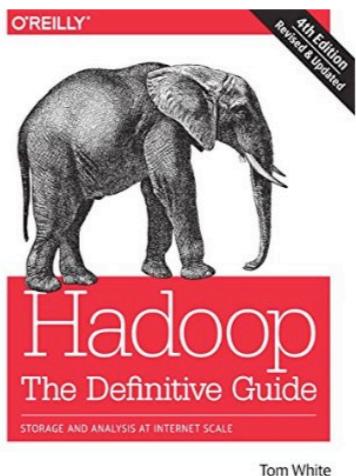
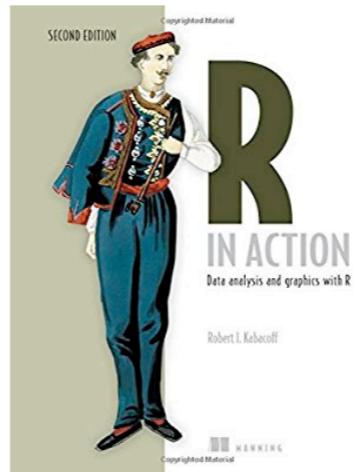
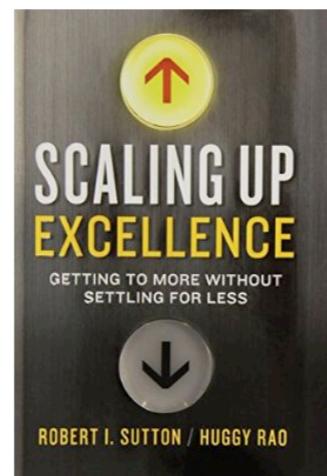
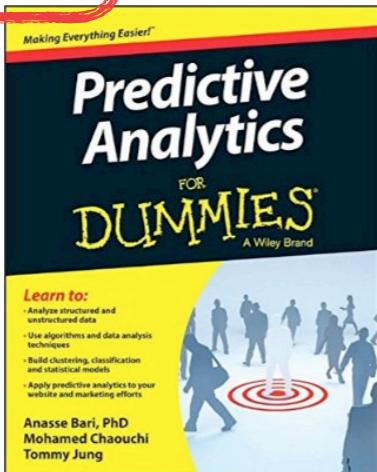
## More top picks for you

### Data Science for Business

What You Need to Know About Data Mining and Data-Analytic Thinking



Foster Provost & Tom Fawcett

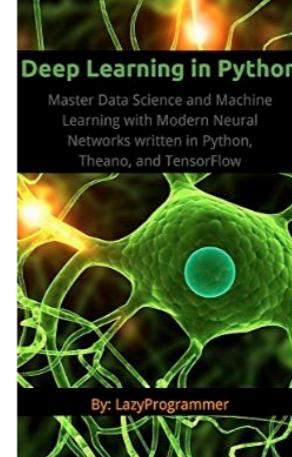
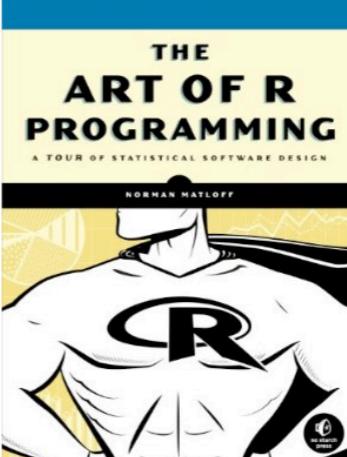
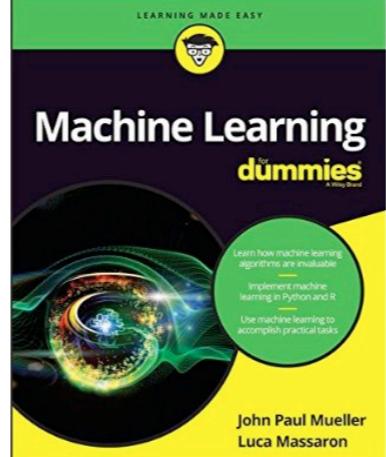
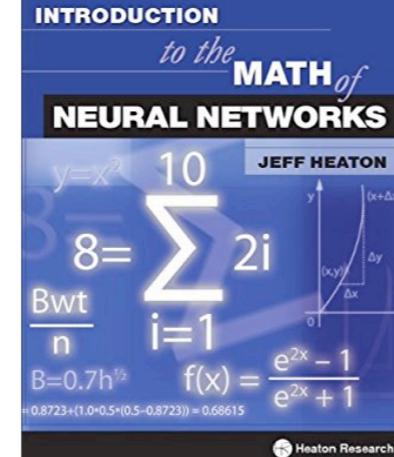
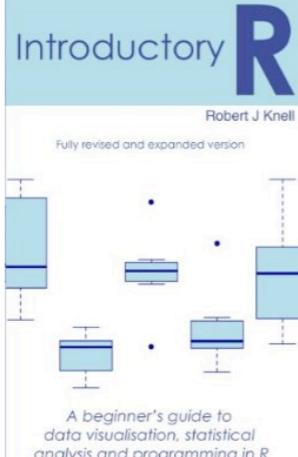


## Recommendations for you in Kindle Store

### Machine Learning With Random Forests And Decision Trees

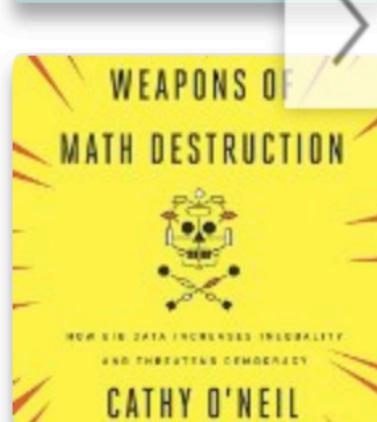
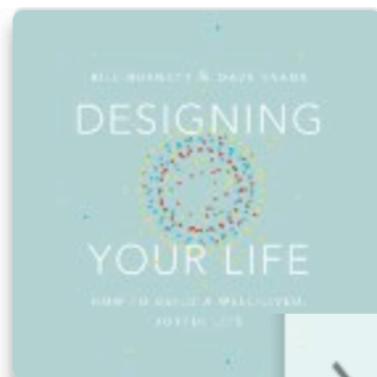
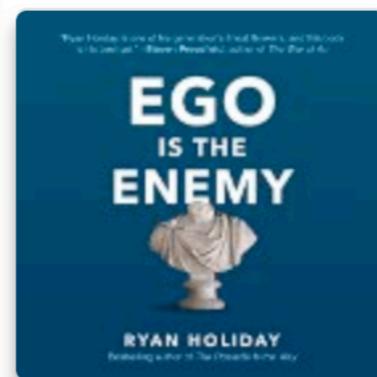
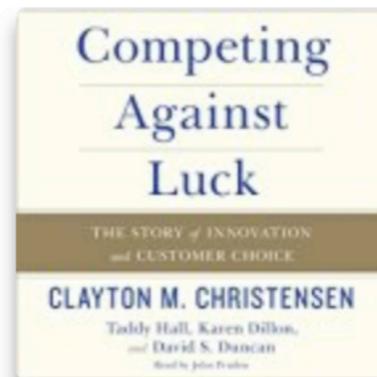
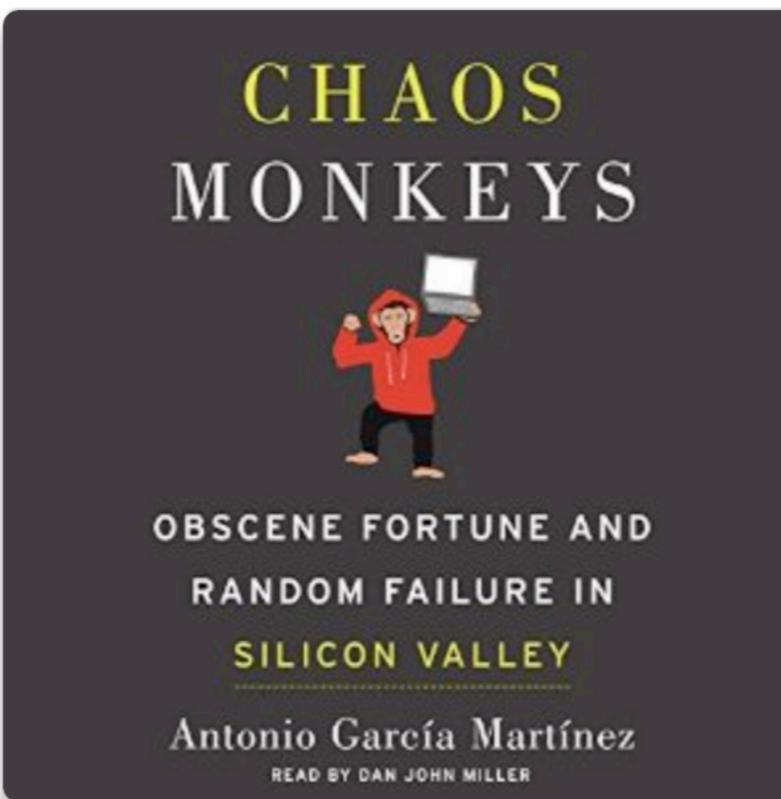
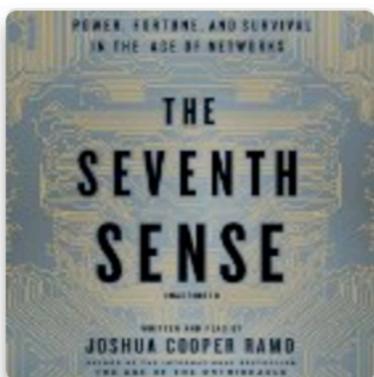
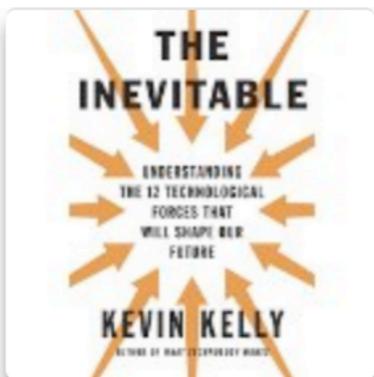
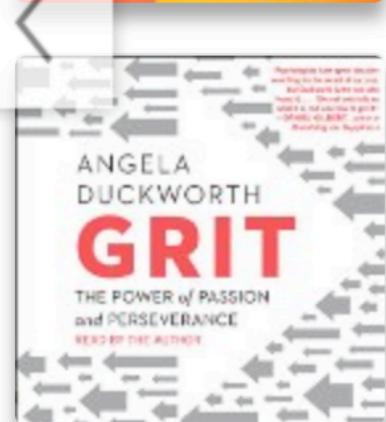
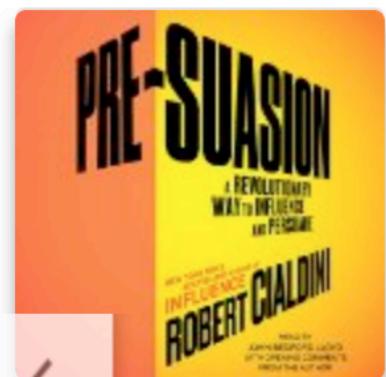


A Mostly Intuitive Guide, But Also Some Python  
SCOTT HARTSHORN



# Audible

Based on Your Past Purchases, We Think You'll Enjoy



# Twitter

Who to follow · Refresh · View all

 **StatsBlogs** @StatsBlogs X

**Follow**

 **Dirk Eddelbuettel** @eddelbu... X

**Follow**

 **ASSA Meeting** @ASSAMeeti... X

**Follow**

 **Find people you know**  
Import your contacts from Gmail

---

**Connect other address books**

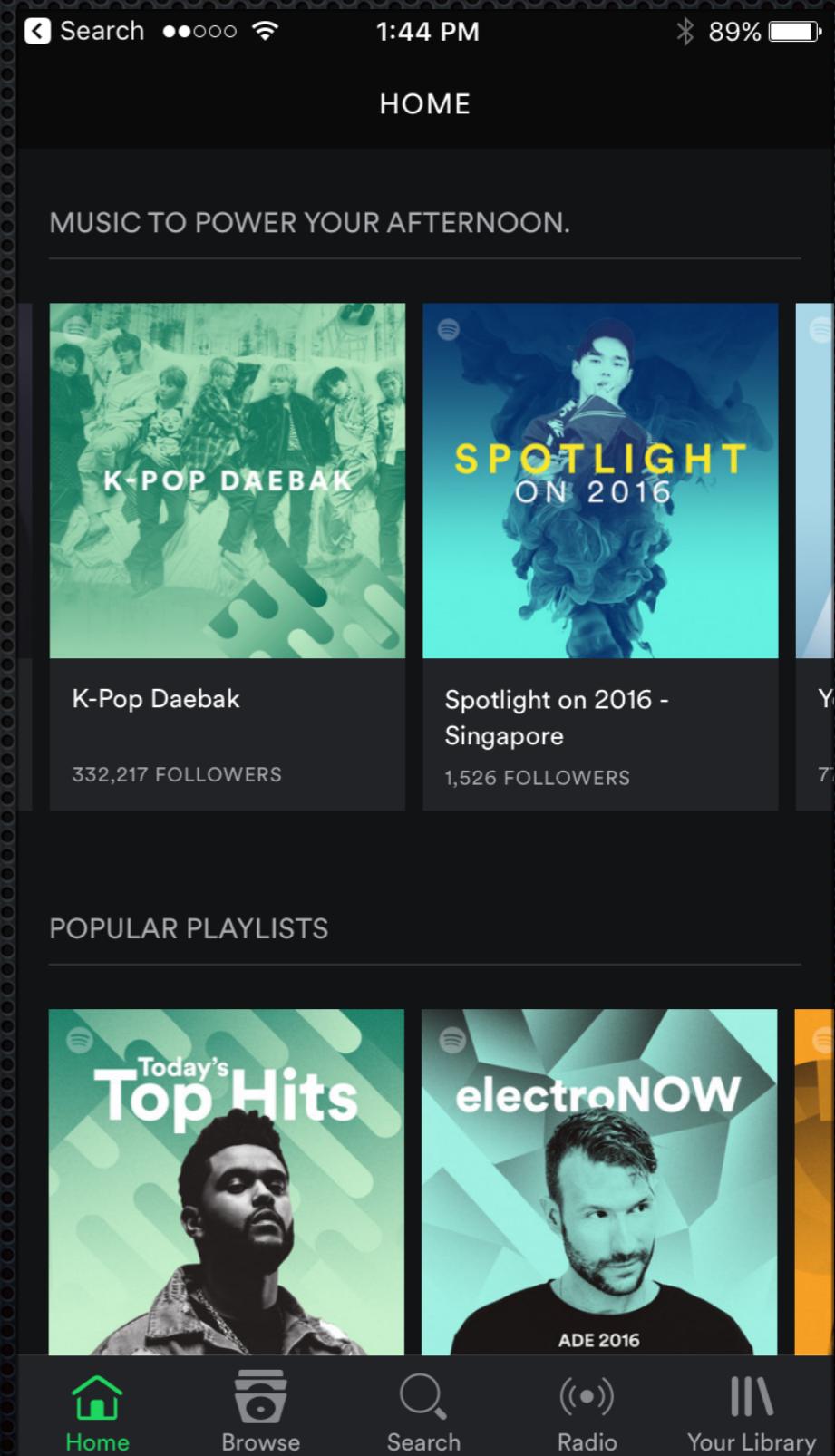
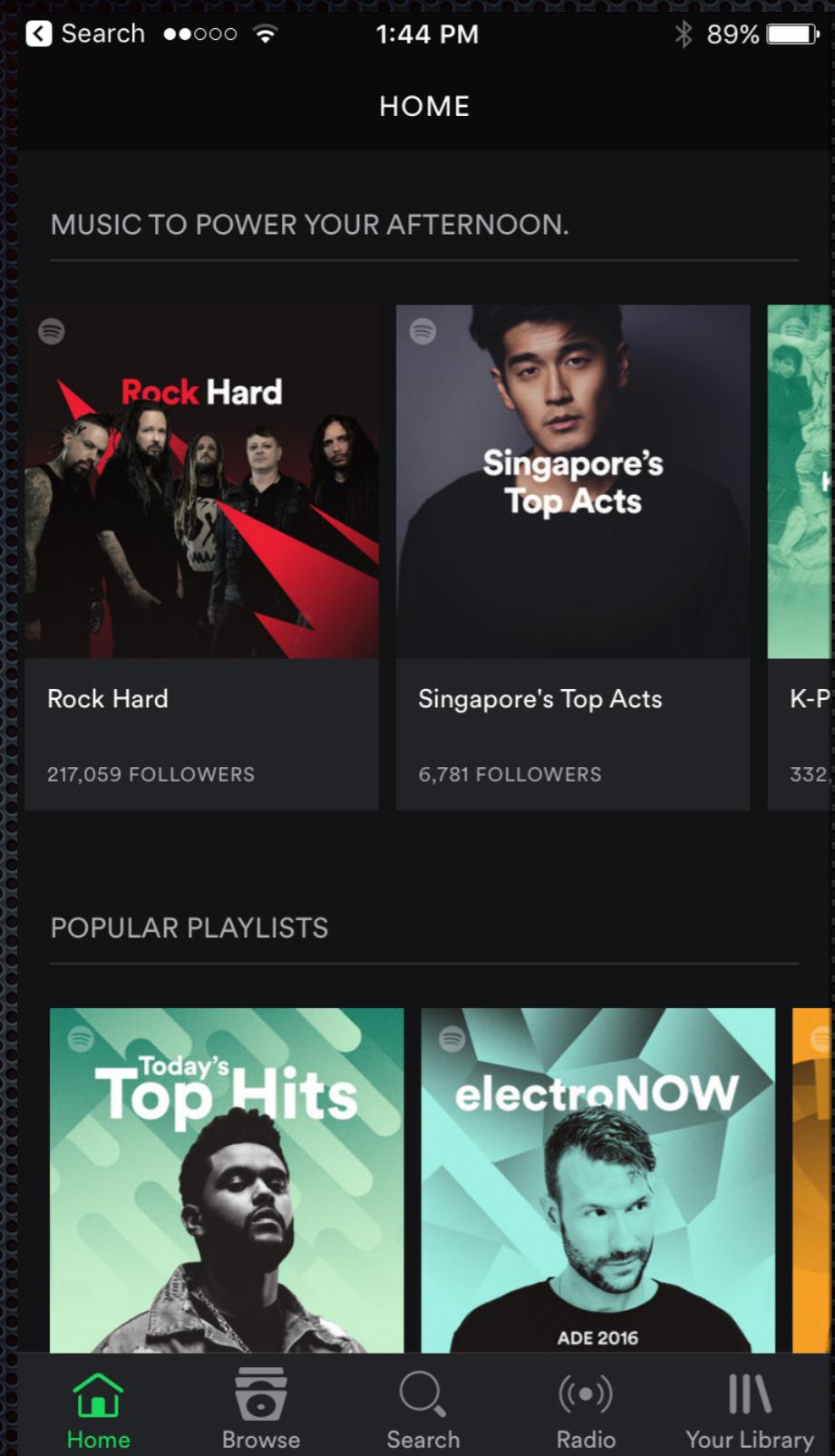
# Google

The screenshot shows a web browser window with a dark background. At the top, there is a horizontal tab bar with several tabs open. From left to right, the visible tab titles are: "r - R", "Dashboard - ...", "https://cran.r... (highlighted)", "https://cran.r... (highlighted)", "R Data Analy...", "tables | OPHI", and "Using R to d...". Below the tab bar is the Google search interface. On the left, the classic multi-colored "Google" logo is visible. To its right is the search input field, which contains the partial query "why is san an". To the right of the input field is a blue search button with a white magnifying glass icon. A dropdown menu is open below the input field, displaying five suggested search terms:

- why is san antonio not a protest town
- why is san antonio so boring
- why is san antonio so hot
- why is san antonio so fat

At the bottom of the search interface, the text "Press Enter to search." is displayed.

# Spotify



# Technologies

- Content-based systems
  - Examine properties of the items recommended
    - Similarity of items is determined by measuring the similarity in their properties
- Collaborative filtering
  - Recommend items based on similarity measures between users and/or items
    - Similarity of items is determined by the similarity of the ratings of those items by the users who have rated both items

# Our Focus

- Our focus is on collaborative filtering
  - The process of identifying similar users and recommending what similar users like is called collaborative filtering
- Content-based systems can be designed based on the item and user attributes. We have done these types of models before
  - e.g., Logistic regression, SVM

# Mechanics

# Utility Matrix

users/ items	i <sub>1</sub>	i <sub>2</sub>	i <sub>3</sub>	...	i <sub>n</sub>
U <sub>1</sub>	r <sub>11</sub>	r <sub>12</sub>	r <sub>13</sub>	...	r <sub>1n</sub>
U <sub>2</sub>	r <sub>21</sub>	r <sub>22</sub>	r <sub>23</sub>	...	r <sub>2n</sub>
U <sub>3</sub>	r <sub>31</sub>	r <sub>32</sub>	r <sub>33</sub>	...	r <sub>3n</sub>
...	...	...	...	...	...
U <sub>m</sub>	r <sub>m1</sub>	r <sub>m2</sub>	r <sub>m3</sub>	...	r <sub>mn</sub>

# Utility Matrix is Sparse

users/ items	$i_1$	$i_2$	$i_3$	...	$i_n$
$U_1$	$r_{11}$		$r_{13}$	...	$r_{1n}$
$U_2$		$r_{22}$		...	
$U_3$	$r_{31}$			...	$r_{3n}$
...	...	...	...	...	...
$U_m$			$r_{m3}$	...	

# Similarity Measures

- Jaccard distance
- Cosine similarity
- Pearson correlation

# Example: Netflix Ratings

	M1	M2	M3	M4	M5	M6	M7
A	4			5	1		
B	5	5	4				
C				2	4	5	
D		3					3

# Jaccard Distance

- Applicable to yes/no ratings
- Ignores the ordinal nature of the rating
- A and B have in common rated only one movie, M1. Collectively they have rated 5 movies.
  - This leads to Jaccard similarity =  $1/5$  and Jaccard distance =  $4/5$
- A and C have in common rated 2 out of 4 movies so Jaccard similarity is  $2/4$  and Jaccard distance is  $2/4$

# Cosine Distance

$$\cos \theta = \frac{A \cdot B}{\|A\| \times \|B\|}$$

Using this formula, cosine similarity between A and B is 0.38 and A and C is 0.322. Thus, A is more similar to B than C.

# Pearson Correlation

- This is centered cosine similarity
  - Center all the ratings row-wise
  - Get cosine similarities

# Types of Collaborative Filtering

- User-based collaborative filtering (UBCF)
  - UBCF is a memory-based algorithm which tries to mimic word-of-mouth by analyzing rating data from many individuals. The assumption is that users with similar preferences will rate items similarly.
- Item-based collaborative filtering (IBCF)
  - IBCF is a model-based approach, which produces recommendations based on the relationship between items inferred from the rating matrix

# Clustering

- We can reduce the number of items and users by clustering
- Usually clustering will happen in sequence
  - Cluster items first then cluster users
  - Keep repeating until we have a reasonable number of clusters