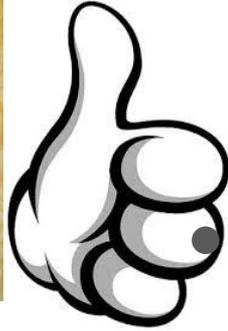


# Recommendation Systems

Stony Brook University  
CSE545, Spring 2019



# Recommendation Systems

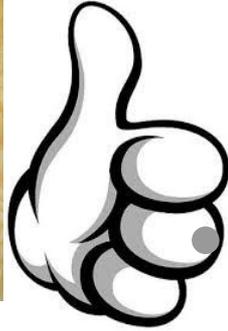


- What other item will this **user** like?  
(based on previously liked items)

How much will user like item X?

?

# Recommendation Systems

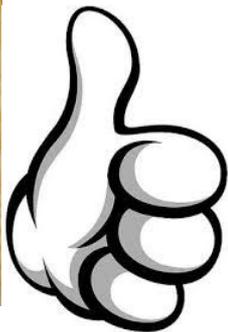


- What other item will this **user** like?  
(based on previously liked items)

How much will user like item X?



# Recommendation Systems



# Recommendation Systems



Past User Ratings

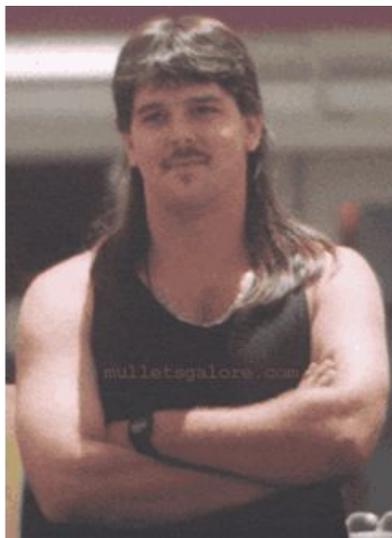


# Recommendation Systems

## Why Big Data?

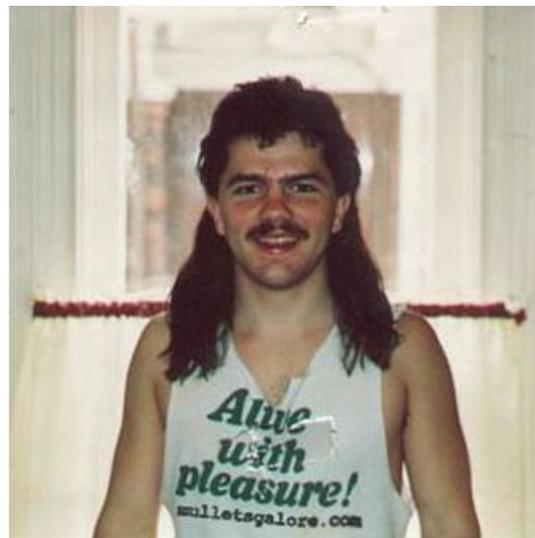
- Data with many potential features (and sometimes observations)
- An application of techniques for finding similar items
  - locality sensitive hashing
  - dimensionality reduction

# Recommendation System: Example



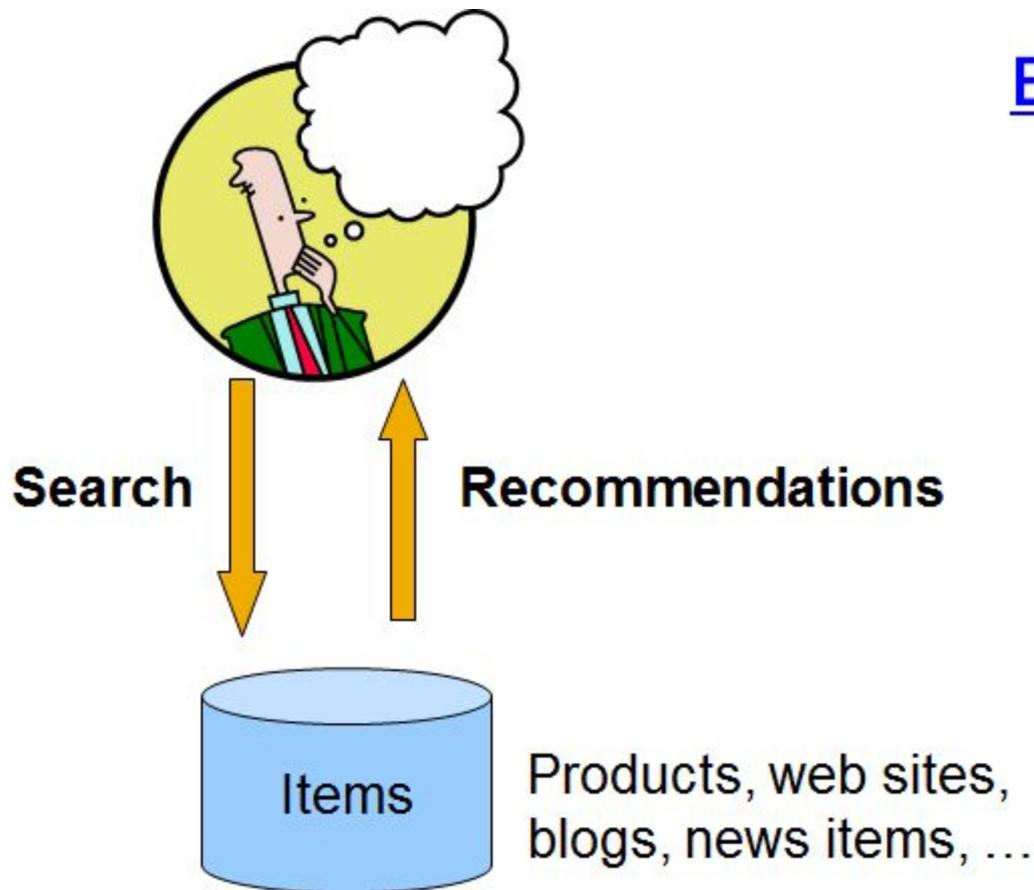
## ■ Customer X

- Buys Metallica CD
- Buys Megadeth CD



## ■ Customer Y

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



## Examples:

amazon.com.



StumbleUpon



del.icio.us



movie lens

helping you find the *right* movies

last.fm  
the social music revolution

Google  
News

YouTube

XBOX  
LIVE

# Enabled by Web Shopping

- Does Wal-Mart have everything you need?

# Enabled by Web Shopping

- Does Wal-Mart have everything you need?



(thelongtail.com)

# Enabled by Web Shopping

- Does Wal-Mart have everything you need?
- A lot of products are only of interest to a small population (i.e. “[long-tail products](#)”).
- However, most people buy many products that are from the long-tail.
- Web shopping enables more choices
  - Harder to search
  - Recommendation engines to the rescue



(thelongtail.com)

# Enabled by Web Shopping

- Does Wal-Mart have everything you need?
- A lot of products are only of interest to a small population (i.e. “[long-tail products](#)”).
- However, most people buy many products that are from
- Web shopping
  - Harder to
  - Recomm



# A Model for Recommendation Systems

Given: *users, items, utility matrix*

# A Model for Recommendation Systems

Given: *users*, *items*, *utility matrix*



<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	3		3
<i>B</i>	5			4	2
<i>C</i>			5	2	

# A Model for Recommendation Systems

Given: *users*, *items*, *utility matrix*



<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	3		3
<i>B</i>	5			4	2
<i>C</i>	?	?	5	2	?

# Recommendation Systems

## Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

# Recommendation Systems

Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

1. Content-based
2. Collaborative
3. Latent Factor

# Recommendation Systems

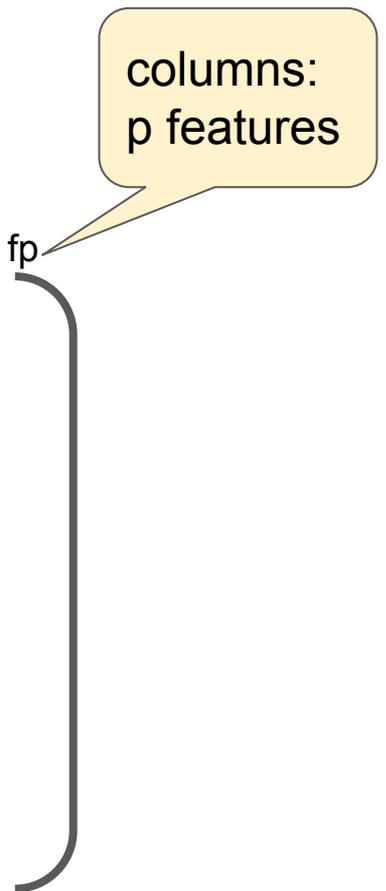
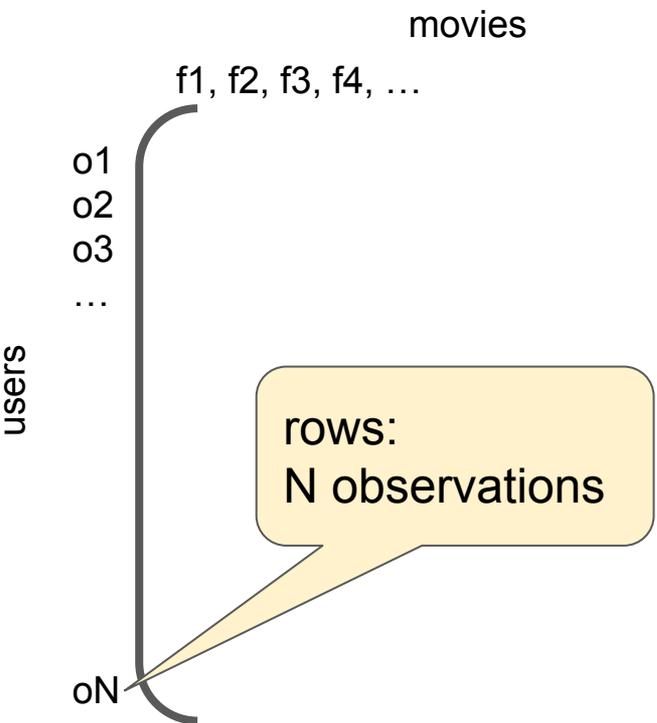
Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

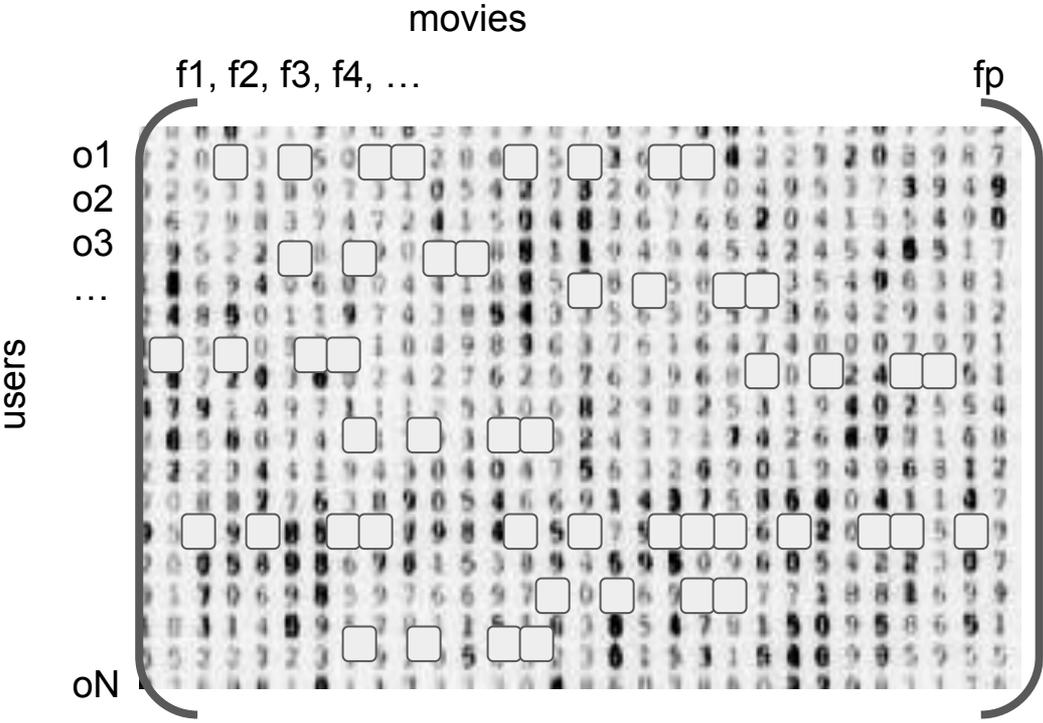
1. Content-based
2. Collaborative
3. **Latent Factor**

# Utility Matrix:

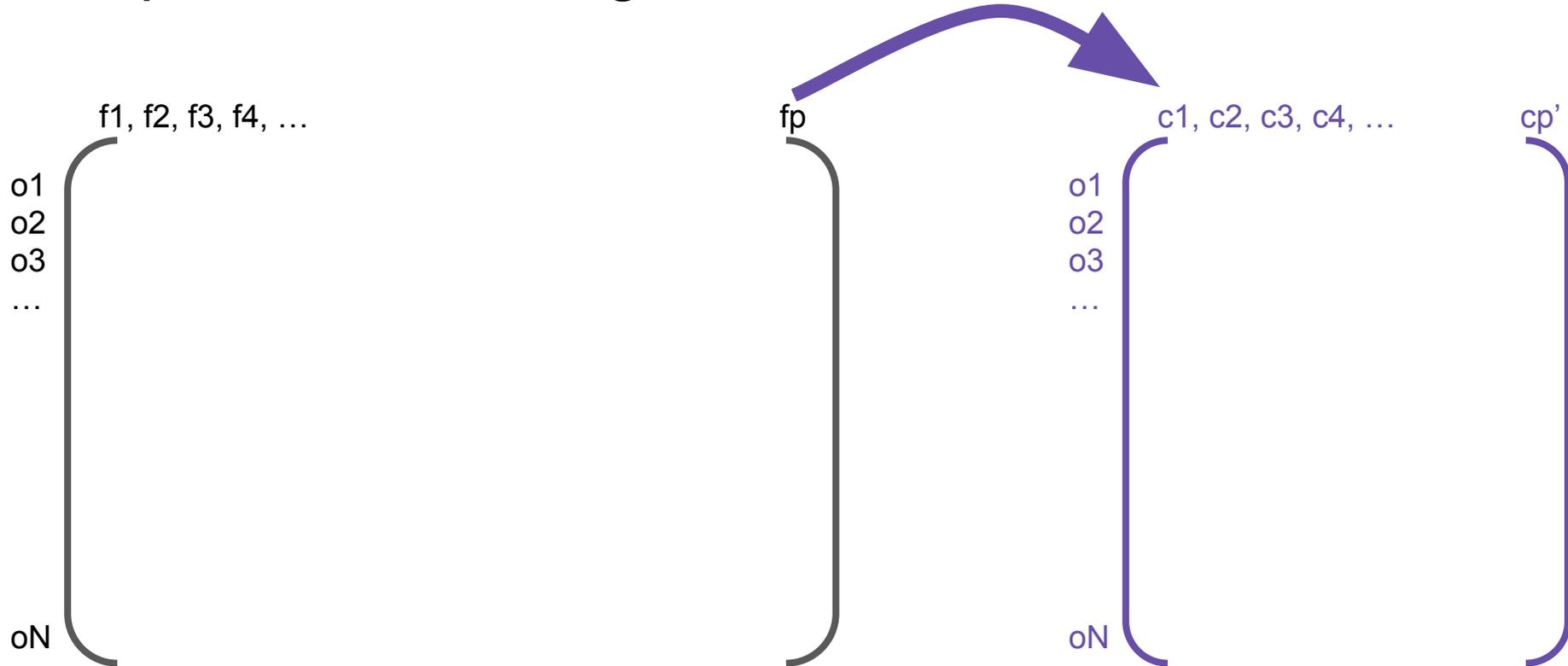




# Problem: Given Incomplete Matrix



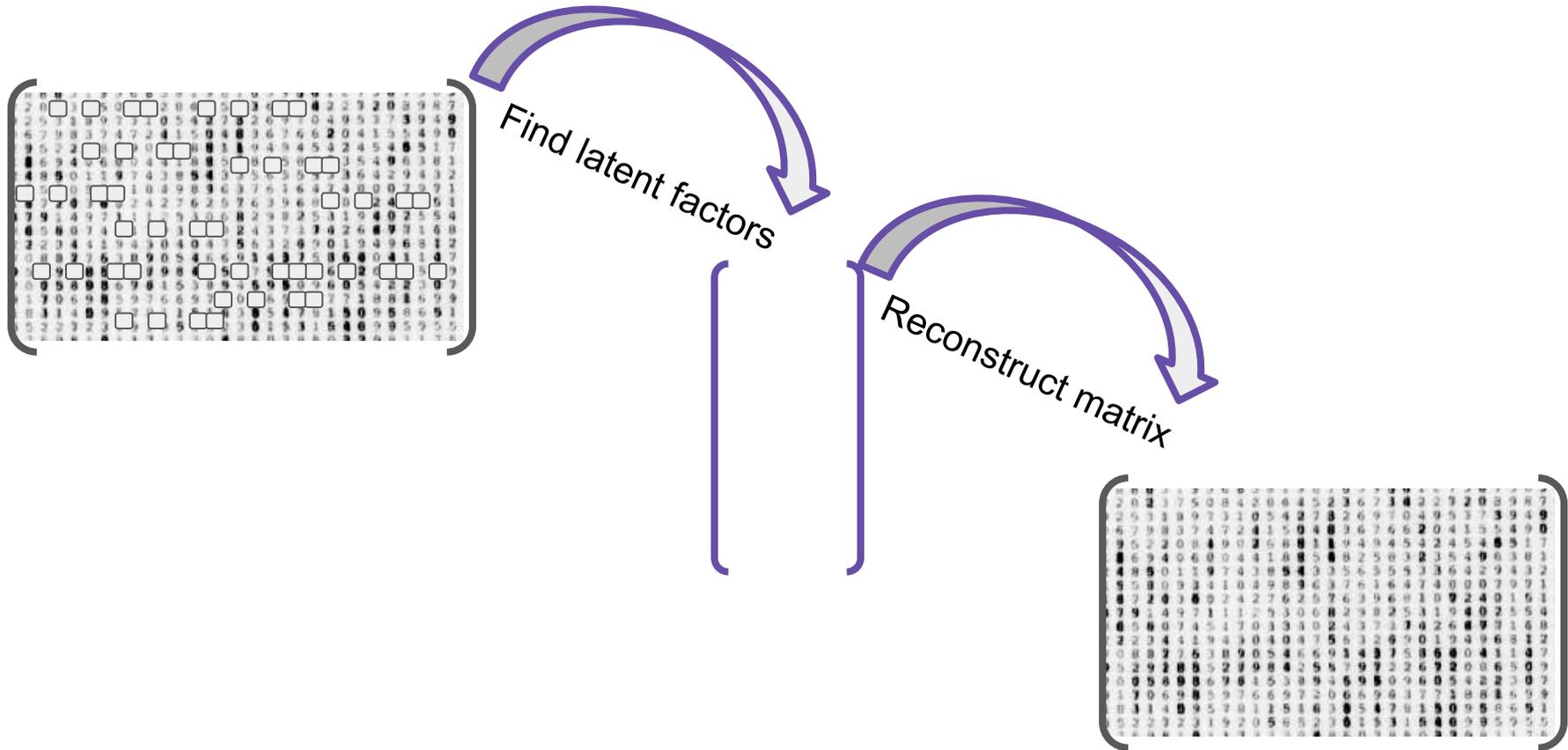
# Complete Matrix using Latent Factors



Dimensionality reduction

Try to best represent but with  $on$   $p'$  columns.

# Complete Matrix using Latent Factors



# Dimensionality Reduction - PCA

Linear approximates of data in  $r$  dimensions.

Found via *Singular Value Decomposition*:

$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

U: “left singular vectors”,

D: “singular values” (diagonal),

V: “right singular vectors”

Projection (dimensionality reduced space) in 3 dimensions:

$$(U_{[n \times 3]} D_{[3 \times 3]} V_{[p \times 3]}^T)$$

To reduce features in new dataset:

$$X_{\text{new}} V = X_{\text{new\_small}}$$

# Dimensionality Reduction - PCA

Linear approximates of data in  $r$  dimensions.

Found via *Singular Value Decomposition*:

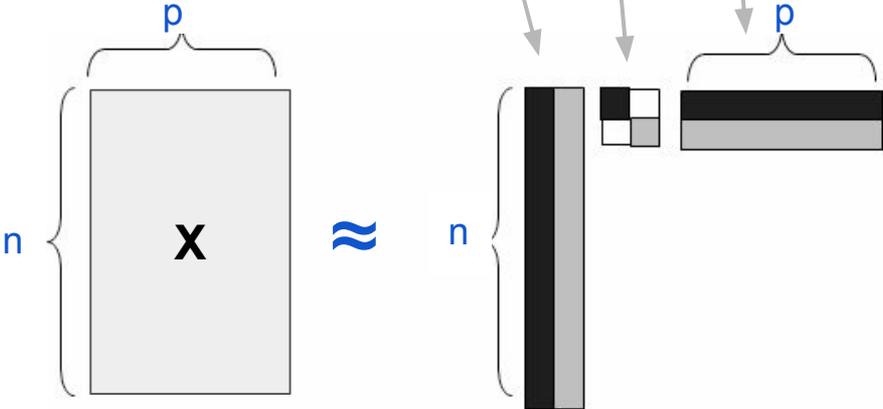
$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

D: “singular values” (diagonal),

U: “left singular vectors”,

V: “right singular vectors”



# Dimensionality Reduction - PCA - Example

$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

<div style="display: flex; flex-direction: column; align-items: center; gap: 10px;"> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">↑</span> </div> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">SciFi</span> </div> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">↓</span> </div> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">↑</span> </div> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">Romnce</span> </div> <div style="display: flex; align-items: center; gap: 5px;"> <span style="color: green;">↓</span> </div> </div>	$=$	<table border="1" style="border-collapse: collapse; text-align: center; margin: 0 auto;"> <thead> <tr> <th style="color: green; writing-mode: vertical-rl; transform: rotate(180deg);">Matrix</th> <th style="color: green; writing-mode: vertical-rl; transform: rotate(180deg);">Alien</th> <th style="color: green; writing-mode: vertical-rl; transform: rotate(180deg);">Serenity</th> <th style="color: green; writing-mode: vertical-rl; transform: rotate(180deg);">Casablanca</th> <th style="color: green; writing-mode: vertical-rl; transform: rotate(180deg);">Amelie</th> </tr> </thead> <tbody> <tr><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td></tr> <tr><td>3</td><td>3</td><td>3</td><td>0</td><td>0</td></tr> <tr><td>4</td><td>4</td><td>4</td><td>0</td><td>0</td></tr> <tr><td>5</td><td>5</td><td>5</td><td>0</td><td>0</td></tr> <tr><td>0</td><td>2</td><td>0</td><td>4</td><td>4</td></tr> <tr><td>0</td><td>0</td><td>0</td><td>5</td><td>5</td></tr> <tr><td>0</td><td>1</td><td>0</td><td>2</td><td>2</td></tr> </tbody> </table>	Matrix	Alien	Serenity	Casablanca	Amelie	1	1	1	0	0	3	3	3	0	0	4	4	4	0	0	5	5	5	0	0	0	2	0	4	4	0	0	0	5	5	0	1	0	2	2	$=$	<table border="1" style="border-collapse: collapse; text-align: center; margin: 0 auto;"> <tbody> <tr><td><b>0.13</b></td><td>0.02</td><td>-0.01</td></tr> <tr><td><b>0.41</b></td><td>0.07</td><td>-0.03</td></tr> <tr><td><b>0.55</b></td><td>0.09</td><td>-0.04</td></tr> <tr><td><b>0.68</b></td><td>0.11</td><td>-0.05</td></tr> <tr><td>0.15</td><td><b>-0.59</b></td><td><b>0.65</b></td></tr> <tr><td>0.07</td><td><b>-0.73</b></td><td><b>-0.67</b></td></tr> <tr><td>0.07</td><td><b>-0.29</b></td><td><b>0.32</b></td></tr> </tbody> </table>	<b>0.13</b>	0.02	-0.01	<b>0.41</b>	0.07	-0.03	<b>0.55</b>	0.09	-0.04	<b>0.68</b>	0.11	-0.05	0.15	<b>-0.59</b>	<b>0.65</b>	0.07	<b>-0.73</b>	<b>-0.67</b>	0.07	<b>-0.29</b>	<b>0.32</b>	$\times$	<table border="1" style="border-collapse: collapse; text-align: center; margin: 0 auto;"> <tbody> <tr><td><b>12.4</b></td><td>0</td><td>0</td></tr> <tr><td>0</td><td><b>9.5</b></td><td>0</td></tr> <tr><td>0</td><td>0</td><td><b>1.3</b></td></tr> </tbody> </table>	<b>12.4</b>	0	0	0	<b>9.5</b>	0	0	0	<b>1.3</b>	$\times$	<table border="1" style="border-collapse: collapse; text-align: center; margin: 0 auto;"> <tbody> <tr><td><b>0.56</b></td><td><b>0.59</b></td><td><b>0.56</b></td><td>0.09</td><td>0.09</td></tr> <tr><td>0.12</td><td>-0.02</td><td>0.12</td><td><b>-0.69</b></td><td><b>-0.69</b></td></tr> <tr><td>0.40</td><td><b>-0.80</b></td><td>0.40</td><td>0.09</td><td>0.09</td></tr> </tbody> </table>	<b>0.56</b>	<b>0.59</b>	<b>0.56</b>	0.09	0.09	0.12	-0.02	0.12	<b>-0.69</b>	<b>-0.69</b>	0.40	<b>-0.80</b>	0.40	0.09	0.09
Matrix	Alien	Serenity	Casablanca	Amelie																																																																																									
1	1	1	0	0																																																																																									
3	3	3	0	0																																																																																									
4	4	4	0	0																																																																																									
5	5	5	0	0																																																																																									
0	2	0	4	4																																																																																									
0	0	0	5	5																																																																																									
0	1	0	2	2																																																																																									
<b>0.13</b>	0.02	-0.01																																																																																											
<b>0.41</b>	0.07	-0.03																																																																																											
<b>0.55</b>	0.09	-0.04																																																																																											
<b>0.68</b>	0.11	-0.05																																																																																											
0.15	<b>-0.59</b>	<b>0.65</b>																																																																																											
0.07	<b>-0.73</b>	<b>-0.67</b>																																																																																											
0.07	<b>-0.29</b>	<b>0.32</b>																																																																																											
<b>12.4</b>	0	0																																																																																											
0	<b>9.5</b>	0																																																																																											
0	0	<b>1.3</b>																																																																																											
<b>0.56</b>	<b>0.59</b>	<b>0.56</b>	0.09	0.09																																																																																									
0.12	-0.02	0.12	<b>-0.69</b>	<b>-0.69</b>																																																																																									
0.40	<b>-0.80</b>	0.40	0.09	0.09																																																																																									

Users to movies matrix

# Dimensionality Reduction - PCA

Linear approximates of data in  $r$  dimensions.

Found via *Singular Value Decomposition*:

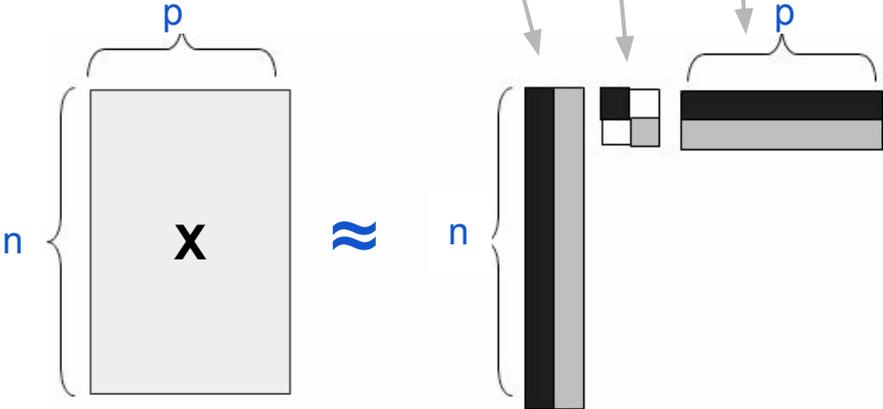
$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$

X: original matrix,

D: “singular values” (diagonal),

U: “left singular vectors”,

V: “right singular vectors”



- **Goal: Minimize the sum of reconstruction errors:**

$$\sum_{i=1}^N \sum_{j=1}^D \|x_{ij} - z_{ij}\|^2$$

- where  $x_{ij}$  are the “old” and  $z_{ij}$  are the “new” coordinates

X: original matrix  
D: “singular values”

s.

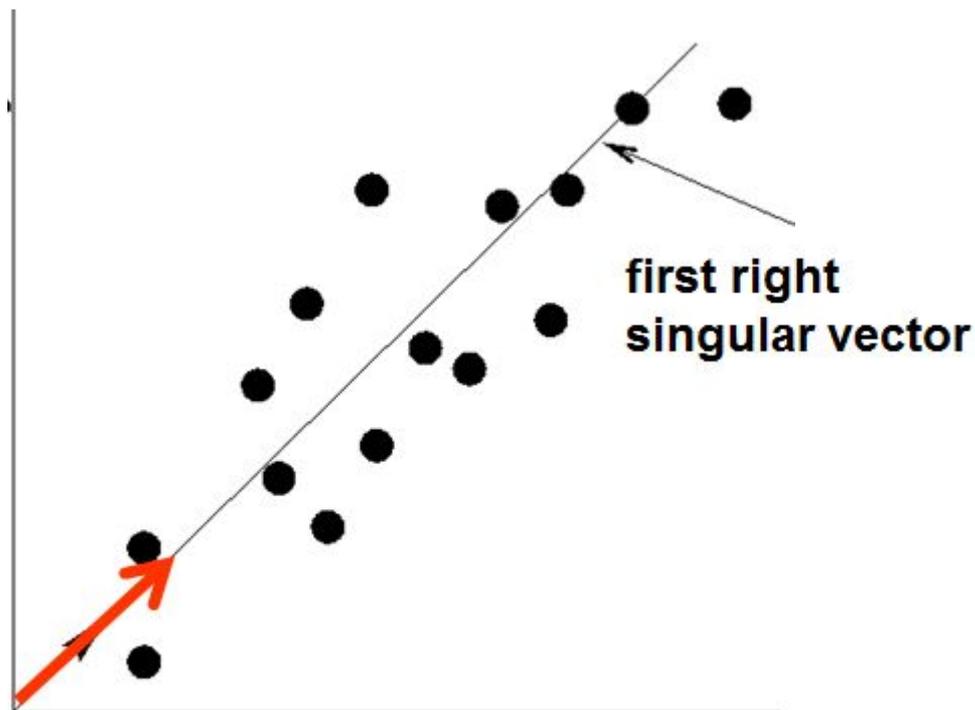
“singular vectors”,  
“singular vectors”

To check how well the original matrix can be reproduced:

$$Z_{[n \times p]} = U D V^T, \text{ How does } Z \text{ compare to original } X?$$

# Dimensionality Reduction - PCA - Example

$$X_{[n \times p]} = U_{[n \times r]} D_{[r \times r]} V_{[p \times r]}^T$$



# Recommendation Systems

Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

1. Content-based
2. Collaborative
3. **Latent Factor**

# Recommendation Systems

Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

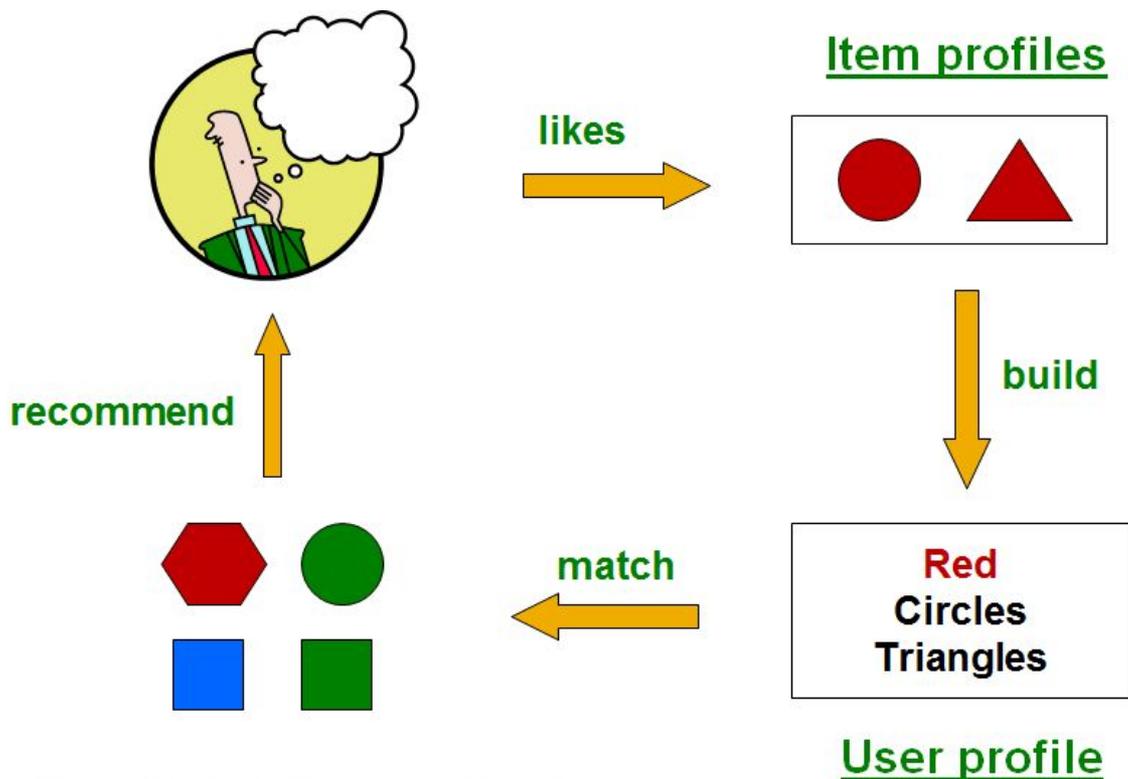
1. **Content-based**
2. Collaborative
3. Latent Factor

# Content-based Rec Systems

Based on similarity of items to past items that they have rated.

# Content-based Rec Systems

Based on similarity of items to past items that they have rated.



# Content-based Rec Systems

Based on similarity of items to past items that they have rated.

1. Build profiles of items (set of features); examples:

*shows*: producer, actors, theme, review

*people*: friends, posts

pick words with tf-idf



# Content-based Rec Systems

Based on similarity of items to past items that they have rated.

1. Build profiles of items (set of features); examples:

*shows*: producer, actors, theme, review

*people*: friends, posts

pick words with tf-idf

2. Construct user profile from item profiles; approach:

average all item profiles

variation: weight by difference from their average

# Content-based Rec Systems

Based on similarity of items to past items that they have rated.

1. Build profiles of items (set of features); examples:

*shows*: producer, actors, theme, review

*people*: friends, posts

pick words with tf-idf

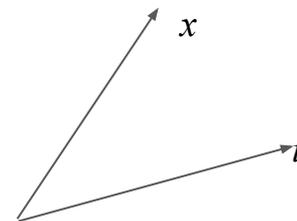
2. Construct user profile from item profiles; approach:

average all item profiles of items they've purchased

variation: weight by difference from their average ratings

3. Predict ratings for new items; approach:

$$utility(user, i) = \cos(x, i) = \frac{x \cdot i}{||x|| \cdot ||i||}$$



# Why Content Based?

- Only need users history
- Captures unique tastes
- Can recommend new items
- Can provide explanations

# Why Content Based?

- Only need users history
- Captures unique tastes
- Can recommend new items
- Can provide explanations
- Need good features
- New users don't have history
- Doesn't venture "outside the box"  
(Overspecialized)

# Why Content Based?

- Only need users history
- Captures unique tastes
- Can recommend new items
- Can provide explanations
- Need good features
- New users don't have history
- Doesn't venture "outside the box"  
(Overspecialized)

(not exploiting other users judgments)

# Collaborative Filtering Rec Systems



(not exploiting other users judgments)

# Recommendation Systems

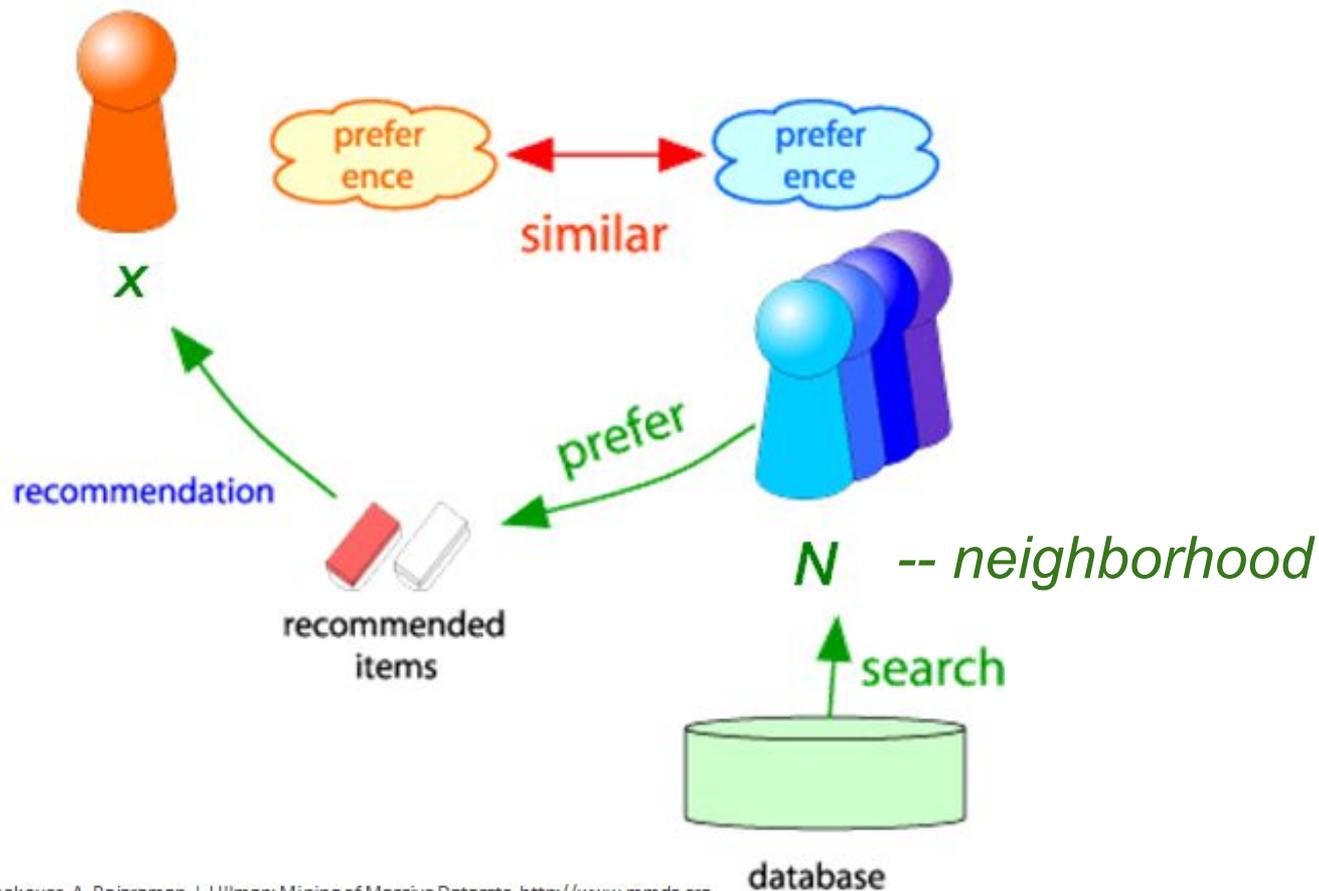
Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

1. Content-based
2. **Collaborative**
3. Latent Factor

# Collaborative Filtering Rec Systems



# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	2		3
<i>B</i>	5			4	2
<i>C</i>			5	2	

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	2		3
<i>B</i>	5			4	2
<i>C</i>			5	2	

*General Idea:*

- 1) Find similar users = "neighborhood"*
- 2) Infer rating based on how similar users rated*

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	2		3
<i>B</i>	5			4	2
<i>C</i>			5	2	

Given: *user, x; item, i; utility matrix, u*

1. Find neighborhood,  $N$  # set of  $k$  users most similar to  $x$  who have also rated  $i$

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
<i>A</i>	4	5	2		3
<i>B</i>	5			4	2
<i>C</i>			5	2	

Given: *user, x; item, i; utility matrix, u*

1. Find neighborhood,  $N$  # set of  $k$  users most similar to  $x$  who have also rated  $i$

*Two Challenges: (1) user bias, (2) missing values*

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
A	4 => 0.5	5 => 1.5	2 => -1.5	=> 0	3 => -0.5
B	5			4	2
C			5	2	

Given: *user, x; item, i; utility matrix, u*

1. Find neighborhood,  $N$  # set of  $k$  users most similar to  $x$  who have also rated  $i$

*Two Challenges: (1) user bias, (2) missing values*

*Solution: subtract user's mean, add zeros for missing*

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
A	4 => 0.5	5 => 1.5	2 => -1.5	=> 0	3 => -0.5
B	5			4	2
C			5	2	

- Given: *user, x; item, i; utility matrix, u*
0. Update *u*: mean center, missing to 0
  1. Find neighborhood, *N* # set of *k* users most similar to *x* who have also rated *i*
    - $\text{sim}(x, \text{other}) = \text{cosine\_sim}(u[x], u[\text{other}])$
    - threshold to top *k* (e.g.  $k = 30$ )

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
A	4 => 0.5	5 => 1.5	2 => -1.5	=> 0	3 => -0.5
B	5			4	2
C			5	2	

- Given: *user, x; item, i; utility matrix, u*
0. Update *u*: mean center, missing to 0
  1. Find neighborhood, *N* # set of *k* users most similar to *x* who have also rated *i*
    - $\text{sim}(x, \text{other}) = \text{cosine\_sim}(u[x], u[\text{other}])$
    - threshold to top *k* (e.g.  $k = 30$ )
  2. Predict utility (rating) of *i* based on *N*

# Collaborative Filtering Rec Systems

<i>user</i>	Game of Thrones	Fargo	Brooklyn Nine-Nine	Silicon Valley	Walking Dead
A	4 => 0.5	5 => 1.5	2 => -1.5	=> 0	3 => -0.5
B	5			4	2
C			5	2	

Given: *user, x; item, i; utility matrix, u*

0. Update *u*: mean center, missing to 0

1. Find neighborhood, *N* # set of *k* users most similar to *x* who have also rated *i*

--  $\text{sim}(x, \text{other}) = \text{cosine\_sim}(u[x], u[\text{other}])$

-- threshold to top *k* (e.g.  $k = 30$ )

2. Predict utility (rating) of *i* based on *N*

-- average, weighted by sim 
$$\text{utility}(x, i) = \frac{\sum_{y \in N} \text{Sim}(x, y) \cdot \text{utility}(y, i)}{\sum_{y \in N} \text{Sim}(x, y)}$$

# Collaborative Filtering Rec Systems

“User-User collaborative filtering”



Given: *user, x; item, i; utility matrix, u*

0. Update *u*: mean center, missing to  $\emptyset$
1. Find neighborhood, *N* # set of *k* users most similar to *x* who have also rated *i*
  - $\text{sim}(x, \text{other}) = \text{cosine\_sim}(u[x], u[\text{other}])$
  - threshold to top *k* (e.g.  $k = 30$ )
2. Predict utility (rating) of *i* based on *N*
  - average, weighted by sim 
$$\text{utility}(x, i) = \frac{\sum_{y \in N} \text{Sim}(x, y) \cdot \text{utility}(y, i)}{\sum_{y \in N} \text{Sim}(x, y)}$$

# Collaborative Filtering Rec Systems

## “User-User collaborative filtering”

### Item-Item:

Flip rows/columns of utility matrix and use same methods.  
(i.e. estimate rating of item  $i$ , by finding similar items,  $j$ )

Given: *user*,  $x$ ; *item*,  $i$ ; *utility matrix*,  $u$

0. Update  $u$ : mean center, missing to  $\emptyset$
1. Find neighborhood,  $N$  # set of  $k$  users most similar to  $x$  who have also rated  $i$ 
  - $\text{sim}(x, \text{other}) = \text{cosine\_sim}(u[x], u[\text{other}])$
  - threshold to top  $k$  (e.g.  $k = 30$ )
2. Predict utility (rating) of  $i$  based on  $N$ 
  - average, weighted by sim

$$\text{utility}(x, i) = \frac{\sum_{y \in N} \text{Sim}(x, y) \cdot \text{utility}(y, i)}{\sum_{y \in N} \text{Sim}(x, y)}$$

# Collaborative Filtering Rec Systems

## “User-User collaborative filtering”

### Item-Item:

Flip rows/columns of utility matrix and use same methods.  
(i.e. estimate rating of item  $i$ , by finding similar items,  $j$ )

Given: *user*,  $x$ ; *item*,  $i$ ; *utility matrix*,  $u$

0. Update  $u$ : mean center, missing to  $\emptyset$

1. Find neighborhood,  $N$  # set of  $k$  **items** most similar to  $i$  also rated by  $x$

--  $\text{sim}(i, \text{other}) = \text{cosine\_sim}(u[i], u[\text{other}])$

-- threshold to top  $k$  (e.g.  $k = 30$ )

2. Predict utility (rating) by  $x$  based on  $N$

-- average, weighted by sim

$$\text{utility}(x, i) = \frac{\sum_{j \in N} \text{Sim}(i, j) \cdot \text{utility}(x, j)}{\sum_{j \in N} \text{Sim}(i, j)}$$

## Item-Item v User-User

**Item-item often works better than user-user. Why?**

Users tend to be more different from each other than items are from other items.

*e.g. Mary likes jazz + rock, Bob likes classical + rock,  
but Mary may still have same rock preferences as Bob*

## Item-Item v User-User

**Item-item often works better than user-user. Why?**

Users tend to be more different from each other than items are from other items.

*e.g. Mary likes jazz + rock, Bob likes classical + rock,  
but Mary may still have same rock preferences as Bob*

***In other words, users span genres but items usually do not.***

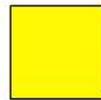
# Item-Item: Example

movies

	1	2	3	4	5	6	7	8	9	10	11	12
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

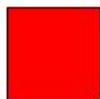


- rating between 1 to 5

# Item-Item: Example

	1	2	3	4	5	6	7	8	9	10	11	12
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	

movies



- estimate rating of movie **1** by user **5**

# Item-Item: Example

	1	2	3	4	5	6	7	8	9	10	11	12	$\text{sim}(1,m)$
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>

movies

Same as cosine sim when subtracting the mean

## Neighbor selection:

Identify movies similar to movie 1, rated by user 5

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: Example

	1	2	3	4	5	6	7	8	9	10	11	12	$\text{sim}(1,m)$
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>

Compute similarity weights:

$$s_{1,3}=0.41, s_{1,6}=0.59$$

# Item-Item: Example

	1	2	3	4	5	6	7	8	9	10	11	12	$\text{sim}(1,m)$
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>

$$\text{utility}(1, 5) = (0.41 \cdot 2 + 0.59 \cdot 3) / (0.41 + 0.59)$$

$$\text{utility}(x, i) = \frac{\sum_{j \in N} \text{Sim}(i, j) \cdot \text{utility}(x, j)}{\sum_{j \in N} \text{Sim}(i, j)}$$

# Recommendation Systems

Problems to tackle:

1. Gathering ratings
2. Extrapolate unknown ratings
  - a. Explicit: based on user ratings and reviews  
(problem: only a few users engage in such tasks)
  - b. Implicit: Learn from actions (e.g. purchases, clicks)  
(problem: hard to learn low ratings)
3. Evaluation

## Common Approaches

1. Content-based
2. Collaborative
3. Latent Factor

# Options for Parallelizing

1. Approximate solutions to PCA (very large speedups with little drawback!):
  - a. **Stochastic Sampling** (also sometimes called "randomized" which is ambiguous):  
Only using a sample rows (i.e. users for recommendation systems)
  - b. **Truncated SVD**: Only optimizing for minimizing reconstruction error based on up to  $r$  dimensions (full SVD solves for up to  $\min(n, p)$  dimensions and then you just truncate the result for the lower rank version). One you do this, by the way, using a smaller sample becomes much less of a problem.
  - c. **Limiting power iterations to a few iterations**: Power iterations from pagerank solves for the first principle component. This can be extended to multiple components.  
(more [here](#).)
2. Distribute the matrix operations. Complex; not as flexible (usually done across processors within node)
3. Data Parallelism: As in other instances stochastic or mini-batch gradient descent.