



Machine Learning with Graphs (MLG)

# Recommender Systems (RecSys)

Basic & Collaborative Filtering

Cheng-Te Li (李政德)

Institute of Data Science  
National Cheng Kung University

[chengte@mail.ncku.edu.tw](mailto:chengte@mail.ncku.edu.tw)





# YouTube

## 推薦影片



GS Warriors vs Houston  
Rockets - Full Game...



Golden State Warriors vs  
Houston Rockets Full Game...



MAYDAY五月天【終於結束的  
起點 Beginning of the End】



于文文《體面》動態歌詞版  
【前任3:再見前任 插曲】



PIGFISH Music Channel  
Top Music KKBOX



柯文哲 我的故事 (全)~柯P來  
輔大

Rapid Highlights  
觀看次數：15萬 · 2 小時前

MLG Highlights  
觀看次數：34萬 · 2 小時前

相信音樂BinMusic  
觀看次數：86萬 · 5 個月前

PIGFISH Music Channel  
觀看次數：3,443萬 · 3 個月前

Top Music KKBOX  
2,002 人正在觀看

fonzae  
觀看次數：7.9萬 · 1 個月前

直播中



田馥甄 - 愛了很久的朋友  
(lyrics)



Lebron James 'Long Live The  
King' Motivational Workout



PROFESSOR VS CONMAN  
Tom "Conman" Connors



Michael Jordan Crazy  
Tomahawk Leaner Dunk vs...



Federer's Historical French  
Open Title - Best Points



海角七號(野玫瑰) 范逸臣&中  
孝介

XIANGLIN HUANG  
觀看次數：12萬 · 1 個月前

basketballprosworkouts  
觀看次數：638萬 · 11 個月前

觀看次數：61萬 · 2 週前

YouToobe  
觀看次數：239萬 · 4 個月前

Valzy  
觀看次數：5.5萬 · 2 天前

De Moore  
觀看次數：69萬 · 9 年前



'Still KD' Kevin Durant  
Workout With Steve Nash...  
basketballprosworkouts  
觀看次數：15萬 · 6 個月前



2018新歌排行榜 (2018最新  
歌曲,歌曲排行榜2018) 201...  
Jean Huang  
觀看次數：478萬 · 4 個月前



政大我在台灣留學的大學  
NCCU My university during...  
Julie Flower  
觀看次數：2.6萬 · 3 週前



Golden State Warriors vs  
Houston Rockets Full Game...  
Ximo Pierto  
觀看次數：320 · 1 小時前



Mayday五月天[步步Step by  
Step]MV官方高音質HD版-電...  
相信音樂BinMusic  
觀看次數：3,127萬 · 4 年前



司马懿劝阻曹叡弑母  
China Zone – 大军师司马懿之...  
觀看次數：2.2萬 · 2 個月前

Watch Instantly

Browse DVDs

Your Queue

Movies You'll ❤️

## Congratulations! Movies we think You will ❤️

Add movies to your Queue, or Rate ones you've seen for even better suggestions.

Spider-Man 3

[Add](#)A rating scale with five stars, where the first four are filled red and the last one is white.  
 Not Interested

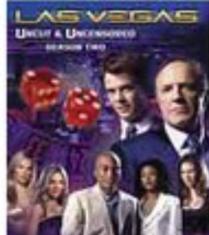
300

[Add](#)A rating scale with five stars, all filled red.  
 Not Interested

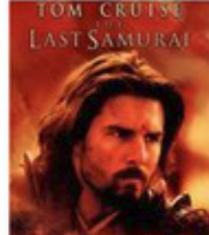
The Rundown

[Add](#)A rating scale with five stars, where the first four are filled red and the last one is white.  
 Not Interested

Bad Boys II

[Add](#)A rating scale with five stars, all filled red.  
 Not InterestedLas Vegas: Season 2  
(6-Disc Series)

The Last Samurai



Star Wars: Episode III

Robot Chicken: Season 3  
(2-Disc Series)



# facebook

You are surrounded by recommendations

Search

News Feed Top News · Most Recent 300+

Share: Status  What's on your mind?

**Rehabber Urban Winery Dream** [www.thefeast.com](http://www.thefeast.com)  
Few venture through the blocks of vacant land, boarded up windows, and overgrown weeds sprawling across the West Side neighborhoods unless for a special visit to Garfield Park Conservatory or one of the handful...  
  
22 minutes ago · Like · Comment · Share

Thanks This is a little bittersweet.  
**The Sad, Beautiful Fact That We're All Going To Miss Almost Everything : Monkey See : NPR** [www.npr.org](http://www.npr.org)  
We take a moment to reflect on being "well-read," and

**塗鴉牆貼文推薦**

commented on photo.

about an hour ago  
[View all 3 comments](#)

You look beautiful! Congratulations.  
about an hour ago

I am so overwhelmed!!! People from all over the country have sent me facebook messages and/or called me to wish me a happy birthday. I want to you know that it's because of YOU that I had such a happy birthday!  
**THANKS SO MUCH!**  
39 minutes ago · Like · Comment  
[2 people like this.](#)

Leng Home 99+

Ah Mol Love's birthday is today  
1 request from Chan Youvireak  
Love@: www.jRoong.com today

RECOMMENDED PAGES

Computerworld Kim Sung and 13 other friends like this. [Like](#)

NBC Connecticut Yareth Mith likes this. [Like](#)

DOL [Like](#)

Marouane Fellaini Piseth Pen and 11 other friends like him. [Like](#)

Manny Pacquiao Sopheap Son and 13 other friends like this. [Like](#)

Polygon 57,799 people like this. [Like](#)

The Next Web Meng Leang Hour and 4 other friends like this. [Like](#)

取消 打卡 回報

Q 搜尋地標

感覺攝影工作室 2F, No.66, Sec. 2, Minsheng E. Rd., Zhongs...

米拉藝視。視覺攝影工作室 台北市民生東路二段61號二樓 - 171 個打卡次...

巴郎子 新疆風味料理 民生東路二段72號 - 2,735 個打卡次 - 0 公尺

這一 12 - 台北市中山區吉林路190號 - 12,722 個打卡次...

達摩動館 台北市吉林路188號 - 1,516 個打卡次 - 40 公尺

史記精緻鴛鴦鍋 民生東路2段62號1樓 - 11,396 個打卡次 - 50...

史記正宗牛肉麵 中山區民生東路二段60號 - 6,670 個打卡次 - 3...

浪漫一生 新北市 - 6,118 個打卡次 - 40 公尺

**粉絲頁推薦**

**打卡地點推薦**

朋友推薦

People You May Know [see all](#)

Kuldip K. Ambastha [Add to Friends](#)

Jason Fang [Add to Friends](#)

Erica Garcia [遊戲推薦](#)

你可能會喜歡的遊戲 [顯示全部](#)

KartRider Dash 100,000 個玩家 [馬上玩](#)

建議的社團 [查看全部社團](#)

台北市暨新北市好吃好玩去處 [+ 加入](#)

**社團推薦**

贊助 [刊登廣告](#)

超人氣影片

請移動一根火柴棒，讓等式成立！謎題揭曉

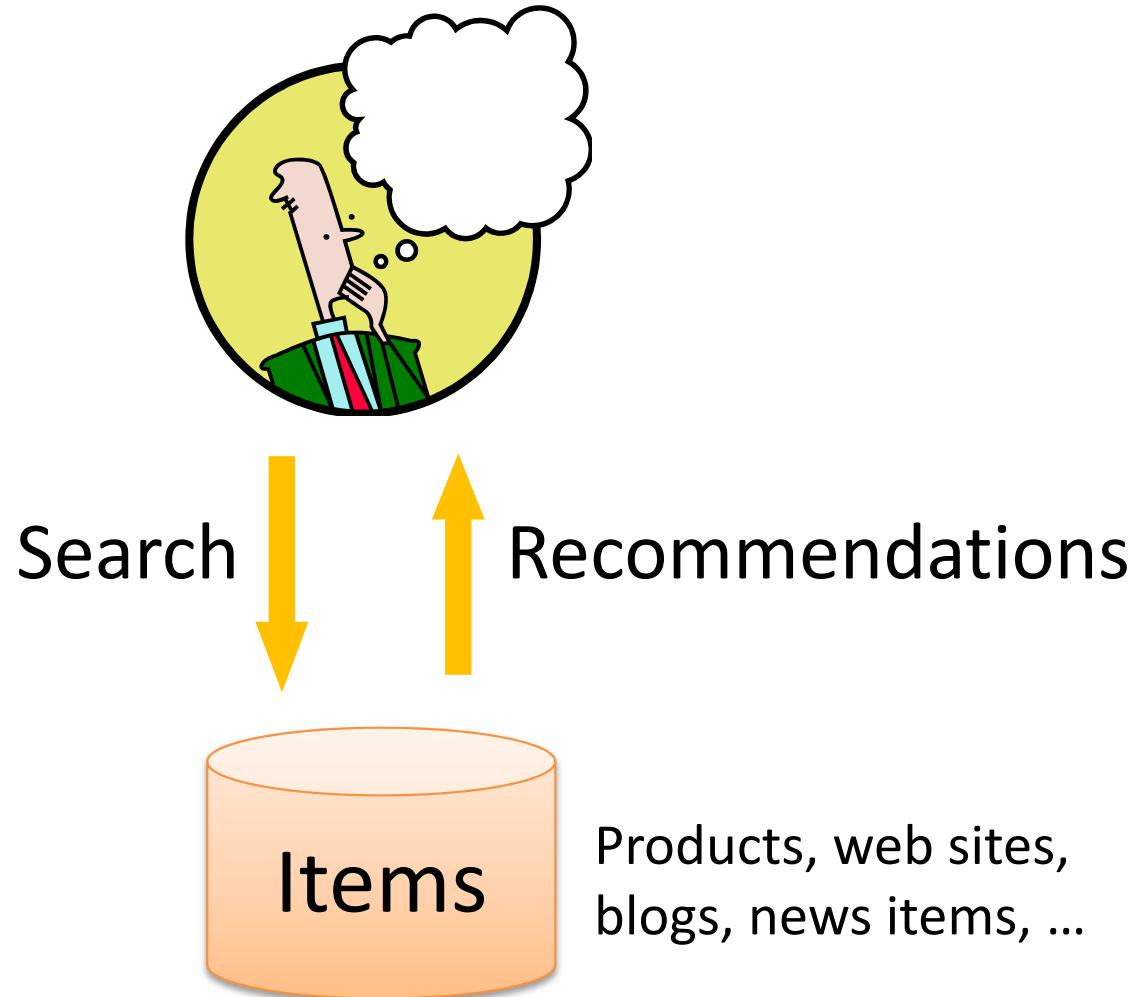
►► <http://113.com.tw/r/S3x>

**廣告推薦**

# Recommender Systems (RecSys)

- RS analyzes **patterns of user interest** in items to provide personalized recommendations
  - Many users view the same movie and each user is likely to view numerous different movies
  - Huge volume of data arise from **user feedbacks** that can be analyzed to provide recommendations
- Goal of RecSys: **predict the rating or preference that user would give to an item**
  - RecSys is very useful for commercial items such as movies, music, books, apps, and TV shows

# 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  Our Target!
  - Amazon, Netflix, YouTube, Instagram, Facebook ...

# Key Problems

- (1) Gathering “known” ratings for matrix
  - How to collect the data about user-item interactions?
- (2) Predicting unknown ratings from known ones
  - Mainly interested in high unknown ratings
    - Not interested in what you don’t like but what you like
- (3) Evaluating extrapolation methods
  - How to measure success/performance of recommendation methods

# Gathering Ratings

- **Explicit**
  - Ask people to rate items
    - Do not work well in practice -- people can't be bothered
  - Crowdsourcing: pay people to label items
  - Incentive (e.g., discounts, coupons) in the websites
- **Implicit**
  - Learn ratings from user actions
    - E.g., purchase implies high ratings
    - E.g., browsing history (page view)
    - E.g., clicks

# Predicting Ratings

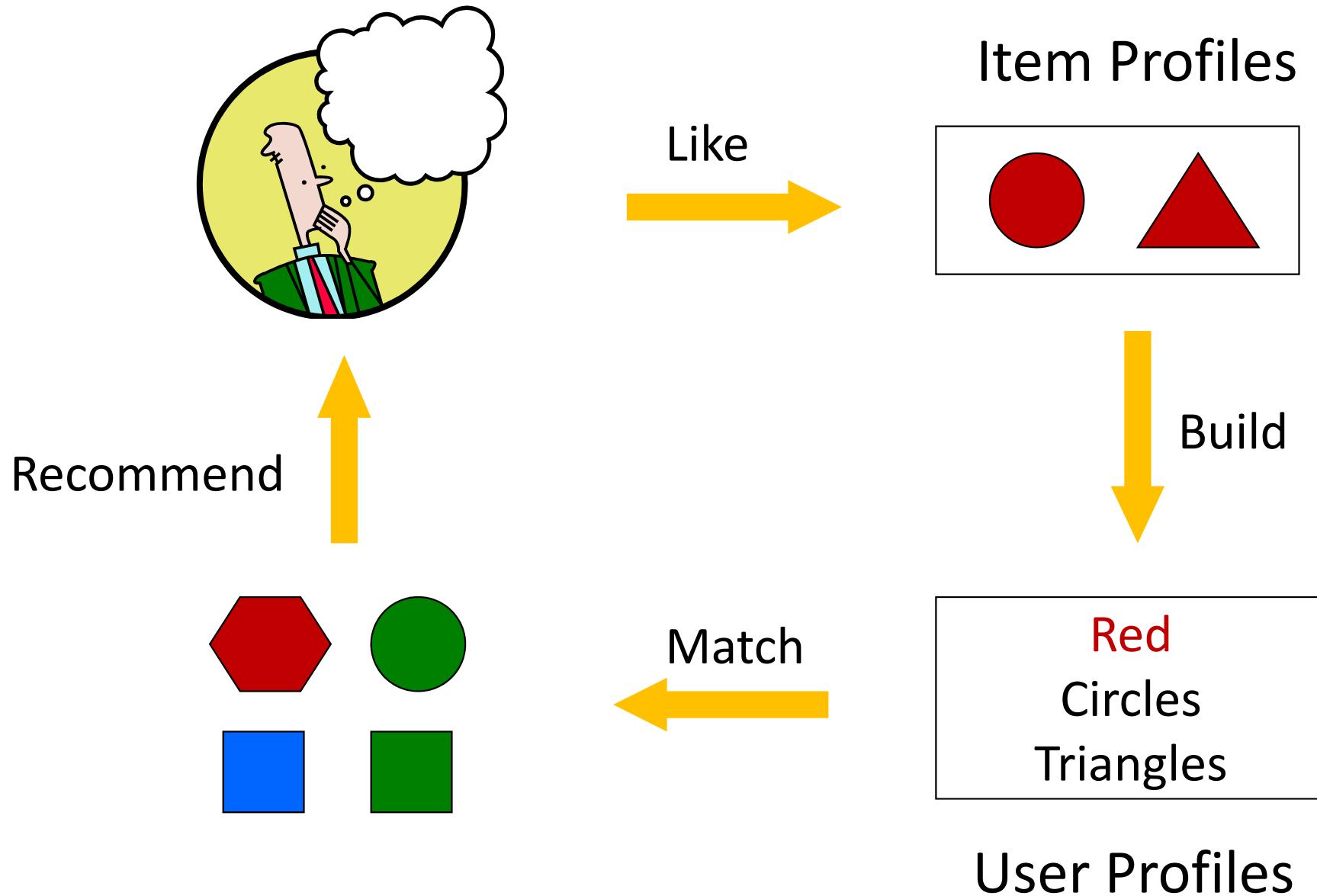
- Key problem: Rating matrix  $R$  is **sparse**
  - Most people have not rated most items
  - **Cold start**
    - New items have no ratings
    - New users have no history
- Three approaches to recommender systems:
  - 1) Content-based Recommendation (CR)
  - 2) Collaborative Filtering (CF)
  - 3) Latent Factor Model (e.g., MF, FM)

# Content-based Recommendation

- Main idea: Recommend items to customer  $u$  similar to previous items rated highly by  $u$
- E.g., Movie recommendations
  - Recommend movies with the same actor(s), director, genre, etc.
- E.g., Websites, blogs, news
  - Recommend other sites with “similar” content



# Plan of Actions



# Item Profiles

- For each item, create an **item profile**
- **Profile is a vector of features**
  - **Movies:** author, title, actor, director,...
  - **Text:** Set of “important” words in document
- **How to pick important features?**
  - Usual heuristic from text mining is **TF-IDF**  
(Term Frequency  $\times$  Inverse Document Frequency)
    - Term  $\rightarrow$  Feature
    - Document  $\rightarrow$  Item

# User Profiles and Prediction

- User profile possibilities:
  - Average of rated item profiles
  - Weighted by difference from average rating for item
- Prediction heuristic:  
**Cosine similarity of user and item profiles**
  - Given user profile  $x_u$  and item profile  $x_i$ , estimate

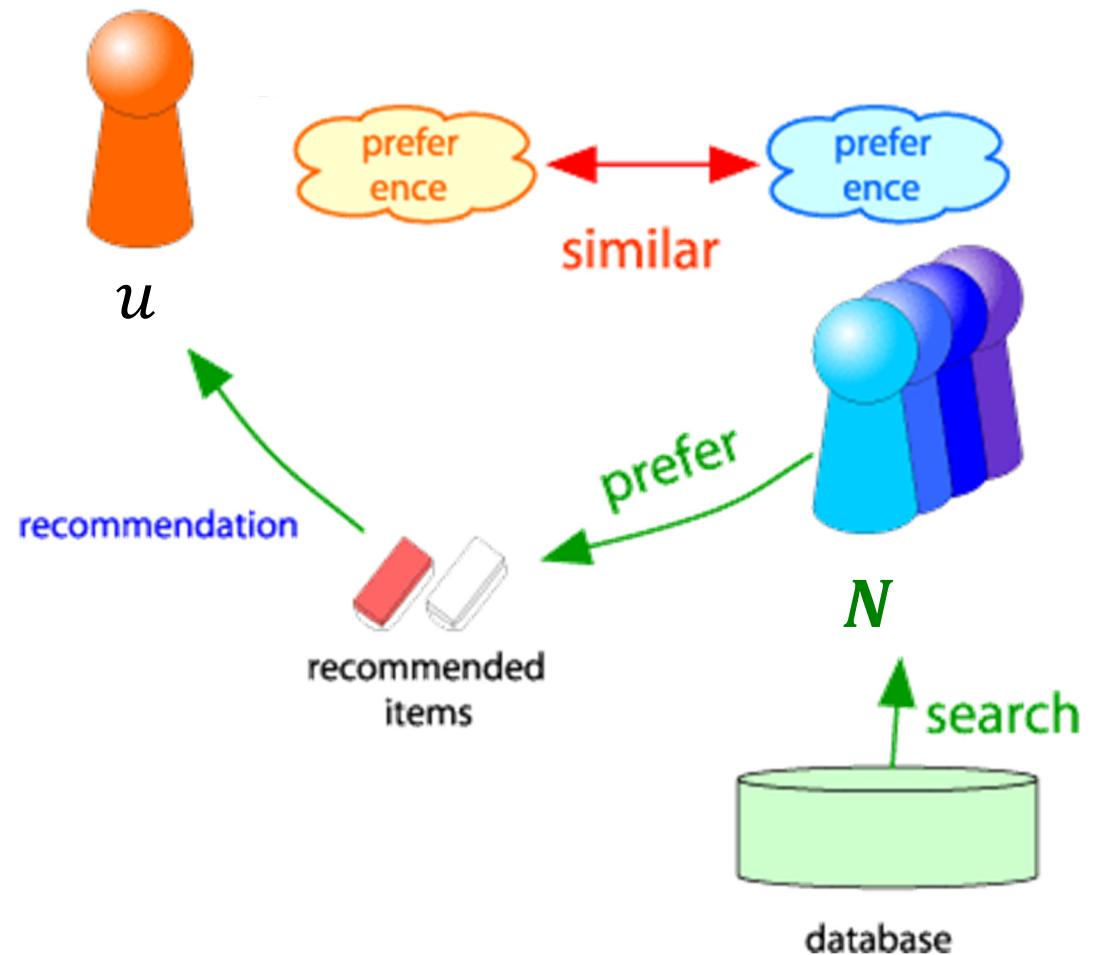
$$\cos(x_u, x_i) = \frac{x_u \cdot x_i}{\|x_u\| \|x_i\|}$$

# Pros/Cons of CR

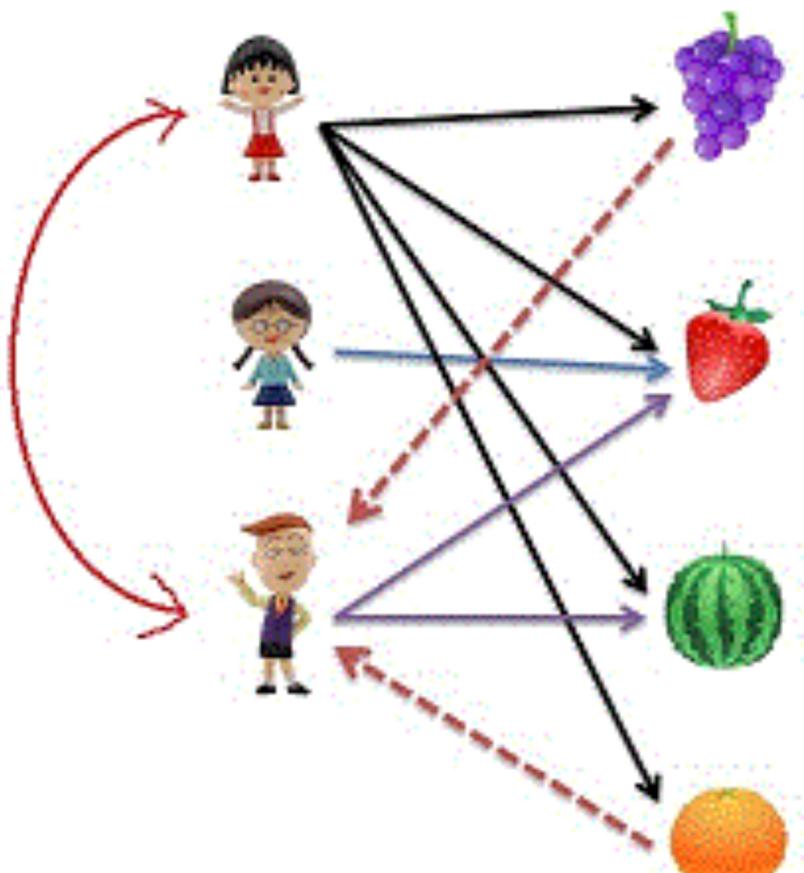
- +: No need for data on other users
  - No cold-start or sparsity problems
- +: Able to recommend to users with unique tastes
- +: Able to recommend new & unpopular items
  - No first-rater problem
- +: Able to provide explanations
  - Can provide explanations of recommended items by listing content-features that caused an item to be recommended
- -: 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

# Collaborative Filtering (CF)

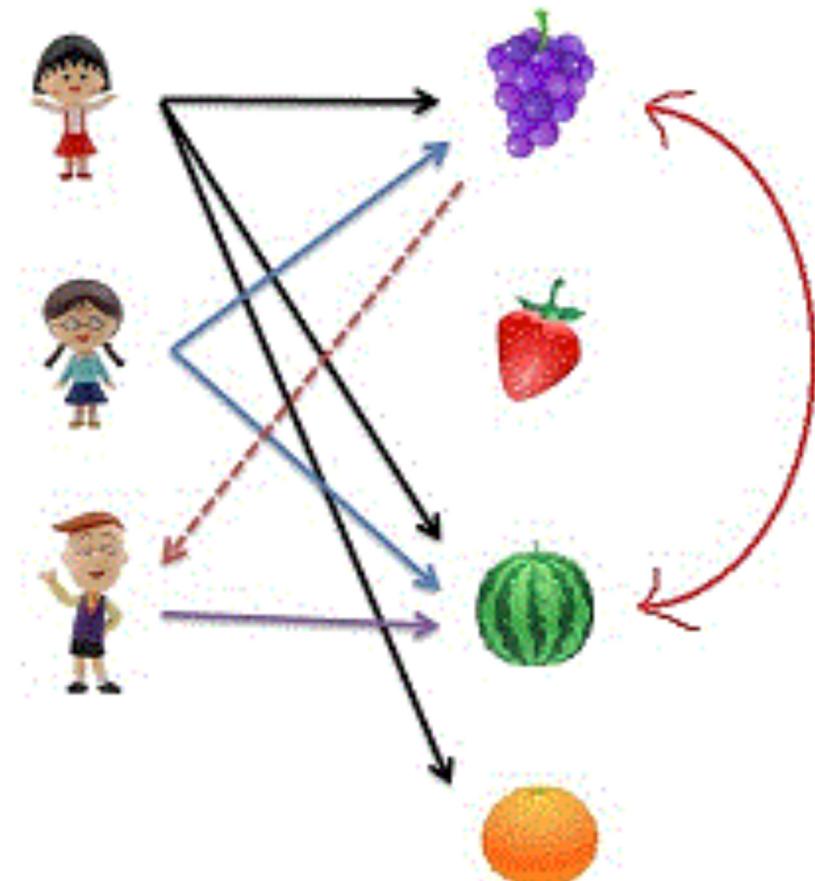
- Consider user  $u$
- Find set  $N$  of other users whose ratings are **similar** to  $u$ 's ratings
- Estimate  $u$ 's ratings based on ratings of users in  $N$



# Collaborative Filtering (CF)



User-based CF



Item-based CF

# User-based CF: Finding Similar Users

- Let  $r_u$  be the vector of user  $u$ 's ratings

- Cosine similarity measure**

- $\text{sim}(u, v) = \cos(r_u, r_v) = \frac{r_u \cdot r_v}{\|r_u\| \|r_v\|}$

$$\begin{aligned} r_u &= \{1, 0, 0, 1, 3\} \\ r_v &= \{1, 0, 2, 2, 0\} \end{aligned}$$

- Problem:** treat missing ratings as “negative” (dissimilar)

- Pearson Correlation Coefficient**

- $S_{uv}$ : items rated by both users  $u$  and  $v$

$$\text{sim}(u, v) = \frac{\sum_{i \in S_{uv}} (r_{ui} - \bar{r}_u)(r_{vi} - \bar{r}_v)}{\sqrt{\sum_{i \in S_{uv}} (r_{ui} - \bar{r}_u)^2} \sqrt{\sum_{i \in S_{uv}} (r_{vi} - \bar{r}_v)^2}}$$

$\bar{r}_u, \bar{r}_v$ : the average ratings of users  $u$  and  $v$

# Cosine 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:  $\text{sim}(A, B) > \text{sim}(A, C)$
- Cosine similarity:**  $0.380 > 0.322$ 
  - Considers missing ratings as “negative”
  - Solution: subtract the (row) mean**

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

	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

# Rating Predictions

From similarity metric to recommendations:

- Let  $r_u$  be the vector of user  $u$ 's ratings
- Let  $N$  be the set of  $k$  users most similar to  $u$  who have rated item  $i$
- Prediction for item  $i$  of user  $u$ :

- $r_{ui} = \frac{1}{k} \sum_{v \in N} r_{vi}$

Shorthand:

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

- Or even better:  $r_{ui} = \frac{\sum_{v \in N} s_{uv} \cdot r_{vi}}{\sum_{v \in N} s_{uv}}$

# Item-Item Collaborative Filtering

- So far: User-User Collaborative Filtering
- Another view: Item-Item CF
  - For item  $i$ , find other similar items
  - Estimate rating for item  $i$  based on ratings for similar items
  - Can use same similarity metrics and prediction functions as in user-user CF model

$$r_{ui} = \frac{\sum_{j \in N(i;u)} s_{ij} \cdot r_{uj}}{\sum_{j \in N(i;u)} s_{ij}}$$

$s_{ij}$ : similarity of items  $i$  and  $j$

$r_{uj}$ : rating of user  $u$  on item  $j$

$N(i;u)$ : set of items rated by  $u$  similar to  $i$

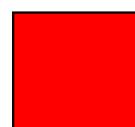
# Item-Item CF ( $|N| = 2$ )

	users												
	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 CF ( $|N| = 2$ )

	users												
	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		



- estimate rating of movie 1 by user 5

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

Compute similarity weights:

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

	users												
	1	2	3	4	5	6	7	8	9	10	11	12	sim(1,m)
movies	1	1		3		?	5			5		4	1.00
2				5	4			4			2	1	-0.18
3	2	4		1	2			3		4	3	5	0.41
4		2	4		5			4			2		-0.10
5			4	3	4	2					2	5	-0.31
6	1		3		3			2			4		0.59

## Neighbor selection:

Identify movies similar to  
movie 1, rated by user 5

Here we use Pearson correlation as similarity:

- Subtract mean rating  $m_i$  from each movie  $i$

$$m_1 = (1+3+5+5+4)/5 = 3.6$$

$$\text{row 1: } [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]$$

- Compute cosine similarities between rows

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

Compute similarity weights:

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

	users												
	1	2	3	4	5	6	7	8	9	10	11	12	sim(1,m)
movies	1	1		3		?	5			5		4	1.00
	2			5	4			4			2	1	3
3	2	4		1	2		3		4	3	5		0.41
4		2	4		5			4			2		-0.10
5			4	3	4	2					2	5	-0.31
6	1		3		3			2			4		0.59

Predict by taking weighted average:

$$r_{1,5} = (0.41*2 + 0.59*3) / (0.41+0.59) = 2.6$$

$$r_{ui} = \frac{\sum_{j \in N(i; x)} s_{ij} \cdot r_{uj}}{\sum_{j \in N(i; x)} s_{ij}}$$

# CF Common Practice

- Define **similarity**  $s_{ij}$  of items  $i$  and  $j$
- Select  $k$  nearest neighbors  $N(i; u)$ 
  - Items most similar to  $i$  that were rated by  $u$
- Estimate rating  $r_{ui}$  as the weighted average:

$$r_{ui} = b_{ui} + \frac{\sum_{j \in N(i; u)} s_{ij} \cdot (r_{uj} - b_{uj})}{\sum_{j \in N(i; u)} s_{ij}}$$

baseline estimate for  $r_{ui}$

$$b_{ui} = \mu + b_u + b_i$$

$$r_{ui} = \frac{\sum_{j \in N(i; u)} s_{ij} \cdot r_{uj}}{\sum_{j \in N(i; u)} s_{ij}}$$

- $\mu$  = overall mean item rating
- $b_x$  = rating deviation of user  $x$   
= (avg. rating of user  $x$ ) –  $\mu$
- $b_i$  = rating deviation of item  $i$   
= (avg. rating of item  $i$ ) –  $\mu$

# Item-based CF vs. User-based CF

CF

	Avatar	LOTR	Matrix	Pirates
Alice	1		0.8	
Bob		0.5		0.3
Carol	0.9		1	0.8
David			1	0.4

- In practice, it has been observed that item-based CF often works better than user-based CF
- Why? Items are simpler, users have multiple tastes

Two similar items (e.g., “Action” movies) tend to have the same group of users

Two similar users could have different tastes (e.g., {Action, Drama, Dramatic}, {Action, Historical, Comedy})

# Pros/Cons of CF

- + Work 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 (極少數的) rated items
- - Popularity bias:
  - Tend to recommend popular items
  - Cannot recommend items to someone with unique taste

# Hybrid Methods

Hybrid Method	Description
Weighted	The <b>scores</b> (or <b>votes</b> ) of several recommendation techniques are combined (e.g. <b>average</b> ) together to produce a single recommendation. Another approach is to use <b>Linear model</b> .
Switching	The system switches between recommendation techniques depending on the current situation (e.g. cold-start)
Mixed	Recommendations from several different recommenders are presented at the same time
Feature Combining	Features from different recommendation data sources are thrown together into a single recommendation algorithm
Cascade	One recommender refines the recommendations given by another
Feature Argumentation	Output from one technique is used as an input feature to another

## USER-BASED COLLABORATIVE FILTERING

	User 1	User 2	User 3	User 4	User 5	User 6
User 1	1.00	0.75	0.63	0.22	0.30	0.00
User 2	0.75	1.00	0.91	0.00	0.00	0.16
User 3	0.63	0.91	1.00	0.00	0.00	0.40
User 4	0.22	0.00	0.00	1.00	0.97	0.64
User 5	0.30	0.00	0.00	0.97	1.00	0.53
User 6	0.00	0.16	0.40	0.64	0.53	1.00

(0.7 x  + 0.6 x ) =  already rated by user  

$$(0.7 \times 4 + 0.6 \times 5) / (0.7 + 0.6) = 4.5$$

 already rated by user  

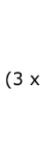
$$(0.6 \times 3) / 0.6 = 3.0$$

## INPUT

	4	3			5	
	5		4		4	
	4		5	3	4	
		3				5
		4				4
			2	4		5

## ITEM-BASED COLLABORATIVE FILTERING

	Item 1	Item 2	Item 3	Item 4	Item 5	Item 6
Item 1	1.00	0.27	0.79	0.32	0.98	0.00
Item 2	0.27	1.00	0.00	0.00	0.34	0.65
Item 3	0.79	0.00	1.00	0.69	0.71	0.18
Item 4	0.32	0.00	0.69	1.00	0.32	0.49
Item 5	0.98	0.34	0.71	0.32	1.00	0.00
Item 6	0.00	0.65	0.18	0.49	0.00	1.00

(4 x  + 3 x  + 5 x ) =  already rated by user  

$$(0.8 \times 4 + 0.7 \times 5) / (0.8 + 0.7) = 4.5$$

 already rated by user  

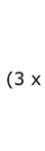
$$(0.7 \times 3) / 0.7 = 3.0$$

## OUTPUTS



## CONTENT-BASED FILTERING

	Item 1	Item 2	Item 3	Item 4	Item 5	Item 6
Item 1	1.00	0.00	0.58	0.00	0.67	0.58
Item 2	0.00	1.00	0.00	0.41	0.00	0.00
Item 3	0.58	0.00	1.00	0.00	0.58	0.75
Item 4	0.00	0.41	0.00	1.00	0.00	0.00
Item 5	0.67	0.00	0.58	0.00	1.00	0.58
Item 6	0.58	0.00	0.75	0.00	0.58	1.00

(4 x  + 3 x  + 5 x ) =  already rated by user  

$$(0.4 \times 3) / 0.4 = 3.0$$

 already rated by user  

$$(0.6 \times 4 + 0.6 \times 5) / (0.6 + 0.6) = 4.5$$

## HYBRID

0.4 x  UB CF + (0.3 x  IB CF) + (0.3 x  CB) =  (0.4 x 4.5 + 0.3 x 4.5) / (0.4 + 0.3) = 4.5  

$$(0.3 \times 3.0 + 0.3 \times 4.5) / (0.3 + 0.3) = 3.8$$

 (0.4 x 3.0 + 0.3 x 3.0) / (0.4 + 0.3) = 3.0