

IEMS 5780 / IERG 4080
Building and Deploying Scalable
Machine Learning Services

Lecture 5 - Recommender Systems

Albert Au Yeung
4th October, 2018

Recommender Systems

Agender

Recommender Systems

- Introduction
- Content-based Recommendation Systems
- Collaborative Filtering
 - User-based Neighbourhood Model
 - Item-based Neighbourhood Model
 - Matrix Factorization
- Recommendation as Classification

Recommender Systems

- We make a lot of **decisions** every day
 - Transportation
 - Restaurant reservation
 - Hotel reservation
 - Music
 - Movies
- We usually rely on some **suggestions** or **recommendations**
 - Family and friends
 - Supervisors, teachers, seniors
 - Experts
 - The general public (word-of-mouth)

Recommender Systems

What are the factors that affect our decisions?

- Our own **preferences** (of the content, the characteristics, the price, etc.)
- What do **most people like**? (which is the blockbuster movie recently?)
- What do **people around us like**? (friends and family)
- What do people **similar to us** like?

Example: Mobile App User Demographics

Gender (Male)				Age (33–100)			
Coef	Share	<i>n</i>	App name	Coef	Share	<i>n</i>	App name
0.81	85 %	150	ESPN	0.53	80 %	42	Great Clips Online Check-in
0.73	80 %	142	Geek - Smarter Shopping	0.48	53 %	1687	Email
0.63	78 %	277	Tinder	0.46	58 %	318	New Words With Friends
0.59	80 %	172	Fallout Shelter	0.44	80 %	65	BINGO Blitz
0.56	86 %	106	WatchESPN	0.43	60 %	380	iHeartRadio - Music & Radio
0.52	72 %	190	Clash of Clans	0.41	54 %	197	Field Agent
0.52	97 %	41	Grindr - Gay chat, meet & date	0.40	55 %	690	Lookout Security & Antivirus
0.49	84 %	96	Yahoo Fantasy Football & More	0.40	92 %	41	DoubleUCasino
Gender (Female)				Age (18–32)			
-1.03	76 %	736	Pinterest	-1.17	78 %	1066	Snapchat
-0.73	84 %	182	Etsy	-0.52	59 %	113	Perk Word Search
-0.61	97 %	79	Period Tracker	-0.49	64 %	88	Summoners War
-0.54	96 %	58	Period Calendar / Tracker	-0.46	59 %	98	Clash of Kings
-0.50	76 %	346	Cartwheel by Target	-0.45	86 %	90	iFunny :)
-0.49	66 %	258	Wish - Shopping Made Fun	-0.45	81 %	158	GroupMe
-0.49	74 %	325	Checkout 51 - Grocery Coupons	-0.42	80 %	68	GIPHY for Messenger
-0.45	74 %	178	Photo Grid - Collage Maker	-0.42	80 %	183	Vine

Ref: [You Are What Apps You Use: Demographic Prediction Based on User's Apps](#)

Example: Mobile App User Demographics

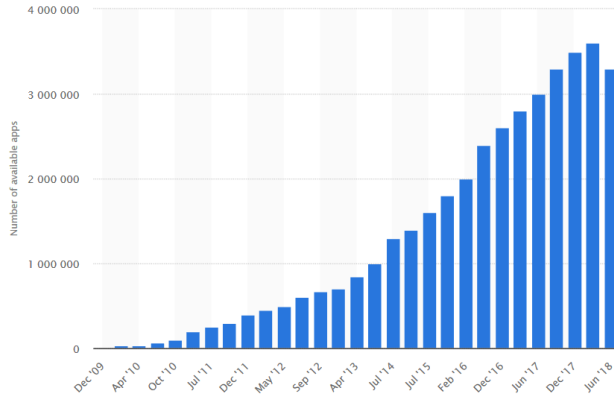
Married (Married)				Income (\geq \$50K)			
Coef	Share	n	App name	Coef	Share	n	App name
0.55	67 %	200	Zillow Real Estate & Rentals	0.58	75 %	141	Fitbit
0.44	67 %	622	Walmart	0.45	66 %	205	LinkedIn
0.44	60 %	823	Pinterest	0.41	65 %	41	com.ws.dm
0.44	74 %	39	Gospel Library	0.37	52 %	141	LG Android QuickMemo+
0.40	59 %	91	USAA Mobile	0.37	58 %	191	Redbox
0.40	80 %	63	ClassDojo	0.36	72 %	22	Like Parent
0.38	60 %	123	ESPN	0.34	66 %	63	Peel Smart Remote
0.37	82 %	28	Deer Hunter 2014	0.34	61 %	220	Yelp
Married (Single)				Income (\leq \$40K)			
-0.89	70 %	810	Snapchat	-0.43	66 %	136	Job Search
-0.78	89 %	114	POF Free Dating App	-0.43	63 %	97	Security policy updates
-0.73	85 %	219	Tinder	-0.37	78 %	23	Solitaire
-0.66	98 %	69	OkCupid Dating	-0.35	67 %	79	Prize Claw 2
-0.48	72 %	269	Tumblr	-0.34	72 %	51	ScreenPay- Get Paid to Unlock
-0.42	72 %	205	SoundCloud - Music & Audio	-0.33	78 %	56	MeetMe
-0.41	65 %	331	Uber	-0.33	62 %	77	Foursquare
-0.41	89 %	69	MeetMe	-0.32	56 %	73	Microsoft Word

Ref: [You Are What Apps You Use: Demographic Prediction Based on User's Apps](#)

Recommender Systems

- There are **many** items out there for us to choose
 - Tens of thousands of movies and songs - More than 1 million apps in the Android and iPhone app stores
 - Millions of books published every year
- We need more efficient way to **filter information**, and identify items **most relevant** to us
- On the other hand, producers also want to **provide consumers things that they really want** (targetted marketing)

Mobile Apps

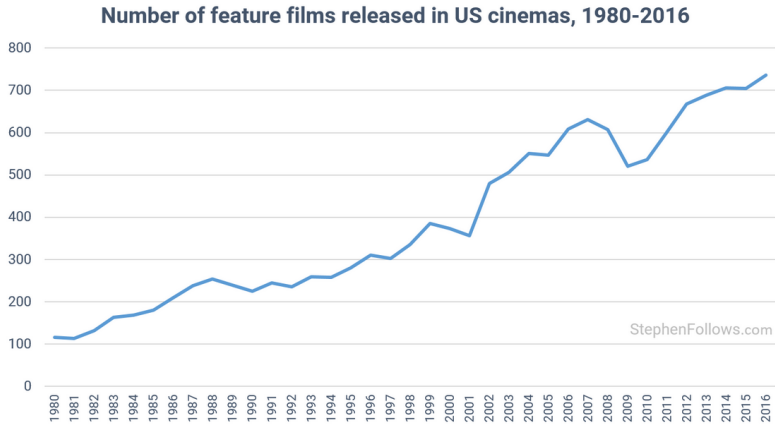


Data visualized by  + a b l e a u

© Statista 2018 

(<https://www.statista.com/statistics/266210/number-of-available-applications-in-the-google-play-store/>)

Movies

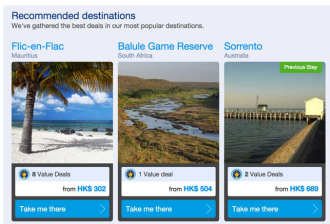
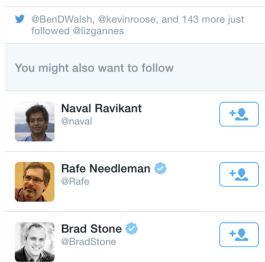
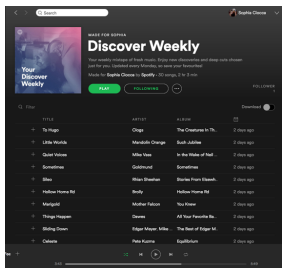


(<https://stephenfollows.com/how-many-films-are-released-each-year/>)

Solution: Recommender Systems

- Use **computers** and **algorithms** to process the huge amount of information and do the **filtering** for us
- Analyse the **tastes** and **preferences** of different people
- Analyse the **characteristics** of different items/products
- Generate **personalized** recommendation based on users' **past activities** and **feedback**
- Systems performing the above tasks are referred to as **recommender systems** / **recommendation systems**

Examples of Recommender Systems

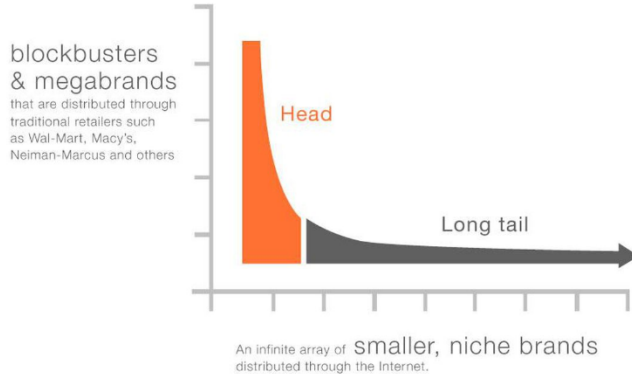


- Music recommendation
- Suggested connections in social networks
- Recommended items in e-commerce Websites
- Recommended travel destinations and accommodations
- ...

Common Strategies

- **Popularity**
 - Recommend items most people like
- **Item Similarity**
 - Recommend **items that are similar** to what the user has already shown interest
- **User Similarity**
 - Recommend items that are preferred by **users who are similar the the target user** in some ways
- **Diversity**
 - Recommend items that are **least known** to the uuser

The Long Tail



<https://www.forbes.com/sites/robinlewis/2016/05/31/the-long-tail-theory-can-be-reality-for-traditional-megabrands/#2b77afbc6372>

Content-based Recommendation

Content-based Recommendation

- **Assumptions**

- Every user has his/her own **interests / tastes / preferences**
- Each user's preferences can be represented as **a summary of what he/she has seen/read/watched/liked in the past**
- A user will prefer something he or she is interested in
- We can compare the **content** or **characteristics** of the items
- Recommend items that are **similar** to what the user has consumed before

Content-based Recommendation

Two Steps

1. Learn user **preferences** (what does the user like?)
2. Find items that **match** these preferences

Problems

1. How do we **learn** user preferences?
2. How do we **represent** user interests?
3. How do we **represent** items?
4. How to measure **similarity**?

Content-based Recommendation

Similarity-based

- **Steps**

- Define features to represent the items (e.g. bag-of-words, author, publish date/time, etc.)
- Construct **user profile** by the items liked by the user (e.g. average of feature vectors)
- Calculate **similarity** between user profile and the new items
- Return a **ranked list** of items

- What is important here is the **user profile**

- How can we **represent** a user?
- A user may have **multiple interests**, or his/her interests may **change over time**

Content-based Recommendation

Limitations of Content-based Methods

- Content (including meta-data) might not be available or enough in some domains
- It is difficult to represent some items by their 'content' (e.g. movies, books, music)
- Content-based methods tend to return very similar items

Collaborative Filtering

Collaborative Filtering

What is Collaborative Filtering (CF)?

- From [Wikipedia](#):

In the newer, narrower sense, collaborative filtering is a method of making automatic predictions (filtering) about the interests of a user by collecting preferences or taste information from many users (collaborating). The underlying assumption of the collaborative filtering approach is that if a person A has the same opinion as a person B on an issue, A is more likely to have B's opinion on a different issue than that of a randomly chosen person.

Collaborative Filtering

Basic Idea

- Instead of relying on the content of the items, we can analyse people's **tastes and preferences**
- Each person is NOT totally different from another
- There are different **kinds** of people, for example:
 - People who like action movies
 - People who read literature
 - People who like spicy food
- By grouping **similar users**, we can recommend similar items to similar users
- This does not require a lot of information about the users and items themselves

Collaborative Filtering

Two Types of Collaborative Filtering

1. **Memory-based**

- Directly use ratings from similar users or items
- Also called **neighbourhood-based** methods

2. **Model-based**

- Mathematical models are used to represent users, items and their relations
- E.g. **matrix factorization**, Bayesian networks, probabilistic models

Collaborative Filtering

User-item Interaction

- In the following discussion, we assume that a user may **rate** an item on a **1 to 5** scale

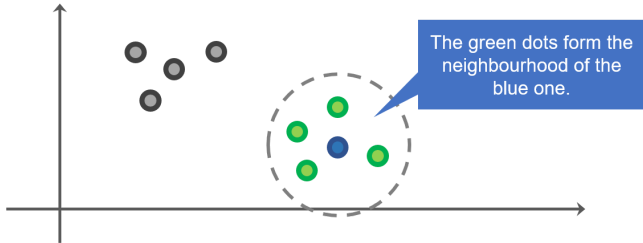
		Items					
		i_0	i_1	i_2	i_3	i_4	i_5
Users	u_0	-	5	-	-	-	2
	u_1	2	-	1	3	-	5
	u_2	-	4	-	4	-	5
	u_3	4	5	-	-	-	-
	u_4	-	-	5	-	1	-
	u_5	2	-	-	4	1	-

Memory-based Collaborative Filtering

Neighbourhood

What is neighbourhood?

- Consider again the **vector space model**
- Users and items can be represented as **vectors** in a high dimensional space
- More similar items/users will appear **closer** to each other



Neighbourhood Models

Two Types of Neighbourhood Models

1. User-based

- We consider **similar users**
- The neighbourhood of a user consists of users who like similar items

2. Item-based

- We consider **similar items**
- The neighbourhood of an item consists of items that are preferred by similar users

User-based Collaborative Filtering

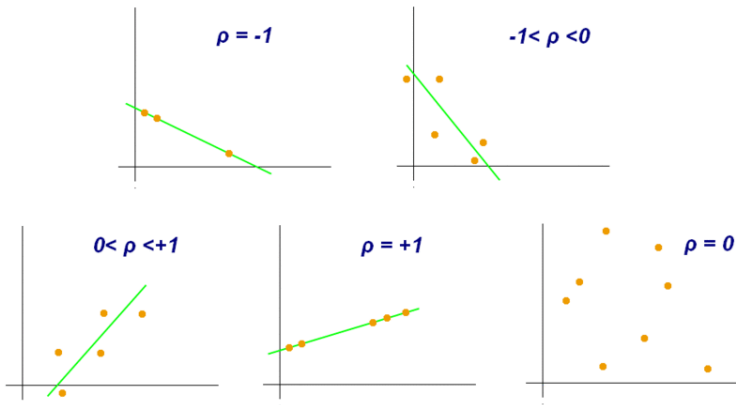
- **Task:** Predict the **extent** to which user u likes item i
- **Procedures:**
 - We have user u , and item i
 - Find a set of users who are similar to u and have rated i
 - Get the **average rating** on i given by this set of users
 - Do this for all items, come up with a **ranking** based on the predict rating

User-based Collaborative Filtering

- In doing **user-based CF**, we made the following **assumptions**:
 - If users had similar tastes in the past, they should have similar tastes in the future
 - User preferences remain **stable and consistent** over time
- Furthermore, we need to define **similarity**
- **Similarity**
 - Should reflect how **close** the tastes and preferences of two users are
 - Similar users should assign **similar ratings** to the same set of items

Similarity - Pearson Correlation

- One way of computing similarity between two users is to use the **pearson correlation**



Similarity - Pearson Correlation

- x, y : users
- $r_{x,i}$: rating assigned to i by x
- I : the set of items rated by both x and y
- $\text{sim}(x, y)$ has a value between -1 and 1

$$\text{sim}(x, y) = \frac{\sum_{i \in I} (r_{x,i} - \bar{r}_x)(r_{y,i} - \bar{r}_y)}{\sqrt{\sum_{i \in I} (r_{x,i} - \bar{r}_x)^2} \sqrt{\sum_{i \in I} (r_{y,i} - \bar{r}_y)^2}}$$

Similarity - Pearson Correlation

- In Python, you can easily compute the correlation using the **scipy** module's [pearsonr function](#)

```
>>> from scipy.stats import pearsonr
>>> a = [5,2,1,4,3,1]
>>> b = [1,5,4,2,1,2]
>>> pearsonr(a,b)
(-0.59628479399994383, 0.21157899460007421)
>>>
>>> a = [5,2,1,4,3,1]
>>> b = [4,3,1,5,4,2]
>>> pearsonr(a,b)
(0.85978530414910515, 0.028111919656069226)
>>>
```


Predicting Ratings

- How to **predict** a user's rating for an item?
- Use the equation below:

$$\text{pred}(x, i) = \overline{r_x} + \frac{\sum_{u \in U} \text{sim}(x, u) \times (r_{u,i} - \overline{r_u})}{\sum_{u \in U} \text{sim}(x, u)}$$

Similarity between x and u

Deviation from mean of ratings of user u

Normalization

- Example: see [l5-user-based-cf.ipynb](#)

Limitations

- **Data sparsity**
 - When there are a lot of users and items, very few overlap between users
- **Does not scale**
 - When there are 10 million users, how can you generate all their neighbourhoods?
- Two users may have similar taste in one domain but very different taste in another

Item-based Collaborative Filtering

Procedures

- We have user u and item i
- Find a set of **items**, which are
 - rated by u
 - given **similar ratings** as i by other users
- Get the **average rating** of this set of items given by u
- Use that as the predicted rating of i by u

Item-based Collaborative Filtering

- In doing **item-based CF**, we made the following **assumptions**:
 - If two items are given similar ratings by the users, they should have **similar characteristics**
- Consider **movies**:
 - Movies of the same genre, by the same director, or feature the same actor/actress will be assigned similar ratings
- Ref: [Amazon.com Recommendations: Item-to-Item Collaborative Filtering](#)

Predicting Ratings

- One way to implement item-based CF is to compute a **weighted combination of ratings** given by u to other items

$$\text{pred}(u, j) = \frac{\sum_{i \in I} \text{sim}(j, i) \times r_{u,i}}{\sum_{i \in I} \text{sim}(j, i)}$$

Similarity between two items

Rating given by the user to other items

Normalization

Neighbourhood-based Collaborative Filtering

- Neighbourhood and recommended items are usually calculated **offline**
- In order to reduce the amount of computation needed, the size of the neighbourhood is usually **limited**
- Previous research shows that a neighbourhood size of **20-50** is quite enough
- Ref: [Empirical Analysis of Design Choices in Neighborhood-based Collaborative Filtering Algorithms](#), Herlocker et al., 2002.

Food for Thought

TED Talk

How Algorithms Shape Our World

TEDGlobal 2011, By Kevin Slavin

http://www.ted.com/talks/kevin_slavin_how_algorithms_shape_our_world.html

Model-based Collaborative Filtering

Model-based Collaborative Filtering

- What we have discussed so far are called **memory-based** collaborative filtering
- We have not **trained** any model that can be used to describe the relationship between inputs and outputs
- Memory-based models are **easy to implement**, but usually are NOT very accurate and NOT **scalable**
- Instead, we can consider a **machine learning approach** to the task of recommendation:
 - Create a model that explains how **ratings are generated**, or how **items are ranked**
 - Use past data to **train/optimize** the parameters of the model

Model-based Collaborative Filtering

We can roughly categorize model-based CF into two types:

- **Matrix Factorization Approach**

- Assume that there are **latent factors** that determine how users rate items
- Decompose the user-item matrix using matrix factorization techniques

- **Classification Approach**

- Consider the task of recommendation as a **classification** problem
- For each given pair of user and item, we determine whether the user will be interested in the item (binary classification)
- Can take into account **contextual information** and **implicit feedbacks**

To be continued

End of Lecture 5