

Recommender Systems

Learning Objectives

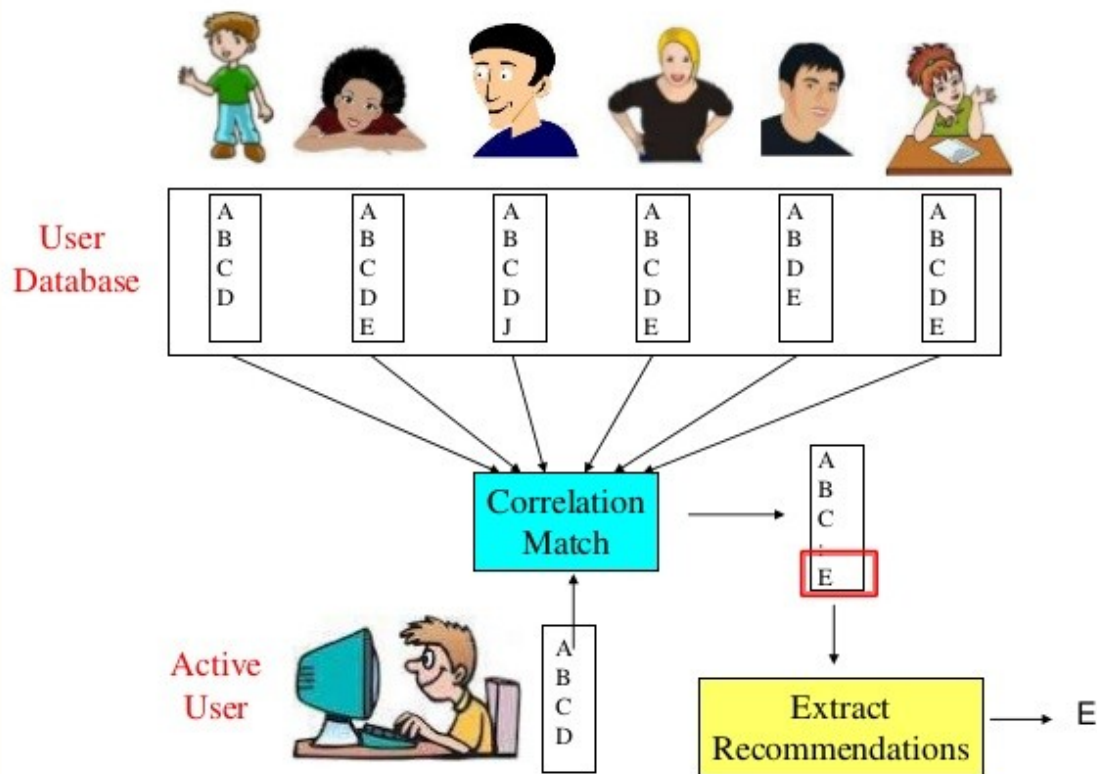
In this lesson, we will cover the following concepts:

- Recommendation system
- The long tail
- A simple popularity-based recommender system
- A collaborative filtering model
- Evaluating a recommendation system

What Is a Recommender System?

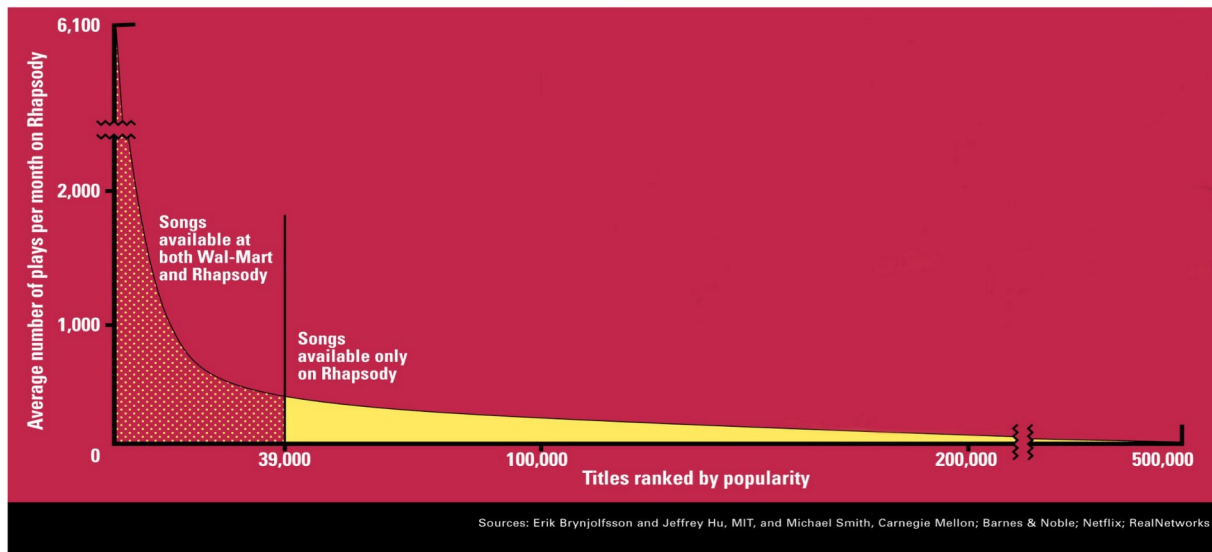
A recommender or recommendation system is a subclass of an information filtering system that seeks to predict the rating or preference that a user would give to an item.

Let's consider the example shown in the figure below. Here, we have a user database, that is, data consisting of items rated by the user. Now, let's suppose that a new user visits and likes five out of ten items on the website. A recommender system recommends the items that the new user might like, based on similarity with other items. We will dive deeper into this concept in the coming sections.



The Theory of Long Tail

- It shows how products in low demand or with low sales volume can collectively make up a market share that exceeds the relatively few current bestsellers and blockbusters but only if the store or distribution channel is large enough.
- The long tail concept looks at less popular goods in lower demand. The use of these goods could increase profitability as consumers navigate away from mainstream markets.
- This can be easily understood by looking at the figure [below](#).



The figure above clearly shows the use of long tail by [Rhapsody] ([https://en.wikipedia.org/wiki/Rhapsody_\(music\)](https://en.wikipedia.org/wiki/Rhapsody_(music))) where they sell music albums both online and off-line. We can clearly observe the following:

- Both Rhapsody and Walmart sell the most popular music albums online, but the former offers 19 times more songs than Walmart. Even though there is a demand for popular music albums, there is also a demand for the less popular online. Recommender systems leverage these less popular items online.

Recommend the Most Popular Items

- Let's consider the movie dataset. We will look carefully at the user ratings and think about what can be done.
- The answer that strikes first is the **most popular item**. This is exactly what we will be doing.
- Technically, this is the fastest method, but it does come with a major drawback, which is a lack of personalization. The dataset has many files; we will be looking at a few of them, mainly the ones that relate to movie ratings.

Popularity-Based Recommender System

- There is a division by section, so the user can look at the section of his or her interest.
- At a time, there are only a few hot topics; there is a high chance that a user wants to read the news which is being read by most others.

Import Libraries

In python, Pandas is used for data manipulation and analysis. NumPy is a package that includes a multidimensional array object and multiple derived objects. Matplotlib is an amazing

visualization library in Python for 2D plots of arrays. Seaborn is an open-source Python library built on top of matplotlib. Mean_squared_error is a library that measures the average of the squares of the errors, which is the average squared difference between the estimated values and the actual value.

These libraries are written with the import keyword.

```
import pandas as pd
import os, io
import numpy as np
from pandas import Series, DataFrame, read_table
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
from sklearn.metrics import mean_squared_error
%matplotlib inline
```

Exporting Dataset from Zip File

Before reading data you need to download "ml-100k.zip" dataset from the resource section and upload it into the Lab. We will use Up arrow icon which is shown in the left side under View icon. Click on the Up arrow icon and upload the file wherever it is downloaded into your system.

After this you will see the downloaded file will be visible on the left side of your lab with all the .ipynb files.

Then, the below snippet will extract the zip dataset to the corresponding folder.

```
import zipfile
with zipfile.ZipFile('ml-100k.zip', 'r') as zip_ref:
    zip_ref.extractall(".")
```

We start to explore the data set of movie ratings and our interest lies particularly in ratings. Let's see how we recommend the most popular (that is, highly rated) movies.

```
#Load the Ratings data
r_cols = ['user_id', 'movie_id', 'rating', 'unix_timestamp']
ratings = read_table('ml-100k//u.data', header=None, sep='\t')
ratings.columns = r_cols

i_cols = ['movie_id', 'movie title', 'release date', 'video release date', 'IMDb URL', 'unknown', 'Action', 'Adventure', 'Animation', 'Children\'s', 'Comedy', 'Crime', 'Documentary', 'Drama', 'Fantasy', 'Film-Noir', 'Horror', 'Musical', 'Mystery', 'Romance', 'Sci-Fi', 'Thriller', 'War', 'Western']
items = read_table('./ml-100k//u.item', sep='|', names=i_cols, encoding='latin-1')
```

```
ratings.head()
```

Ratings is a variable that stores all the columns from the ml-100k dataset in the u.data file. The head() function displays the first five rows from ratings.

Let's Build a Popularity-Based Recommender System

With our initial exploration, we decided that ideal data would be the one where we could also have the movie ratings with us. Let's see how we are able to do this.

We will use the pd.merge function that is used to combine data on common columns or indices.

```
new_data = pd.merge(items, ratings, on='movie_id')
new_data = new_data[['movie_id', 'movie title', 'user_id', 'rating']]
new_data.head()
```

New_data is a variable that stores data read by the pd.merge function. It consists of items and ratings. The head() function displays the first five rows from new_data.

Before proceeding to build the recommender system, we will observe the following steps to recommend movies:

- Find unique users
- Count the number of times the movie has been seen
- [Rank](#) the scores (counts)

```
def popularity(train, title, ids):
    train_data_grouped = train.groupby([title])
    [ids].count().reset_index() #user_id #movie title
    train_data_grouped.rename(columns = {ids: 'score'}, inplace=True)

    train_data_sort = train_data_grouped.sort_values(['score', title],
    ascending = [0, 1])
    train_data_sort['Rank'] =
train_data_sort['score'].rank(ascending=0, method='first')
    popularity_recommendations = train_data_sort.head(10)
    return popularity_recommendations

popularity(new_data, 'movie title', 'user_id')
```

Drawback

Having recommended the movies, we can immediately conclude that the major drawback of such a system would be the **lack of personalization**.

Collaborative Filtering

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 person A has the same opinion as 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.

Types of Collaborative Filtering

User-Based Collaborative Filtering

In this type, we find look-alike customers (based on similarity) and offer products that the first customer's look-alike chose in the past. This algorithm is very effective but takes a lot of time and resources. It computes every customer pair information, which takes time. Therefore, for big base platforms, this algorithm is hard to implement without a very strong parallelizing system.

1. Build a matrix of things each user bought or viewed or rated
2. Compute similarity scores between users
3. Find users similar to you
4. Recommend stuff they bought or viewed or rated that you haven't yet

Problems

1. People are fickle, so their tastes tend to change
2. There are usually more people than things

Item-Based Collaborative Filtering

It is quite similar to the previous algorithm, but instead of finding customer look-alikes, it tries to find items that look alike. Once we have an item look-alike matrix, we can easily recommend similar items to customers who have purchased an item from the store. This algorithm is far less resource-consuming than user-based collaborative filtering.

1. Find every pair of movies that were watched by the same person
2. Measure the similarity of rating across all the users who watched both
3. Sort movies by the similarity strength

Interesting fact

Item-based collaboration is extensively used in Amazon, and they came out with it in great detail. You can read more at [Amazon](#)

Let's get started with building our item-based collaborative recommender system. For convenience, let's split this into two parts.

- To find similarities between items
- To recommend them to users

Item-based collaborative filtering would be the most feasible solution, as the number of items is always lesser than the number of users and it improves the computational speed.

Leverage the Pandas

- To begin with, we will use the pandas pivot table to look at relationships between movies and we will use the pivot table in pandas. Pivot table in pandas is an excellent tool to summarize one or more numeric variable based on two other categorical variables.
- We start building a utility matrix (matrix consisting of movies and ratings)

```
movie_ratings = new_data.pivot_table(index=['user_id'],columns=['movie
title'],values='rating')
movie_ratings.head()
```

The above table gives information about the rating given by each user against the movie title. There are many NaN as it is not necessary for each user to review each movie. Let's start by looking at the geeks' most favorite, Star Wars, and see how it correlates pairwise with other movies in the table.

Similarity Function

To decide the similarity between two items in the dataset, let's briefly look at the popular similarity functions.

Terminology

- Let r_x denote the rating of the item x given by the user and r_y be the rating of item y. To find the similarity pairwise between two items the following metrics can be used:

cosine Index

$$\text{sim}(r_x, r_y) = \cos(r_x, r_y) = \frac{r_x r_y}{(|r_x|)(|r_y|)}$$

The major problem is that it treats missing values as negative.

Pearson Index

S_{xy} = Items x and y both have ratings

$$\text{sim}(r_x, r_y) = \frac{\sum_{x \in S} (r_{xs} - r_{xm})(r_{ys} - r_{ym})}{\sqrt{\sum_{x \in S} (r_{xs} - r_{xm})^2} \sqrt{\sum_{x \in S} (r_{ys} - r_{ym})^2}}$$

Jaccard Index

$$Jaccard\ Index = \frac{Number\ in\ both\ sets}{Number\ in\ either\ set}$$

Let's start with the Pearson Index in this case. Now that we have understood how similar products can be found, let's start with the movie, Star Wars.

```
StarWarsRatings = movie_ratings['Star Wars (1977)']
StarWarsRatings.head()
```

Now, let's use the **corrwith()** function to check the pairwise correlation of Star Wars's user rating with other films in the column.

```
similarmovies = movie_ratings.corrwith(StarWarsRatings)
similarmovies = similarmovies.dropna()
df = pd.DataFrame(similarmovies)
df.head()
```

If we look at the data closely, we will find something incorrect.

The potential reason here is that a handful of people who have seen obscure films are messing up our movies. We want to get rid of the movies that only a few people have watched that show incorrect results.

We have used groupby function that involves some combination of splitting the object, applying a function, and combining the results and sort_values function that sorts by the values along either axis.

```
movie_stats = new_data.groupby('movie title').agg({'rating':
[ np.size, np.mean ]})

check = movie_stats.sort_values([ ('rating', 'mean') ], ascending=False)
check.head()
```

Now, we can clearly observe that there are movies that have very few rating counts (size). Therefore, we set a threshold of the movie count to have at least 100 ratings.

```
popularmovies = movie_stats['rating']['size']>=100

movie_stats[popularmovies].sort_values([ ('rating', 'mean') ], ascending=False)[:10]

df =
movie_stats[popularmovies].join(DataFrame(similarmovies, columns=['similarity']))
df.sort_values('similarity', ascending=False)[:20]
```


Building an End-to-End Recommender System

We will list points that need to be followed to recommend a movie based on what we did till now :

- Compute the correlation score for every pair in the matrix
- Choose a user and find his or her movies of interest
- Recommend movies to him or her
- Improve on the recommendation

The pandas method **corr()** will compute the correlation score for every pair in the matrix. This gives a correlation score between every pair of movies in turn creating a sparse matrix. Let's see how this looks.

```
corrMatrix = movie_ratings.corr(method='pearson',min_periods=100)
corrMatrix.head()
```

Now, we want to recommend movies to a friend, so let's have a look at the movies our friend has rated.

```
friend_ratings = movie_ratings.loc[1].dropna()[1:4]
friend_ratings

simcandidates= pd.Series()
for i in range(0,len(friend_ratings.index)):
    print('Adding similars to ', friend_ratings.index[i])

    print('-----')
    sims = corrMatrix[friend_ratings.index[i]].dropna()
    sims = sims.map(lambda x: x*friend_ratings[i]) # Assigning lower
weights to movies with lower ratings.
    simcandidates = simcandidates.append(sims)

    print('sorting')

    simcandidates.sort_values(inplace=True,ascending=False)

    print(simcandidates.head(10))
```

Some movies come up more than once, because they are very similar to the ones that the user has rated. Let's eliminate them.

```
simcandidates = simcandidates.groupby(simcandidates.index).sum()
simcandidates.sort_values(inplace=True,ascending=False)
simcandidates.head(10)
```

Having done all the computations using pandas, we can see that it is computationally intensive. We have a Python module that does that for us.

Using the Surprise Module

[Python Surprise](#) is an easy-to-use Python scikit for recommender systems. Let's see how to build a recommender system using the surprise module and focus on the model inspired by K-Nearest Neighbors (KNN).

Common Practice

1. Define Similarity S_{ij} in terms of i and j
2. Select K nearest neighbors $N(i;X)$
 - Items most similar to i that were rated by X
3. Estimate rating r_{xi} as the weighted average

$$r_{x_i} = b_{x_i} + \frac{\sum_{j \in N(i;X)} S_{ij} (r_{x_j} - b_{x_j})}{\sum_{j \in N(i;X)} S_{ij}}$$

Here, the term b_{x_i} is the baseline estimator for the rating comprising three terms: the overall mean movie rating, rating deviation of user x , and rating deviation of the movie i .

Evaluation Metrics

Comparing Predictions with Known Ratings

RMSE

- Root Mean Square Error (RMSE)
 - $\sqrt{\frac{1}{N} \sum_{x_i} (\text{predicted } r_{x_i} - \text{actual } r_{x_i}^i)^2}$ here r_{x_i} is the predicted rating and $r_{x_i}^i$ is the actual rating
- Precision at top 10
 - % of those in top 10

Note: In this lesson, we saw the use of the recommender systems.