

## Task 1: Introduction



For this project we are going to create a recommendation engine for movies for users based on there past behaviour.

---

We will focus on the **collaborative filtering** approach, that is:

The user is recommended items that people with similar tastes and preferences liked in the past. In another word, this method predicts unknown ratings by using the similarities between users.

---

Note: This notebook uses `python 3` and these packages: `pandas` , `numpy` , `matplotlib` and `scikit-surprise`

We can install them using:

```
pip3 install pandas matplotlib numpy scikit-surprise
```

### 1.1: Installing Libraries

```
In [1]: print('>> Installing Libraries')
!pip3 install pandas matplotlib numpy scikit-surprise

print('>> Libraries Installed')

>> Installing Libraries
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: pandas in /home/rhyme/.local/lib/python3.6/site-packages (1.1.0)
Requirement already satisfied: matplotlib in /home/rhyme/.local/lib/python3.6/site-packages (3.3.0)
Requirement already satisfied: numpy in /home/rhyme/.local/lib/python3.6/site-packages (1.19.1)
Requirement already satisfied: scikit-surprise in /home/rhyme/.local/lib/python3.6/site-packages (1.1.1)
Requirement already satisfied: python-dateutil>=2.7.3 in /home/rhyme/.local/lib/python3.6/site-packages (from pandas) (2.8.1)
Requirement already satisfied: pytz>=2017.2 in /home/rhyme/.local/lib/python3.6/site-packages (from pandas) (2020.1)
Requirement already satisfied: pyparsing!=2.0.4,!=2.1.2,!=2.1.6,>=2.0.3 in /home/rhyme/.local/lib/python3.6/site-packages (from matplotlib) (2.4.7)
Requirement already satisfied: cycler>=0.10 in /home/rhyme/.local/lib/python3.6/site-packages (from matplotlib) (0.10.0)
Requirement already satisfied: kiwisolver>=1.0.1 in /home/rhyme/.local/lib/python3.6/site-packages (from matplotlib) (1.2.0)
Requirement already satisfied: pillow>=6.2.0 in /home/rhyme/.local/lib/python3.6/site-packages (from matplotlib) (7.2.0)
Requirement already satisfied: six>=1.10.0 in /usr/lib/python3/dist-packages (from scikit-surprise) (1.11.0)
Requirement already satisfied: joblib>=0.11 in /home/rhyme/.local/lib/python3.6/site-packages (from scikit-surprise) (0.16.0)
Requirement already satisfied: scipy>=1.0.0 in /home/rhyme/.local/lib/python3.6/site-packages (from scikit-surprise) (1.5.2)
WARNING: You are using pip version 20.2.1; however, version 20.2.4 is available.
You should consider upgrading via the '/usr/bin/python3 -m pip install --upgrade pip' command.
>> Libraries Installed
```

## 1.2: Importing Libraries

First of all, we will need to import some libraries. This includes surprise which we will use to create the recommendation system.

```
In [2]: print('>> Importing Libraries')

import pandas as pd

from surprise import Reader, Dataset, SVD

from surprise.accuracy import rmse, mae
from surprise.model_selection import cross_validate

print('>> Libraries imported.')

>> Importing Libraries
>> Libraries imported.
```

## Task 2: Importing Data

We will use open-source dataset from GroupLens Research ([movielens.org](http://movielens.org) (<http://movielens.org>))

### 2.1: Importing the Data

The dataset is saved in a `ratings.csv` file. We will use pandas to take a look at some of the rows.

```
In [3]: df=pd.read_csv('ratings.csv')
df.head()
```

```
Out[3]:
```

	userId	movieId	rating	timestamp
0	1	1	4.0	964982703
1	1	3	4.0	964981247
2	1	6	4.0	964982224
3	1	47	5.0	964983815
4	1	50	5.0	964982931

### 2.2 Dropping timestamp

We won't be using the timestamp when user gave the particular rating. So we will drop that column.

```
In [4]: df.drop('timestamp',axis=1,inplace=True)#axis=1 means dropping column
df.head()
```

```
Out[4]:
```

	userId	movieId	rating
0	1	1	4.0
1	1	3	4.0
2	1	6	4.0
3	1	47	5.0
4	1	50	5.0

### 2.3 Check for Missing Data

It's a good practice to check if the data has any missing values. In real world data, this is quite common and must be taken care of before any data pre-processing or model training.

```
In [5]: df.isna().sum()
```

```
Out[5]:
```

userId	0
movieId	0
rating	0
dtype:	int64

## Task 3: EDA (Exploratory data analysis)

In statistics, exploratory data analysis is an approach to analyzing data sets to summarize their main characteristics.

### 3.1 Number of movies/users

```
In [6]: n_movies=df["movieId"].nunique()
n_users=df["userId"].nunique()
print(f"Number of unique movies: {n_movies}")
print(f"Number of unique users: {n_users}")
```

```
Number of unique movies: 9724
Number of unique users: 610
```

### 3.2 Sparsity of our data

Sparsity (%) = (No of missing values / (Total Values))X100

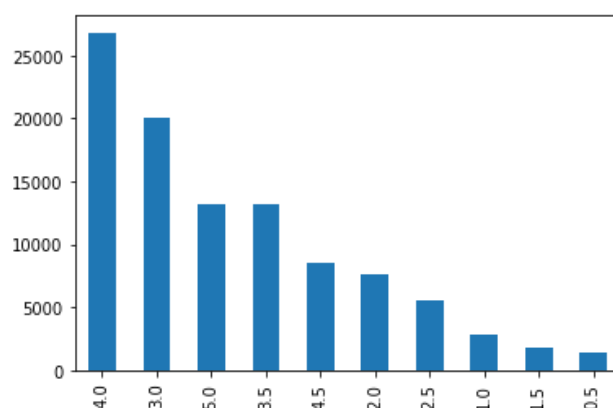
```
In [7]: available_ratings=df["rating"].count()
total_ratings=n_movies*n_users
missing_ratings=total_ratings-available_ratings
sparsity=(missing_ratings/total_ratings)*100
print(f"Sparsity: {sparsity}")
print(missing_ratings,total_ratings,available_ratings)
```

```
Sparsity: 98.30003169443864
5830804 5931640 100836
```

### 3.3 Ratings Distribution

```
In [8]: df["rating"].value_counts().plot(kind="bar")
```

```
Out[8]: <AxesSubplot:>
```



## Task 4: Dimensionality Reduction

To reduce the dimensionality of the dataset, we will filter out rarely rated movies and rarely rating users

## 4.1 Filter movies with less than 3 ratings

```
In [9]: filter_movies=df["movieId"].value_counts() > 3
filter_movies=filter_movies[filter_movies].index.tolist()
#print(filter_movies)
```

## 4.2 Filter users with less than 3 movies rated

```
In [10]: filter_users=df["userId"].value_counts() > 3
filter_users=filter_users[filter_users].index.tolist()
#print(filter_users)
```

## 4.3 Remove rarely rated movies and rarely rating users

```
In [11]: print(f"Original Shape: {df.shape}")
df=df[(df["movieId"].isin(filter_movies)) & (df["userId"].isin(filter_user
s))]
print(f"New Shape: {df.shape}")

Original Shape: (100836, 3)
New Shape: (92394, 3)
```

# Task 5: Create Training and Test Sets

## 5.1 Columns used for training

```
In [12]: cols=["userId","movieId","rating"]
```

## 5.2 Create surprise dataset

```
In [13]: reader=Reader(rating_scale=(0.5,5))
data=Dataset.load_from_df(df[cols],reader)
```

## 5.3 Create Train-set and Prediction-set

```
In [14]: trainset=data.build_full_trainset()#to build model
antiset=trainset.build_anti_testset()#combination of all users and movies th
at do not have training yet
```

# Task 6: Creating and training the model

## 6.1 Creating the model

### SVD (Singular Value Decomposition)

Interaction Matrix =  $A \times B \times C$

```
In [15]: algo=SVD(n_epochs=25,verbose=True)#epochs means go over the data n times and  
try to reduce the error
```

## 6.2 Training the model

**Mean Absolute Error (MAE):** MAE measures the average magnitude of the errors in a set of predictions, without considering their direction.

**Root mean squared error (RMSE):** RMSE is the square root of the average of squared differences between prediction and actual observation.

```
In [16]: cross_validate(algo,data,measures=['RMSE','MAE'],cv=5,verbose=True)#verbose
         is true means we see the model as it is being trained
         print('>> Training Done')
```

Processing epoch 0  
Processing epoch 1  
Processing epoch 2  
Processing epoch 3  
Processing epoch 4  
Processing epoch 5  
Processing epoch 6  
Processing epoch 7  
Processing epoch 8  
Processing epoch 9  
Processing epoch 10  
Processing epoch 11  
Processing epoch 12  
Processing epoch 13  
Processing epoch 14  
Processing epoch 15  
Processing epoch 16  
Processing epoch 17  
Processing epoch 18  
Processing epoch 19  
Processing epoch 20  
Processing epoch 21  
Processing epoch 22  
Processing epoch 23  
Processing epoch 24  
Processing epoch 0  
Processing epoch 1  
Processing epoch 2  
Processing epoch 3  
Processing epoch 4  
Processing epoch 5  
Processing epoch 6  
Processing epoch 7  
Processing epoch 8  
Processing epoch 9  
Processing epoch 10  
Processing epoch 11  
Processing epoch 12  
Processing epoch 13  
Processing epoch 14  
Processing epoch 15  
Processing epoch 16  
Processing epoch 17  
Processing epoch 18  
Processing epoch 19  
Processing epoch 20  
Processing epoch 21  
Processing epoch 22  
Processing epoch 23  
Processing epoch 24  
Processing epoch 0  
Processing epoch 1  
Processing epoch 2  
Processing epoch 3  
Processing epoch 4  
Processing epoch 5  
Processing epoch 6  
Processing epoch 7  
Processing epoch 8  
Processing epoch 9  
Processing epoch 10  
Processing epoch 11  
Processing epoch 12  
Processing epoch 13  
Processing epoch 14  
Processing epoch 15  
Processing epoch 16  
Processing epoch 17



## Task 7: Predictions

### 7.1 Predict ratings for all pairs (user, items) that are NOT in the training set.

```
In [17]: predictions=algo.test(antiset)
```

```
In [18]: predictions[0]
```

```
Out[18]: Prediction(uid=1, iid=318, r_ui=3.529119856267723, est=5, details={'was_impos  
sible': False})
```

### 7.2 Recommending top 3 movies movies based on predictions

```
In [20]: from collections import defaultdict
def get_top_n(predictions,n):
    top_n=defaultdict(list)
    for uid,iid,_,est,_ in predictions:
        top_n[uid].append((iid,est))

    for uid,user_ratings in top_n.items():
        user_ratings.sort(key=lambda x: x[1],reverse=True)#x[0]=iid,x[1]=est
        top_n[uid]=user_ratings[:n]
    return top_n
pass
top_n=get_top_n(predictions,n=3)
for uid,user_ratings in top_n.items():
    print(uid,[iid for (iid,rating) in user_ratings])
```

```
1 [318, 1704, 58559]
2 [750, 898, 1204]
3 [1617, 4226, 1204]
4 [3435, 356, 3275]
5 [912, 1204, 1217]
6 [260, 1204, 2115]
7 [1198, 1217, 1262]
8 [1208, 858, 1198]
9 [1222, 741, 7153]
10 [1204, 2788, 55118]
11 [3451, 912, 866]
12 [50, 110, 527]
13 [48516, 356, 1217]
14 [5690, 912, 3275]
15 [1172, 1204, 1221]
16 [3275, 475, 1204]
17 [106100, 898, 2700]
18 [1204, 475, 750]
19 [527, 5618, 1204]
20 [1215, 1204, 260]
21 [1204, 741, 903]
22 [1204, 608, 1197]
23 [527, 2324, 7153]
24 [3836, 750, 1215]
25 [110, 593, 923]
26 [1204, 2959, 48516]
27 [3504, 5690, 177593]
28 [1204, 750, 2067]
29 [741, 898, 608]
30 [593, 1136, 1197]
31 [969, 1272, 933]
32 [4973, 1208, 1200]
33 [1204, 750, 898]
34 [1617, 2858, 44195]
35 [1204, 914, 111]
36 [750, 177593, 1204]
37 [1204, 912, 527]
38 [1204, 1228, 541]
39 [4973, 750, 48516]
40 [5690, 1223, 778]
41 [1246, 912, 5690]
42 [1256, 1104, 1148]
43 [50, 260, 333]
44 [177593, 608, 2959]
45 [527, 3836, 916]
46 [912, 1208, 1204]
47 [527, 858, 898]
48 [4993, 318, 527]
49 [898, 750, 912]
50 [912, 2028, 50]
51 [1222, 2716, 2959]
52 [527, 1208, 1214]
53 [6, 47, 50]
54 [7361, 1215, 750]
55 [1204, 260, 56367]
56 [1204, 741, 750]
57 [318, 1732, 2329]
58 [2324, 5952, 4973]
59 [899, 2788, 741]
60 [750, 2959, 1276]
61 [527, 1204, 750]
62 [2324, 1172, 899]
63 [175, 54997, 81834]
64 [750, 933, 2788]
65 [7361, 110, 750]
66 [898, 112552, 4973]
67 [1204, 50, 608]
68 [1204, 2788, 3066]
```

In [ ]: