

# Movie Recommendation System

Welcome to the code notebook for Recommender Systems with Python. In this lecture we will develop basic recommendation systems using Python and pandas. In this notebook, we will focus on providing a basic recommendation system by suggesting items that are most similar to a particular item, in this case, movies. Keep in mind, this is not a true robust recommendation system, to describe it more accurately, it just tells you what movies/items are most similar to your movie choice. There is no project for this topic, instead you have the option to work through the advanced lecture version of this notebook (totally optional!).

Let's get started!

## Import Libraries

```
In [1]: import numpy as np
import pandas as pd
```

## Get the Data

```
In [2]: column_names = ['user_id', 'item_id', 'rating', 'timestamp']
df = pd.read_csv('u.data', names=column_names, sep='\t')
```

```
In [3]: df.head()
```

```
Out[3]:
```

	user_id	item_id	rating	timestamp
0	0	50	5	881250949
1	0	172	5	881250949
2	0	133	1	881250949
3	196	242	3	881250949
4	186	302	3	891717742

Now let's get the movie titles:

```
In [4]: movie_titles = pd.read_csv('Movie_Id_Titles')
movie_titles.head()
```

```
Out[4]:
```

	item_id	title
0	1	Toy Story (1995)
1	2	GoldenEye (1995)
2	3	Four Rooms (1995)
3	4	Get Shorty (1995)
4	5	Copycat (1995)

We can merge them together

```
In [5]: df = pd.merge(df, movie_titles, on='item_id')
```

```
In [6]: df.head()
```

Out[6]:

	user_id	item_id	rating	timestamp	title
0	0	50	5	881250949	Star Wars (1977)
1	290	50	5	880473582	Star Wars (1977)
2	79	50	4	891271545	Star Wars (1977)
3	2	50	5	888552084	Star Wars (1977)
4	8	50	5	879362124	Star Wars (1977)

In [7]:

```
pd.merge(df,movie_titles, on='item_id').head(20)
```

Out[7]:

	user_id	item_id	rating	timestamp	title_x	title_y
0	0	50	5	881250949	Star Wars (1977)	Star Wars (1977)
1	290	50	5	880473582	Star Wars (1977)	Star Wars (1977)
2	79	50	4	891271545	Star Wars (1977)	Star Wars (1977)
3	2	50	5	888552084	Star Wars (1977)	Star Wars (1977)
4	8	50	5	879362124	Star Wars (1977)	Star Wars (1977)
5	274	50	5	878944679	Star Wars (1977)	Star Wars (1977)
6	227	50	4	879035347	Star Wars (1977)	Star Wars (1977)
7	99	50	5	885679998	Star Wars (1977)	Star Wars (1977)
8	305	50	5	886321799	Star Wars (1977)	Star Wars (1977)
9	108	50	4	879879739	Star Wars (1977)	Star Wars (1977)
10	63	50	4	875747292	Star Wars (1977)	Star Wars (1977)
11	234	50	4	892079237	Star Wars (1977)	Star Wars (1977)
12	97	50	5	884239471	Star Wars (1977)	Star Wars (1977)
13	117	50	5	880126022	Star Wars (1977)	Star Wars (1977)
14	70	50	4	884064188	Star Wars (1977)	Star Wars (1977)
15	318	50	2	884495696	Star Wars (1977)	Star Wars (1977)
16	145	50	5	885557660	Star Wars (1977)	Star Wars (1977)
17	124	50	3	890287508	Star Wars (1977)	Star Wars (1977)
18	253	50	4	891628518	Star Wars (1977)	Star Wars (1977)
19	271	50	5	885848640	Star Wars (1977)	Star Wars (1977)

## EDA

Let's Explore the data a bit and get a look at some of the best rated movies.

## Visualization Imports

In [8]:

```
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style('white')
%matplotlib inline
```

Let's create a ratings dataframe with average rating and number of rating:

In [9]:

```
df.groupby('title')['rating'].mean().sort_values(ascending=False).head()
```

Out[9]:

title	
They Made Me a Criminal (1939)	5.0
Marlene Dietrich: Shadow and Light (1996)	5.0
Saint of Fort Washington, The (1993)	5.0
Someone Else's America (1995)	5.0
Star Kid (1997)	5.0
Name: rating, dtype: float64	

```
In [10]: df.groupby('title')['rating'].count().sort_values(ascending=False).head()
```

```
Out[10]:
title
Star Wars (1977)      584
Contact (1997)        509
 Fargo (1996)         508
Return of the Jedi (1983)  507
Liar Liar (1997)       485
Name: rating, dtype: int64
```

```
In [11]: ratings = pd.DataFrame(df.groupby('title')['rating'].mean())
```

```
In [12]: ratings.head()
```

```
Out[12]:
           rating
title
'Til There Was You (1997)  2.333333
1-900 (1994)              2.600000
101 Dalmatians (1996)     2.908257
12 Angry Men (1957)       4.344000
187 (1997)                3.024390
```

Now set the number of rating column:

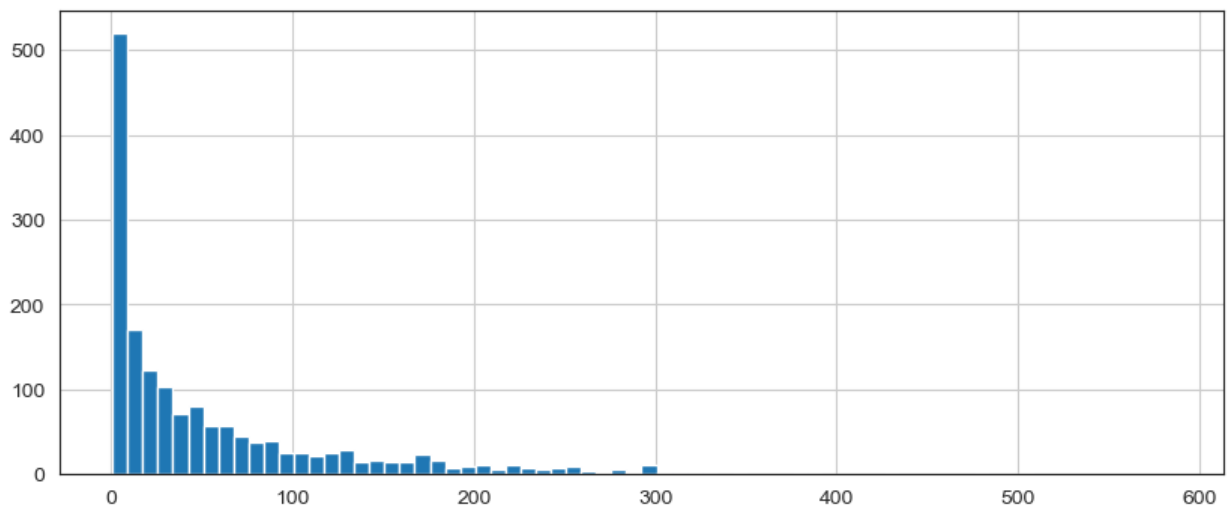
```
In [13]: ratings['num of ratings'] = pd.DataFrame(df.groupby('title')['rating'].count())
```

```
In [14]: ratings.head()
```

```
Out[14]:
           rating  num of ratings
title
'Til There Was You (1997)  2.333333          9
1-900 (1994)              2.600000          5
101 Dalmatians (1996)     2.908257        109
12 Angry Men (1957)       4.344000       125
187 (1997)                3.024390         41
```

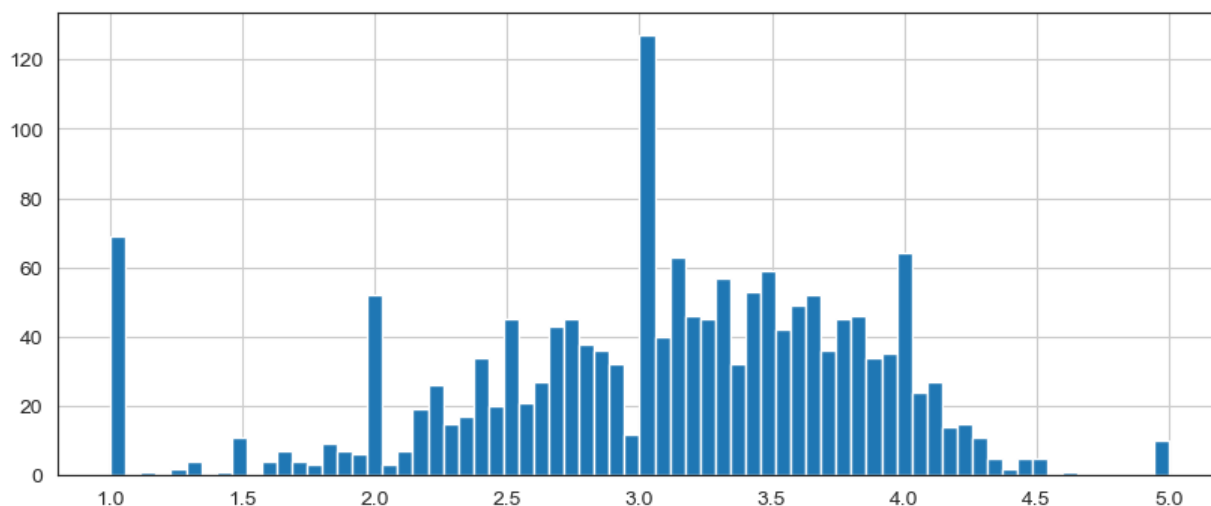
Now a few histogram

```
In [15]: plt.figure(figsize=(10,4))
ratings['num of ratings'].hist(bins=70);
```



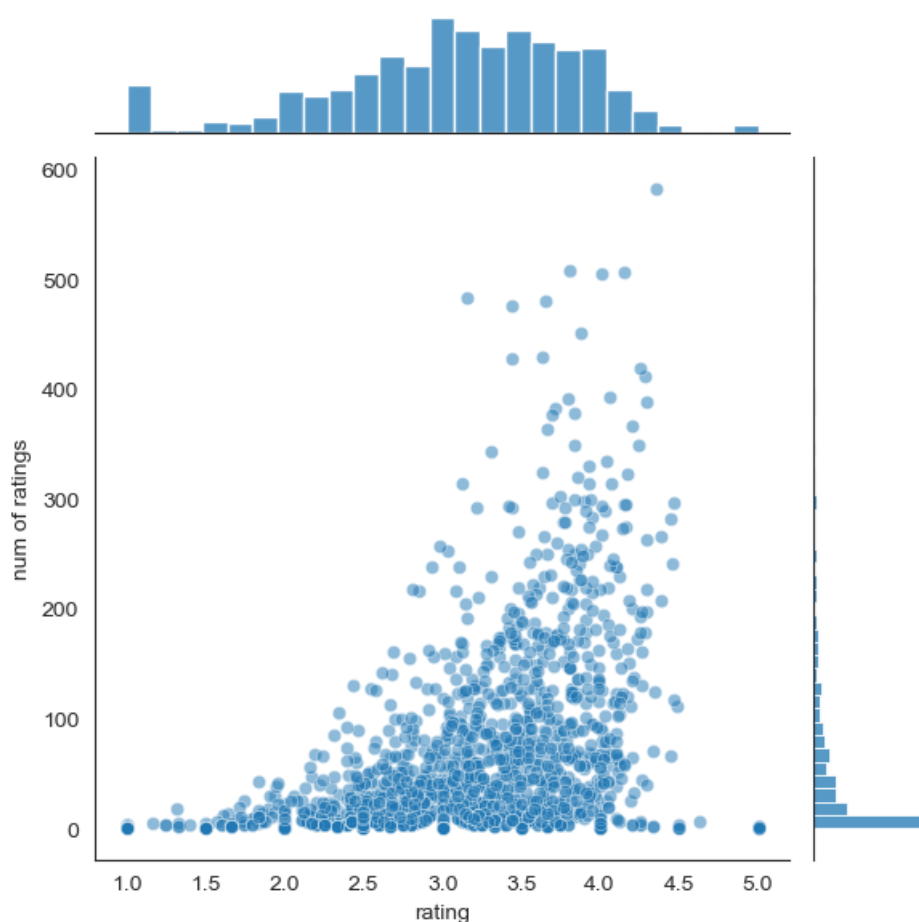
```
In [16]: plt.figure(figsize=(10,4))
ratings['rating'].hist(bins=70)
```

```
Out[16]: <Axes: >
```



```
In [17]: sns.jointplot(x='rating', y='num of ratings', data=ratings,alpha=0.5)
```

```
Out[17]: <seaborn.axisgrid.JointGrid at 0x27ff90aa110>
```



Okay! Now that we have a general idea of what the data look like, let's move on to creating a simple recommendation system:

## Recommending Similar Movies

Now let's create a matrix that has the user ids on one axis and the movie title on another axis. Each cell will then consist of the rating the user gave to that movie. Note there will be a lot of NaN values, because most people have not seen most of the movies.

```
In [18]: moviemat = df.pivot_table(index='user_id', columns='title', values='rating')
```

```
In [19]: moviemat
```

Out[19]:

	title	'Til There Was You (1997)	1-900 (1994)	101 Dalmatians (1996)	12 Angry Men (1957)	187 (1997)	2 Days in the Valley (1996)	20,000 Leagues Under the Sea (1954)	2001: A Space Odyssey (1968)	3 Ninjas: High Noon At Mega Mountain (1998)	39 Steps, The (1935)	...	Yankee Zulu (1994)	Year of the Horse (1997)	You So Crazy (1994)
user_id															
0		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
1		NaN	NaN	2.0	5.0	NaN	NaN	3.0	4.0	NaN	NaN	...	NaN	NaN	NaN
2		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	1.0	NaN	...	NaN	NaN	NaN
3		NaN	NaN	NaN	NaN	2.0	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
4		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
...		...	...	...	...	...	...	...	...	...	...	...	...	...	...
939		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
940		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
941		NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
942		NaN	NaN	NaN	NaN	NaN	NaN	NaN	3.0	NaN	3.0	...	NaN	NaN	NaN
943		NaN	NaN	NaN	NaN	NaN	2.0	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

944 rows × 1664 columns



Most related movies:

```
In [20]: ratings.sort_values('num of ratings', ascending=False).head(10)
```

	rating	num of ratings
title		
Star Wars (1977)	4.359589	584
Contact (1997)	3.803536	509
Fargo (1996)	4.155512	508
Return of the Jedi (1983)	4.007890	507
Liar Liar (1997)	3.156701	485
English Patient, The (1996)	3.656965	481
Scream (1996)	3.441423	478
Toy Story (1995)	3.878319	452
Air Force One (1997)	3.631090	431
Independence Day (ID4) (1996)	3.438228	429

Let's choose two movies: Starwars, a sci-fi movies. And Liar Liar, a comedy.

```
In [21]: starwars_user_ratings = moviemat['Star Wars (1977)']
liarliar_user_ratings = moviemat['Liar Liar (1997)']
```

```
In [22]: starwars_user_ratings.head()
```

user_id	
0	5.0
1	5.0
2	5.0
3	NaN
4	5.0

Name: Star Wars (1977), dtype: float64

We can then use corrwith() method to get correlations between two pandas series:

```
In [23]: similar_to_starwars = moviemat.corrwith(starwars_user_ratings)
similar_to_liarliar = moviemat.corrwith(liarliar_user_ratings)
```

```
C:\Users\RISHABH\anaconda3\lib\site-packages\numpy\lib\function_base.py:2845: RuntimeWarning: Degrees of freedom <= 0 for slice
  c = cov(x, y, rowvar, dtype=dtype)
C:\Users\RISHABH\anaconda3\lib\site-packages\numpy\lib\function_base.py:2704: RuntimeWarning: divide by zero encountered in divide
  c *= np.true_divide(1, fact)
C:\Users\RISHABH\anaconda3\lib\site-packages\numpy\lib\function_base.py:2845: RuntimeWarning: Degrees of freedom <= 0 for slice
  c = cov(x, y, rowvar, dtype=dtype)
C:\Users\RISHABH\anaconda3\lib\site-packages\numpy\lib\function_base.py:2704: RuntimeWarning: divide by zero encountered in divide
  c *= np.true_divide(1, fact)
```

In [24]: `similar_to_starwars`

```
Out[24]:
title
'Til There Was You (1997)          0.872872
1-900 (1994)                      -0.645497
101 Dalmatians (1996)             0.211132
12 Angry Men (1957)              0.184289
187 (1997)                       0.027398
...
Young Guns II (1990)              0.228615
Young Poisoner's Handbook, The (1995) -0.007374
Zeus and Roxanne (1997)           0.818182
unknown                           0.723123
Á köldum klaka (Cold Fever) (1994) NaN
Length: 1664, dtype: float64
```

In [25]: `corr_starwars = pd.DataFrame(similar_to_starwars, columns=['Correlation'])`

In [26]: `corr_starwars.dropna(inplace=True)`  
`corr_starwars.head()`

Out[26]:

	Correlation
title	
'Til There Was You (1997)	0.872872
1-900 (1994)	-0.645497
101 Dalmatians (1996)	0.211132
12 Angry Men (1957)	0.184289
187 (1997)	0.027398

Now if we sort the dataframe by correlation, we should get the most similar movies, however note that we get some results that don't really make sense. This is because there are a lot of movies only watched once by users who also watched star wars (it was the most popular movie).

In [27]: `corr_starwars.sort_values('Correlation', ascending=False).head(10)`

Out[27]:

	Correlation
title	
Commandments (1997)	1.0
Cosi (1996)	1.0
No Escape (1994)	1.0
Stripes (1981)	1.0
Man of the Year (1995)	1.0
Hollow Reed (1996)	1.0
Beans of Egypt, Maine, The (1994)	1.0
Good Man in Africa, A (1994)	1.0
Old Lady Who Walked in the Sea, The (Vieille qui marchait dans la mer, La) (1991)	1.0
Outlaw, The (1943)	1.0

Let's fix this by filtering out movies that have less than 100 reviews (this value was chosen based off the histogram from earlier).

```
In [28]: corr_starwars = corr_starwars.join(ratings['num of ratings'])
corr_starwars.head()
```

```
Out[28]:
```

	Correlation	num of ratings
title		
'Til There Was You (1997)	0.872872	9
1-900 (1994)	-0.645497	5
101 Dalmatians (1996)	0.211132	109
12 Angry Men (1957)	0.184289	125
187 (1997)	0.027398	41

Now sort the values and notice how the titles make a lot more sense:

```
In [29]: corr_starwars[corr_starwars['num of ratings']>100].sort_values('Correlation',ascending=False).head(10)
```

```
Out[29]:
```

	Correlation	num of ratings
title		
Star Wars (1977)	1.000000	584
Empire Strikes Back, The (1980)	0.748353	368
Return of the Jedi (1983)	0.672556	507
Raiders of the Lost Ark (1981)	0.536117	420
Austin Powers: International Man of Mystery (1997)	0.377433	130
Sting, The (1973)	0.367538	241
Indiana Jones and the Last Crusade (1989)	0.350107	331
Pinocchio (1940)	0.347868	101
Frighteners, The (1996)	0.332729	115
L.A. Confidential (1997)	0.319065	297

Now the same for the comedy Liar Liar:

```
In [30]: corr_liarliar = pd.DataFrame(similar_to_liarliar, columns=['Correlation'])
corr_liarliar.dropna(inplace=True)
corr_liarliar = corr_liarliar.join(ratings['num of ratings'])
corr_liarliar[corr_liarliar['num of ratings']>100].sort_values('Correlation',ascending=False).head(10)
```

```
Out[30]:
```

	Correlation	num of ratings
title		
Liar Liar (1997)	1.000000	485
Batman Forever (1995)	0.516968	114
Mask, The (1994)	0.484650	129
Down Periscope (1996)	0.472681	101
Con Air (1997)	0.469828	137
Pretty Woman (1990)	0.469790	164
101 Dalmatians (1996)	0.469765	109
Michael (1996)	0.442022	119
Waterworld (1995)	0.438405	102
Indiana Jones and the Last Crusade (1989)	0.414427	331

You can find this project on [GitHub](#).