

Import the Library

```
In [13]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

load the dataset

```
In [14]: df_user=pd.read_csv('users.dat',sep=";",names=['UserID','Gender','Age','Occupation','Zip Code'],engine='python')
```

Out[15]:

	UserID	Gender	Age	Occupation	Zip Code
0	1	F	1	10	48067
1	2	M	56	16	70072
2	3	M	25	15	55117
3	4	M	45	7	02400
4	5	M	25	20	56455
...
6035	6036	F	25	15	32603
6036	6037	F	45	1	76006
6037	6038	F	56	1	14706
6038	6039	F	45	0	01060
6039	6040	M	25	6	11106

6040 rows x 5 columns

```
In [14]: df_user['Gender'].value_counts().idxmax()
```

Out[15]:

```
'M'
```

```
In [15]: df_user.shape
```

Out[15]: (6040, 5)

Out[16]:

```
In [15]: df_user.isna().sum().any
<bound method Series.any of UserID      0
Age      0
Gender    0
Occupation 0
Zip Code  0
dtype: int64>
```

```
In [17]: df_movies=pd.read_csv('movies.dat',sep=";",names=['MovieID','Title','Genres'],engine='python')
```

Out[18]:

	MovieID	Title	Genres
0	1	Toy Story (1995)	Animation Children's Comedy
1	2	Jurassic Park (1993)	Adventure Children's Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama
4	5	Father of the Bride Part II (1995)	Comedy

In [19]: df_movies.shape

Out[19]: (3683, 3)

Out[10]:

```
In [19]: df_movies.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3683 entries, 0 to 3682
Data columns (total 3 columns):
# Column Non-Null Count  Dtype
---  ---  ---
0 MovieID 3683 non-null   int64
1 Title   3683 non-null   object
2 genres  3683 non-null   object
memory usage: 91.1+ KB
```

```
In [11]: df_movies.isna().sum()
```

Out[11]: MovieID 0
Title 0
Genres 0
dtype: int64

```
In [12]: df_ratings=pd.read_csv('ratings.dat',sep=";",names=['UserID','MovieID','Title','Timestamp'],engine='python')
```

Out[13]:

Out[13]: (1088209, 4)

In [14]: df_ratings.isna().sum()

Out[14]:

```
UserID      0
MovieID      0
Title       0
Timestamp    0
dtype: int64
```

In [15]: df_ratings.head()

Out[15]:

	UserID	MovieID	Title	Timestamp
0	1	1193	5	978530760
1	1	661	3	9783021209
2	1	914	4	9783019568
3	1	3408	4	9783002075
4	1	2355	5	978524291

Now merge the dataset movies and ratings and then merge the resultant dataset with user dataset

```
In [16]: df_MovieRatings= df_movies.merge(df_ratings,on='MovieID',how='inner')
```

Out[17]: df_MovieRatings.head()

	MovieID	Title_x	Genres	UserID	Title_y	Timestamp
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978524298
1	1	Toy Story (1995)	Animation Children's Comedy	6	5	978237008
2	1	Toy Story (1995)	Animation Children's Comedy	8	4	978232496
3	1	Toy Story (1995)	Animation Children's Comedy	9	5	978252952
4	1	Toy Story (1995)	Animation Children's Comedy	10	5	97825474

In [18]: df_MovieRatings.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1088209 entries, 0 to 1088208
Data columns (total 6 columns):
# Column Non-Null Count  Dtype
---  ---  ---
0 MovieID 1088209 non-null   int64
1 Title_x  1088209 non-null   object
2 genres   1088209 non-null   object
3 UserID   1088209 non-null   int64
4 Title_y  1088209 non-null   int64
5 Timestamp 1088209 non-null   int64
memory usage: 53.4+ MB
```

In [19]: df_MovieRatings.shape

Out[19]: (1088209, 6)

In [20]: df_MovieRatings.isna().sum()

Out[20]:

```
MovieID      0
Title_x      0
Title_y      0
UserID      0
Genre       0
Timestamp    0
dtype: int64
```

```
In [21]: df_master=df_MovieRatings.merge(df_user,on='UserID',how='inner')
```

Out[21]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1088209 entries, 0 to 1088208
Data columns (total 18 columns):
# Column Non-Null Count  Dtype
---  ---  ---
0 MovieID 1088209 non-null   int64
1 Title_x  1088209 non-null   object
2 genres   1088209 non-null   object
3 UserID   1088209 non-null   int64
4 Title_y  1088209 non-null   int64
5 Timestamp 1088209 non-null   int64
6 Age       1088209 non-null   int64
7 Zip Code  1088209 non-null   object
memory usage: 10.9+ MB
```

Out[21]: False

```
In [22]: # to csv file
df_master.to_csv('master_data.csv')
```

Out[22]:

	MovieID	Title_x	Genres	UserID	Title_y	Timestamp	Gender	Age	Occupation	Zip Code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978524298	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978524261	F	1	10	48067
2	150	Apollo 13 (1995)	Drama	1	5	9785031777	F	1	10	48067
3	260	Star Wars: Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300700	F	1	10	48067
4	527	Schindler's List (1993)	Drama War	1	5	978524195	F	1	10	48067

Explore the data using EDA

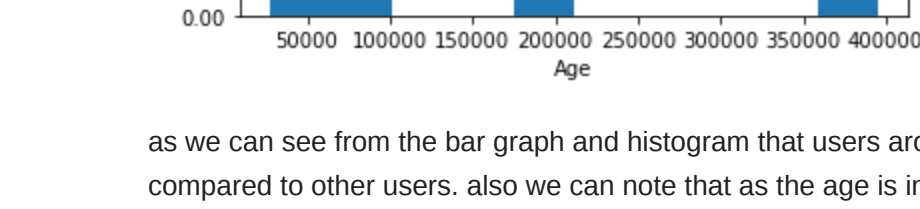
User Age Distribution

```
In [25]: df_master['Age'].value_counts().plot(kind='bar')
```

Out[25]:



```
In [26]: df_master['Age'].value_counts().plot(kind='hist')
plt.xlabel('Age')
plt.ylabel('User Count')
plt.title('User Age Distribution')
plt.show()
```



as we can see from the bar graph and histogram that users around the age of 25 are the one who are more involved in giving ratings than other users this also means that they are the one who are watching more movies as compared to other users, also we can note that as the age is increasing, less user are involved in ratings

User Ratings of the movie 'Toy Story'

```
In [27]: #Extract Toy Story
toystory=df_master[df_master['Title_x'].str.contains('Toy Story')==True]
```

Out[27]:

	MovieID	Title_x	Genres	UserID	Title_y	Timestamp	Gender	Age	Occupation	Zip Code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978524298	F	1	10	48067
1	3114	Toy Story 2 (1999)	Animation Children's Comedy	4	4	97802314	F	1	10	48067
53	1	Toy Story (1995)	Animation Children's Comedy	6	4	978327008	F	50	9	55117
124	1	Toy Story (1995)	Animation Children's Comedy	8	4	978232496	M	25	12	11413
263	1	Toy Story (1995)	Animation Children's Comedy	9	5	978252952	M	25	17	61654

...
99888	3114	Toy Story 2 (1999)	Animation Children's Comedy	3023	4	97041248	F	25	7	82108
99907	3114	Toy Story 2 (1999)	Animation Children's Comedy	5600	2	96805206	M	26	18	90968
99946	3114	Toy Story 2 (1999)	Animation Children's Comedy	2189	4	974607816	M	1	10	60148
99969	3114	Toy Story 2 (1999)	Animation Children's Comedy	159	4	96956644	F	45	0	37922
1000192	3114	Toy Story 2 (1999)	Animation Children's Comedy	5727	5	956492554	M	25	4	92843

3662 rows x 10 columns

```
In [28]: toystory.groupby(['Title_x','Title_y']).size()
```

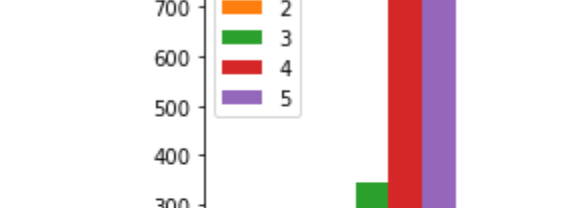
Out[28]:

```
Toy Story (1995)
Title_y
1      16
2      61
3      345
4      835
5      820
```

Toy Story 2 (1999) 1 25
2 44
3 214
4 578
5 724

```
In [30]: toystory.groupby(['Title_x','Title_y']).size().unstack().plot(kind='barh',legend=True)
```

Out[30]:



```
In [31]: toystory.groupby(['Title_x','Title_y']).size().unstack().plot(kind='bar',legend=True)
```

Out[31]:



The above graph indicates that, Toy Story 2 has been rated 5 by the most users where in the case of Toy Story 1 ratings 4 has a little edge over ratings 5, so most of the users have given 4 or 5 ratings which is positive note for the movie.

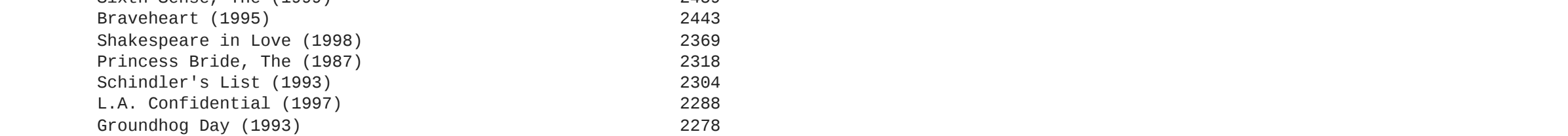
Top 25 movies by viewership rating

```
In [32]: dfTop25=df_master.groupby('Title_x').size().sort_values(ascending=False)[1:25]
dfTop25
```

	Title_x	size
0	American Beauty (1999)	3428
1	Star Wars: Episode V - The Empire Strikes Back (1980)	2998
2	Star Wars: Episode VI - Return of the Jedi (1983)	2883
3	Jurassic Park (1993)	2653
4	Saving Private Ryan (1998)	2653
5	Terminator 2: Judgment Day (1991)	2649
6	Matrix, The (1999)	2589
7	Back to the Future (1985)	2578
8	Silence of the Lambs, The (1991)	2578
9	Men in Black (1997)	2578
10	Raiders of the Lost Ark (1981)	2514
11	Peter Dinklage (1995)	2513
12	Sixth Sense, The (1999)	2459
13	Braveheart (1995)	2443
14	Shakespeare in Love (1998)	2389
15	Princess Bride, The (1987)	2319
16	Schindler's List (1993)	2284
17	L.A. Confidential (1997)	2227
18	Groundhog Day (1993)	2278
19	E.T. the Extra-Terrestrial (1982)	2269
20	Star Wars: Episode I - The Phantom Menace (1999)	2245
21	Being John Malkovich (1999)	2241
22	Shogun: Redemption, The (1994)	2227
23	Godfather, The (1972)	2223
24

Out[33]:

```
<AxesSubplot: xlabel='Title_x'>
```



here are the top25 movies by viewership ratings, the insights that we can take from the above plot is that American Beauty is the top most movies in the viewership rating

Find the ratings for all the movies reviewed by a particular user of user id = 2696

```
In [34]: user_2696 = df_master[df_master.UserID==2696]
user_2696
```

Out[34]:

	MovieID	Title_x	Genres	UserID	Title_y	Timestamp	Gender	Age	Occupation	Zip Code
91035	350	Client, The (1994)	Drama Mystery Thriller	2696	3	97308886	M	25	7	24210
91036	1002	Love Story (1996)	Drama Mystery	2696	5	97308842	F	1	10	48067
91037	892	Basic Instinct (1992)	Mystery Thriller	2696	4	97308886	M	25	7	24210
91038	1007	E.T. the Extra-Terrestrial (1982)	Children's Drama Fantasy Sci-Fi	2696	3	97308890	M	25	7	24210
91039	1258	Swing, The (1997)	Comedy	2696	4	97308710	M	25	7	24210
91040	1270	Back to the Future (1985)	Comedy Sci-Fi	2696	2	97308676	M	25	7	24210
91041	1589	Cop Land (1997)	Crime Drama Mystery	2696	4	97308865	M	25	7	24210
91042	1617	L.A. Confidential (1997)	Crime Film-Noir Mystery Thriller	2696	4	97308842	M	25	7	24210
91043	1625	Game, The (1997)	Mystery Thriller	2696	4	97308842	M	25	7	24210
91044	1644	I Know What You Did Last Summer (1997)	Horror Mystery Thriller	2696	2	97308920	M	25	7	24210
91045	1645	Twelve Monkeys (1995)	Science Fiction Thriller	2696	4	97308904	M	25	7	24210
91046	1711	Man on the Moon (1999)	Comedy Crime Drama Mystery	2696	4	97308904	M	25	7	24210
91047	1783	Palmiro (1998)	Film-Noir Mystery Thriller	2696	4	97308865	M	25	7	24210
91048	1805	Wild Things (1998)	Crime Drama Mystery Thriller	2696	4	97308886	M	25	7	24210
91049	1892	Perfect Murder: A (1998)	Mystery Thriller	2696	4	97308804	M	25	7	24210
91050	2338	I Still Know What You Did Last Summer (1998)	Horror Mystery Thriller	2696	2	97308920	M	25	7	24210
91051	2369	Psycho (1960)	Crime Horror Thriller	2696	4	97308710	M	25	7	24210
91052	2713	Law Paced (1998)	Horror Thriller	2696	1	97308710	M	25	7	24210
91053	3116	Taken: Mr. Ripley, The (1999)	Drama Mystery Thriller	2696	4	97308865	M	25	7	24210
91054	3386	JFK (1991)	Drama Mystery	2696	1	97308842	M	25	7	24210

Feature Engineering:

Use column genres: Find out all the unique genres (Here split the data in column genre making a list and then process the data to find out only the unique categories of genres)

Out[35]:

```
df_master['Genres']
```

Out[35]:

```
0 Animation|Children's|Comedy
1 Animation|Children's|Musical|Romance
2 Animation|Adventure|Fantasy|Sci-Fi
3 Drama|War
...
```

```
1088204 3513 American Beauty (1999)
1088205 3535 Keeping the Faith
1088206 3555 U-571 (2000)
1088207 3578 Gladiator (2000)
Name: Genres, Length: 1088209, dtype: object
```

```
In [36]: dfGenere=df_master['Genres'].str.split('|')
```

Out[37]:

```
dfGenre
```

Out[37]:

```
0 [Animation, Children's, Comedy]
1 [Animation, Children's, Musical, Romance]
2 [Animation, Children's, Comedy, Crime]
3 [Action, Adventure, Fantasy, Sci-Fi]
4 [Drama, War]
...
```

```
1088204 [Drama, Thriller]
1088205 [Comedy, Horror, Thriller]
1088206 [Comedy, Romance]
1088207 [Action, Thriller]
1088208 [Action, Drama]
Name: Genres, Length: 1088209, dtype: object
```

```
In [38]: listgenres=set()
for genre in dfGenre:
    listgenres=listgenres.union(set(genre))
```

Out[38]:

```
listgenres
```

Out[39]:

```
['Action',
'Adventure',
'Animation',
'Children's',
'Comedy',
'Crime',
'Documentary',
'Drama',
'Film-Noir',
'Horror',
'Musical',
'Mystery',
'Romance',
'Sci-Fi',
'Thriller',
'War',
'Western']
```

Out[40]:

```
len(listgenres)
18
```

Create a separate column for each genre category with a one-hot encoding (1 and 0) whether or not the movie belongs to that genre