

```
# importing libra for pandas and numpy
import pandas as pd
import numpy as np
```

Import the three datasets

```
movies = pd.read_csv('movies.dat', sep=';', header=None, engine='python', names=['MovieID', 'Title', 'Genres'])
# MovieID:Title:Genres
```

```
movies.head()
```

	MovieID	Title	Genres
0	1	Toy Story (1995)	Animation Children's Comedy
1	2	Junany (1995)	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

```
rating = pd.read_csv('ratings.dat', sep=';', header=None, engine='python', names=['UserID', 'MovieID', 'Rating', 'Timestamp'])
# UserID:MovieID:Rating:Timestamp
```

```
rating.head()
```

	UserID	MovieID	Rating	Timestamp
0	1	1193	5	978300760
1	1	661	3	978302109
2	1	914	3	978301968
3	1	3408	4	978300275
4	1	2355	5	978824291

```
users = pd.read_csv('users.dat', sep=';', header=None, engine='python', names=['UserID', 'Gender', 'Age', 'Occupat
# UserID:Gender:Age:Occupation:Zip-code
```

```
users.head()
```

	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	02460
4	5	M	25	20	55455

```
# shape of the datasets
```

```
print(f'Movie: {movies.shape}, Rating: {rating.shape}, Users: {users.shape}')
```

```
Movie: (3883, 3), Rating: (100209, 4), Users: (6040, 5)
```

Checking NaN is datasets

```
movies.isnull().sum()
```

```
MovieID 0
Title 0
Genres 0
dtype: int64
```

```
rating.isnull().sum()
```

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

```
users.isnull().sum()
```

```
UserID 0
Gender 0
Age 0
Occupation 0
Zip-code 0
dtype: int64
```

```
# checking for duplicated values
```

```
print(f'Movie duplicates: {movies.duplicated().sum()}, Rating duplicates: {rating.duplicated().sum()}, User du
Movie duplicates: 0, Rating duplicates: 0, User duplicates: 0
```

Create a new dataset

- (i) Merge two tables at a time
- (ii) Merge the tables using two primary keys MovieID & UserID)

```
one = pd.merge(movies, rating, on='MovieID')
```

```
two = pd.merge(one, users, on='UserID')
```

```
two.head()
```

	MovieID	Title	Genres	UserID	Rating	Timestamp	Gender	Age	Occupation	Zip-code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978824268	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978824351	F	1	10	48067
2	150	Apollo 13 (1995)	Drama	1	5	978301777	F	1	10	48067
3	260	Star Wars Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300760	F	1	10	48067
4	527	Schindler's List (1993)	Drama War	1	5	978824195	F	1	10	48067

```
two.shape
```

```
(1000209, 10)
```

```
master_data = two.copy()
```

Master Data after merging two datasets

```
master_data.head()
```

	MovieID	Title	Genres	UserID	Rating	Timestamp	Gender	Age	Occupation	Zip-code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978824268	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978824351	F	1	10	48067
2	150	Apollo 13 (1995)	Drama	1	5	978301777	F	1	10	48067
3	260	Star Wars Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300760	F	1	10	48067
4	527	Schindler's List (1993)	Drama War	1	5	978824195	F	1	10	48067

Explore the datasets using visual representations (graphs or tables), also include your comments on the following:

- User Age Distribution
- User rating of the movie "Toy Story"
- Top 25 movies by viewership rating
- Find the ratings for all the movies reviewed by a particular user of user id = 2696

```
# changing the columns into lower case
```

```
master_data.columns = master_data.columns.str.lower()
```

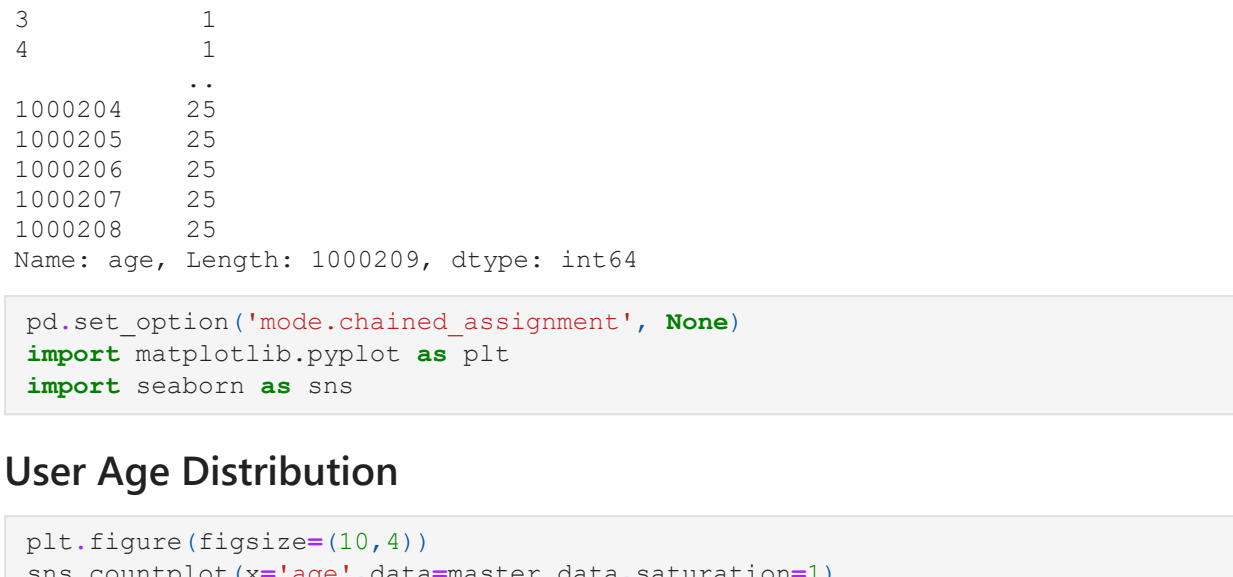
```
master_data['age']
```

```
0 1
1 1
2 1
3 1
4 1
..
1000204 25
1000205 25
1000206 25
1000207 25
1000208 25
Name: age, Length: 1000209, dtype: int64
```

```
pd.set_option('mode.chained_assignment', None)
import matplotlib.pyplot as plt
import seaborn as sns
```

User Age Distribution

```
plt.figure(figsize=(10,4))
sns.countplot(x='age', data=master_data, saturation=1)
plt.show()
```



Insights:

- From the countplot, it is seen that users of age around 25 has seen many movies

User rating of the movie "Toy Story"

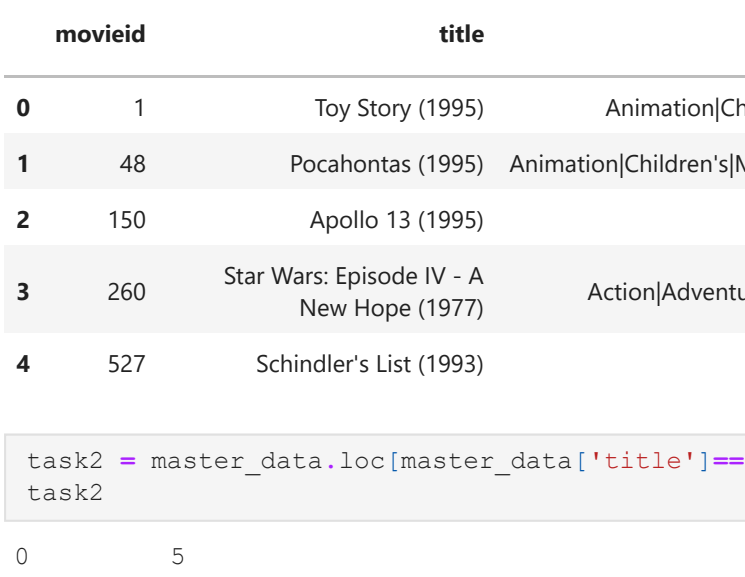
```
master_data.head()
```

	movieid	title	genres	userid	rating	timestamp	gender	age	occupation	zip-code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978824268	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978824351	F	1	10	48067
2	150	Apollo 13 (1995)	Drama	1	5	978301777	F	1	10	48067
3	260	Star Wars Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300760	F	1	10	48067
4	527	Schindler's List (1993)	Drama War	1	5	978824195	F	1	10	48067

```
task2 = master_data.loc[master_data['title']=='Toy Story (1995)', 'rating']
```

```
task2
0 5
53 4
124 4
263 5
369 5
..
575166 5
575214 5
575485 4
575589 4
575653 3
Name: rating, Length: 2077, dtype: int64
```

```
sns.distplot(task2, bins=10)
plt.xlabel('Ratings')
plt.ylabel('Count')
plt.title('TOY STORY- User Rating')
plt.show()
```



Comments:

- 'Histogram' has been used above.
- 'Rating 4' is most given by user for the movie ToyStory

Top 25 movies by viewership rating

```
master_data.head()
```

	movieid	title	genres	userid	rating	timestamp	gender	age	occupation	zip-code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978824268	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978824351	F	1	10	48067
2	150	Apollo 13 (1995)	Drama	1	5	978301777	F	1	10	48067
3	260	Star Wars Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300760	F	1	10	48067
4	527	Schindler's List (1993)	Drama War	1	5	978824195	F	1	10	48067

```
task3 = master_data.groupby('title')['rating'].mean()
task3.head()
```

```
title
$1,000,000 Duck (1971) 3.027027
Night Mother (1986) 4.371429
'Til There Was You (1997) 2.692308
Dorothy (The 1959) 2.910891
...And Justice for All (1979) 3.713568
Name: rating, dtype: float64
```

```
task3_25 = task3.sort_values(ascending=False)
task3_frame = task3_25.to_frame()[0:25]
task3_frame
```

	rating
title	
Gate of Heavenly Peace, The (1995)	5.000000
Lured (1947)	5.000000
Ulysses (Ulyse) (1954)	5.000000
Smashing Time (1967)	5.000000
Follow the Birch (1998)	5.000000
Song of Freedom (1936)	5.000000
Bittersweet Motel (2000)	4.554558
Baby, The Indian (1973)	5.000000
One Little Indian (1973)	5.000000
Schlafes Bruder (Brother of Sleep) (1995)	5.000000
I Am Cuba (Soy Cuba/Ya Cuba) (1964)	4.800000
Lamerica (1979)	4.750000
Apple, The (Sib) (1998)	4.666667
Sanjuro (1962)	4.608696
Seven Samurai (The Magnificent Seven) (Shichinin no samurai) (1954)	4.560510
Shawshank Redemption, The (1994)	4.554558
Godfather, The (1972)	4.524966
Close Shave, A (1995)	4.520548
Unsuspected, The (1995)	4.517106
Schindler's List (1993)	4.510417
Wrong Trousers, The (1993)	4.507937
Dangerous Game (1993)	4.500000
Mamma Roma (1962)	4.500000
Inheritors, The (Die Siebtebaumen) (1998)	4.500000
Hour of the Pig, The (1993)	4.500000

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

```
task4 = master_data.loc[master_data['userid']==2696, ['title', 'rating']]
task4.head(25)
```

	title	rating
991035	Client, The (1954)	3
991036	Lone Star (1996)	5
991037	Basic Instinct (1992)	4
991038	ET, the Extra-Terrestrial (1982)	3
991039	Shining, The (1980)	4
991040	Back to the Future (1985)	2
991041	Cop Land (1997)	3
991042	LA Confidential (1997)	4
991043	Game, The (1997)	4
991044	I Know What You Did Last Summer (1997)	2
991045	Devil's Advocate, The (1997)	4
991046	Midnight in the Garden of Good and Evil (1997)	4
991047	Palmetto (1998)	4
991048	Wild Things (1998)	4
991049	Perfect Murder, A (1998)	4
991050	I Still Know What You Did Last Summer (1998)	2
991051	Lake (1998)	4
991052	Psyche Placid (1999)	1
991053	Talented Mr. Ripley, The (1999)	4
991054	JFK (1991)	1

Find out all the unique genres

```
un_genres = master_data['genres'].copy()
un_genres.duplicated().sum()
```

```
999908
```

```
un_genres.drop_duplicates(keep='first')
task5
0 Animation|Children's|Comedy
1 Animation|Children's|Musical|Romance
2 Drama
3 Action|Adventure|Fantasy|Sci-Fi
4 Drama|War
..
79242 Romance|Western
80659 Drama|Romance|Western
236412 Comedy|Film-Noir|Thriller
303913 Film-Noir|Horror
920103 Fantasy
Name: genres, Length: 301, dtype: object
```

```
from itertools import combinations
from collections import Counter
```

```
count = Counter()
```

```
for row in un_genres:
    row_list = row.split('|')
    un_count.update(combinations(row_list,1))
print(count)
```

```
{('Comedy',): 356580, ('Drama',): 354529, ('Action',): 257457, ('Thriller',): 189680, ('Sci-Fi',): 1
57294, ('Romance',): 147523, ('Adventure',): 133953, ('Crime',): 78341, ('Horror',): 76386, ('Children's',)
72186, ('War',): 68527, ('Animation',): 43293, ('Musical',): 41533, ('Mystery',): 40178, ('Fantasy',): 3630
1, ('Western',): 20683, ('Film-Noir',): 18261, ('Documentary',): 7910}
```

```
m = count.most_common()
```

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

Create a separate column for each genre category with a one-hot encoding

```
un_genres = master_data['genres'].str.get_dummies()
dummy
```

	Action	Adventure	Animation	Children's	Comedy	Crime	Documentary	Drama	Fantasy	Film-Noir	Horror	Musical	Mystery	Rom
0	0	0	1	1	1	0	0	0	0	0	0	0	0	0
1	0	0	1	1	0	0	0	0	0	0	0	0	1	0
2	0	0	0	0	0	0	0	1	0	0	0	0	0	0
3	1	1	0	0	0	0	0	0	0	1	0	0	0	0
4	0	1	0	0	0	0	0	0	1	0	0	0	0	0
...
1000204	0	0	0	0	0	0	0	0	1	0	0	0	0	0
1000205	0	0	0	0	0	1	0	0	0	0	0	1	0	0
1000206	0	0	0	0	0	1	0	0	0	0	0	0	0	0
1000207	1	0	0	0	0	0	0	0	0	0	0	0	0	0
1000208	1	0	0	0	0	0	0	0	1	0	0	0	0	0
1000209 rows x 18 columns														

Determine the features affecting the ratings of any particular movie.

```
master_data.head()
```

	movieid	title	genres	userid	rating	timestamp	gender	age	occupation	zip-code
0	1	Toy Story (1995)	Animation Children's Comedy	1	5	978824268	F	1	10	48067
1	48	Pocahontas (1995)	Animation Children's Musical Romance	1	5	978824351	F	1	10	48067
2	2	Apollo 13 (1995)	Drama	1	5	978301777	F	1	10	48067
3	3	Star Wars: Episode IV - A New Hope (1977)	Action Adventure Fantasy Sci-Fi	1	4	978300760	F	1	10	48067
4	4	Schindler's List (1993)	Drama War	1	5	978824195	F	1	10	48067
...
1000204	1000204	Rules of Engagement (2000)	Drama Thriller	5727	4	958489970	M	25	4	...
1000205	1000205	American Psycho (2000)	Comedy Horror Thriller	5727	2	958489970	M	25	4	...
1000206	1000206	Keeping the Faith (2000)	Comedy Romance	5727	5	958489992	M	25	4	...
1000207	1000207	U-571 (2000)	Action Thriller	5727	3	958490069	M	25	4	...
1000208	1000208	Gladiator (2000)	Action Drama	5727	5	958490171	M	25	4	...
1000209 rows x 29 columns										

```
new_df.drop(columns='key_0', inplace=True)
```

```
new_df.columns
```

```
Index(['movieid', 'title', 'genres', 'userid', 'rating', 'timestamp', 'gender', 'age', 'occupation', 'zip-code', 'Fantasy', 'Film-Noir', 'Horror', 'Musical', 'Mystery', 'Romance', 'Sci-Fi', 'Thriller', 'War', 'Western'], dtype='object')
```

```
new_df.dtypes
```

```
movieid int64
title object
genres object
userid int64
rating int64
timestamp int64
gender object
age int64
occupation int64
zip-code object
Action int64
Adventure int64
Animation int64
Children's int64
Comedy int64
Crime int64
Documentary int64
Drama int64
Fantasy int64
Film-Noir int64
Horror int64
Musical int64
Mystery int64
Romance int64
Sci-Fi int64
Thriller int64
War int64
Western object
dtype: object
```

```
feature_affecting = new_df.copy()
```

```
feature_affecting.drop(columns=['title', 'genres', 'zip-code'], inplace=True)
```

```
feature_affecting['gender'].replace('F', 0, inplace=True)
```

```
feature_affecting['gender'].replace('M', 1, inplace=True)
```

```
feature_affecting.head()
```

	movieid	rating	timestamp	gender	age	occupation	Action	Adventure	Animation
--	---------	--------	-----------	--------	-----	------------	--------	-----------	-----------

