

Python project

By: Chirag Malhotra

MovieLens 100k Dataset

Data Source: Kaggle

MovieLens data sets were collected by the GroupLens Research Project at the University of Minnesota.

This data set consists of:

- 100,000 ratings (1-5) from 943 users on 1682 movies.
- Each user has rated at least 20 movies.
- Simple demographic info for the users (age, gender, occupation, zip)

The data was collected through the MovieLens web site (movielens.umn.edu) during the seven-month period from September 19th, 1997 through April 22nd, 1998.

Importing the data

```
In [15]: 1 import pandas as pd
          2 import numpy as np
          3 columns = ['User_id', 'Movie_id', 'Rating', 'timestamp']
          4 file_path = r"D:\AIDTM\term 2\EDA\archive\ml-100k\u.data"
          5 movie_data = pd.read_csv(file_path, sep='\t', names=columns)
          6 movie_data
```

```
Out[15]:
```

	User_id	Movie_id	Rating	timestamp
0	196	242	3	881250949
1	186	302	3	891717742
2	22	377	1	878887116
3	244	51	2	880606923
4	166	346	1	886397596
...
99995	880	476	3	880175444
99996	716	204	5	879795543
99997	276	1090	1	874795795
99998	13	225	2	882399156
99999	12	203	3	879959583

100000 rows × 4 columns

Part 1 : Data Cleaning

1). Identify and handle missing values in the dataset. Are there any rows or columns you should remove? Justify your decisions.

```
In [31]: 1 movie_data.isnull().sum()
```

```
Out[31]: User_id          0
Movie_id          0
Rating            0
timestamp         0
day_of_week       0
hour_of_day       0
time_of_day_category 0
dtype: int64
```

Interpretation:

- There are no missing values in the dataset, so there is no need to remove the rows or columns in the dataset.
- If there are more no of values missing in the columns than it is always advised to remove the column, cause that may skew the result.
- Skewed results may give us wrong insights, which may result in bad business decisions.

2). Explore the "timestamp" column. Can you extract any useful information like day of the week or time of day for rentals? Clean and format the timestamp data accordingly.

```
In [32]: 1 import pandas as pd
2 movie_data['timestamp'] = pd.to_datetime(movie_data['timestamp'], unit='s')
3 movie_data['rental_day_of_week'] = movie_data['timestamp'].dt.day_name()
4 movie_data['rental_time_of_day'] = movie_data['timestamp'].dt.time
5
6 print(movie_data[['timestamp', 'rental_day_of_week', 'rental_time_of_d']])
```

	timestamp	rental_day_of_week	rental_time_of_day
0	1970-01-01 00:00:00.881250949	Thursday	00:00:00.881250
1	1970-01-01 00:00:00.891717742	Thursday	00:00:00.891717
2	1970-01-01 00:00:00.878887116	Thursday	00:00:00.878887
3	1970-01-01 00:00:00.880606923	Thursday	00:00:00.880606
4	1970-01-01 00:00:00.886397596	Thursday	00:00:00.886397
5	1970-01-01 00:00:00.884182806	Thursday	00:00:00.884182
6	1970-01-01 00:00:00.881171488	Thursday	00:00:00.881171
7	1970-01-01 00:00:00.891628467	Thursday	00:00:00.891628
8	1970-01-01 00:00:00.886324817	Thursday	00:00:00.886324
9	1970-01-01 00:00:00.883603013	Thursday	00:00:00.883603

Interpretation:

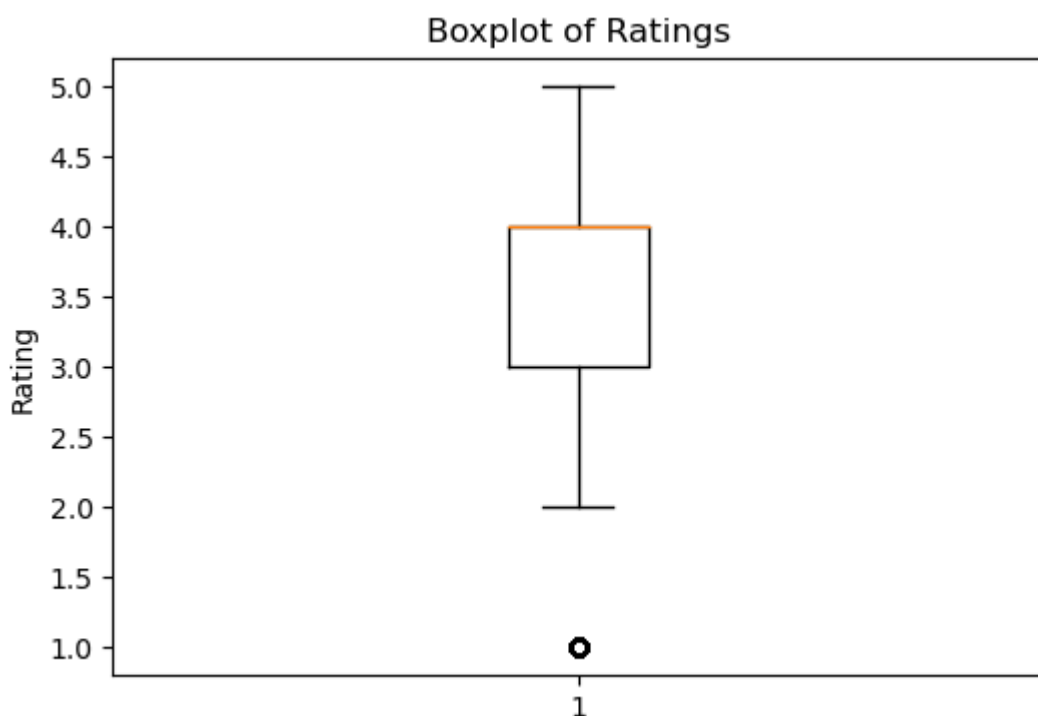
- With this timestamp column we have got the day of the week and also the time of the day.

- By visualizing and analyzing the day of the week and time of day for rentals, we can know the trends that can be useful for decision-making, such as optimizing rental operations, identifying popular time slots, or adjusting pricing strategies based on demand during certain hours or days

3). Look for outliers and inconsistencies in the "rating" column. How would you handle outliers? Justify your chosen method.

```
In [34]: 1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 print(movie_data['Rating'].describe())
6 plt.figure(figsize=(6, 4))
7 plt.boxplot(movie_data['Rating'])
8 plt.title('Boxplot of Ratings')
9 plt.ylabel('Rating')
10 plt.show()
```

```
count    100000.000000
mean         3.529860
std         1.125674
min         1.000000
25%         3.000000
50%         4.000000
75%         4.000000
max         5.000000
Name: Rating, dtype: float64
```



Interpretation:

- The ratings are skewed to the right, meaning that most of the ratings are higher than the median.
- This indicates that the majority of the reviewers were satisfied with the product or service being rated.

- The outlier at a rating of 1 is an extreme value that deviates significantly from the rest of the ratings.

4). Check for duplicate entries in the dataset.Remove duplicates if necessary.

In [36]:

```
1 duplicate_rows = movie_data.duplicated()
2 print("Number of duplicate rows:", duplicate_rows.sum())
3
4 duplicate_data = movie_data[duplicate_rows]
5 print("Duplicate rows:")
6 print(duplicate_data)
```

Number of duplicate rows: 0

Duplicate rows:

Empty DataFrame

Columns: [User_id, Movie_id, Rating, timestamp, day_of_week, hour_of_day, time_of_day_category, rental_day_of_week, rental_time_of_day]

Index: []

Interpretation:

- There are no duplicate values in the data so there is nothing to remove.
- Duplicate values can increase the size and complexity of the data, making it harder to store, process, and analyze.
- Duplicate values can skew the distribution and statistics of the data, leading to inaccurate and misleading results.

Part - 2 : Data Visualization

1). Create a histogram or boxplot to visualize the distribution of movie ratings.Are there any interesting trends?

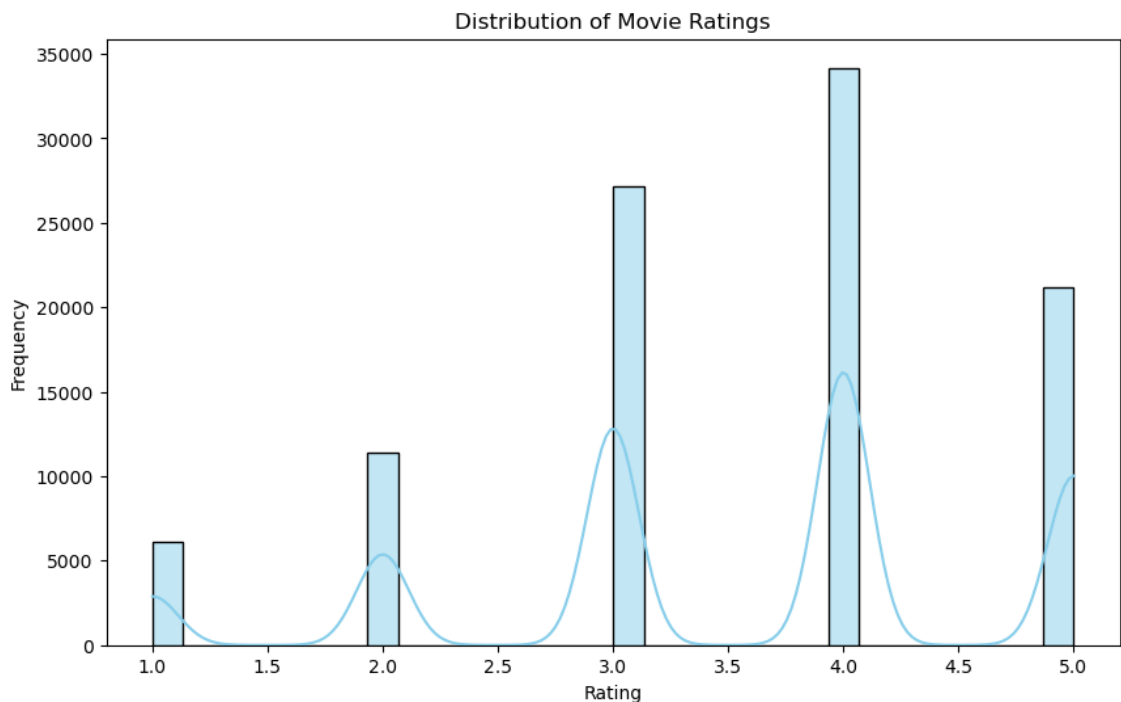
```
In [38]: 1 plt.figure(figsize=(10, 6))
2 sns.histplot(movie_data['Rating'], bins=30, kde=True, color='skyblue')
3 plt.title('Distribution of Movie Ratings')
4 plt.xlabel('Rating')
5 plt.ylabel('Frequency')
6 plt.show()
7
8 plt.figure(figsize=(8, 6))
9 sns.boxplot(x=movie_data['Rating'], color='lightcoral')
10 plt.title('Box Plot of Movie Ratings')
11 plt.xlabel('Rating')
12 plt.show()
```

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):

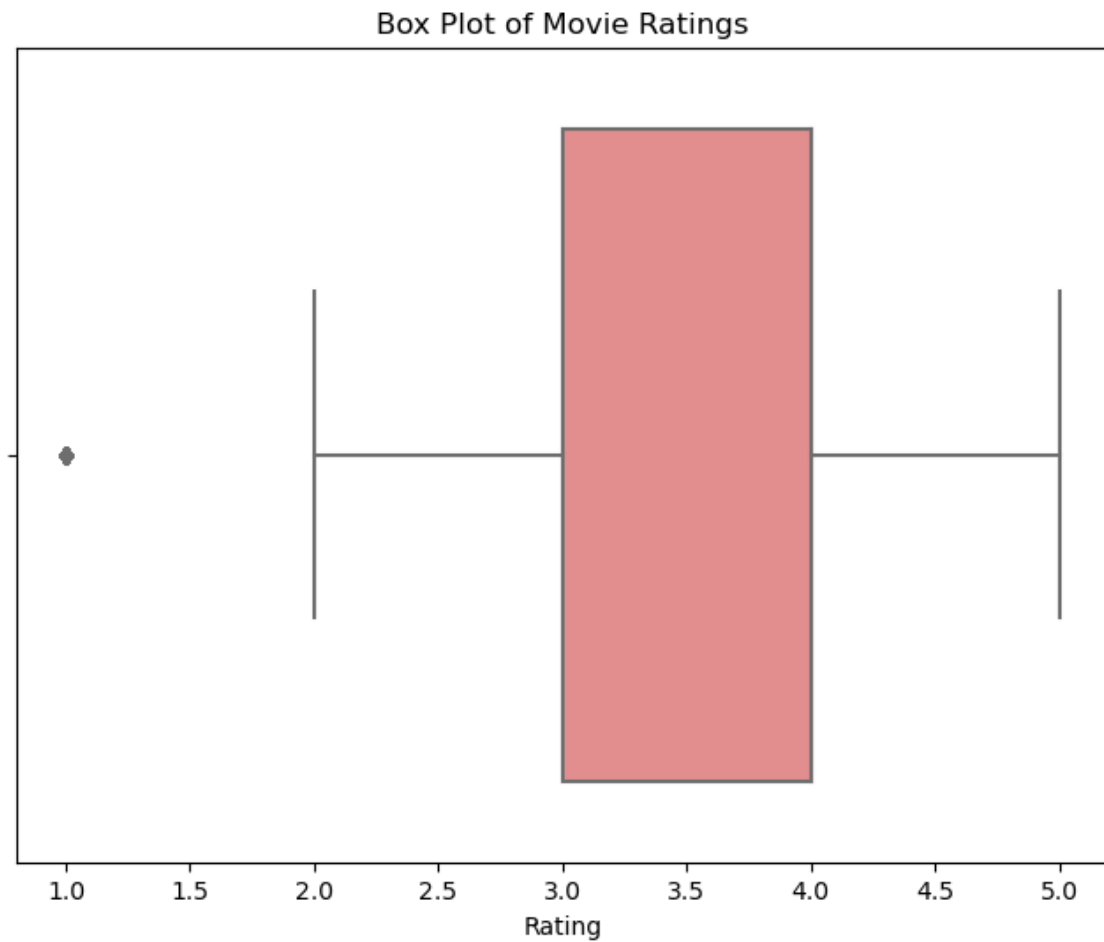
C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1119: FutureWarning: use_inf_as_na option is deprecated and will be removed in a future version. Convert inf values to NaN before operating instead.

with pd.option_context('mode.use_inf_as_na', True):



C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):

**Interpretation:**

- In the histogram, you can observe the distribution of ratings, whether they are concentrated in a specific range or spread out across various values.
- The boxplot provides information about the median, quartiles, and potential outliers in the rating distribution.

3). Use a bar chart or heatmap to represent the most popular movies (top 10 or 20) among users.

```
In [48]: 1 ratings_count = movie_data.groupby('Movie_id')['Rating'].count()
2 top_movies = ratings_count.sort_values(ascending=False).head(10)
3 plt.figure(figsize=(12, 8))
4 sns.barplot(x=top_movies.values, y=top_movies.index, palette='viridis')
5 plt.title('Top 10 Most Popular Movies')
6 plt.xlabel('Number of Ratings')
7 plt.ylabel('Movie Title')
8 plt.show()
```

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

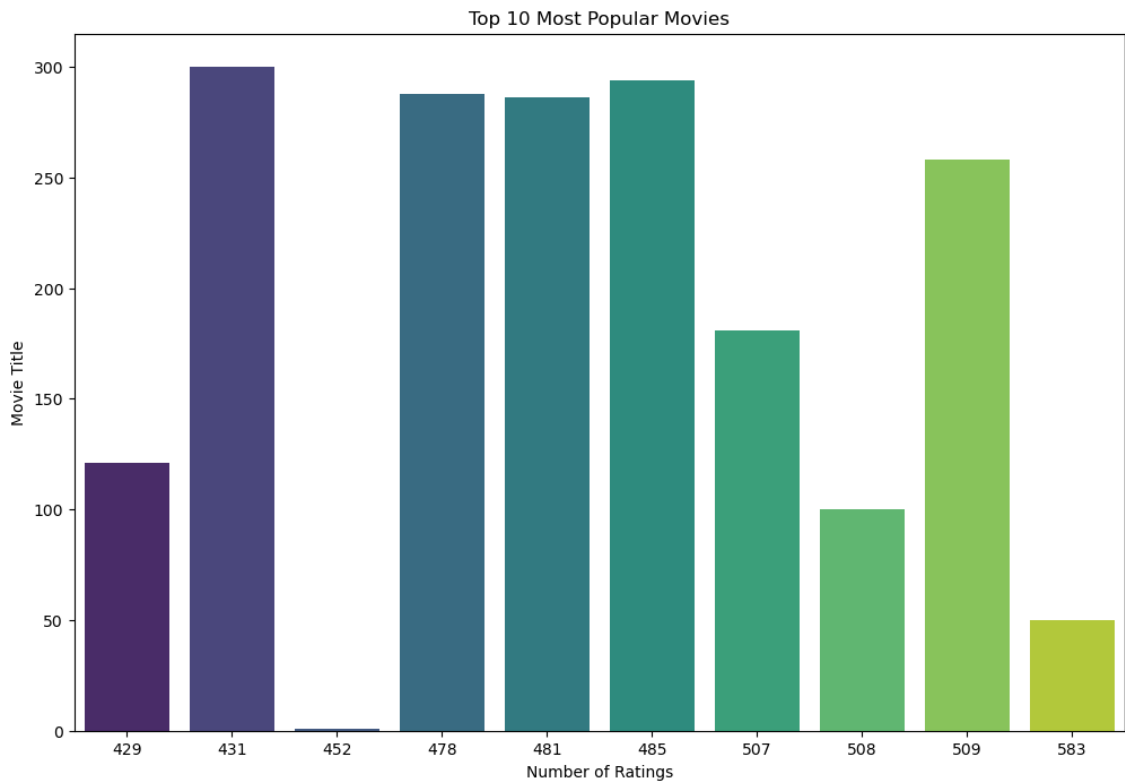
if pd.api.types.is_categorical_dtype(vector):

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):



Interpretation:

- The resulting bar chart provides a visual representation of the popularity of movies based on the number of user ratings.
- Movies with higher bars have received more ratings, suggesting a higher level of user engagement.

4). Visualize the distribution of rentals by day of the week or time of day. Are there any preferred rental times?

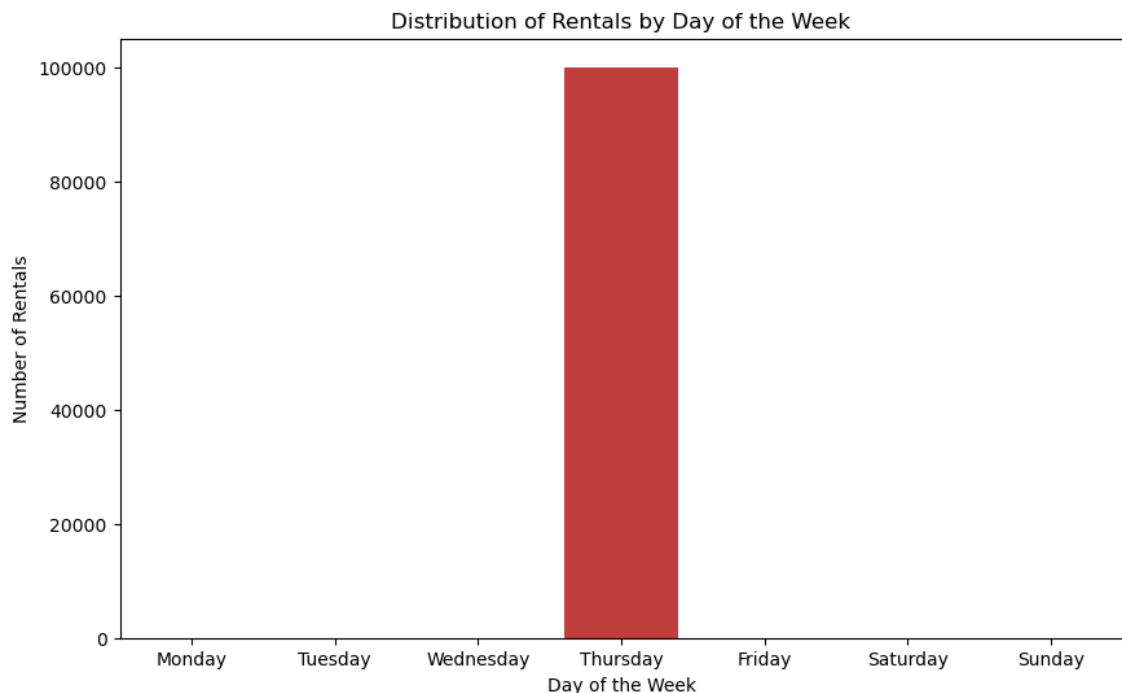
```
In [52]: 1 movie_data['timestamp'] = pd.to_datetime(movie_data['timestamp'])
2 movie_data['day_of_week'] = movie_data['timestamp'].dt.day_name()
3 plt.figure(figsize=(10, 6))
4 sns.countplot(x='day_of_week', data=movie_data, order=['Monday', 'Tues
5 plt.title('Distribution of Rentals by Day of the Week')
6 plt.xlabel('Day of the Week')
7 plt.ylabel('Number of Rentals')
8 plt.show()
```

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):

C:\Users\pf3zw\anaconda3\Lib\site-packages\seaborn_oldcore.py:1498: FutureWarning: is_categorical_dtype is deprecated and will be removed in a future version. Use isinstance(dtype, CategoricalDtype) instead

if pd.api.types.is_categorical_dtype(vector):



Interpretation:

- The number of rentals is extremely high on Thursday, compared to the other days of the week.
- This could indicate that Thursday is a popular day for renting products or services, or that there is a special promotion or event on that day.

Part 3: Data Summarization

1). Calculate basic descriptive statistics for the "rating" column (mean, median, standard deviation).


```
In [54]: 1 rating_statistics = movie_data['Rating'].describe()
2 print(rating_statistics)
```

```
count    100000.000000
mean         3.529860
std         1.125674
min         1.000000
25%         3.000000
50%         4.000000
75%         4.000000
max         5.000000
Name: Rating, dtype: float64
```

Interpretation:

- The average rating is about 3.53, which means that most ratings are above the midpoint of 3.
- The standard deviation is about 1.13, which means that there is some variability in the ratings, but not too much. Most ratings are within one standard deviation of the mean, or between 2.4 and 4.6.
- The median (or 50th percentile) is 4, which means that half of the ratings are above 4 and half are below 4.

2). Summarize the number of unique users and movies in the dataset.

```
In [56]: 1 unique_users = movie_data['User_id'].nunique()
2 unique_movies = movie_data['Movie_id'].nunique()
3
4 print("Number of Unique Users:", unique_users)
5 print("Number of Unique Movies:", unique_movies)
```

```
Number of Unique Users: 943
Number of Unique Movies: 1682
```

Interpretation:

- Number of Unique Users: This represents the count of distinct user IDs in the dataset which is 943
- Number of Unique Movies: This represents the count of distinct movie IDs in the dataset which is 1682

3). Find the most active users (top 10 or 20) based on the number of movie rentals.

```
In [58]: 1 user_rental_counts = movie_data.groupby('User_id')['Movie_id'].count()
2 top_users = user_rental_counts.sort_values(ascending=False).head(10)
3 print("Top 10 Most Active Users:")
4 print(top_users)
```

Top 10 Most Active Users:

User_id

405 737

655 685

13 636

450 540

276 518

416 493

537 490

303 484

234 480

393 448

Name: Movie_id, dtype: int64

Interpretation:

- The resulting top_users Series contains user IDs as the index and the corresponding counts of movie rentals.
- This information helps identify the users who have rented the most movies, indicating their level of activity or engagement with the rental service.