

Question 1

```
In [1]: #pip install pyppeteer

In [2]: #import libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [3]: df=pd.read_csv('dataset/imdb_top_1000.csv')

In [4]: df.head()
```

```
Out[4]:
```

	Poster_Link	Series_Title	Released_Year	Certificate	Runtime	Genre	IMDB_Rating
0	https://media-amazon.com/images/M/MV5BMDFkYT...	The Shawshank Redemption	1994	A	142 min	Drama	9.3
1	https://media-amazon.com/images/M/MV5BMjE2MjY...	The Godfather	1972	A	175 min	Crime, Drama	9.2
2	https://media-amazon.com/images/M/MV5BMjE2MjY...	The Dark Knight	2008	UA	152 min	Action, Crime, Drama	9.1
3	https://media-amazon.com/images/M/MV5BMWwMjY...	The Godfather: Part II	1974	A	202 min	Crime, Drama	9.1
4	https://media-amazon.com/images/M/MV5BMWU4Mz...	12 Angry Men	1957	U	96 min	Crime, Drama	9.1

```
In [5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 16 columns):
 # Column                Non-Null Count  Dtype
---  --
 0 Poster_Link            1000 non-null   object
 1 Series_Title            1000 non-null   object
 2 Released_Year           1000 non-null   object
 3 Certificate              899 non-null    object
 4 Runtime                 1000 non-null   object
 5 Genre                   1000 non-null   object
 6 IMDB_Rating              1000 non-null   float64
 7 Overview                1000 non-null   object
 8 Meta_score               843 non-null    float64
 9 Director                1000 non-null   object
10 Star1                   1000 non-null   object
11 Star2                   1000 non-null   object
12 Star3                   1000 non-null   object
13 Star4                   1000 non-null   object
14 No_of_Votes             1000 non-null   int64
15 Gross                   831 non-null    object
dtypes: float64(2), int64(1), object(13)
memory usage: 125.1+ KB
```

```
In [6]: #lets see what is the percentage of null values in total dataset
df.isnull().sum()/len(df)
```

```
Out[6]:
Poster_Link      0.000
Series_Title      0.000
Released_Year     0.000
Certificate        0.101
Runtime           0.000
Genre             0.000
IMDB_Rating       0.000
Overview          0.000
Meta_score        0.157
Director          0.000
Star1             0.000
Star2             0.000
Star3             0.000
Star4             0.000
No_of_Votes       0.000
Gross             0.169
dtype: float64
```

As the null values are max 17% in total I will replace with other values instead of dropping them

```
In [7]: #fill the null values with the meta_score mean
df['Meta_score']=df['Meta_score'].fillna(df['Meta_score'].mean())

In [8]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 16 columns):
 # Column                Non-Null Count  Dtype
---  --
 0 Poster_Link            1000 non-null   object
 1 Series_Title            1000 non-null   object
 2 Released_Year           1000 non-null   object
 3 Certificate              899 non-null    object
 4 Runtime                 1000 non-null   object
 5 Genre                   1000 non-null   object
 6 IMDB_Rating              1000 non-null   float64
 7 Overview                1000 non-null   object
 8 Meta_score               1000 non-null   float64
 9 Director                1000 non-null   object
10 Star1                   1000 non-null   object
11 Star2                   1000 non-null   object
12 Star3                   1000 non-null   object
13 Star4                   1000 non-null   object
14 No_of_Votes             1000 non-null   int64
15 Gross                   831 non-null    object
dtypes: float64(2), int64(1), object(13)
memory usage: 125.1+ KB
```

```
In [9]: #change Gross from object to numeric numbers
df['Gross']=df['Gross'].str.replace(',','')
df['Gross']=df['Gross'].apply(pd.to_numeric)
```

```
In [10]: #fill the missing values with the mean of the Gross column
df['Gross'].fillna(df['Gross'].mean(),inplace=True)
```

```
In [11]: df['Gross'].head()
```

```
Out[11]:
0    28341469.0
1    134966411.0
2    534858444.0
3    57300000.0
4    4360000.0
Name: Gross, dtype: float64
```

```
In [12]: df.describe()
```

```
Out[12]:
```

	IMDB_Rating	Meta_score	No. of Votes	Gross
count	1000.000000	1000.000000	1.000000e+03	1.000000e+03
mean	7.949300	77.97153	2.73692e+05	6.80347e+08
std	0.725491	11.36206	3.27372e+05	1.00037e+08
min	7.600000	72.00000	5.55262e+04	1.30500e+03
25%	7.700000	78.00000	2.00000e+05	5.01291e+06
50%	7.900000	77.97153	1.38548e+05	4.23894e+07
75%	8.100000	85.25000	3.74161e+05	6.80347e+08
max	9.300000	100.00000	2.34311e+06	9.36662e+07

```
In [13]: df['Certificate'].is_unique
```

```
Out[13]: False
```

```
In [14]: df['Certificate'].nunique()
```

```
Out[14]: 16
```

```
In [15]: #Drop unnecessary columns
df.drop(columns=['Poster_Link','Series_Title','Certificate','Overview'],inplace=True)
```

```
In [16]: #Remaining columns
df.columns
```

```
Out[16]: Index(['Released_Year', 'Runtime', 'Genre', 'IMDB_Rating', 'Meta_score',
        'Director', 'Star1', 'Star2', 'Star3', 'Star4', 'No_of_Votes', 'Gross'],
      dtype='object')
```

```
In [17]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 12 columns):
 # Column                Non-Null Count  Dtype
---  --
 0 Released_Year           999 non-null   int32
 1 Runtime                 1000 non-null   object
 2 Genre                   999 non-null   object
 3 IMDB_Rating              999 non-null   float64
 4 Meta_score               1000 non-null   float64
 5 Director                1000 non-null   object
 6 Star1                   1000 non-null   object
 7 Star2                   1000 non-null   object
 8 Star3                   1000 non-null   object
 9 Star4                   1000 non-null   object
10 No_of_Votes             999 non-null   int64
11 Gross                   1000 non-null   float64
dtypes: float64(3), int32(1), int64(1), object(8)
memory usage: 93.9+ KB
```

```
In [18]: #Drop PG value from Released_Year
df=df.drop(df[df['Released_Year']== 'PG'].index)
```

```
In [19]: #Change Released_Year datatype to integer
df['Released_Year']=df['Released_Year'].astype('int')
```

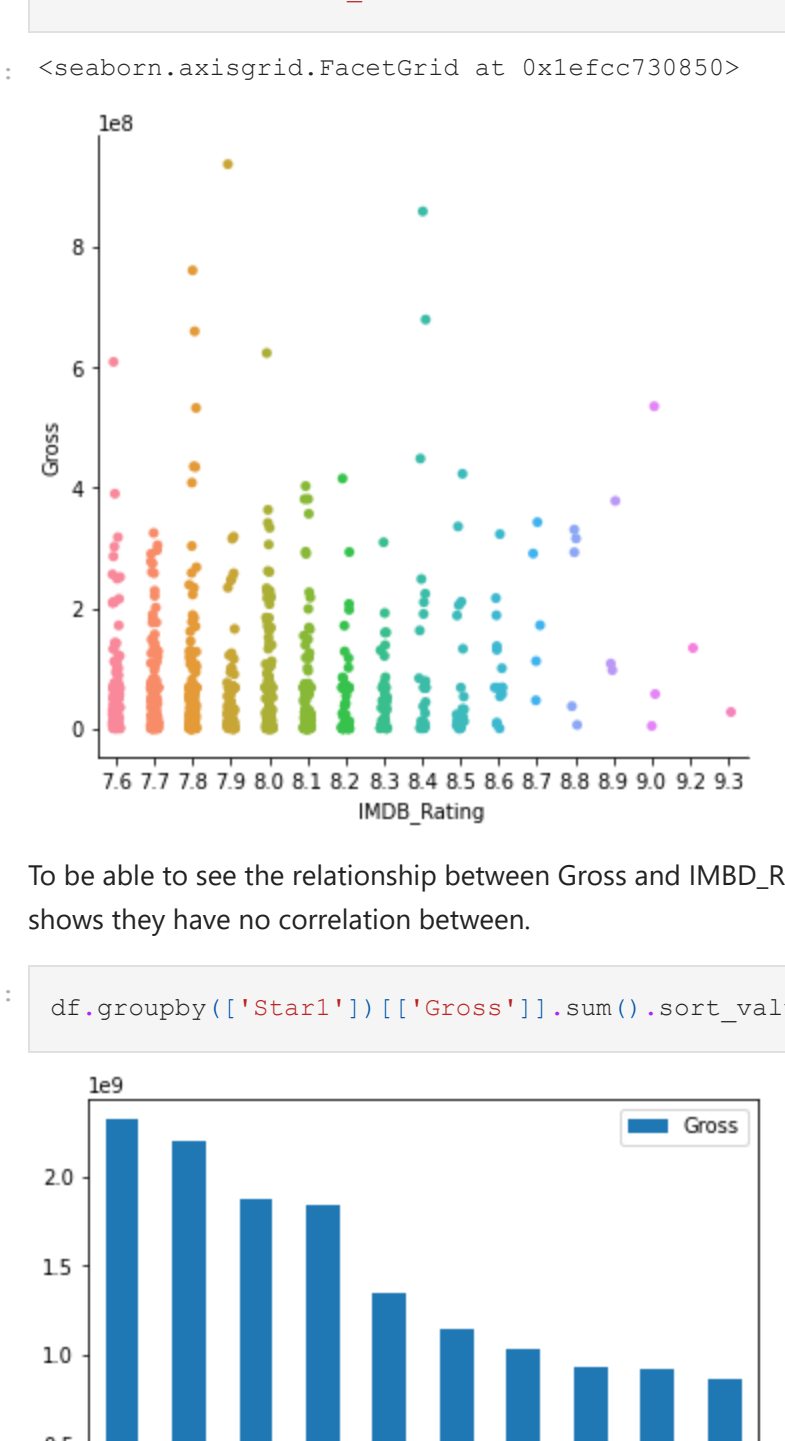
```
In [20]: #Remove commas and change Runtime feature datatype to integer
df['Runtime']=df['Runtime'].str.replace(',','')
df['Runtime']=df['Runtime'].astype('int')
```

```
In [21]: #After cleaning dataset there are no null values left
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 999 entries, 0 to 999
Data columns (total 12 columns):
 # Column                Non-Null Count  Dtype
---  --
 0 Released_Year           999 non-null   int32
 1 Runtime                 999 non-null   int32
 2 Genre                   999 non-null   object
 3 IMDB_Rating              999 non-null   float64
 4 Meta_score               999 non-null   float64
 5 Director                999 non-null   object
 6 Star1                   999 non-null   object
 7 Star2                   999 non-null   object
 8 Star3                   999 non-null   object
 9 Star4                   999 non-null   object
10 No_of_Votes             999 non-null   int64
11 Gross                   999 non-null   float64
dtypes: float64(3), int32(2), int64(1), object(6)
memory usage: 93.7+ KB
```

```
In [22]: #some visualization techniques- distribution of Runtime
sns.distplot(df['Runtime'])
```

```
Out[22]: DisEmNullah\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is deprecated and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:xlabel='Runtime', ylabel='Density'>
```



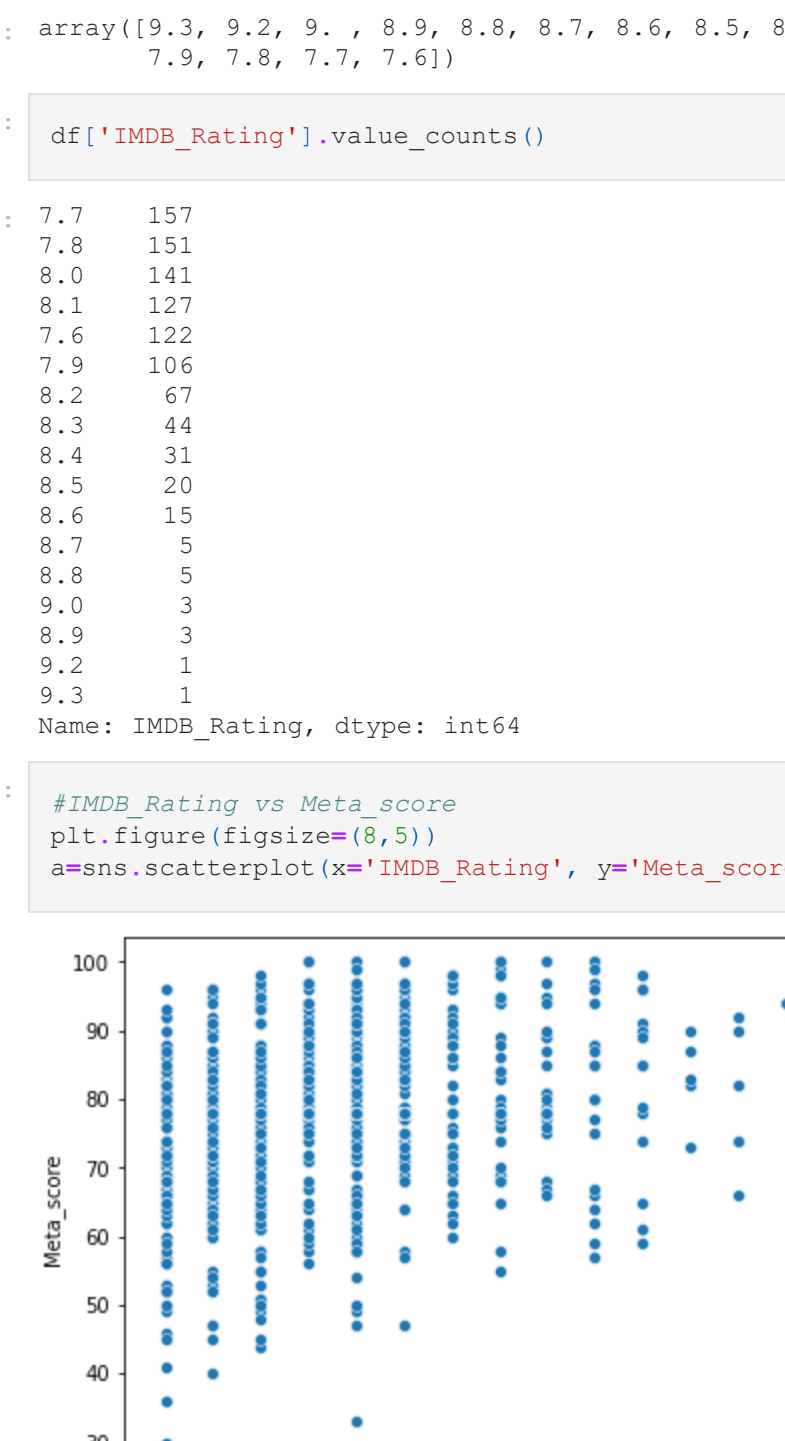
To be able to see the distribution of Runtime of movies distplot is used. The Runtime has Gaussian distribution and the films between 100min and 150m have the highest Density.

```
In [23]: # to see when earlier and latest movies are produced
df['Released_Year'].min(),df['Released_Year'].max()
```

```
Out[23]: (1920, 2020)
```

```
In [24]: sns.distplot(df['Released_Year'])
```

```
Out[24]: DisEmNullah\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is deprecated and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:xlabel='Released_Year', ylabel='Density'>
```



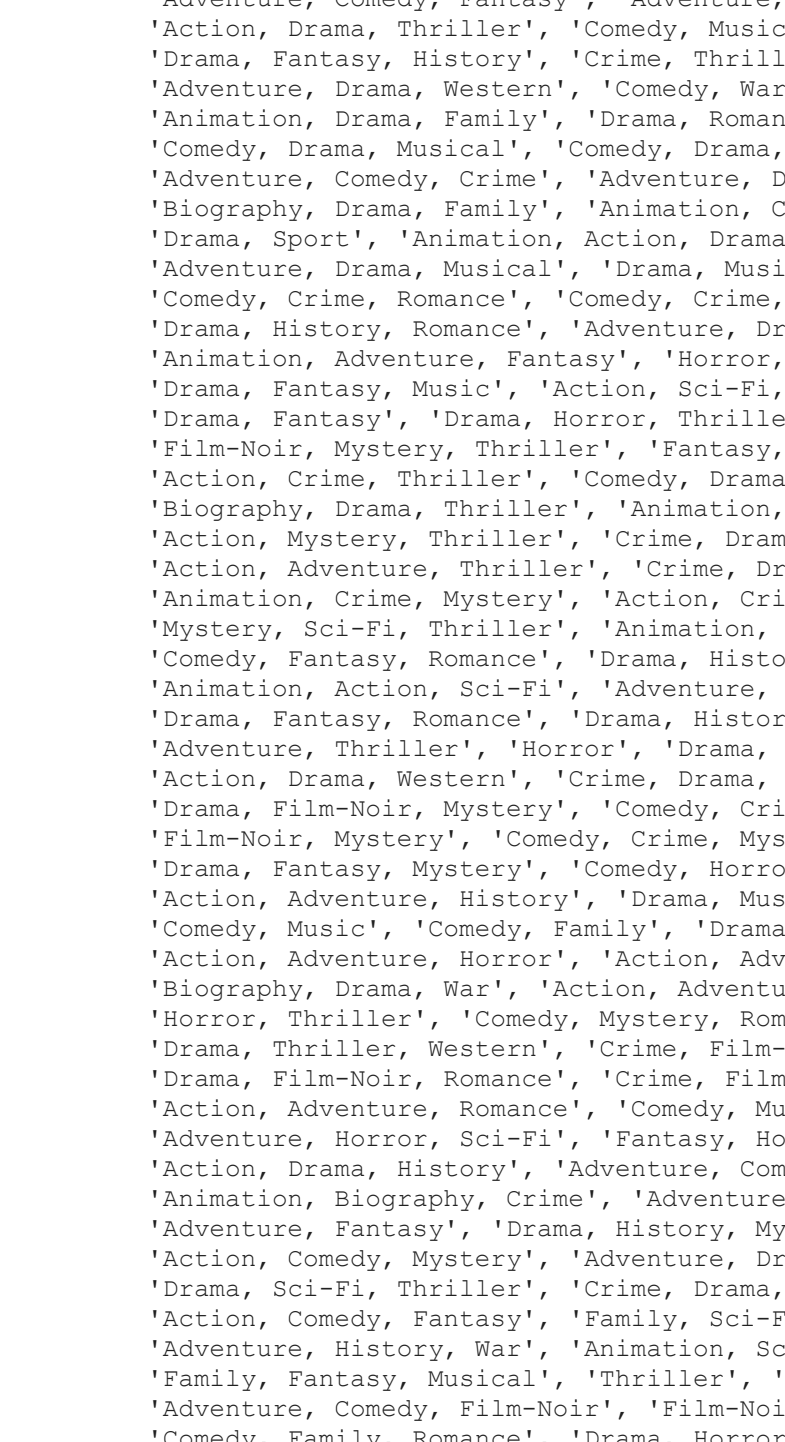
Released_Year of movies has left skewed distribution and also the number of movies are increasing gradually

```
In [25]: df['Gross'].head()
```

```
Out[25]:
0    28341469.0
1    134966411.0
2    534858444.0
3    57300000.0
4    4360000.0
Name: Gross, dtype: float64
```

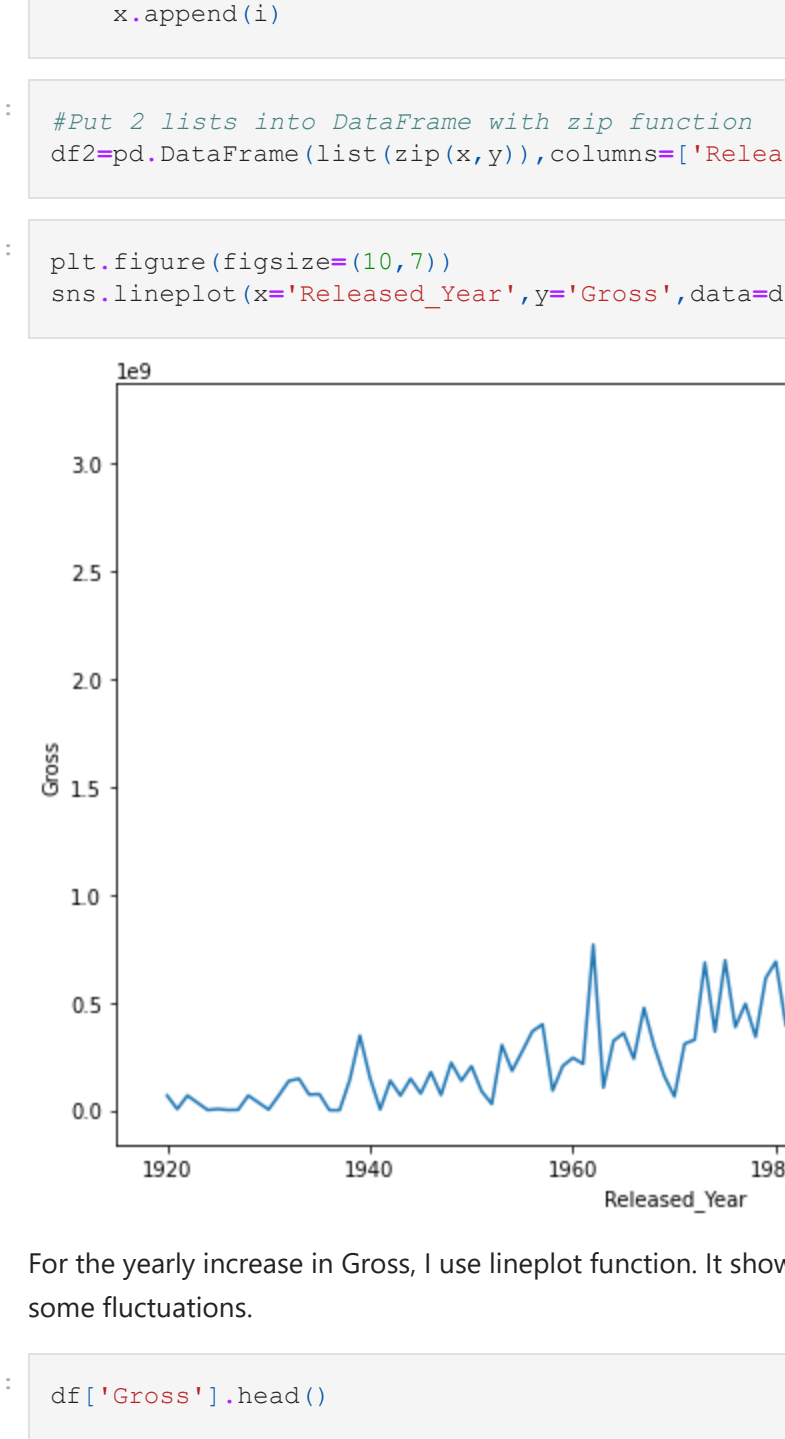
```
In [26]: sns.distplot(df['Gross'])
```

```
Out[26]: DisEmNullah\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is deprecated and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)
<AxesSubplot:xlabel='Gross', ylabel='Density'>
```



```
In [27]: fig, (ax1, ax2) = plt.subplots(2,1)
df['Gross'].hist(ax=ax1, bins=50)
ax1.tick_params(labelsize=14)
ax1.set_xlabel('Gross', fontsize=14)
ax1.set_ylabel('Occurrence', fontsize=14)
np.log(df['Gross']).hist(ax=ax2, bins=50)
ax2.tick_params(labelsize=14)
ax2.set_xlabel('log10(Gross)', fontsize=14)
ax2.set_ylabel('Occurrence', fontsize=14)
```


```
Out[27]: Text(0, 0.5, 'Occurrence')
```



To be able to see the Gross value distribution for movies the log transformation is applied and the highest revenue is shown to be between 17K and 20K.

```
In [28]: sns.catplot(x='IMDB_Rating',y='Gross',data=df)
```

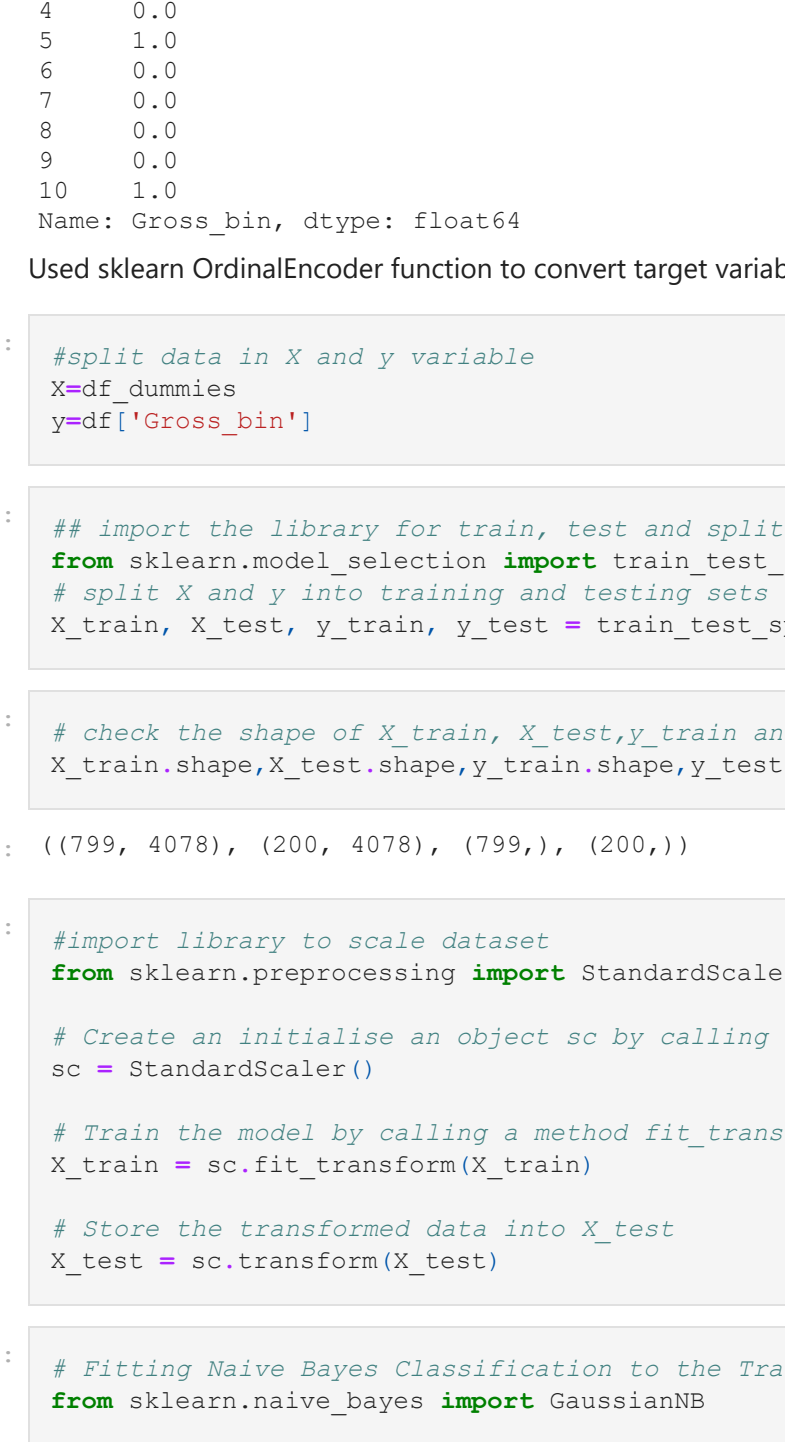
```
Out[28]: <seaborn.axisgrid.FacetGrid at 0x1efcc730850>
```



To be able to see the relationship between Gross and IMDB_Rating features the catplot is applied, the plot shows they have no correlation between.

```
In [29]: df.groupby(['Star1'])[['Gross']].sum().sort_values(by='Gross',ascending=False)[:10].plot()
```

```
Out[29]:
```



```
In [30]: df.groupby(['Director'])[['Gross']].sum().sort_values(by='Gross',ascending=False)[:10].plot()
```

```
Out[30]:
```


To be able to see which top 10 Star1 and directors brought most revenue I use groupby function to create datasets and then apply bar plot. The highest revenue for Star1 is Tom Hanks and for director is Steven Spielberg.

```
In [31]: df['IMDB_Rating'].unique()
```

```
Out[31]: array([9.3, 9.2, 9. , 8.9, 8.8, 8.7, 8.6, 8.5, 8.4, 8.3, 8.2, 8.1, 8. ,
       7.9, 7.8, 7.7, 7.6])
```

```
In [32]: df['IMDB_Rating'].value_counts()
```

```
Out[32]:
7.7    157
7.8    151
8.0    141
8.1    127
7.6    122
7.9    106
8.2     67
8.3     44
8.4     31
8.5     20
8.6     15
8.7      5
9.2      3
9.0      3
9.3      1
Name: IMDB_Rating, dtype: int64
```

```
In [33]: #IMDB_Rating vs Meta_score
plt.figure(figsize=(8,5))
sns.scatterplot(x='IMDB_Rating',y='Meta_score',data=df, c='None')
```

```
Out[33]:
```


Also there is no relationship between IMDB_Rating and Meta_score

```
In [34]: #See Genre column unique values
df['Genre'].unique()
```

```
Out[34]: array(['Drama', 'Crime, Drama', 'Action, Crime, Drama', 'Action, Crime, Drama, Biography, Drama',
        'Action, Adventure, Drama', 'Biography, Drama, History', 'Action, Adventure, Fantasy', 'Comedy, Drama, Thriller',
        'Adventure, Drama, Sci-Fi', 'Animation, Adventure, Family', 'Drama, War', 'Crime, Drama, Fantasy', 'Crime, Drama, Romance',
        'Crime, Drama, Mystery', 'Crime, Drama, Thriller', 'Crime, Drama, Mystery', 'Drama, Music', 'Action, Drama, Mystery', 'Drama, History, Mystery',
        'Biography, Drama, Music', 'Crime, Mystery, Thriller', 'Biography, Drama, Mystery', 'Crime, Mystery, Thriller', 'Animation, Adventure, Drama',
        'Adventure, Comedy, Sci-Fi', 'Horror, Mystery, Thriller', 'Drama, Romance, War', 'Comedy, Drama, Family', 'Animation, Drama, Fantasy',
        'Action, Biography, Drama', 'Animation, Adventure, Fantasy', 'Horror, Mystery, Thriller', 'Drama, Fantasy, Thriller', 'Comedy, Drama, Fantasy',
        'Action, Adventure', 'Comedy, Drama', 'Drama, Fantasy', 'Drama, Horror', 'Drama, Thriller, War', 'Drama, Fantasy, Horror',
        'Crime, Drama, Musical', 'Adventure, Thriller', 'Crime, Mystery', 'Drama, Music', 'Comedy', 'Drama, Film-Noir', 'Comedy, Drama, War',
        'Drama, Thriller, War', 'Drama, Fantasy, Horror', 'Crime, Drama, Musical', 'Adventure, Thriller', 'Comedy, Crime', 'Drama, Family, Sport',
        'Animation, Adventure, Comedy', 'Adventure, Drama, Thriller', 'Comedy, Crime, Drama', 'Crime, Drama, Sci-Fi', 'Adventure, Sci-Fi',
        'Adventure, Biography, Drama', 'Adventure, Mystery, Thriller', 'Comedy, Musical, Romance', 'Crime, Drama, Film-Noir', 'Drama, History, Mystery',
        'Action, Drama, War', 'Action, Drama', 'Adventure, Comedy, Drama', 'Biography, Drama, Sport', 'Action, Comedy, Crime',
        'Action, Biography, Crime', 'Drama, Mystery, Thriller', 'Action, Adventure, Fantasy', 'Horror, Mystery, Thriller', 'Drama, Fantasy, Music',
        'Action, Thriller', 'Action, Adventure, Comedy', 'Adventure, Comedy, Fantasy', 'Adventure, Drama, History', 'Action, Drama, Thriller',
        'Comedy, Musical, Romance', 'Drama, Fantasy, Romance', 'Drama, History, Mystery', 'Drama, Film-Noir, Thriller', 'Horror, Sci-Fi',
        'Adventure, Horror, Sci-Fi', 'Fantasy, Horror', 'Action, Crime, Thriller', 'Comedy, Drama, Music', 'Adventure, Comedy, Drama',
        'Biography, Drama, Sport', 'Action, Comedy, Crime', 'Action, Mystery, Thriller', 'Crime, Drama, Musical', 'Comedy, Western',
        'Adventure, Horror, Sci-Fi', 'Fantasy, Horror', 'Action, Drama, History', 'Adventure, Comedy, Family',
        'Animation, Biography, Crime', 'Adventure, Biography, Crime', 'Drama, Comedy, Romance', 'Action, Comedy, Fantasy', 'Drama, Fantasy, Music',
        'Action, Thriller', 'Action, Adventure, Comedy', 'Adventure, Comedy, Fantasy', 'Adventure, Drama, History', 'Action, Drama, Thriller',
        'Comedy, Musical, Romance', 'Drama, Fantasy, Romance', 'Drama, History, War', 'Adventure, Thriller', 'Horror, Sci-Fi',
        'Adventure, Horror, Sci-Fi', 'Fantasy, Horror', 'Action, Drama, Western', 'Crime, Drama, Horror', 'Drama, Film-Noir, Mystery',
        'Comedy, Crime, Thriller', 'Film-Noir, Mystery, Thriller', 'Comedy, Crime, Thriller', 'Animation, Comedy, Fantasy',
        'Adventure, Horror, Sci-Fi', 'Comedy, Western', 'Drama, Fantasy, Mystery', 'Comedy, Horror', 'Action, Adventure, History',
        'Drama, Music, Mystery', 'Comedy, Music', 'Comedy, Family', 'Comedy, Drama, Music', 'Action, Adventure, Horror',
        'Action, Adventure, Biography', 'Biography, Drama, War', 'Action, Adventure, Western', 'Horror, Thriller', 'Comedy, Mystery, Romance',
        'Drama, Thriller, Western', 'Crime, Film-Noir, Thriller', 'Drama, Film-Noir, Romance', 'Crime, Film-Noir, Mystery',
        'Action, Adventure, Romance', 'Comedy, Musical, Musical', 'Adventure, Horror, Sci-Fi', 'Fantasy, Horror', 'Action, Drama, History',
        'Adventure, Comedy, Family', 'Animation, Biography, Crime', 'Adventure, Biography, Crime', 'Drama, Comedy, Romance',
        'Action, Comedy, Thriller', 'Crime, Drama, Fantasy', 'Family, Fantasy, Musical', 'Adventure, Comedy, Film-Noir',
        'Film-Noir, Thriller', 'Comedy, Family, Romance', 'Comedy, Musical, War', 'Biography, Drama, Romance', 'Drama, History, Music',
        'Animation, Action, Fantasy', 'Animation, Comedy, Fantasy', 'Animation, Adventure, War', 'Drama, Horror, Mystery',
        'Animation, Comedy, Crime', 'Adventure, Adventure, Crime', 'Action, Adventure, Mystery', 'Action, Adventure, Family',
        'Action, Crime, Mystery', 'Animation, Drama, Romance', 'Drama, War, Western', 'Adventure, Comedy, War'],
      dtype=object)
```

```
In [35]: #Disitribution of top 10 Genres in barplot
df['Genre'].value_counts()[:10].plot.bar()
```

```
Out[35]:
```



```
In [36]: #Find unique values for Released_Year with Gross and append to the lists
y=[]
x=[]
for i in df['Released_Year'].unique():
    y.append(i)
    x.append(df[df['Released_Year']==i]['Gross'].sum())
```

```
In [37]: #Put x,y lists into DataFrame with zip function
df2=pd.DataFrame(list(zip(x,y)),columns=['Released_Year','Gross']).sort_values('Gross')
```

```
In [38]: #Plot figure (figsize=(10,7))
sns.lineplot(x='Released_Year',y='Gross',data=df2)
```

```
Out[38]:
```


For the yearly increase in Gross, I use lineplot function. It shows Revenue is increasing gradually despite some fluctuations.

```
In [39]: df['Gross'].head()
```

```
Out[39]:
0    28341469.0
1    134966411.0
2    534858444.0
3    57300000.0
4    4360000.0
Name: Gross, dtype: float64
```

```
In [40]: min_value = df['Gross'].min()
max_value = df['Gross'].max()
print(min_value)
print(max_value)
```

```
Out[40]:
1305.0
936662225.0
```

```
In [41]: import numpy as np
bins = np.linspace(min_value,max_value,4)
```

```
Out[41]: array([1.13050000e+03, 3.12221612e+08, 6.24441918e+08, 9.36662225e+08])
```

```
In [42]: labels = ['Average', 'High', 'Very High']
```

```
In [43]: df['Gross_bin'] = pd.cut(df['Gross'], bins=bins, labels=labels, include_lowest=True)
```

```
In [44]: plt.hist(df['Gross_bin'], bins=3)
plt.show()
```

```
Out[44]:
```


I use Gross column to convert continuous data to ordinary categorical features for classification model I use Min() and Max() of the data and I use 'pandas.cut' function. Then after that I plot the new categories and it shows that most of the movies' Gross fall into Average category.

```
In [45]: df.columns
```

```
Out[45]: Index(['Released_Year', 'Runtime', 'Genre', 'IMDB_Rating', 'Meta_score',
        'Director', 'Star1', 'Star2', 'Star3', 'Star4', 'No_of_Votes', 'Gross',
        'Gross_bin'],
      dtype='object')
```

```
In [46]: obj=['Genre','Director','Star1','Star2','Star3','Star4']
```

```
In [47]: cols_to_std = ['Runtime','Released_Year','IMDB_Rating','Meta_score','No_of_Votes']
```

```
In [48]: dummies=pd.get_dummies(df[obj],drop_first=True)
df_dummies = pd.concat([df[cols_to_std],dummies],axis=1)
```

```
Out[48]:
```

	Runtime	Released_Year	IMDB_Rating	Meta_score	No. of Votes	Genre_Action_Biography	Genre_Action_Adventure	Genre_Action_Crime
0	142	1994	9.3	80.0	2343110	0	0	0
1	175	1972	9.2	108.0	1620367	0	0	0
2	152	2008	9.0	84.0	2303232	0	0	0
3	202	1974	9.0	90.0	1129952	0	0	0
4	96	1957	9.0	96.0	689845	0	0	0

5 rows x 4078 columns

```
In [50]: df_dummies.columns
```

```
Out[50]: Index(['Runtime', 'Released_Year', 'IMDB_Rating', 'Meta_score', 'No. of Votes',
        'Genre_Action_Biography', 'Genre_Action_Adventure', 'Genre_Action_Crime',
        'Genre_Action_Drama', 'Genre_Action_Adventure_Drama',
        'Genre_Action_Adventure_Family',
        'Star4_Yōnosuke Itō', 'Star4_Zac Mattoon O'Brien', 'Star4_Zach Grenier',
        'Star4_Zarah Jane McKenzie', 'Star4_Zeppo Marx', 'Star4_Ziyi Zhang',
        'Star4_Zoe Saldana', 'Star4_Zoe Kravitz', 'Star4_Álvaro Guerrero',
        'Star4_Emile Vallée'],
      dtype='object', length=4078)
```

I create obj list with non numeric columns and use Dummies function to convert them to numeric values and then concatenate with cols_to_std list which already have numeric values and after that I have 4078 columns.

```
In [51]: from sklearn.preprocessing import OrdinalEncoder
ord_enc = OrdinalEncoder()
df['Gross_bin'] = ord_enc.fit_transform(df['Gross_bin'])
df['Gross_bin'].head(11)
```

```
Out[51]:
0    0.0
1    1.0
2    1.0
3    0.0
4    0.0
5    1.0
6    0.0
7    0.0
8    0.0
9    1.0
10   1.0
Name: Gross_bin, dtype: float64
```

Used sklearn OrdinalEncoder function to convert target variable into categorical data

```
In [52]: #to encode X and y variable
x=df[['Gross_bin']]
y=df['Gross_bin']
```

```
In [53]: #I import the library for train, test and splitting of teh data set
from sklearn.model_selection import train_test_split
# split X and y into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X,y,test_size=0.20)
```

```
In [54]: # check the shape of X_train, X_test, y_train and y_test
X_train.shape, X_train.head(), X_train.shape, y_train.shape, y_test.shape
```

```
Out[54]: ((799, 4078), (200, 4078), (799, 1), (200, 1))
```

```
In [55]: #import library to scale dataset
from sklearn.preprocessing import StandardScaler
# Create an initialise an object sc by calling a method StandardScaler()
sc = StandardScaler()
# Train the model by calling a method fit_transform
X_train = sc.fit_transform(X_train)
# Scale the transformed data into X_test
X_test = sc.transform(X_test)
```

```
In [56]: # Fitting Naive Bayes Classification to the Training set with linear kernel
from sklearn.naive_bayes import GaussianNB
# Create an initialise an object 'nvcclassifier' by calling a method 'GaussianNB()'
nvcclassifier = GaussianNB()
# Call fit() method for training the dataset
nvcclassifier.fit(X_train, y_train)
```

```
Out[56]: GaussianNB()
```


[illegible]