

im2_db

```
In [83]: 1 import os
          2 # set working directory
          3 working_directory = r"C:/Users/PC/Documents/moringa/Project Phase
          4 os.chdir(working_directory)
```

Objectives

1. what are the popular genres in the Movie Industry?
2. What is the relationship between average rating and number of votes?
3. What is the trend of the number of votes over the years
4. What is the trend of the average rating over the years
5. Why average rating and number of votes are moving in different direction with respect to time

```
In [84]: 1 # importing relevant libraries
          2 import pandas as pd
          3 import numpy as np
          4 import matplotlib.pyplot as plt
          5 %matplotlib inline
          6 import seaborn as sns
          7 import sqlite3
          8 import csv
          9 import os
         10 import zipfile
```

```
In [85]: 1 im=os.path.join('im2.db')
          2 conn=sqlite3.connect(im)
          3 cursor=conn.cursor()
```

```
In [45]: 1 table_name_query="""SELECT name
          2                AS 'Table Names'
          3                From sqlite_master
          4                WHERE type='table';"""
          5 pd.read_sql(table_name_query,conn)
```

Out[45]:

	Table Names
0	movie_basics
1	directors
2	known_for
3	movie_akas
4	movie_ratings
5	persons
6	principals
7	writers

```
In [46]: 1 # Query the relevant tables
2 movie_basics_query = """
3 SELECT *
4 FROM movie_basics;
5 """
6
7 movie_ratings_query = """
8 SELECT *
9 FROM movie_ratings;
10 """
11 df_movie_basics_query = pd.read_sql(movie_basics_query, conn)
12 df_movie_ratings_query = pd.read_sql(movie_ratings_query, conn)
13
14
15
```

```
In [47]: 1 # Display the tables above each other
2 df_movie_basics = pd.read_sql(movie_basics_query, conn)
3 df_movie_ratings = pd.read_sql(movie_ratings_query, conn)
4
```

```
In [48]: 1 df_movie_basics = pd.DataFrame(df_movie_basics)
2 df_movie_basics.columns
```

```
Out[48]: Index(['movie_id', 'primary_title', 'original_title', 'start_year',
               'runtime_minutes', 'genres'],
              dtype='object')
```

```
In [49]: 1 print(df_movie_basics.shape)
2 print(df_movie_ratings.shape)
```

```
(146144, 6)
(73856, 3)
```

```
In [50]: 1 df_movie_ratings = pd.DataFrame(df_movie_ratings)
2 df_movie_ratings.columns
```

```
Out[50]: Index(['movie_id', 'averagerating', 'numvotes'], dtype='object')
```

```
In [51]: 1 df_combined = pd.concat([df_movie_basics, df_movie_ratings], axis = 1
2 # print(df_combined)
3 df_combined = pd.DataFrame(df_combined)
4 df_combined.tail()
```

Out[51]:

	movie_id	primary_title	original_title	start_year	runtime_minutes	genres
146139	tt9916538	Kuambil Lagi Hatiku	Kuambil Lagi Hatiku	2019	123.0	Drama
146140	tt9916622	Rodolpho Teóphilo - O Legado de um Pioneiro	Rodolpho Teóphilo - O Legado de um Pioneiro	2015	NaN	Documentary
146141	tt9916706	Dankyavar Danka	Dankyavar Danka	2013	NaN	Comedy
146142	tt9916730	6 Gunn	6 Gunn	2017	116.0	None
146143	tt9916754	Chico Albuquerque - Revelações	Chico Albuquerque - Revelações	2013	NaN	Documentary

```
In [52]: 1 pd.merge(df_movie_basics, df_movie_ratings, how = 'left', on = 'movie_id')
```

Out[52]:

	movie_id	primary_title	original_title	start_year	runtime_minutes	genres
0	tt0063540	Sunghursh	Sunghursh	2013	175.0	Action
1	tt0066787	One Day Before the Rainy Season	Ashad Ka Ek Din	2019	114.0	Bic
2	tt0069049	The Other Side of the Wind	The Other Side of the Wind	2018	122.0	
3	tt0069204	Sabse Bada Sukh	Sabse Bada Sukh	2018	NaN	C
4	tt0100275	The Wandering Soap Opera	La Telenovela Errante	2017	80.0	Comedy, I
...	

In [53]:

```
1 print(df_movie_basics.info())
2 print(df_movie_ratings.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 146144 entries, 0 to 146143
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   movie_id              146144 non-null object
1   primary_title         146144 non-null object
2   original_title        146123 non-null object
3   start_year            146144 non-null int64
4   runtime_minutes       114405 non-null float64
5   genres                140736 non-null object
dtypes: float64(1), int64(1), object(4)
memory usage: 6.7+ MB
None
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73856 entries, 0 to 73855
Data columns (total 3 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   movie_id              73856 non-null object
1   averagerating         73856 non-null float64
2   numvotes              73856 non-null int64
dtypes: float64(1), int64(1), object(1)
memory usage: 1.7+ MB
None
```

In [54]:

```
1 # check missing values
2 print(df_movie_basics.isnull().sum())
3 print(df_movie_ratings.isnull().sum())
```

```
movie_id          0
primary_title     0
original_title    21
start_year        0
runtime_minutes   31739
genres            5408
dtype: int64
movie_id          0
averagerating     0
numvotes          0
dtype: int64
```

```
In [55]: 1 # handle missing values
2 df_movie_basics = df_movie_basics.dropna()
3 df_movie_ratings = df_movie_ratings.dropna()
4 df_movie_basics
5 df_movie_ratings
```

Out[55]:

	movie_id	averagerating	numvotes
0	tt10356526	8.3	31
1	tt10384606	8.9	559
2	tt1042974	6.4	20
3	tt1043726	4.2	50352
4	tt1060240	6.5	21
...
73851	tt9805820	8.1	25
73852	tt9844256	7.5	24
73853	tt9851050	4.7	14
73854	tt9886934	7.0	5
73855	tt9894098	6.3	128

73856 rows × 3 columns

```
In [56]: 1 df_movie_basics.describe()
2
```

Out[56]:

	start_year	runtime_minutes
count	112232.000000	112232.000000
mean	2014.402078	86.261556
std	2.639042	167.896646
min	2010.000000	1.000000
25%	2012.000000	70.000000
50%	2014.000000	87.000000
75%	2017.000000	99.000000
max	2022.000000	51420.000000

In [57]: `df_movie_ratings.describe()`

Out[57]:

	averagerating	numvotes
count	73856.000000	7.385600e+04
mean	6.332729	3.523662e+03
std	1.474978	3.029402e+04
min	1.000000	5.000000e+00
25%	5.500000	1.400000e+01
50%	6.500000	4.900000e+01
75%	7.400000	2.820000e+02
max	10.000000	1.841066e+06

```

1 merged_df = pd.merge(df_movie_basics, df_movie_ratings, on='movie_id')
2 # identify top 10 genres
3 all_genres = merged_df['genres'].str.split(',', expand = True).stack()
4 top_genres = all_genres.value_counts().head(5)
5 # display the list of the first top 10 genres
6 print("top 5 Genres:")
7 print(top_genres.index.tolist())
8

```

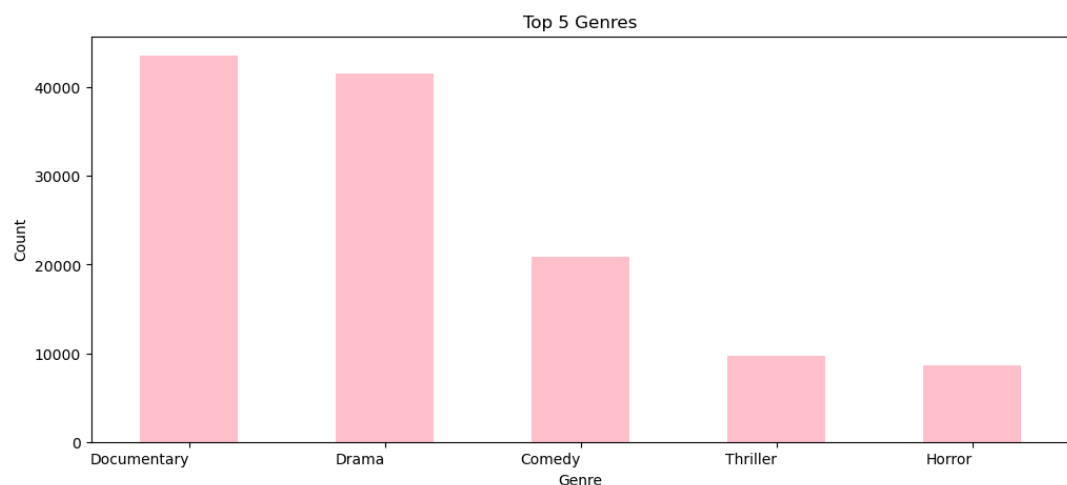
top 5 Genres:

['Documentary', 'Drama', 'Comedy', 'Thriller', 'Horror']

```

1 # visualize the most common genres
2 plt.figure(figsize = (12, 5))
3 ax = top_genres.plot(kind = 'bar', color = 'pink')
4 ax.set_xticklabels(ax.get_xticklabels(), rotation = 0, ha = 'right')
5 plt.title('Top 5 Genres')
6 plt.xlabel('Genre')
7 plt.ylabel('Count')
8 plt.savefig('fig1.png')
9 plt.show()

```

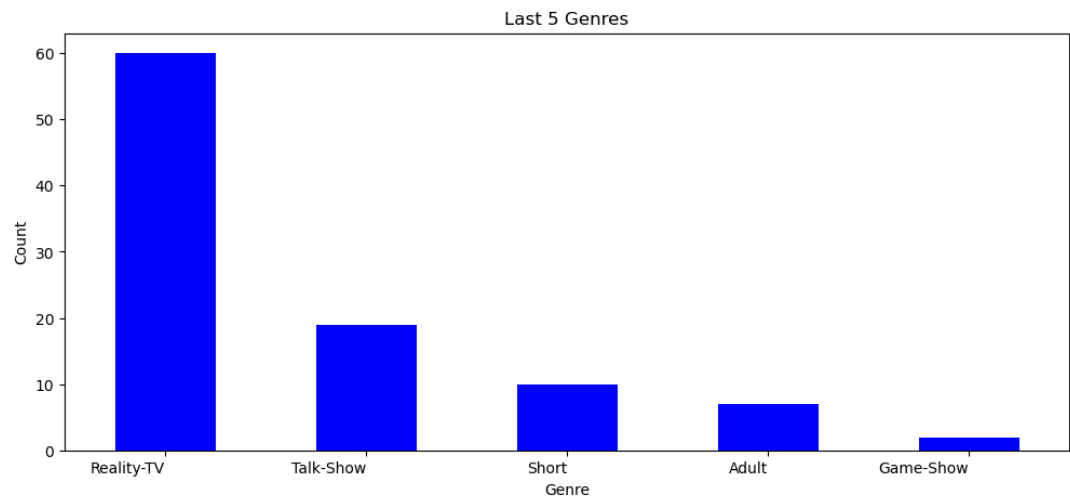


```
In [60]: 1 merged_df = pd.merge(df_movie_basics, df_movie_ratings, on='movie_id')
2 # identify last 5 genres
3 all_genres = merged_df['genres'].str.split(',', expand = True).stack()
4 last_genres = all_genres.value_counts().tail(5)
5 # display the list of the first last 5 genres
6 print("last 5 Genres:")
7 print(last_genres.index.tolist())
8
```

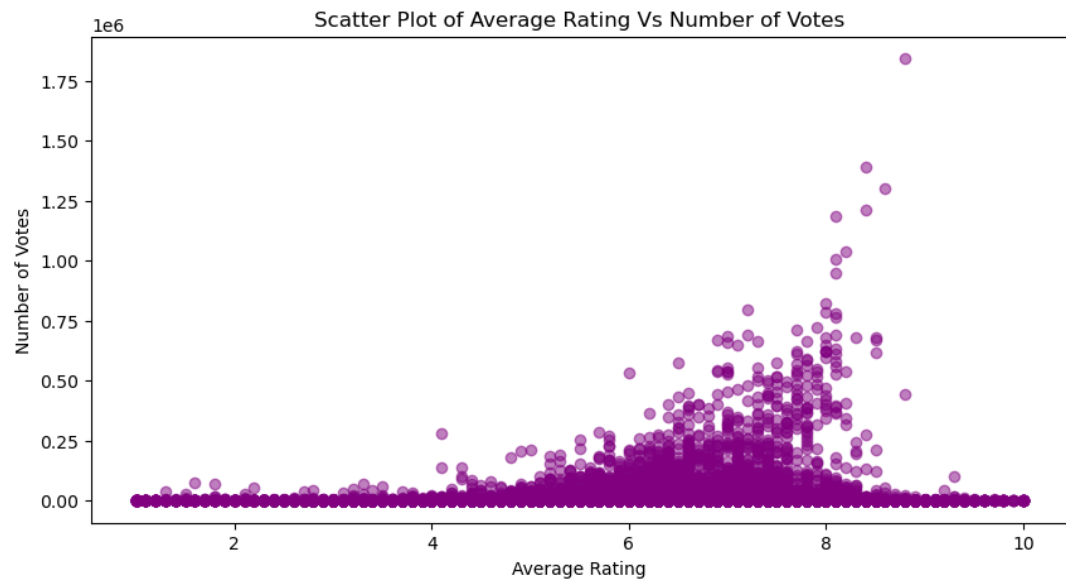
last 5 Genres:

['Reality-TV', 'Talk-Show', 'Short', 'Adult', 'Game-Show']

```
In [61]: 1 # visualize the last genres
2 plt.figure(figsize = (12, 5))
3 ax = last_genres.plot(kind = 'bar', color = 'blue')
4 ax.set_xticklabels(ax.get_xticklabels(), rotation = 0, ha = 'right')
5 plt.title('Last 5 Genres')
6 plt.xlabel('Genre')
7 plt.ylabel('Count')
8 plt.savefig('fig2.png')
9 plt.show()
```



```
In [62]: ▶ 1 # select relevant columns for the scatter plot
2 scatter_data = merged_df[['averagerating', 'numvotes']]
3 # create a scatter plot
4 plt.figure(figsize = (10, 5))
5 plt.scatter(scatter_data['averagerating'], scatter_data['numvotes'])
6 plt.title('Scatter Plot of Average Rating Vs Number of Votes')
7 plt.xlabel('Average Rating')
8 plt.ylabel('Number of Votes')
9 plt.savefig('fig3.png')
10 plt.show()
11
12
```



```
In [63]: ▶ 1 print(df_movie_basics.genres.unique())
2 # print(df_movie_ratings.genres.unique())

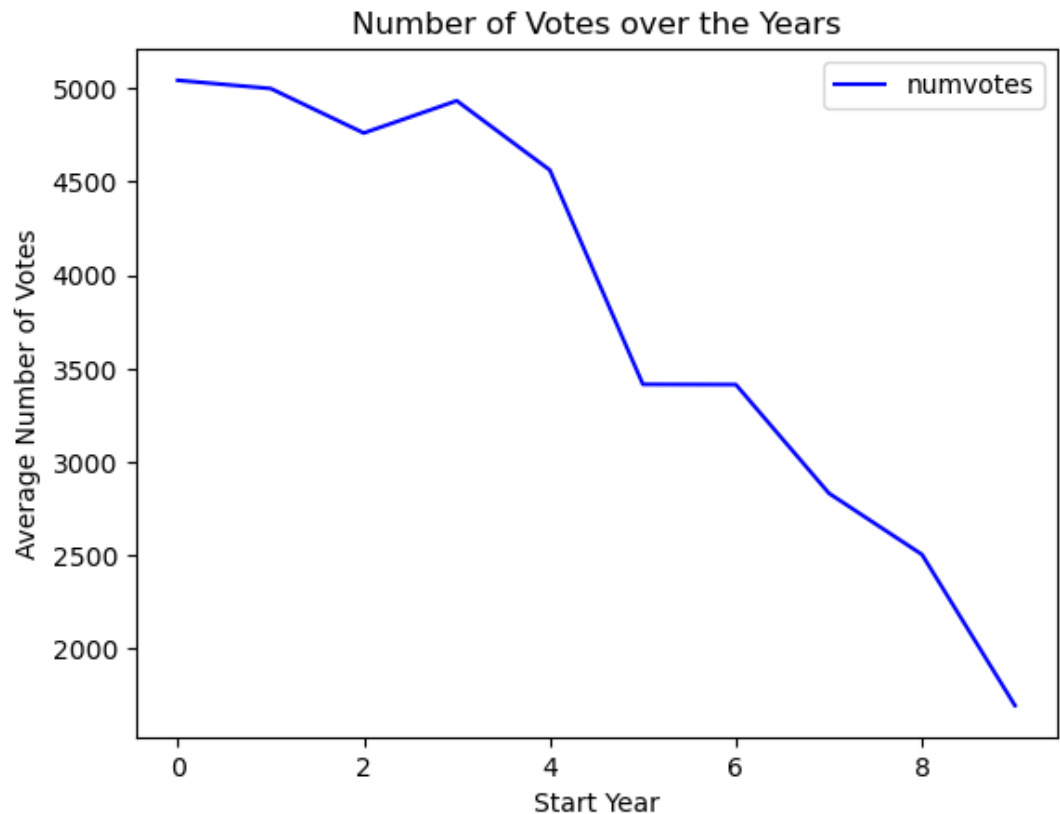
['Action, Crime, Drama' 'Biography, Drama' 'Drama' ...
 'Mystery, Reality-TV, Thriller' 'Music, Musical, Reality-TV' 'Family, Wa
r']
```



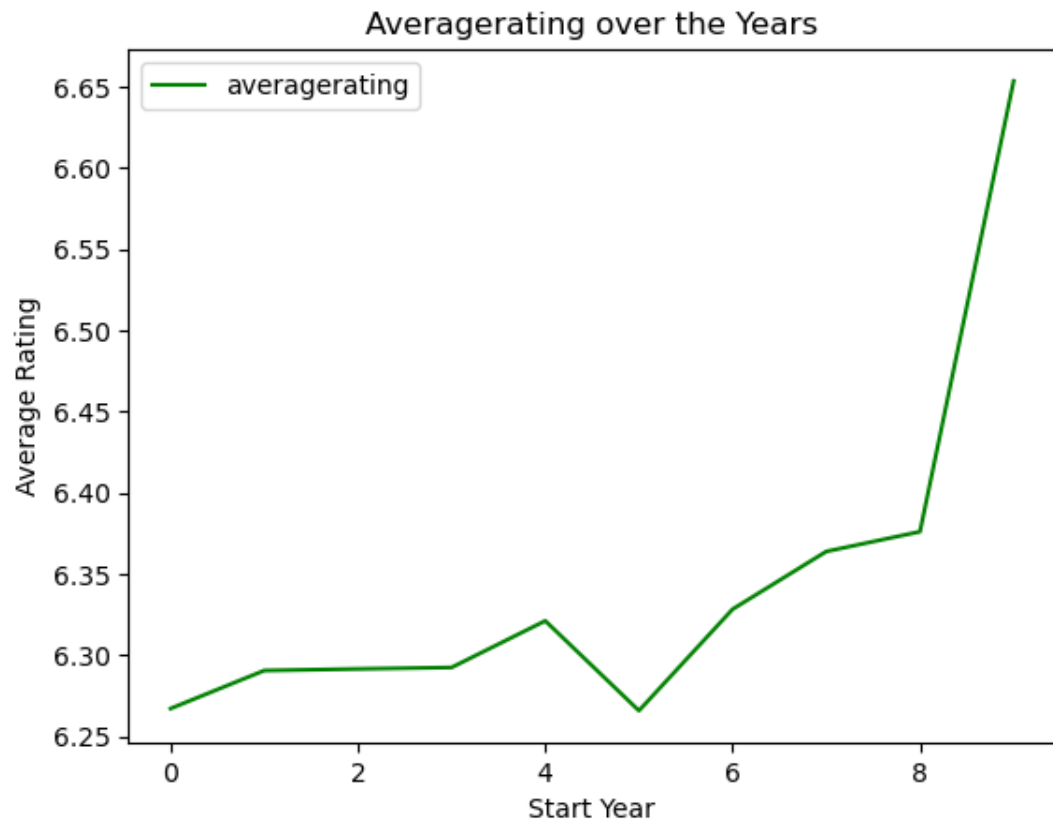
```
In [64]: 1 # Group the data by start_year and calculating the average rating of
2 grouped_data = merged_df.groupby('start_year').agg({'averagerating':
3 # Sorting the data by start_year in ascending order
4 grouped_data = grouped_data.sort_values('start_year')
5 print(grouped_data)
```

start_year	averagerating	numvotes
2010	6.267108	5043.330407
2011	6.290523	4999.282024
2012	6.291509	4761.220652
2013	6.292406	4933.784507
2014	6.321108	4562.855340
2015	6.265817	3415.804183
2016	6.328285	3414.603597
2017	6.363825	2832.789971
2018	6.375993	2505.006998
2019	6.653492	1696.149113
2020	NaN	NaN
2021	NaN	NaN
2022	NaN	NaN

```
In [65]: 1
2 C=grouped_data
3 plt.plot(np.arange(0,len(C)),C["numvotes"],color="blue",label="num
4 # plt a line graph
5 plt.title("Number of Votes over the Years")
6 plt.xlabel("Start Year")
7 plt.ylabel("Average Number of Votes")
8 plt.legend()
9 plt.savefig('fig4.png')
10 plt.show()
```



```
In [66]: ▶ 1 C=grouped_data
2 plt.plot(np.arange(0,len(C)),C["averagerating"],color="green",label="averagerating")
3 # plt a line graph
4 plt.title("Averagerating over the Years")
5 plt.xlabel("Start Year")
6 plt.ylabel("Average Rating")
7 plt.legend()
8 plt.savefig('fig5.png')
9 plt.show()
```



bom_movie_gross_csv

Objectives for bom_movie_gross_csv

1. what is the relationship between domestic gross and foreign gross?
2. What is the trend for domestic gross and foreign gross over the years
3. Is there any outliers in the gross revenue?
4. Is there any significant difference between the domestic and foreign revenue over the years?

```
In [67]: 1 # Load the data
2 df =pd.read_csv('bom.movie1_gross.csv')
3 df
```

Out[67]:

	title	studio	domestic_gross	foreign_gross	year
0	Toy Story 3	BV	415000000.0	652000000	2010
1	Alice in Wonderland (2010)	BV	334200000.0	691300000	2010
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000	2010
3	Inception	WB	292600000.0	535700000	2010
4	Shrek Forever After	P/DW	238700000.0	513900000	2010
...
3382	The Quake	Magn.	6200.0	NaN	2018
3383	Edward II (2018 re-release)	FM	4800.0	NaN	2018
3384	El Pacto	Sony	2500.0	NaN	2018
3385	The Swan	Synergetic	2400.0	NaN	2018
3386	An Actor Prepares	Grav.	1700.0	NaN	2018

3387 rows × 5 columns

```
In [68]: 1 df.shape
```

Out[68]: (3387, 5)

```
In [69]: 1 df.columns
```

Out[69]: Index(['title', 'studio', 'domestic_gross', 'foreign_gross', 'year'], dtype='object')

```
In [70]: 1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3387 entries, 0 to 3386
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   title                 3387 non-null   object
1   studio                3382 non-null   object
2   domestic_gross        3359 non-null   float64
3   foreign_gross         2037 non-null   object
4   year                  3387 non-null   int64
dtypes: float64(1), int64(1), object(3)
memory usage: 132.4+ KB
```

In [71]: 1 df.dtypes

Out[71]: title object
studio object
domestic_gross float64
foreign_gross object
year int64
dtype: object

In [72]: 1 df.dtypes

Out[72]: title object
studio object
domestic_gross float64
foreign_gross object
year int64
dtype: object

In [73]: 1 # check missing values
2 print(df.isnull().sum())

title 0
studio 5
domestic_gross 28
foreign_gross 1350
year 0
dtype: int64

In [74]: 1 # handle missing values
2 df = df.dropna()
3 df

Out[74]:

	title	studio	domestic_gross	foreign_gross	year
0	Toy Story 3	BV	415000000.0	652000000	2010
1	Alice in Wonderland (2010)	BV	334200000.0	691300000	2010
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000	2010
3	Inception	WB	292600000.0	535700000	2010
4	Shrek Forever After	P/DW	238700000.0	513900000	2010
...
3275	I Still See You	LGF	1400.0	1500000	2018
3286	The Catcher Was a Spy	IFC	725000.0	229000	2018
3309	Time Freak	Grindstone	10000.0	256000	2018
3342	Reign of Judges: Title of Liberty - Concept Short	Darin Southa	93200.0	5200	2018
3353	Antonio Lopez 1970: Sex Fashion & Disco	FM	43200.0	30000	2018

2007 rows × 5 columns

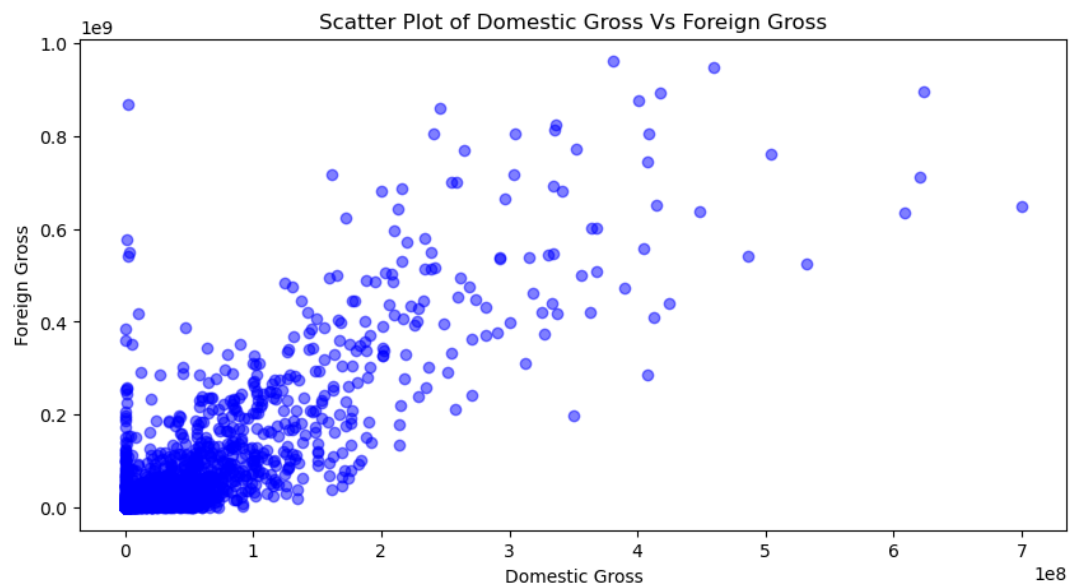
```
In [75]: 1 import warnings
        2 warnings.filterwarnings('ignore')
```

```
In [76]: 1 df.describe()
```

Out[76]:

	domestic_gross	year
count	2.007000e+03	2007.000000
mean	4.701984e+07	2013.506228
std	8.162689e+07	2.597997
min	4.000000e+02	2010.000000
25%	6.700000e+05	2011.000000
50%	1.670000e+07	2013.000000
75%	5.605000e+07	2016.000000
max	9.367000e+08	2018.000000

```
In [77]: 1 df['domestic_gross'] = pd.to_numeric(df['domestic_gross'], errors =
        2 df['foreign_gross'] = pd.to_numeric(df['foreign_gross'], errors =
        3 # create a scatter plot
        4 plt.figure(figsize = (10, 5))
        5 plt.scatter(df['domestic_gross'], df['foreign_gross'], alpha=0.5,
        6 plt.title('Scatter Plot of Domestic Gross Vs Foreign Gross')
        7 plt.xlabel('Domestic Gross')
        8 plt.ylabel('Foreign Gross')
        9 plt.savefig('fig3.png')
       10 plt.show()
```



```
In [78]: 1 # DataFrame for domestic_gross and foreign_gross
2 domestic_data = df.domestic_gross
3 foreign_data = df.foreign_gross
4 # display the data for the domestic_gross column
5 print("Domestic Gross Data:")
6 print(domestic_data)
7 # display the data for foreign_gross column
8 print("\nForeign Gross Data:")
9 print(foreign_data)
10
```

Domestic Gross Data:

```
0      415000000.0
1      334200000.0
2      296000000.0
3      292600000.0
4      238700000.0
```

...

```
3275      1400.0
3286      725000.0
3309      10000.0
3342      93200.0
3353      43200.0
```

Name: domestic_gross, Length: 2007, dtype: float64

Foreign Gross Data:

```
0      652000000.0
1      691300000.0
2      664300000.0
3      535700000.0
4      513900000.0
```

...

```
3275      1500000.0
3286      229000.0
3309      256000.0
3342       5200.0
3353      30000.0
```

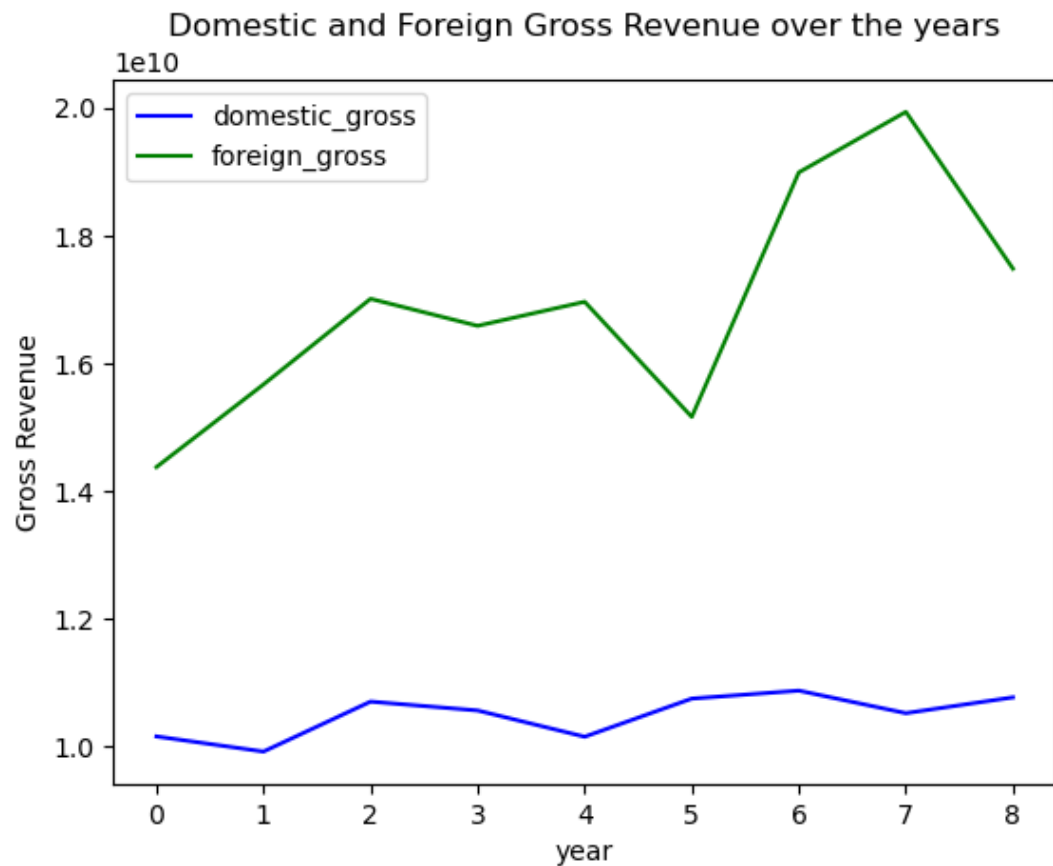
Name: foreign_gross, Length: 2007, dtype: float64

```
In [79]: 1 df1 = df.groupby(['year'])[['domestic_gross', 'foreign_gross']].sum()
2 df1.head()
```

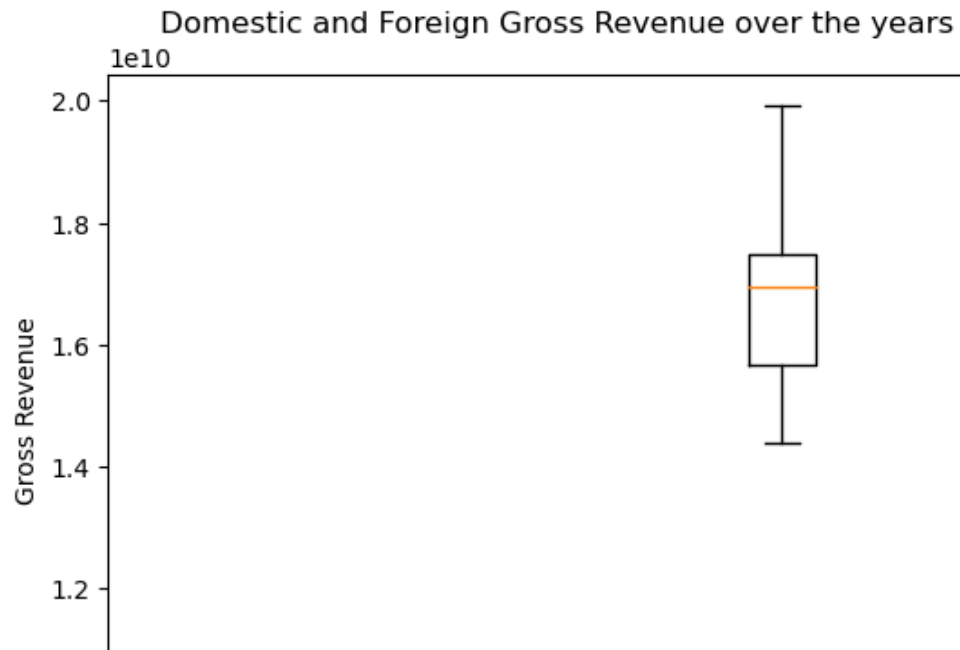
Out[79]:

	year	domestic_gross	foreign_gross
0	2010	1.015274e+10	1.436937e+10
1	2011	9.915690e+09	1.566287e+10
2	2012	1.069786e+10	1.700298e+10
3	2013	1.055885e+10	1.658024e+10
4	2014	1.014798e+10	1.695667e+10

```
In [80]: 1 C=df1
2 plt.plot(np.arange(0,len(C)),C["domestic_gross"],color="blue",label="domestic_gross")
3 plt.plot(np.arange(0,len(C)),C["foreign_gross"],color="green",label="foreign_gross")
4 # plt a line graph
5 plt.title("Domestic and Foreign Gross Revenue over the years")
6 plt.xlabel("year")
7 plt.ylabel("Gross Revenue")
8 plt.legend()
9 plt.savefig('fig7.png')
10 plt.show()
```



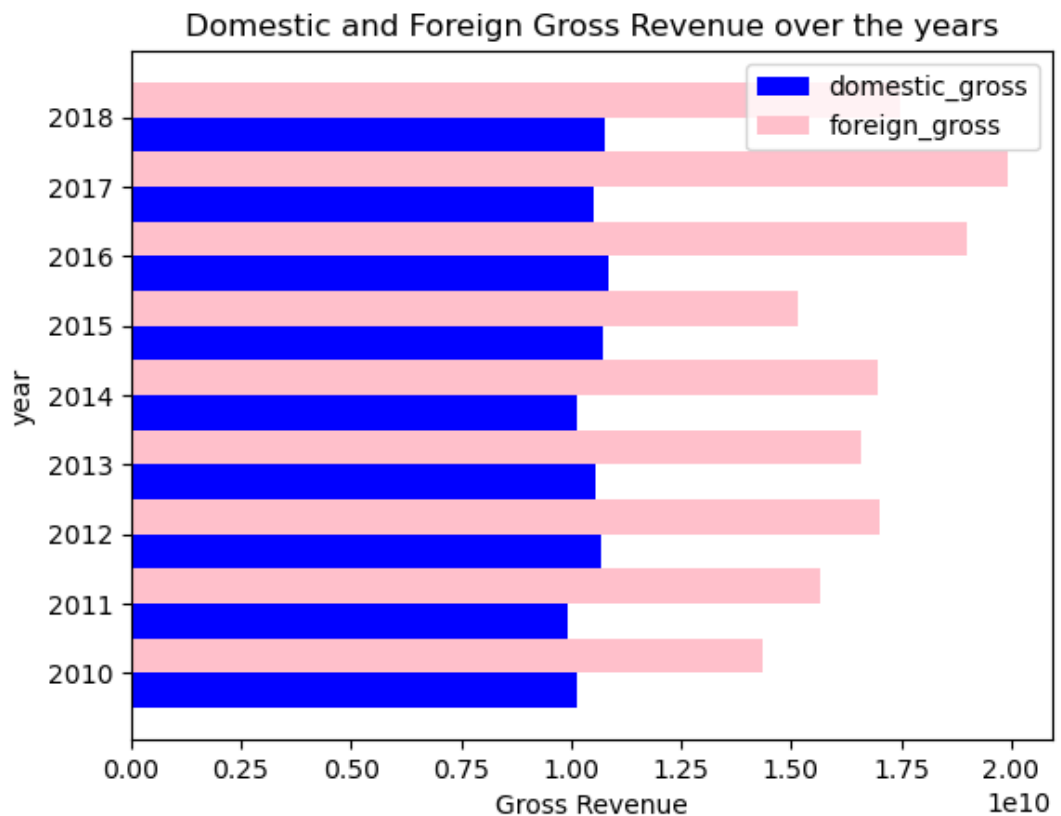
```
In [81]: ▶ 1 C=df1
2 plt.boxplot([df1["domestic_gross"], df1["foreign_gross"]], labels :
3 # plt.boxplot
4 plt.title("Domestic and Foreign Gross Revenue over the years")
5 plt.xlabel("Variable")
6 plt.ylabel("Gross Revenue")
7 plt.savefig('fig8.png')
8 plt.show()
```




```

In [82]: 1 df1 = df.groupby(['year'])[['domestic_gross', 'foreign_gross']].sum
2 C=df1
3 # convert the "year" column to datetime and extract the year
4 bar_width = 0.5
5 plt.barh(np.arange(len(df1)), df1["domestic_gross"], height = bar_w
6 plt.barh(np.arange(len(df1)) + bar_width, df1["foreign_gross"], he
7 # set y-axis ticks and labels
8 plt.yticks(np.arange(len(df1)) + bar_width / 2, df1["year"])
9 plt.title("Domestic and Foreign Gross Revenue over the years")
10 plt.xlabel("Gross Revenue")
11 plt.ylabel("year")
12 plt.legend()
13 plt.savefig('fig9.png')
14 plt.show()

```



```

In [ ]: 1

```