

Final Project Submission

Please fill out:

- Student name: Charles Ndegwa
- Student pace: self paced / part time
- Scheduled project review date/time:
- Instructor name:

```
In [302]: 1 # Import standard packages
          2 import pandas as pd
          3 import numpy as np
          4 import matplotlib.pyplot as plt
          5 import seaborn as sns
          6 import sqlite3
          7 %matplotlib inline
```

```
In [303]: 1 filepath = '/Users/ndegwa/Documents/GitHub/Phase/data'
```

```
In [304]: 1 #start by opening a connection to the database with sqlite3.conn
          2 conn = sqlite3.connect(f'{filepath}/im.db')
```

```
In [305]: 1 #cursor object to execute SQL commands
          2 cur = conn.cursor()
```

```
In [306]: 1 # special query for finding the table names
          2 cur.execute("""SELECT name FROM sqlite_master WHERE type = 'tabl
          3 # Fetch the result and store it in table_names
          4 table_names = cur.fetchall()
          5 table_names
```

```
Out[306]: [('movie_basics',),
           ('directors',),
           ('known_for',),
           ('movie_akas',),
           ('movie_ratings',),
           ('persons',),
           ('principals',),
           ('writers',)]
```

```
In [307]: 1 #show all information about the movie_basics table
          2 cur.execute("""SELECT * FROM movie_basics;""").fetchall()
```

```
Out[307]: [('tt0063540', 'Sunghursh', 'Sunghursh', 2013, 175.0, 'Action,Crime,Drama'),
            ('tt0066787',
             'One Day Before the Rainy Season',
             'Ashad Ka Ek Din',
             2019,
             114.0,
             'Biography,Drama'),
            ('tt0069049',
             'The Other Side of the Wind',
             'The Other Side of the Wind',
             2018,
             122.0,
             'Drama'),
            ('tt0069204',
             'Sabse Bada Sukh',
             'Sabse Bada Sukh',
             2018,
             None,
             'Comedy,Drama')]
```

```
In [308]: 1 #only run after above code to get column names of select table
          2 cur.description
```

```
Out[308]: (('movie_id', None, None, None, None, None, None),
            ('primary_title', None, None, None, None, None, None),
            ('original_title', None, None, None, None, None, None),
            ('start_year', None, None, None, None, None, None),
            ('runtime_minutes', None, None, None, None, None, None),
            ('genres', None, None, None, None, None, None))
```

```
In [309]: 1 #show all information about the movie_ratings table
          2 cur.execute("""SELECT * FROM movie_ratings;""").fetchall()
```

```
Out[309]: [('tt10356526', 8.3, 31),
            ('tt10384606', 8.9, 559),
            ('tt1042974', 6.4, 20),
            ('tt1043726', 4.2, 50352),
            ('tt1060240', 6.5, 21),
            ('tt1069246', 6.2, 326),
            ('tt1094666', 7.0, 1613),
            ('tt1130982', 6.4, 571),
            ('tt1156528', 7.2, 265),
            ('tt1161457', 4.2, 148),
            ('tt1171222', 5.1, 8296),
            ('tt1174693', 5.8, 2381),
            ('tt1181840', 7.0, 5494),
            ('tt1193623', 8.0, 5),
            ('tt1199588', 5.5, 74),
            ('tt1204784', 5.8, 6),
            ('tt1210166', 7.6, 326657),
            ('tt1212419', 6.5, 87288),
            ('tt1220911', 5.0, 941),
            ('tt1220920', 7.4, 120142)]
```

```
In [310]: 1 #only run after above code to get column names of select table
          2 cur.description
```

```
Out[310]: (('movie_id', None, None, None, None, None, None),
          ('averagerating', None, None, None, None, None, None),
          ('numvotes', None, None, None, None, None, None))
```

```
In [311]: 1 #express the result of sql query in a DataFrame
          2 pd.DataFrame(
          3     # Execute the SQL query to select rows from movie_ratings wh
          4     cur.execute("""SELECT * FROM movie_ratings WHERE numvotes >
          5     columns=[x[0] for x in cur.description])
```

```
Out[311]:
```

	movie_id	averagerating	numvotes
0	tt1210166	7.6	326657
1	tt1229238	7.4	428142
2	tt1232829	7.2	477771
3	tt1403981	7.1	129443
4	tt1535109	7.8	387402
...
869	tt2205697	7.2	78903
870	tt2234003	7.4	52266
871	tt2386490	7.6	60769
872	tt2404461	7.8	41191
873	tt7048622	7.7	11168

874 rows × 3 columns

```
In [312]: 1 # Execute the SQL query to select rows from movie_ratings where
          2 query = """
          3 SELECT movie_ratings.*, movie_basics.primary_title, movie_basics
          4         movie_basics.start_year, movie_basics.runtime_minutes, mo
          5 FROM movie_ratings
          6 INNER JOIN movie_basics ON movie_ratings.movie_id = movie_basics
          7 WHERE movie_ratings.numvotes > 10000 AND movie_ratings.averagera
          8 """
```

```
In [313]: 1 # Fetch the data and create a DataFrame
          2 df = pd.DataFrame(cur.execute(query).fetchall(), columns=[x[0] f
          3 # View only the columns from 'movie_basics'
          4 df_movie_basics = df[['primary_title', 'original_title', 'start_
          5 df_movie_basics
```

Out [313]:

	primary_title	original_title	start_year	runtime_minutes	genres
0	Moneyball	Moneyball	2011	133.0	Biography,Drama,Sport
1	Mission: Impossible - Ghost Protocol	Mission: Impossible - Ghost Protocol	2011	132.0	Action,Adventure,Thriller
2	21 Jump Street	21 Jump Street	2012	109.0	Action,Comedy,Crime
3	Remember Me	Remember Me	2010	113.0	Drama,Romance
4	Captain Phillips	Captain Phillips	2013	134.0	Biography,Drama,Thriller
...
869	Stuck in Love.	Stuck in Love.	2012	97.0	Comedy,Drama,Romance
870	Calvary	Calvary	2014	102.0	Drama
871	How to Train Your Dragon: The Hidden World	How to Train Your Dragon: The Hidden World	2019	104.0	Action,Adventure,Animation
872	The Past	Le passé	2013	130.0	Drama,Mystery
873	The Insult	L'insulte	2017	113.0	Crime,Drama,Thriller

874 rows × 5 columns

```
In [314]: 1 # Create an empty list to store the new rows
          2 new_rows = []
          3
          4 # Iterate over each row in the DataFrame
          5 for index, row in df_movie_basics.iterrows():
          6     # Split the 'genres' string by comma and iterate over the re
          7     for genre in row['genres'].split(','):
          8         # Create a new row with the details from other columns a
          9         new_row = row.copy()
         10         new_row['genres'] = genre.strip() # Strip leading/trail
         11         new_rows.append(new_row)
         12
         13 # Create a new DataFrame from the list of new rows
         14 df_individual_genres = pd.DataFrame(new_rows)
         15
         16 # Print the new DataFrame
         17 print(df_individual_genres)
```

```
primary_title \
0 Moneyball
0 Moneyball
0 Moneyball
1 Mission: Impossible – Ghost Protocol
1 Mission: Impossible – Ghost Protocol
..
872 The Past
872 The Past
873 The Insult
873 The Insult
873 The Insult

original_title start_year runtime_minu
tes \
0 Moneyball 2011 13
3.0
0 Moneyball 2011 13
3.0
0 Moneyball 2011 13
3.0
1 Mission: Impossible – Ghost Protocol 2011 13
2.0
1 Mission: Impossible – Ghost Protocol 2011 13
2.0
.. ...
...
872 Le passé 2013 13
0.0
872 Le passé 2013 13
0.0
873 L'insulte 2017 11
3.0
873 L'insulte 2017 11
3.0
873 L'insulte 2017 11
3.0

genres
0 Biography
0 Drama
0 Sport
1 Action
1 Adventure
..
872 Drama
872 Mystery
873 Crime
873 Drama
873 Thriller
```

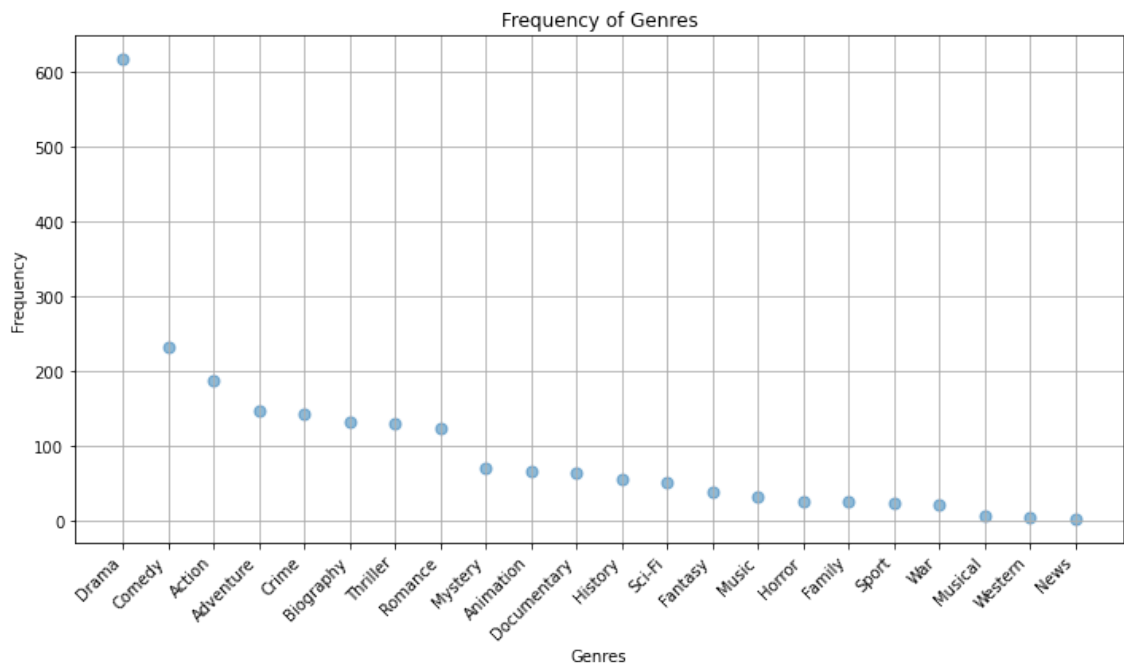
[2190 rows x 5 columns]

In [315]: 1 `print(df_individual_genres.info())`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2190 entries, 0 to 873
Data columns (total 5 columns):
#   Column              Non-Null Count  Dtype
---  -
0   primary_title        2190 non-null   object
1   original_title       2190 non-null   object
2   start_year           2190 non-null   int64
3   runtime_minutes      2190 non-null   float64
4   genres               2190 non-null   object
dtypes: float64(1), int64(1), object(3)
memory usage: 102.7+ KB
None
```

In [316]:

```
1 # Count the occurrences of each genre
2 genre_counts = df_individual_genres['genres'].value_counts()
3
4 # Plot the scatter plot
5 plt.figure(figsize=(10, 6))
6 plt.scatter(genre_counts.index, genre_counts.values, s=50, alpha=0.5)
7 plt.title('Frequency of Genres')
8 plt.xlabel('Genres')
9 plt.ylabel('Frequency')
10 plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for readability
11 plt.grid(True)
12 plt.tight_layout()
13 plt.show()
```



```
In [317]: 1 #path to data files
          2 filepath = "/Users/ndegwa/Documents/GitHub/Phase/data"
          3
          4 #read all dataset files
          5 bom_movie_gross = pd.read_csv(f'{filepath}/bom.movie_gross.csv')
          6 tmdb_movies = pd.read_csv(f'{filepath}/tmdb.movies.csv')
          7 tn_movie_budgets = pd.read_csv(f'{filepath}/tn.movie_budgets.csv')
          8 rt_movie_info = pd.read_csv(f'{filepath}/rt.movie_info.tsv', sep
          9 rt_reviews = pd.read_csv(f'{filepath}/rt.reviews.tsv', sep='\t',
          10
```

```
In [318]: 1 bom_movie_gross.columns
```

```
Out[318]: Index(['title', 'studio', 'domestic_gross', 'foreign_gross', 'year'], dtype='object')
```

```
In [319]: 1 tmdb_movies.columns
```

```
Out[319]: Index(['Unnamed: 0', 'genre_ids', 'id', 'original_language', 'original_title',
                 'popularity', 'release_date', 'title', 'vote_average', 'vote_count'],
                 dtype='object')
```


In [320]:

1	tmdb_movies.info
---	------------------

```
Out[320]: <bound method DataFrame.info of Unnamed: 0 genre_
ids id original_language \
0 0 [12, 14, 10751] 12444 en
1 1 [14, 12, 16, 10751] 10191 en
2 2 [12, 28, 878] 10138 en
3 3 [16, 35, 10751] 862 en
4 4 [28, 878, 12] 27205 en
...
26512 26512 [27, 18] 488143 en
26513 26513 [18, 53] 485975 en
26514 26514 [14, 28, 12] 381231 en
26515 26515 [10751, 12, 28] 366854 en
26516 26516 [53, 27] 309885 en
```

```
original_title popularity rel
ease_date \
0 Harry Potter and the Deathly Hallows: Part 1 33.533 2
010-11-19
1 How to Train Your Dragon 28.734 2
010-03-26
2 Iron Man 2 28.515 2
010-05-07
3 Toy Story 28.005 1
995-11-22
4 Inception 27.920 2
010-07-16
...
...
...
26512 Laboratory Conditions 0.600 2
018-10-13
26513 _EXHIBIT_84xxx_ 0.600 2
018-05-01
26514 The Last One 0.600 2
018-10-01
26515 Trailer Made 0.600 2
018-06-22
26516 The Church 0.600 2
018-10-05
```

```
title vote_average
vote_count
0 Harry Potter and the Deathly Hallows: Part 1 7.7
10788
1 How to Train Your Dragon 7.7
7610
2 Iron Man 2 6.8
12368
3 Toy Story 7.9
10174
4 Inception 8.3
22186
...
...
...
26512 Laboratory Conditions 0.0
1
26513 _EXHIBIT_84xxx_ 0.0
1
26514 The Last One 0.0
1
26515 Trailer Made 0.0
1
```

26516
1

[26517 rows x 10 columns]>

In [321]: 1 tn_movie_budgets.columns

Out[321]: Index(['id', 'release_date', 'movie', 'production_budget', 'domestic_gross',
 'worldwide_gross'],
 dtype='object')

In [322]: 1 rt_movie_info.columns

Out[322]: Index(['id', 'synopsis', 'rating', 'genre', 'director', 'writer',
 'theater_date', 'dvd_date', 'currency', 'box_office', 'runtime',
 'studio'],
 dtype='object')

In [323]: 1 rt_reviews.columns

Out[323]: Index(['id', 'review', 'rating', 'fresh', 'critic', 'top_critic',
 'publisher',
 'date'],
 dtype='object')

In [324]: 1 *# Displaying basic information about the selected DataFrame*
 2 print("Basic Info:")
 3 tn_movie_budgets.info()
 4

Basic Info:

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 5782 entries, 0 to 5781

Data columns (total 6 columns):

#	Column	Non-Null Count	Dtype
0	id	5782 non-null	int64
1	release_date	5782 non-null	object
2	movie	5782 non-null	object
3	production_budget	5782 non-null	object
4	domestic_gross	5782 non-null	object
5	worldwide_gross	5782 non-null	object

dtypes: int64(1), object(5)

memory usage: 271.2+ KB

```
In [325]: 1 # Handling missing values
2 print("\nHandling Missing Values:")
3 tn_movie_budgets.isnull().sum()
4
5 # Handling duplicate rows
6 print("\nHandling Duplicate Rows:")
7 tn_movie_budgets = tn_movie_budgets.drop_duplicates()
8
9 # Summary statistics
10 print("\nSummary Statistics:")
11 tn_movie_budgets.describe()
```

Handling Missing Values:

Handling Duplicate Rows:

Summary Statistics:

```
Out [325]:
```

	id
count	5782.000000
mean	50.372363
std	28.821076
min	1.000000
25%	25.000000
50%	50.000000
75%	75.000000
max	100.000000

```
In [326]: 1 # Convert 'production_budget', 'domestic_gross', and 'worldwide_gross' to float
2 tn_movie_budgets['production_budget'] = tn_movie_budgets['production_budget'].astype(float)
3 tn_movie_budgets['domestic_gross'] = tn_movie_budgets['domestic_gross'].astype(float)
4 tn_movie_budgets['worldwide_gross'] = tn_movie_budgets['worldwide_gross'].astype(float)
```

```
In [327]: 1 # Convert 'release_date' column to datetime
2 tn_movie_budgets['release_date'] = pd.to_datetime(tn_movie_budgets['release_date'])
```

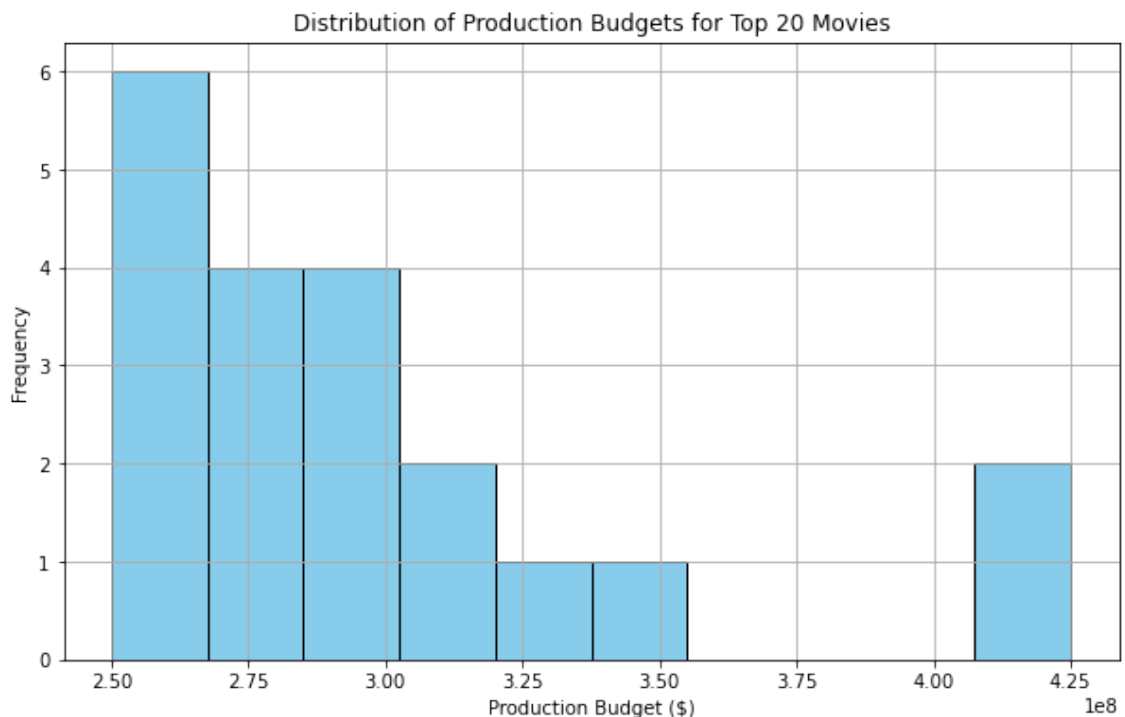
1. Production Budgets Analysis

```
In [328]: 1
          2 # Summary statistics of production budgets
          3 print("Production Budgets Summary Statistics:")
          4 tn_movie_budgets['production_budget'].describe()
```

Production Budgets Summary Statistics:

```
Out [328]: count      5.782000e+03
           mean      3.158776e+07
           std       4.181208e+07
           min       1.100000e+03
           25%       5.000000e+06
           50%       1.700000e+07
           75%       4.000000e+07
           max       4.250000e+08
           Name: production_budget, dtype: float64
```

```
In [329]: 1 # Sort the DataFrame by production budget in descending order and
          2 top_20_movies = tn_movie_budgets.nlargest(20, 'production_budget')
          3
          4 # Plot a histogram of the production budgets for the top 20 movies
          5 plt.figure(figsize=(10, 6))
          6 plt.hist(top_20_movies['production_budget'], bins=10, color='skyblue')
          7 plt.title('Distribution of Production Budgets for Top 20 Movies')
          8 plt.xlabel('Production Budget ($)')
          9 plt.ylabel('Frequency')
         10 plt.grid(True)
         11 plt.show()
         12
```



```
In [330]: 1 # Sort the DataFrame by production budget in descending order and
2 top_20_movies = tn_movie_budgets.nlargest(20, 'production_budget')
3
4 # Plot a bar chart of the production budgets for the top 20 movies
5 plt.figure(figsize=(14, 8))
6 plt.bar(top_20_movies['movie'], top_20_movies['production_budget'])
7 plt.title('Production Budget for Top 20 Movies')
8 plt.xlabel('Movie Title')
9 plt.ylabel('Production Budget ($)')
10 plt.xticks(rotation=90, fontsize=10) # Rotate movie titles for
11 plt.grid(True)
12 plt.tight_layout()
13
14 # Set font family to a commonly available one that supports a wide
15 plt.rcParams['font.family'] = 'DejaVu Sans'
16
17 plt.show()
```

/Applications/anaconda3/envs/learn-env/lib/python3.8/site-packages/matplotlib/backends/backend_agg.py:238: RuntimeWarning: Glyph 128 missing from current font.

font.set_text(s, 0.0, flags=flags)

/Applications/anaconda3/envs/learn-env/lib/python3.8/site-packages/matplotlib/backends/backend_agg.py:238: RuntimeWarning: Glyph 153 missing from current font.

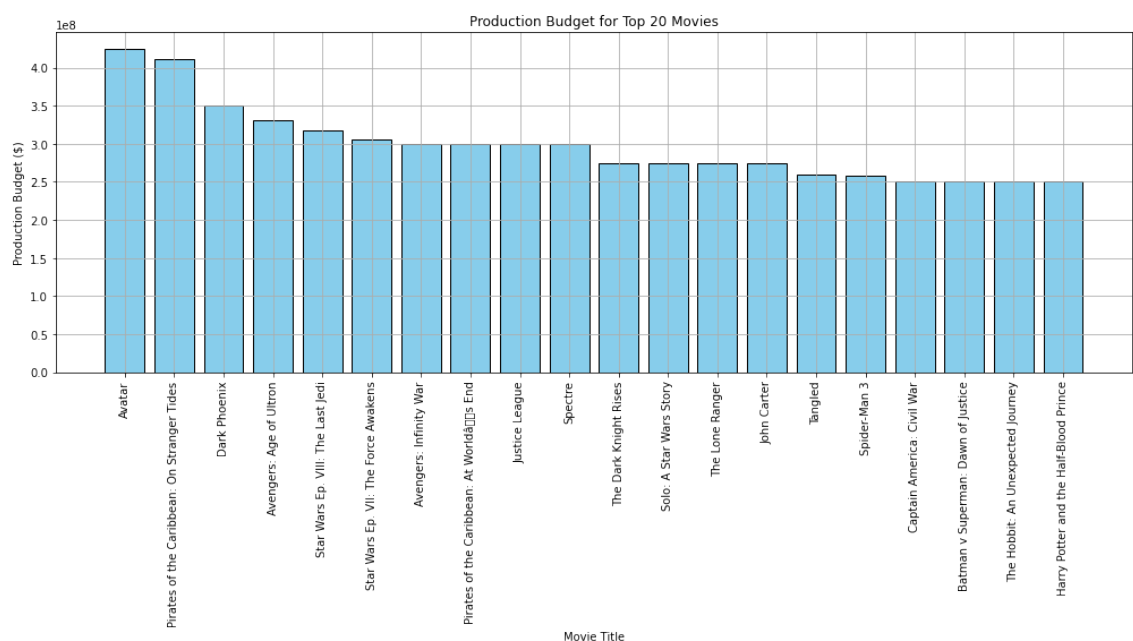
font.set_text(s, 0.0, flags=flags)

/Applications/anaconda3/envs/learn-env/lib/python3.8/site-packages/matplotlib/backends/backend_agg.py:201: RuntimeWarning: Glyph 128 missing from current font.

font.set_text(s, 0, flags=flags)

/Applications/anaconda3/envs/learn-env/lib/python3.8/site-packages/matplotlib/backends/backend_agg.py:201: RuntimeWarning: Glyph 153 missing from current font.

font.set_text(s, 0, flags=flags)



2.Domestic and Worldwide Gross Revenues Analysis

```
In [331]: 1 # Summary statistics of domestic gross revenues
          2 print("\nDomestic Gross Summary Statistics:")
          3 tn_movie_budgets['domestic_gross'].describe()
          4
```

Domestic Gross Summary Statistics:

```
Out[331]: count      5.782000e+03
          mean      4.187333e+07
          std       6.824060e+07
          min       0.000000e+00
          25%      1.429534e+06
          50%      1.722594e+07
          75%      5.234866e+07
          max      9.366622e+08
          Name: domestic_gross, dtype: float64
```

```
In [332]: 1 #Summary statistics of Worldwide gross revenues
          2 print("\nWorldwide Gross Summary Statistics:")
          3 tn_movie_budgets['worldwide_gross'].describe()
```

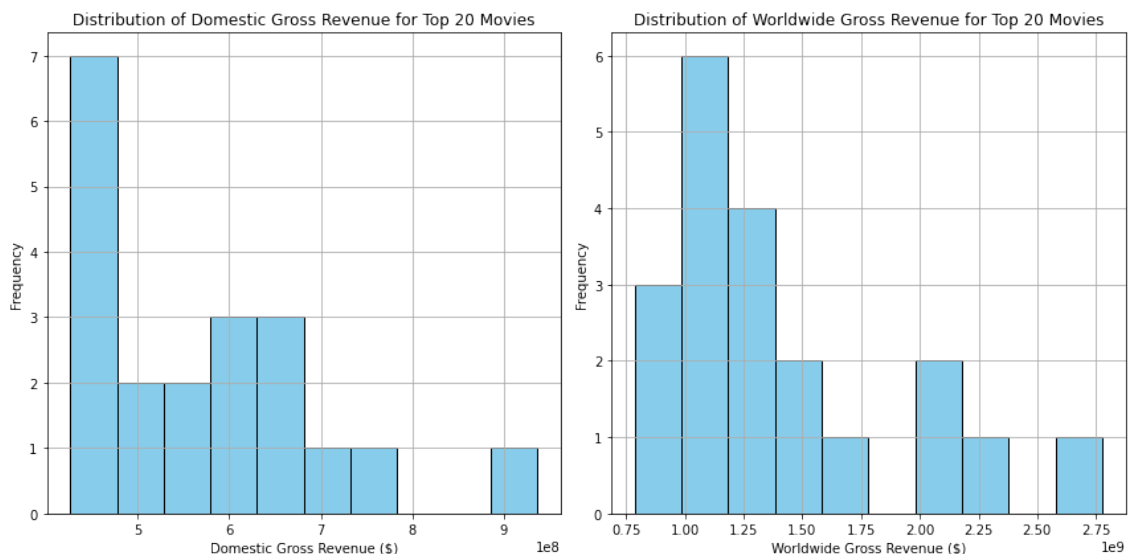
Worldwide Gross Summary Statistics:

```
Out[332]: count      5.782000e+03
          mean      9.148746e+07
          std      1.747200e+08
          min      0.000000e+00
          25%      4.125415e+06
          50%      2.798445e+07
          75%      9.764584e+07
          max      2.776345e+09
          Name: worldwide_gross, dtype: float64
```

```

In [333]: 1 # Sort the DataFrame by domestic and worldwide gross revenues in
2 top_20_movies = tn_movie_budgets.nlargest(20, ['domestic_gross',
3
4 # Plot histograms of domestic and worldwide gross revenues for t
5 plt.figure(figsize=(12, 6))
6
7 # Histogram of domestic gross revenue
8 plt.subplot(1, 2, 1)
9 plt.hist(top_20_movies['domestic_gross'], bins=10, color='skyblue')
10 plt.title('Distribution of Domestic Gross Revenue for Top 20 Mov
11 plt.xlabel('Domestic Gross Revenue ($)')
12 plt.ylabel('Frequency')
13 plt.grid(True)
14
15 # Histogram of worldwide gross revenue
16 plt.subplot(1, 2, 2)
17 plt.hist(top_20_movies['worldwide_gross'], bins=10, color='skyblue')
18 plt.title('Distribution of Worldwide Gross Revenue for Top 20 Mo
19 plt.xlabel('Worldwide Gross Revenue ($)')
20 plt.ylabel('Frequency')
21 plt.grid(True)
22
23 plt.tight_layout()
24 plt.show()

```



3. Release Dates Analysis

```

In [334]: 1 # Extracting year and month from release_date
2 tn_movie_budgets['release_year'] = pd.to_datetime(tn_movie_budg
3 tn_movie_budgets['release_month'] = pd.to_datetime(tn_movie_budg

```

```

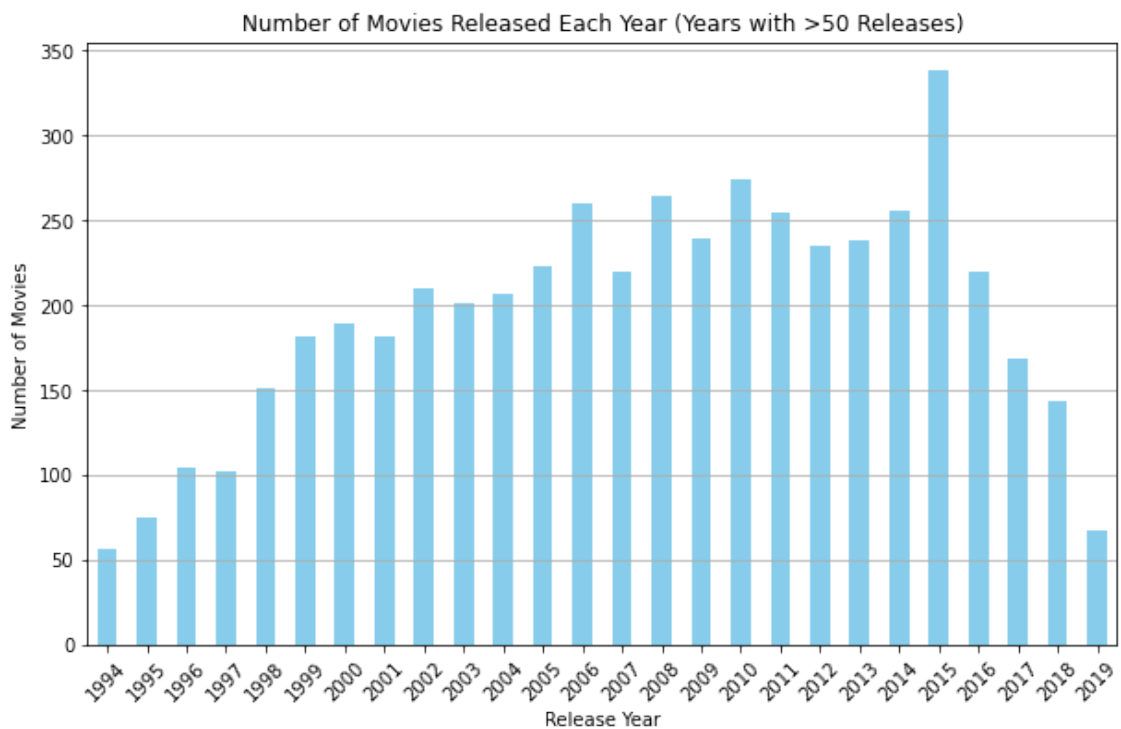
In [335]: 1 # Convert release_date to datetime and extract release year
2 tn_movie_budgets['release_date'] = pd.to_datetime(tn_movie_budg
3 tn_movie_budgets['release_year'] = tn_movie_budgets['release_dat
4

```

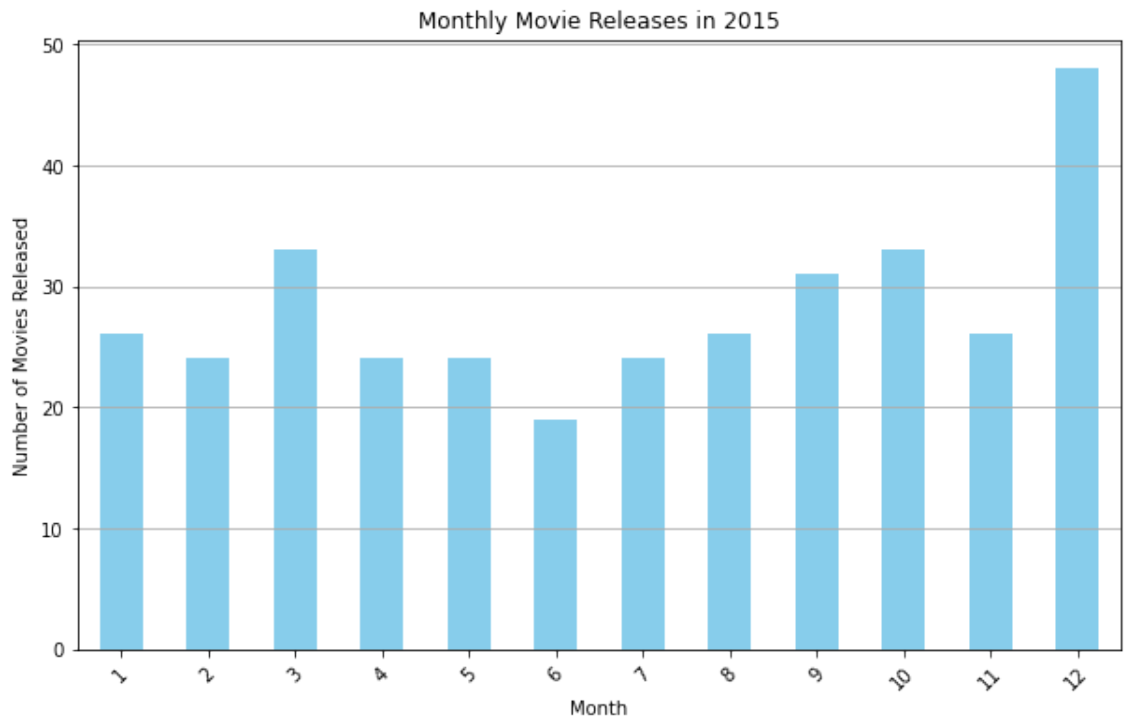


```
In [336]: 1 # Counting the number of movies released each year
2 movies_per_year = tn_movie_budgets['release_year'].value_counts()
3
4 # Filtering years with more than 50 releases
5 years_with_more_than_50_releases = movies_per_year[movies_per_year > 50]
6
7 # Filtering the DataFrame to include only those years
8 filtered_df = tn_movie_budgets[tn_movie_budgets['release_year'].isin(years_with_more_than_50_releases.index)]
```

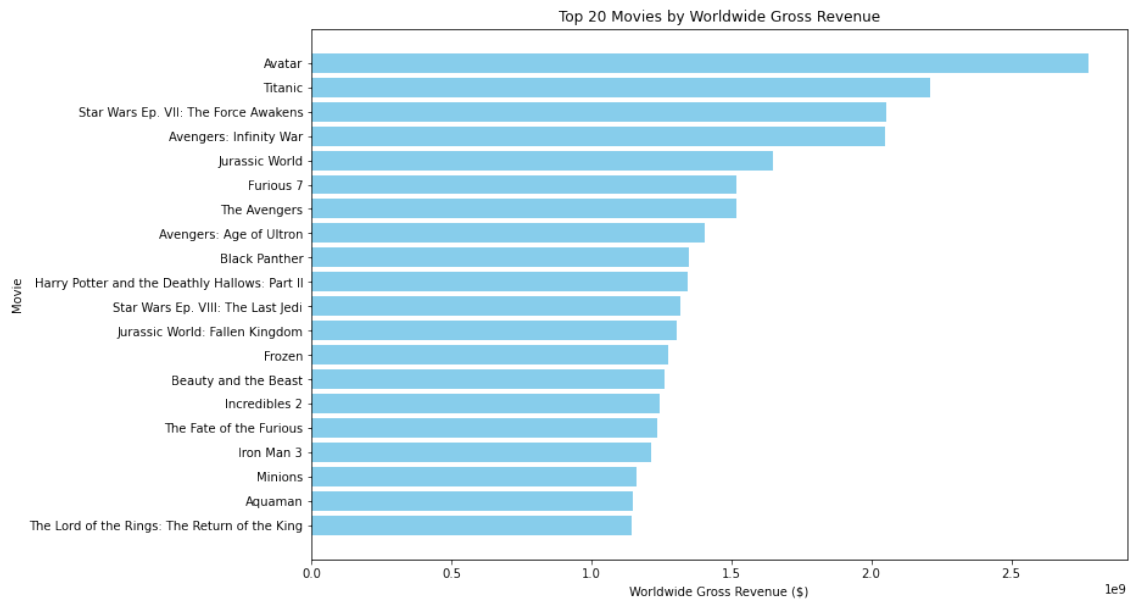
```
In [337]: 1
2 # Plotting number of movies released each year for years with more than 50 releases
3 plt.figure(figsize=(10, 6))
4 filtered_df['release_year'].value_counts().sort_index().plot(kind='bar')
5 plt.title('Number of Movies Released Each Year (Years with >50 Releases)')
6 plt.xlabel('Release Year')
7 plt.ylabel('Number of Movies')
8 plt.xticks(rotation=45)
9 plt.grid(axis='y')
10 plt.show()
```



```
In [341]: 1
2 # Find the year with the most releases
3 most_releases_year = monthly_counts.sum(axis=1).idxmax()
4
5 # Plotting monthly movie releases for the year with the most releases
6 plt.figure(figsize=(10, 6))
7 monthly_counts.loc[most_releases_year].plot(kind='bar', color='steelblue')
8 plt.title(f'Monthly Movie Releases in {most_releases_year}')
9 plt.xlabel('Month')
10 plt.ylabel('Number of Movies Released')
11 plt.xticks(rotation=45)
12 plt.grid(axis='y')
13 plt.show()
```



```
In [ ]: 1 # Sorting the DataFrame by worldwide gross revenue in descending
2 top_20_movies = tn_movie_budgets.nlargest(20, 'worldwide_gross')
3
4 # Creating a bar plot for the top 20 movies
5 plt.figure(figsize=(12, 8))
6 plt.barh(top_20_movies['movie'], top_20_movies['worldwide_gross'])
7 plt.xlabel('Worldwide Gross Revenue ($)')
8 plt.ylabel('Movie')
9 plt.title('Top 20 Movies by Worldwide Gross Revenue')
10 plt.gca().invert_yaxis() # Invert y-axis to display the highest
11 plt.show()
```



From the above analysis the Drama genre is a huge contributor to the movie production industry having a stake in almost all combination of movies produced.

As of 2015 the results show that movie have increased in production with most releases being after the 2nd half of the year.

From the World Gross Revenue it was noted that fiction movies have greatest impact and should be highly considered.