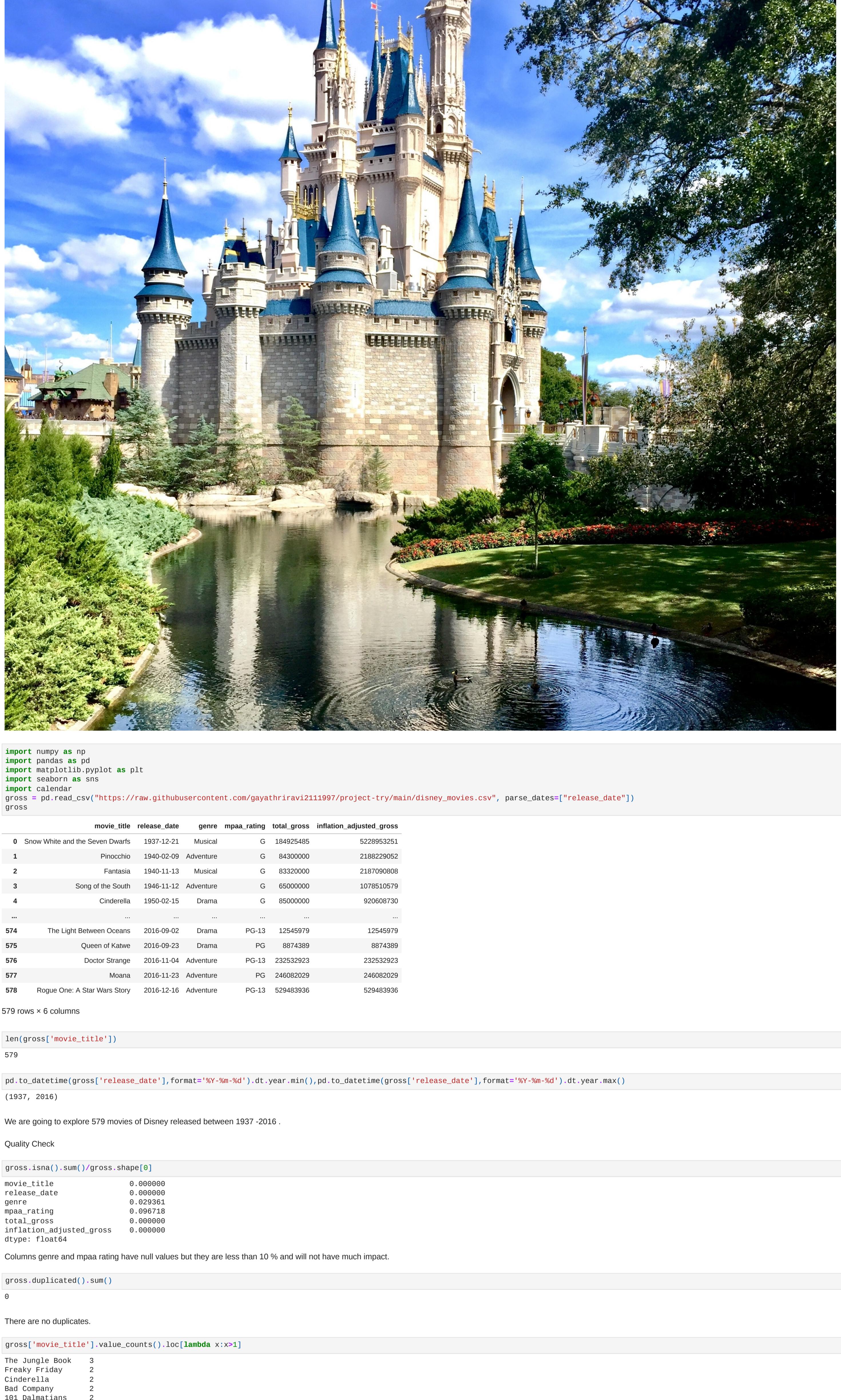


Disney movies and Box Office Success

1. The dataset

Walt Disney Studios is the foundation on which The Walt Disney Company was built. The Studios has produced more than 600 films since their debut film, Snow White and the Seven Dwarfs in 1937. While many of its films were big hits, some of them were not. In this notebook, we will explore a dataset of Disney movies and analyze what contributes to the success of Disney movies.



```
In [3]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import calendar
gross = pd.read_csv("https://raw.githubusercontent.com/gayathrirav1211997/project-try/main/disney_movies.csv", parse_dates=[“release_date”])
```

Out[3]: movie_title release_date genre mpaa_rating total_gross inflation_adjusted_gross

	movie_title	release_date	genre	mpaa_rating	total_gross	inflation_adjusted_gross
0	Snow White and the Seven Dwarfs	1937-12-21	Musical	G	184925485	5228953251
1	Pinocchio	1940-02-09	Adventure	G	84300000	2188229052
2	Fantasia	1940-11-13	Musical	G	83320000	218709008
3	Song of the South	1946-11-12	Adventure	G	65000000	1078510579
4	Cinderella	1950-02-15	Drama	G	85000000	920608730
...
574	The Light Between Oceans	2016-09-20	Drama	PG-13	12545979	12545979
575	Queen of Katwe	2016-09-23	Drama	PG	8874389	8874389
576	Doctor Strange	2016-11-04	Adventure	PG-13	232532923	232532923
577	Moana	2016-11-23	Adventure	PG	246082029	246082029
578	Rogue One: A Star Wars Story	2016-12-16	Adventure	PG-13	529483936	529483936

579 rows × 6 columns

```
In [5]: len(gross[“movie_title”])
```

579

```
In [8]: pd.to_datetime(gross[“release_date”],format=“%Y-%m-%d”).dt.year.min(),pd.to_datetime(gross[“release_date”],format=“%Y-%m-%d”).dt.year.max()
```

Out[8]: (1937, 2016)

We are going to explore 579 movies of Disney released between 1937-2016.

Quality Check

```
In [9]: gross.isna().sum()/gross.shape[0]
```

movie_title	0.000000
release_date	0.000000
mpaa_rating	0.000000
total_gross	0.000000
inflation_adjusted_gross	0.000000

dtype: float64

Columns genre and mpaa rating have null values but they are less than 10 % and will not have much impact.

```
In [11]: gross.duplicated().sum()
```

0

There are no duplicates.

```
In [12]: gross[“movie_title”].value_counts().loc[lambda x:x>1]
```

movie_title	count
The Jungle Book	3
Freaky Friday	2
Cinderella	2
Bedeviled	2
101 Dalmatians	2

Name: movie_title, dtype: int64

There are 5 movie's which have multiple entries. We need to be careful while summarizing.

Create Features

```
In [13]: gross[“release_year”]=pd.to_datetime(gross[“release_date”]).dt.year
gross[“release_month”]=pd.to_datetime(gross[“release_date”]).dt.month.apply(lambda x:calendar.month_abbr[x])
```

```
In [14]: gross.head()
```

Out[14]: movie_title release_date genre mpaa_rating total_gross inflation_adjusted_gross release_year release_month

	movie_title	release_date	genre	mpaa_rating	total_gross	inflation_adjusted_gross	release_year	release_month
0	Snow White and the Seven Dwarfs	1937-12-21	Musical	G	184925485	5228953251	1937	Dec
1	Pinocchio	1940-02-09	Adventure	G	84300000	2188229052	1940	Feb
2	Fantasia	1940-11-13	Musical	G	83320000	218709008	1940	Nov
3	Song of the South	1946-11-12	Adventure	G	65000000	1078510579	1946	Nov
4	Cinderella	1950-02-15	Drama	G	85000000	920608730	1950	Feb

1. Top ten movies at the box office

Let's started by exploring the data. We will check which are the 10 Disney movies that have earned the most at the box office. We can do this by sorting movies by their inflation-adjusted gross (we will call it adjusted gross from this point onward).

```
In [16]: gross.sort_values(by=“inflation_adjusted_gross”, ascending=False).head(10)
```

Out[16]: movie_title release_date genre mpaa_rating total_gross inflation_adjusted_gross release_year release_month

	movie_title	release_date	genre	mpaa_rating	total_gross	inflation_adjusted_gross	release_year	release_month
0	Snow White and the Seven Dwarfs	1937-12-21	Musical	G	184925485	5228953251	1937	Dec
1	Pinocchio	1940-02-09	Adventure	G	84300000	2188229052	1940	Feb
2	Fantasia	1940-11-13	Musical	G	83320000	218709008	1940	Nov
3	Song of the South	1946-11-12	Adventure	G	65000000	1078510579	1946	Nov
4	Cinderella	1950-02-15	Drama	G	85000000	920608730	1950	Feb

From the top 10 movies above, it seems that some genres are more popular than others. So, we will check which genres are growing stronger in popularity. To do this, we will group movies by genre and then by year to see the adjusted gross of each genre in each year.

```
In [17]: gross[“release_year”]=pd.datetimeIndex(gross[“release_date”]).year
group=gross.groupby([“genre”, “release_year”]).mean()
```

genre_yearly = group.reset_index()

genre_yearly.head(10)

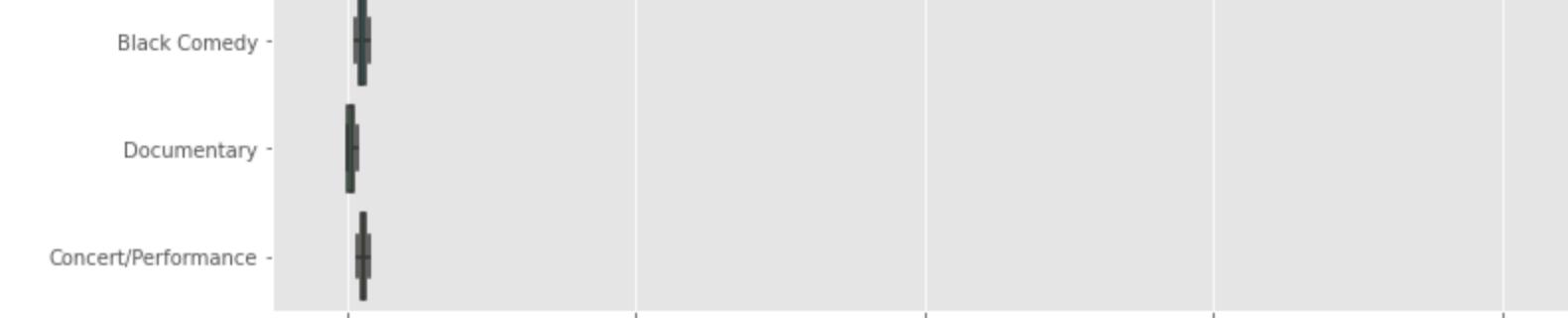
```
Out[17]: genre release_year total_gross inflation_adjusted_gross
```

genre	release_year	total_gross	inflation_adjusted_gross
Action	1981	0.0	0.0
Action	1982	20918576.0	7718495.0
Action	1983	17577696.0	36053517.0
Action	1984	21943553.5	44682157.0
Action	1985	19180582.0	39545796.0
Action	1986	63037553.5	122162426.5
Action	1987	135291096.0	257755262.5

1. Visualize the genre popularity trend

We will make a plot out of these means of groups to better see how box office revenues have changed over time.

```
In [18]: import matplotlib.pyplot as plt
import seaborn as sns
fig, ax=plt.subplots(figsize=(14,10))
plt.style.use(‘ggplot’)
sns.lineplot(x=“release_year”, y=“inflation_adjusted_gross”, data=genre_yearly, hue=“genre”)
plt.xlabel(“Release Year”)
plt.ylabel(“Frequency”)
plt.show()
```



1. Data transformation

The line plot supports our belief that some genres are growing faster in popularity than others. For Disney movies, Action and Adventure genres are growing the fastest. Next, we will build a linear regression model to understand the relationship between genre and box office gross. Since linear regression requires numerical variables and the genre variable is a categorical variable, we'll use a technique called one-hot encoding to convert the categorical variables to numerical. This technique transforms each category value into a new column and assigns a 1 or 0 to the column. For this dataset, there will be 11 dummy variables, one for each genre except the action genre which we will use as a baseline. For example, if a movie is an adventure movie, like The Lion King, the adventure variable will be 1 and other dummy variables will be 0. Since the action genre is our baseline, if a movie is an action movie, such as The Avengers, all dummy variables will be 0.

```
In [19]: genre_dummies = pd.get_dummies(gross[“genre”], dropFirst=True)
genre_dummies.head()
```

Out[19]: Adventure Black Comedy Comedy ConcertPerformance Documentary Drama Horror Musical Romantic Comedy Thriller/Suspense Western

0 0 0 0 0 0 0 0 0 0 0

1 1 0 0 0 0 0 0 0 0 0

2 0 0 0 0 0 0 0 0 0 0

3 1 0 0 0 0 0 0 0 0 0

4 0 0 0 0 0 0 0 0 0 0

5 0 0 0 0 0 0 0 0 0 0

6 0 0 0 0 0 0 0 0 0 0

7 0 0 0 0 0 0 0 0 0 0

8 0 0 0 0 0 0 0 0 0 0

9 0 0 0 0 0 0 0 0 0 0

10 0 0 0 0 0 0 0 0 0 0

11 0 0 0 0 0 0 0 0 0 0

12 0 0 0 0 0 0 0 0 0 0

13 0 0 0 0 0 0 0 0 0 0

14 0 0 0 0 0 0 0 0 0 0

15 0 0 0 0 0 0 0 0 0 0

16 0 0 0 0 0 0 0 0 0 0

17 0 0 0 0 0 0 0 0 0 0

18 0 0 0 0 0 0 0 0 0 0

19 0 0 0 0 0 0 0 0 0 0

20 0 0 0 0 0 0 0 0 0 0

21 0 0 0 0 0 0 0 0 0 0

22 0 0 0 0 0 0 0 0 0 0

23 0 0 0 0 0 0 0 0 0 0

24 0 0 0 0 0 0 0 0 0 0

25 0 0 0 0 0 0 0 0 0 0

26 0 0 0 0 0 0 0 0 0 0

27 0 0 0 0 0 0 0 0 0 0

28 0 0 0 0 0 0 0 0 0 0

29 0 0 0 0 0 0 0 0 0 0

30 0 0 0 0 0 0 0 0 0 0

31 0 0 0 0 0 0 0 0 0 0

32 0 0 0 0 0 0 0 0 0 0

33 0 0 0 0 0 0 0 0 0 0

34 0 0 0 0 0 0 0 0 0 0

35 0 0 0 0 0 0 0 0 0 0

36 0 0 0 0 0 0 0 0 0 0

37 0 0 0 0 0 0 0 0 0 0

38 0 0 0 0 0 0 0 0 0 0

39 0 0 0 0 0 0 0 0 0 0

40 0 0 0 0 0 0 0 0 0 0

41 0 0 0 0 0 0 0 0 0 0

42 0 0 0 0 0 0 0 0 0 0

43 0 0 0 0 0 0 0 0 0 0

44 0