

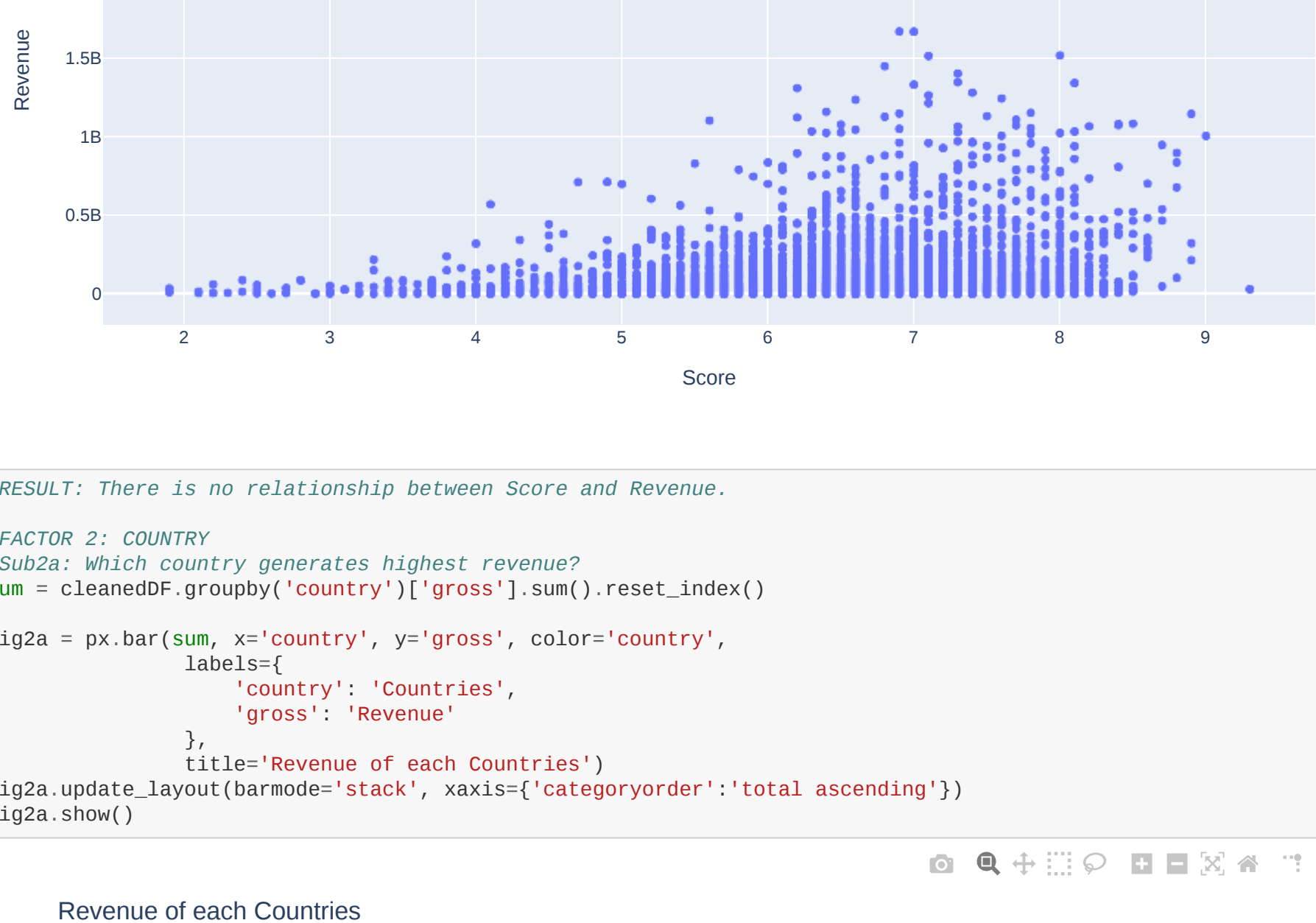
```
In [1]: #https://www.kaggle.com/datasets/danielgrijalvas/movies
#Project #3, LinhPhan, Nov 28
import pandas as pd
import plotly.express as px

#import csv file
rawdf = pd.read_csv('movies.csv')

#clean data
cleanedDF = rawdf.dropna()

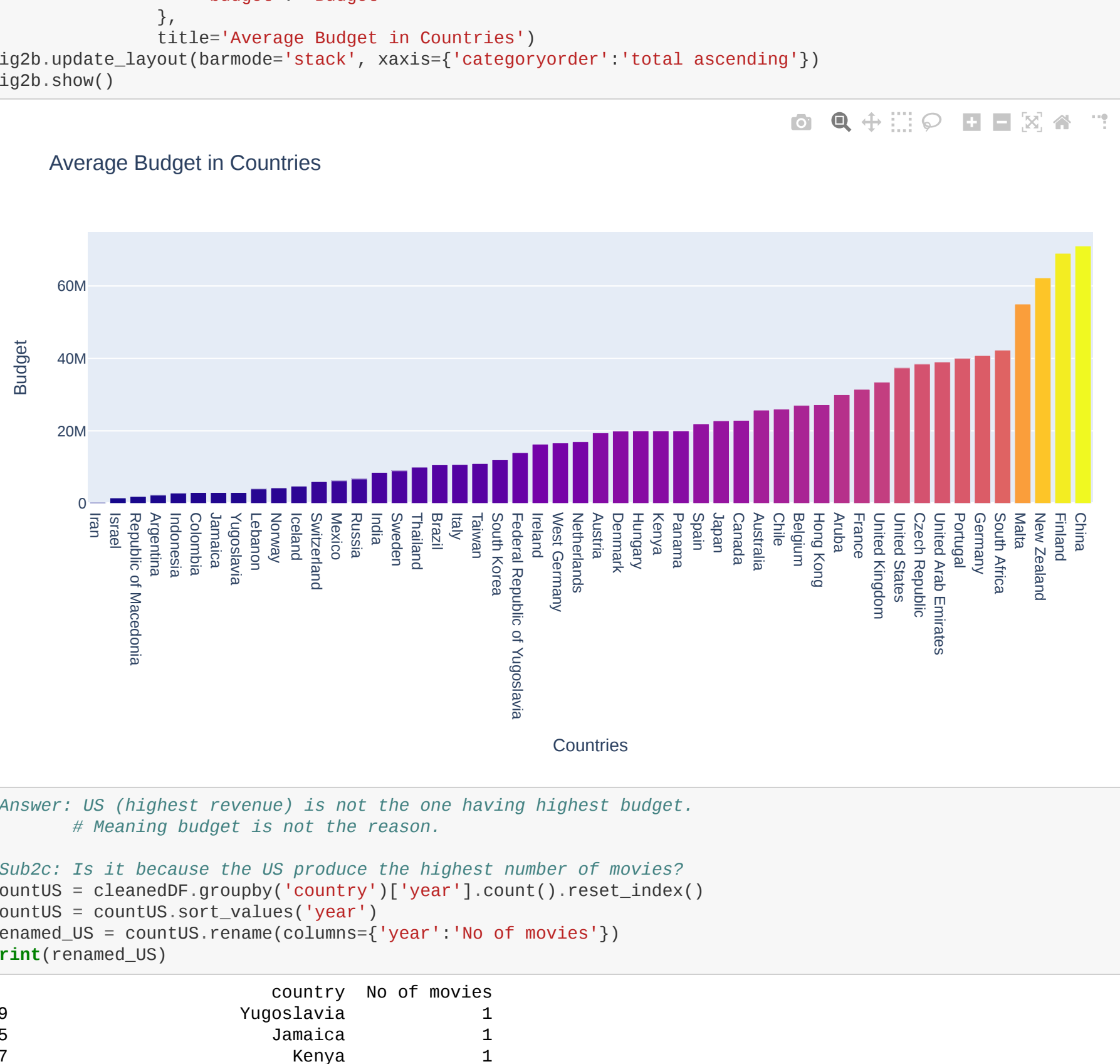
#MAIN QUESTION: How do Score, Country, Budget, or Genre affect Revenue?

#FACTOR 1: SCORE
#Sub1: Does movies with higher score have higher revenue?
fig1 = px.scatter(cleanedDF, x='score', y='gross',
                 labels={
                     'score': 'Score',
                     'gross': 'Revenue'
                 },
                 title='Gross Revenue versus Score')
fig1.show()
```



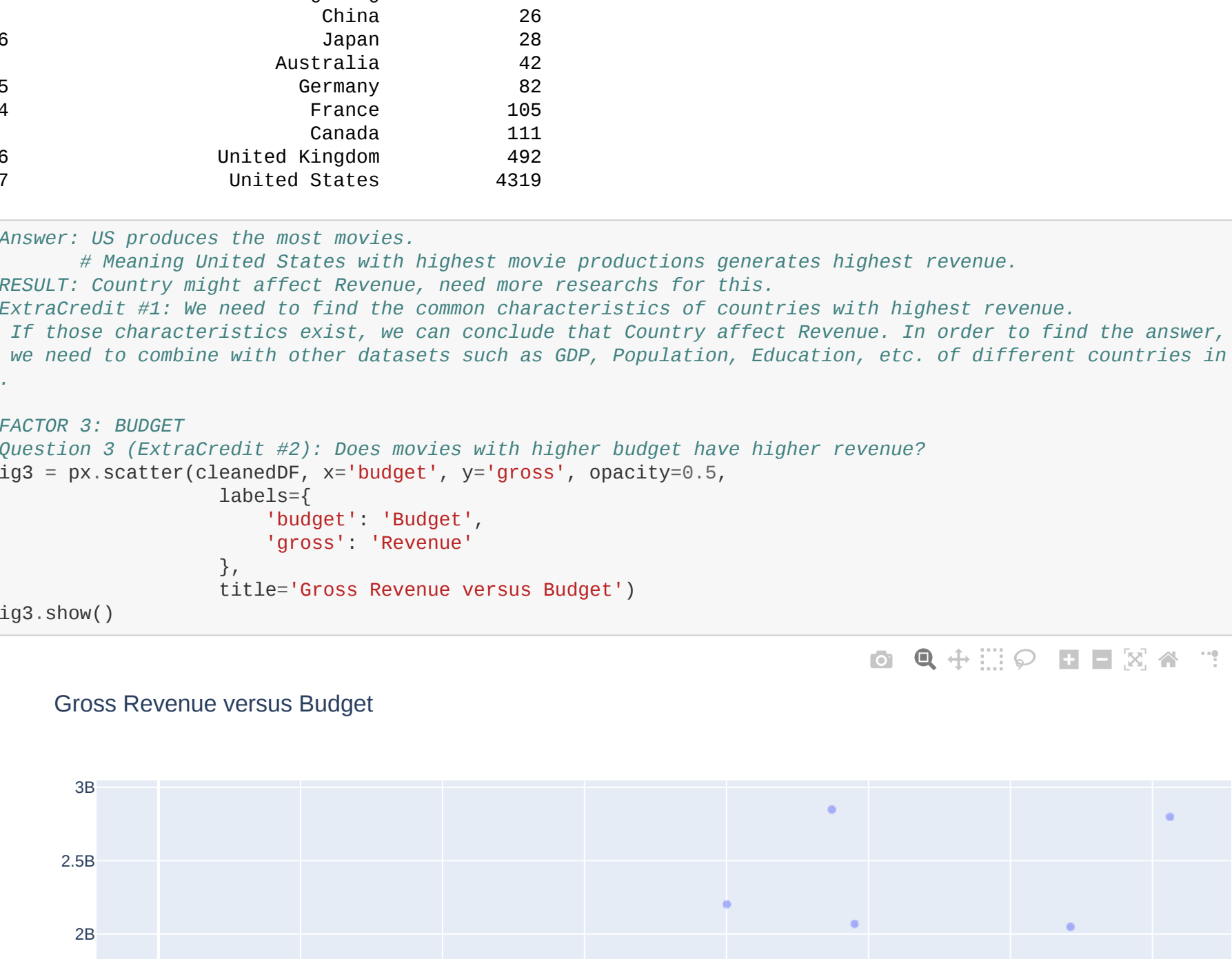
```
In [2]: #RESULT: There is no relationship between Score and Revenue.

#FACTOR 2: COUNTRY
#Sub2a: Which country generates highest revenue?
sum = cleanedDF.groupby('country')['gross'].sum().reset_index()
fig2a = px.bar(sum, x='country', y='gross', color='country',
              labels={
                  'country': 'Countries',
                  'gross': 'Revenue'
              },
              title='Revenue of each Countries')
fig2a.update_layout(barmode='stack', xaxis='categoryorder':'total ascending')
fig2a.show()
```



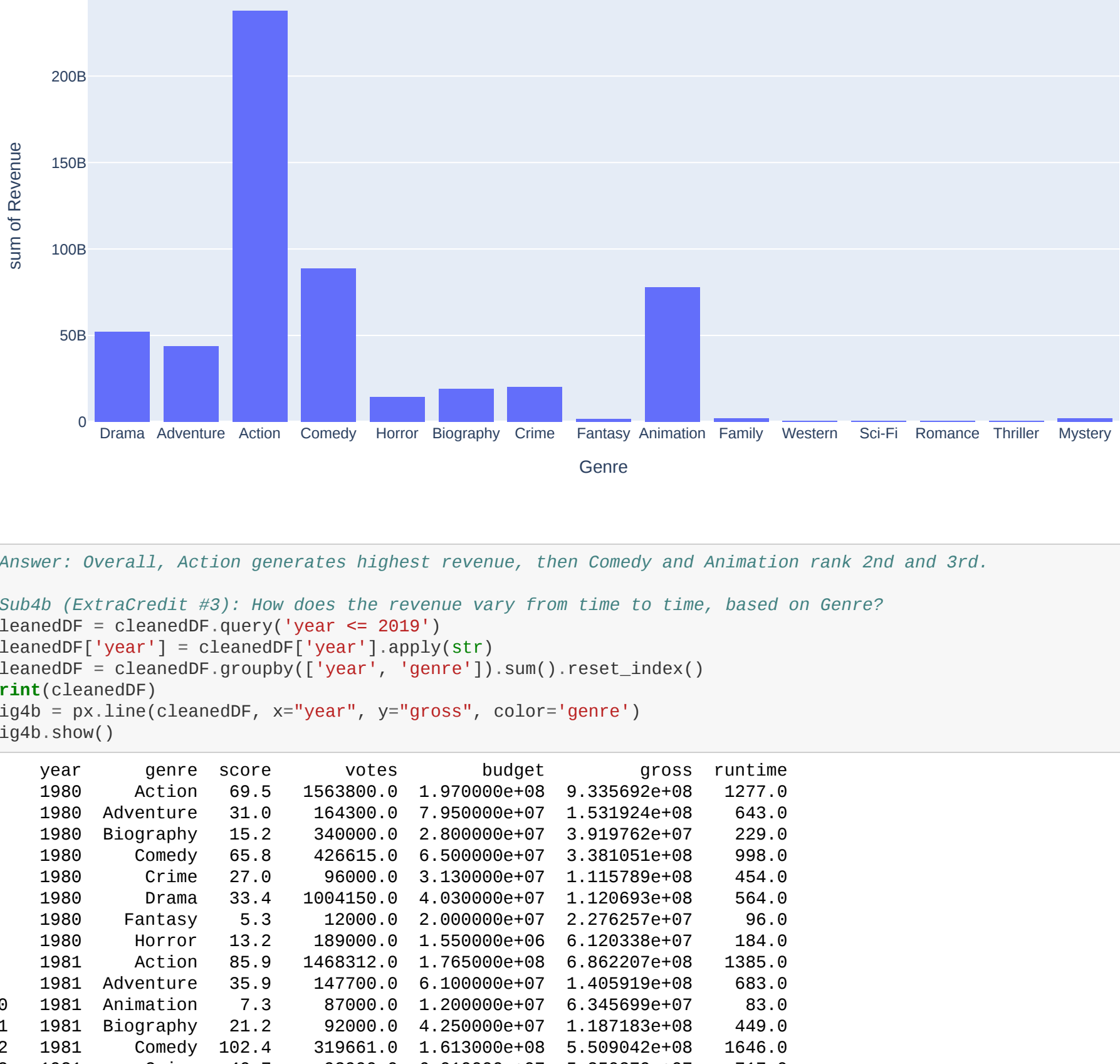
```
In [3]: #Answer: United States generates highest revenue.

#Sub2b: Is it because the US having the highest budget?
average = cleanedDF.groupby('country').mean().reset_index()
average = cleanedDF.groupby('country')['budget'].mean().reset_index()
fig2b = px.bar(average, x='country', y='budget', color='budget',
              labels={
                  'country': 'Countries',
                  'budget': 'Budget'
              },
              title='Average Budget in Countries')
fig2b.update_layout(barmode='stack', xaxis='categoryorder':'total ascending')
fig2b.show()
```



```
In [4]: #Answer: US (highest revenue) is not the one having highest budget.
# Meaning budget is not the reason.

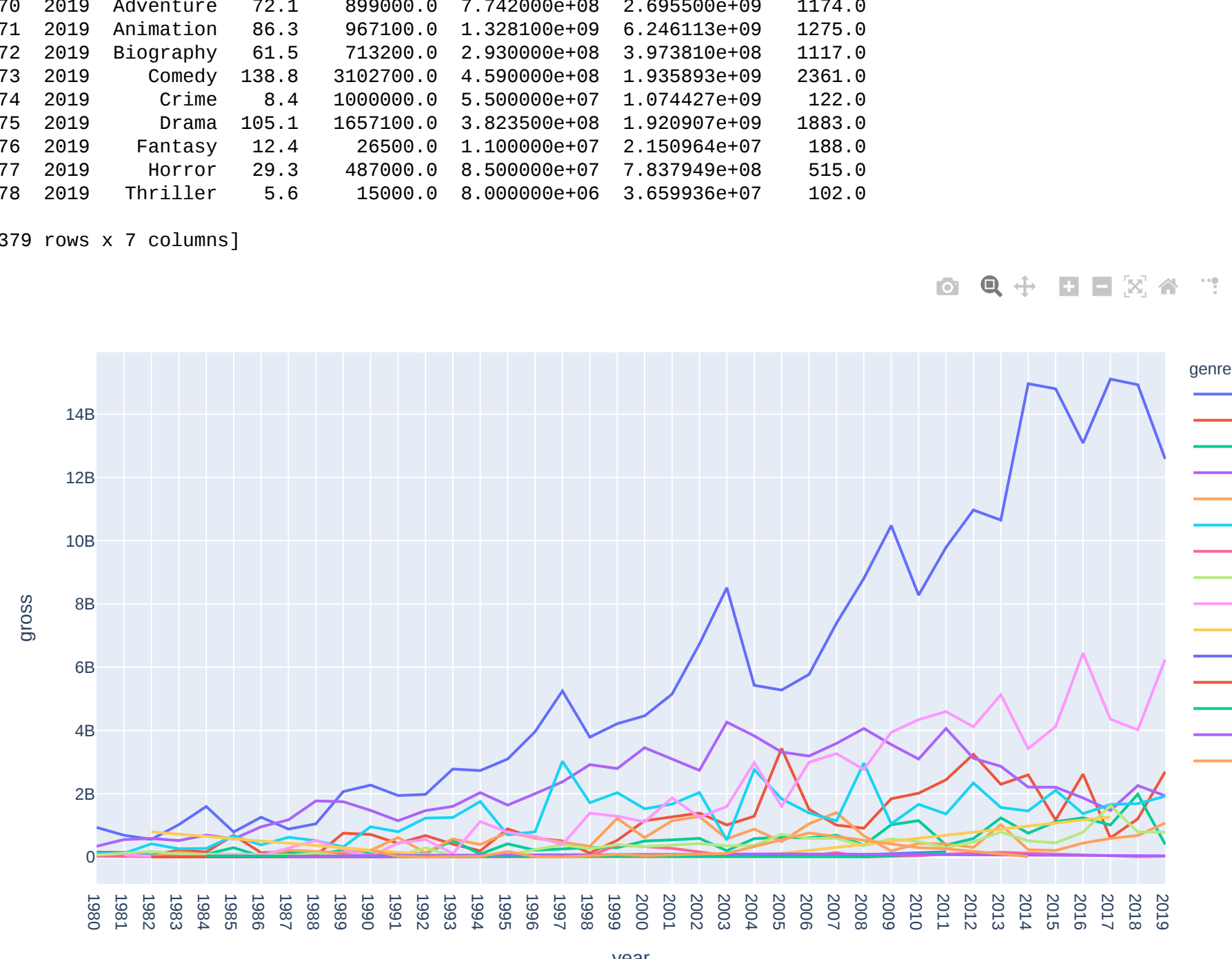
#Sub2c: Is it because the US produce the highest number of movies?
countUS = cleanedDF.groupby('country')['year'].count().reset_index()
countUS = countUS.sort_values('year')
renamed_US = countUS.rename(columns={'year':'No of movies'})
print(renamed_US)
```



```
In [5]: #Answer: US produces the most movies.
# # Meaning United States with highest movie productions generates highest revenue.

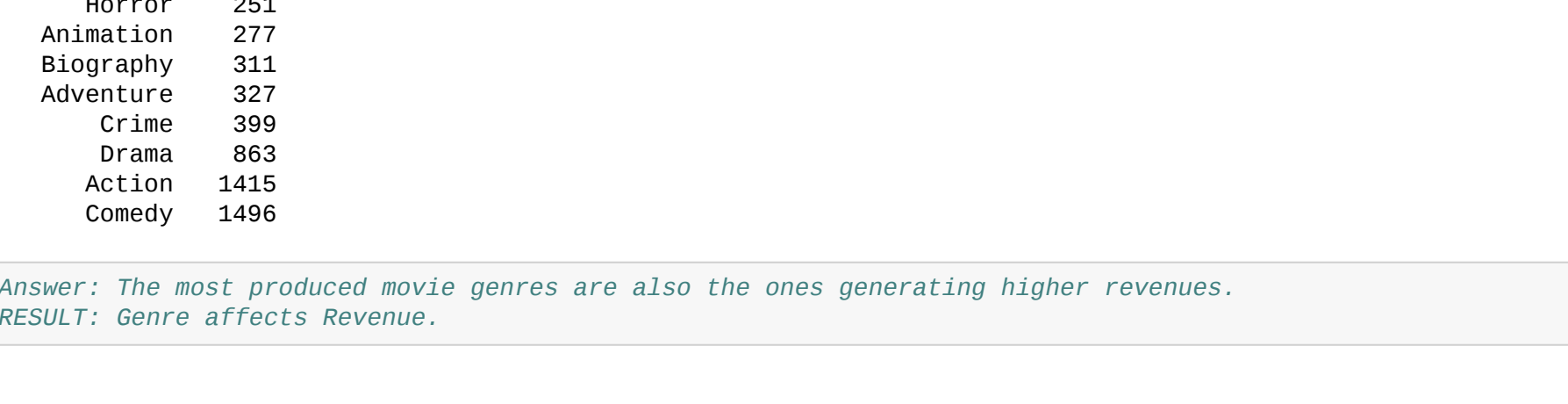
#RESULT: Country might affect Revenue, need more researchs for this.
#ExtraCredit #1: We need to find the common characteristics of countries with highest revenue.
# If those characteristics exist, we can conclude that Country affect Revenue. In order to find the answer,
# we need to combine with other datasets such as GDP, Population, Education, etc. of different countries in the world.

#FACTOR 3: BUDGET
#Question 3 (ExtraCredit #2): Does movies with higher budget have higher revenue?
fig3 = px.scatter(cleanedDF, x='budget', y='gross', opacity=0.5,
                 labels={
                     'budget': 'Budget',
                     'gross': 'Revenue'
                 },
                 title='Gross Revenue versus Budget')
fig3.show()
```



```
In [6]: #Answer: Movies with higher budget generate higher revenue while movies with lower budget generate less revenue,
# With a few exceptions (Outliers).
#RESULT: Budget affects Revenue.

#FACTOR 4: GENRE
#Sub4a: Gross revenue for each genre of movies?
fig4a = px.histogram(cleanedDF, x='genre', y='gross',
                    labels={
                        'genre': 'Genre',
                        'gross': 'Revenue'
                    },
                    title='Gross Revenue based on Genre of movies')
fig4a.show()
```



```
In [8]: #Answer: Overall, Action generates highest revenue, then Comedy and Animation rank 2nd and 3rd.

#Sub4b (ExtraCredit #3): How does the revenue vary from time to time, based on Genre?
cleanedDF = cleanedDF.query('year <= 2019')
cleanedDF['year'] = cleanedDF['year'].apply(str)
cleanedDF = cleanedDF.groupby(['year', 'genre']).sum().reset_index()
print(cleanedDF)
fig4b = px.line(cleanedDF, x='year', y='gross', color='genre')
fig4b.show()
```



```
In [7]: #Answer: Revenue of each Genre fluctuates from 1980 to 2019, with Action generates highest revenue most of the year
# followed by Animation and Comedy.

#Sub2c: How many movies produced in each kind of genres?
countGE = cleanedDF.groupby('genre')['year'].count().reset_index()
countGE = countGE.sort_values('year')
renamed_GE = countGE.rename(columns={'year':'count'})
print(renamed_GE)
```



```
In [ ]: #Answer: The most produced movie genres are also the ones generating higher revenues.
#RESULT: Genre affects Revenue.
```