



RECOMMENDING MICROSOFT MOVIES TO MAKE

Authors: DANIEL WAHOME

Overview

A one-paragraph overview of the project, including the business problem, data, methods, results and recommendations.

Microsoft is planning to create a movie production company and would like to get recommendations on the type of films that the type of films to produce according to the gross trending movies in the box office. Using Python and the provided box office data I will separate the works per studio see what they have done and what has worked for them the methods I have used are data exploration, dealt with missing values, explored categorical data, identified outliers, performed correlation, and visualized the data

the results I found are that films that people tend to love according to gross sales

Boyhood- Coming of Age story

Heartbreaker- comedy drama

Jurassic World Universe -Sci-fi Adventure Supported by a sequel to the first Jurassic Park boosting sales in the same year

Dark Knight Rises- Action

Harry Potter and the Deathly Hallows Part 2 - Fantasy Adventure also boosted by a prequel the part one of it

Business Problem

Microsoft is planning to create a movie production studio and would like to get recommendations on the type of films that the type of films to produce. They have no clue of where to begin and would like to utilize data from the gross trending movies in the box office to influence the kind of movies they should make

I chose to ask the following question; according to the top studio houses : which film has the highest international grossing

domestic grossing

which film has the highest

belong to

which genre do these films

released

what year were these films

this will help them come up with genres that they can start producing on. as they are new to the business and want a large amount of people to view their films for commercial success, following the mass produced films and genres will put them in the billboards

Data Understanding

Describe the data being used for this project.

the data being used is the box_office data as it will help us to see what the domestic and international viewers like to see, then from there we can emulate

the data consists of the title, studio name, domestic gross ammount, foreign gross ammount and the year the film premiered

i intend to categorise the data accoring to the studios and find out their properties

```
In [1]: # Import standard packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import os
import sqlite3
%matplotlib inline
```

```
In [2]: #Loading all csv files to use for analysis
box_office = pd.read_csv(r'C:\Users\DANIEL\Desktop\PROJECT 1 ASSETS\DATA\bom.movie_gross.csv')
tmdb = pd.read_csv(r'C:\Users\DANIEL\Desktop\PROJECT 1 ASSETS\DATA\tmdb.movies.csv')
the_numbers = pd.read_csv(r'C:\Users\DANIEL\Desktop\PROJECT 1 ASSETS\DATA\tn.movie_budgets.csv')
```

Begin with the Box Office data to understand it clean it and get it ready for manipulation and visualization

```
In [3]: box_office
```

```
Out[3]:
```

	title	studio	domestic_gross	foreign_gross	year
0	Toy Story 3	BV	415000000.0	652000000.0	2010
1	Alice in Wonderland (2010)	BV	334200000.0	691300000.0	2010
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000.0	2010
3	Inception	WB	292600000.0	535700000.0	2010
4	Shrek Forever After	P/DW	238700000.0	513900000.0	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 [4]: #checking out the different data types in Box Office csv file
box_office.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  float64
4  year            3387 non-null  int64
dtypes: float64(2), int64(1), object(2)
memory usage: 132.4+ KB

```

```

In [5]: #a brief statistical description of the data
        box_office.describe()

```

```

Out[5]:

```

	domestic_gross	foreign_gross	year
count	3.359000e+03	2.037000e+03	3387.000000
mean	2.874585e+07	7.503550e+07	2013.958075
std	6.698250e+07	1.373874e+08	2.478141
min	1.000000e+02	6.000000e+02	2010.000000
25%	1.200000e+05	3.700000e+06	2012.000000
50%	1.400000e+06	1.900000e+07	2014.000000
75%	2.790000e+07	7.550000e+07	2016.000000
max	9.367000e+08	9.605000e+08	2018.000000

```

In [6]: #getting info on any missing data. this way i can identify how to make it uniform and usable
        #then i can use it to identify outliers and do the necessary changes according to the findings
        box_office.isnull().sum()

```

```

Out[6]: title          0
        studio         5
        domestic_gross  28
        foreign_gross   1350
        year           0
        dtype: int64

```

DROPPING AND FILLING IN MISSING VALUES

```

In [7]: box_office = box_office.dropna(subset=['studio'])
        box_office
        #i have chosen to drop the rows in the studio column that were missing
        #this is because the data in there cannot be extropolated using mean or mode as it is not an int
        #if i physically added new data(random) it wouldnt affect the analysis greatly
        #so i chose to frop them

```

```

Out[7]:

```

	title	studio	domestic_gross	foreign_gross	year
0	Toy Story 3	BV	415000000.0	652000000.0	2010
1	Alice in Wonderland (2010)	BV	334200000.0	691300000.0	2010
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000.0	2010
3	Inception	WB	292600000.0	535700000.0	2010
4	Shrek Forever After	P/DW	238700000.0	513900000.0	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

3382 rows × 5 columns

TREATING MISSING VALUES USING MEDIAN AS THE DATA COULD CONTAIN OUTLIERS MAKING IT SKEWED

```
In [8]: box_office['domestic_gross'] = box_office['domestic_gross'].fillna(box_office['domestic_gross'].median())
box_office
box_office['foreign_gross'] = box_office['foreign_gross'].fillna(box_office['foreign_gross'].median())
box_office
```

<ipython-input-8-e1a4cc905a53>:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
box_office['domestic_gross'] = box_office['domestic_gross'].fillna(box_office['domestic_gross'].median())
```

<ipython-input-8-e1a4cc905a53>:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

```
box_office['foreign_gross'] = box_office['foreign_gross'].fillna(box_office['foreign_gross'].median())
```

```
Out[8]:
```

	title	studio	domestic_gross	foreign_gross	year
0	Toy Story 3	BV	415000000.0	652000000.0	2010
1	Alice in Wonderland (2010)	BV	334200000.0	691300000.0	2010
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000.0	2010
3	Inception	WB	292600000.0	535700000.0	2010
4	Shrek Forever After	P/DW	238700000.0	513900000.0	2010
...
3382	The Quake	Magn.	6200.0	19000000.0	2018
3383	Edward II (2018 re-release)	FM	4800.0	19000000.0	2018
3384	El Pacto	Sony	2500.0	19000000.0	2018
3385	The Swan	Synergetic	2400.0	19000000.0	2018
3386	An Actor Prepares	Grav.	1700.0	19000000.0	2018

3382 rows × 5 columns

```
In [9]: box_office.isnull().sum()
```

```
Out[9]: title          0
studio          0
domestic_gross  0
foreign_gross   0
year           0
dtype: int64
```

```
In [10]: box_office.duplicated().sum()
```

```
Out[10]: 0
```

```
In [11]: box_office['studio'].value_counts().head(15)
```

```
Out[11]: IFC          166
Uni.            147
WB              140
Magn.           136
Fox             136
SPC             123
```

```

Sony      110
BV        106
LGF       103
Par.      101
Eros      89
Wein.     77
CL        74
Strand    68
FoxS      67
Name: studio, dtype: int64

```

```
In [12]: box_office['studio'].value_counts().tail(10)
```

```

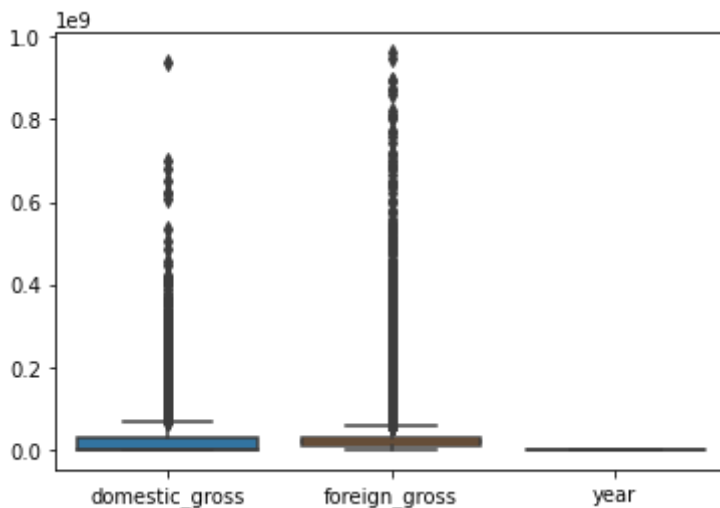
Out[12]: Arth.      1
Triu      1
RLJ       1
SMod      1
PBS       1
ParC      1
ELS       1
MPFT      1
MOM       1
DF        1
Name: studio, dtype: int64

```

USING SCATTER PLOT TO VISUALIZE OUTLIERS

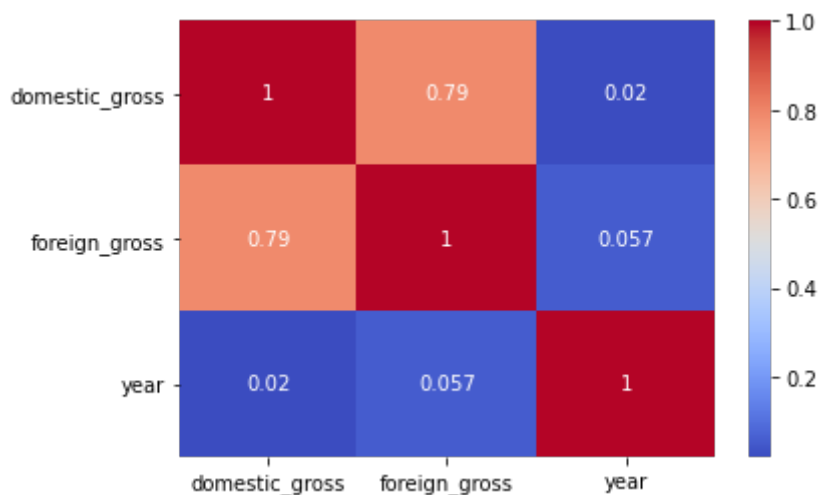
```
In [13]: sns.boxplot(data=box_office)
#i decided not to drop the outliers as nothing is out proportional range
```

```
Out[13]: <AxesSubplot:>
```



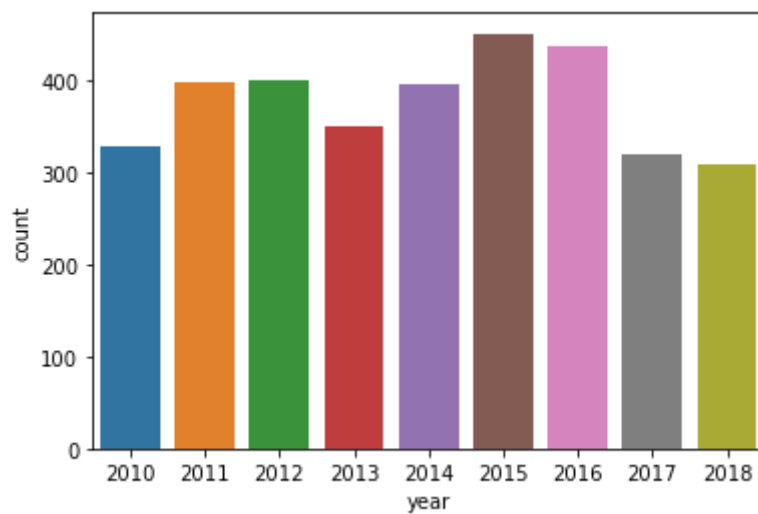
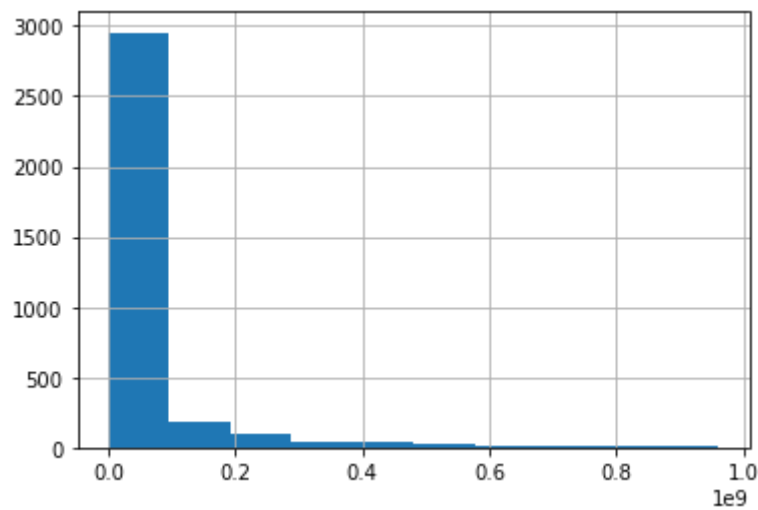
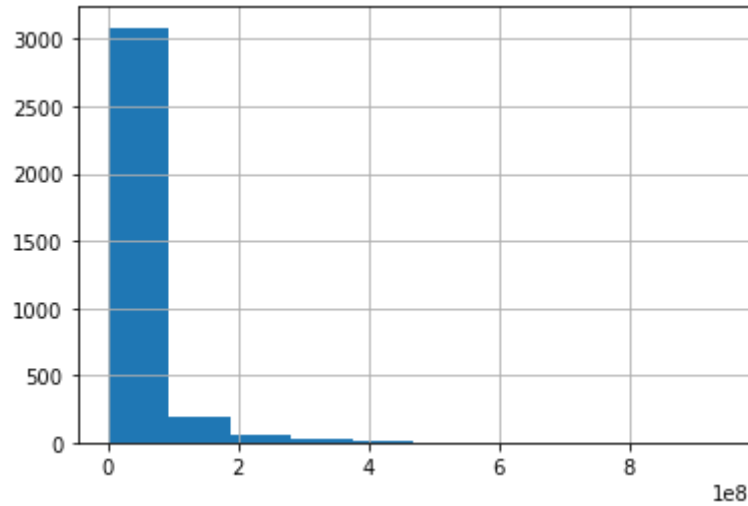
```
In [14]: correlation_matrix = box_office.corr()

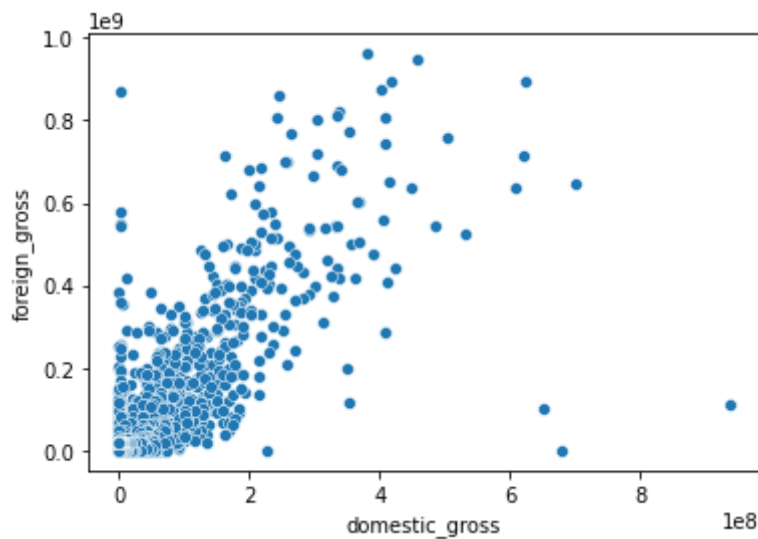
# Visualize correlation matrix
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
plt.show()
```



THERE IS A POSITIVE CORRELATION BETWEEN DOMESTIC AND FOREIGN GROSS AMOUNTS OF 0.79 THERE IS LITTLE TO NO CORRELATION BETWEEN YEAR AND

```
In [15]: box_office['domestic_gross'].hist()  
plt.show()  
box_office['foreign_gross'].hist()  
plt.show()  
sns.countplot(x='year', data=box_office)  
plt.show()  
sns.scatterplot(x='domestic_gross', y='foreign_gross', data=box_office)  
plt.show()
```





FOR THIS DATA SET: CATEGORIZATION BY STUDIO

1. IFC STUDIO

TITLES BY STUDIO IFC

THERE ARE 166 TITLES BY IFC STUDIO

THE HIGHEST DOMESTIC GROSSING FILM BY THEM IS "BOYHOOD"

GROSSING AT: 25400000.0

PREMIERED IN 2014

ACCORDING TO THE BAR GRAPH 2014 WAS THE YEAR WITH THE MOST DOMESTIC GROSS INCOME

THE HIGHEST FOREIGN GROSSING FILM IS "HEARTBREAKER"

GROSSING AT: 46900000.0

PREMIERED IN 2010

ACCORDING TO THE BAR GRAPH 2010 WAS THE YEAR WITH THE MOST FOREIGN GROSS INCOME

```
In [16]: studio1 = box_office.loc[box_office['studio']=='IFC']
studio1
```

```
Out[16]:
```

	title	studio	domestic_gross	foreign_gross	year
110	Heartbreaker	IFC	504000.0	46900000.0	2010
112	The Good, the Bad, the Weird	IFC	128000.0	44100000.0	2010
151	Soul Kitchen	IFC	277000.0	17600000.0	2010
166	Looking for Eric	IFC	55800.0	11500000.0	2010
190	Vincere	IFC	619000.0	5100000.0	2010
...
3324	Ghost Stories	IFC	149000.0	19000000.0	2018
3335	Mary Shelley	IFC	109000.0	19000000.0	2018

	title	studio	domestic_gross	foreign_gross	year
3344	The House That Jack Built	IFC	88000.0	19000000.0	2018
3361	A Ciambra	IFC	41900.0	19000000.0	2018
3374	The Escape	IFC	14000.0	19000000.0	2018

166 rows × 5 columns

```
In [17]: highest_gross = studio1['domestic_gross'].max()
print('the highest domestic grossing film by IFC studio is:', highest_gross)
```

the highest domestic grossing film by IFC studio is: 25400000.0

```
In [18]: target = 25400000.0
row = studio1[studio1['domestic_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

25400000.0 Boyhood

```
In [19]: target = 'Boyhood'
row = studio1.loc[studio1['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

Boyhood 2014

```
In [20]: type(studio1)
```

Out[20]: pandas.core.frame.DataFrame

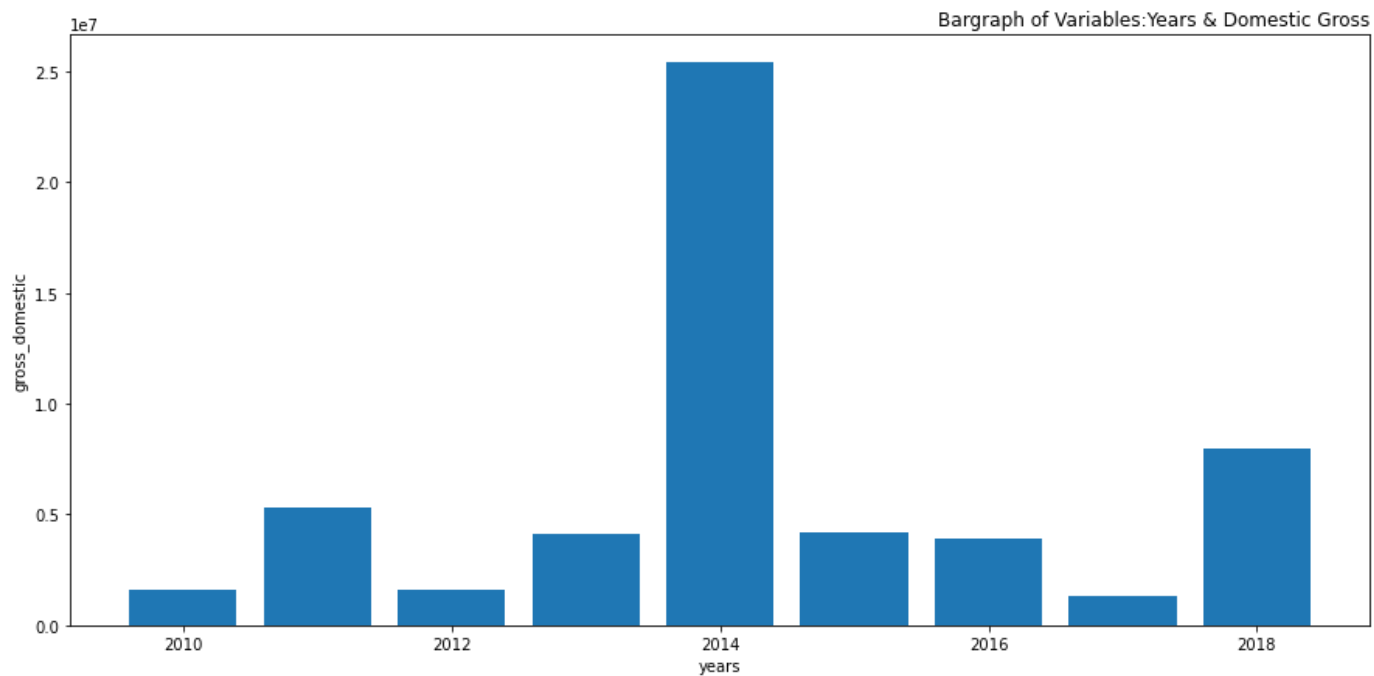
```
In [21]: studio1 = pd.DataFrame(studio1)

years = studio1['year']
grossd = studio1['domestic_gross']
grossf = studio1['foreign_gross']

fig = plt.figure(figsize= (15,7))

plt.bar(years, grossd)
plt.xlabel ('years')
plt.ylabel ('gross_domestic')
plt.title('Bargraph of Variables:Years & Domestic Gross', loc = "right")

plt.show()
```

In [22]:

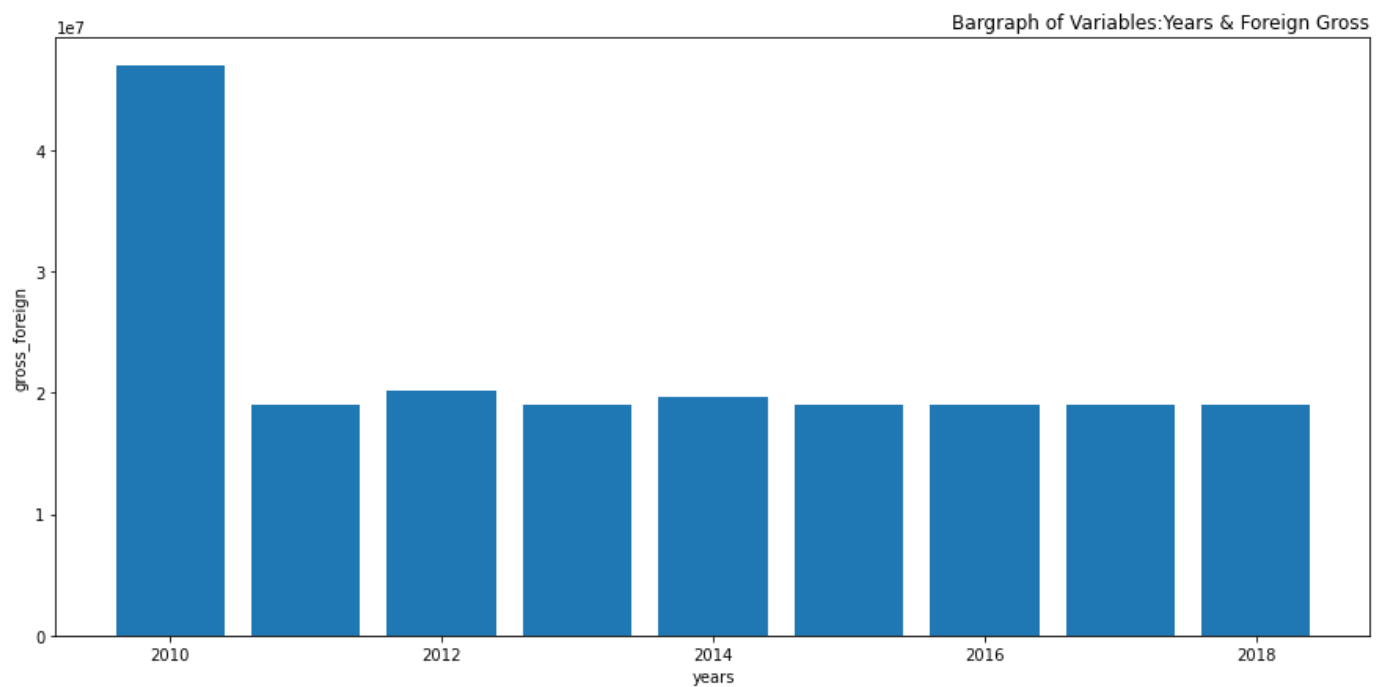
```
studio1 = pd.DataFrame(studio1)

years = studio1['year']
grossd = studio1['domestic_gross']
grossf = studio1['foreign_gross']

fig = plt.figure(figsize= (15,7))

plt.bar(years, grossf)
plt.xlabel ('years')
plt.ylabel('gross_foreign')
plt.title('Bargraph of Variables:Years & Foreign Gross', loc = "right")

plt.show()
```



In [23]:

```
highest_gross = studio1['foreign_gross'].max()
print('the highest domestic grossing film by IFC studio is:', highest_gross)
```

the highest domestic grossing film by IFC studio is: 46900000.0

```
In [24]: target = 46900000.0
row = studio1[studio1['foreign_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

46900000.0 Heartbreaker

```
In [25]: target = 'Heartbreaker'
row = studio1.loc[studio1['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

Heartbreaker 2010

2. UNIVERSAL STUDIOS

STUDIO UNIVERSAL

THERE ARE TITLES BY UNIVERSAL STUDIOS

THE HIGHEST DOMESTIC GROSSING FILM BY THEM IS "JURASSIC WORLD"

GROSSING AT: 652300000.0

PREMIERED IN 2015

ACCORDING TO THE BAR GRAPH 2015 WAS THE YEAR WITH THE MOST DOMESTIC GROSS INCOME

THE HIGHEST FOREIGN GROSSING FILM IS "JURASSIC WORLD: FALLEN KINGDOM"

GROSSING AT: 891800000.0

PREMIERED IN 2018

ACCORDING TO THE BAR GRAPH 2018 WAS THE YEAR WITH THE MOST FOREIGN GROSS INCOME

```
In [26]: studio2 = box_office.loc[box_office['studio']=='Uni.']
studio2
```

Out[26]:

	title	studio	domestic_gross	foreign_gross	year
8	Despicable Me	Uni.	251500000.0	291600000.0	2010
18	Robin Hood	Uni.	105300000.0	216400000.0	2010
20	Little Fockers	Uni.	148400000.0	162200000.0	2010
49	The Wolfman	Uni.	62000000.0	77800000.0	2010
66	Green Zone	Uni.	35100000.0	59800000.0	2010
...
3148	Mortal Engines	Uni.	16000000.0	67700000.0	2018
3172	Breaking In (2018)	Uni.	46800000.0	4600000.0	2018
3219	Welcome to Marwen	Uni.	10800000.0	2100000.0	2018

	title	studio	domestic_gross	foreign_gross	year
3289	Schindler's List (2018 re-release)	Uni.	833000.0	19000000.0	2018
3369	Loving Pablo	Uni.	22000.0	19000000.0	2018

147 rows × 5 columns

```
In [27]: highest_gross = studio2['domestic_gross'].max()
print('the highest domestic grossing film by Universal studio is:', highest_gross)
```

the highest domestic grossing film by Universal studio is: 652300000.0

```
In [28]: target = 652300000.0
row = studio2[studio2['domestic_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

652300000.0 Jurassic World

```
In [29]: target = 'Jurassic World'
row = studio2.loc[studio2['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

Jurassic World 2015

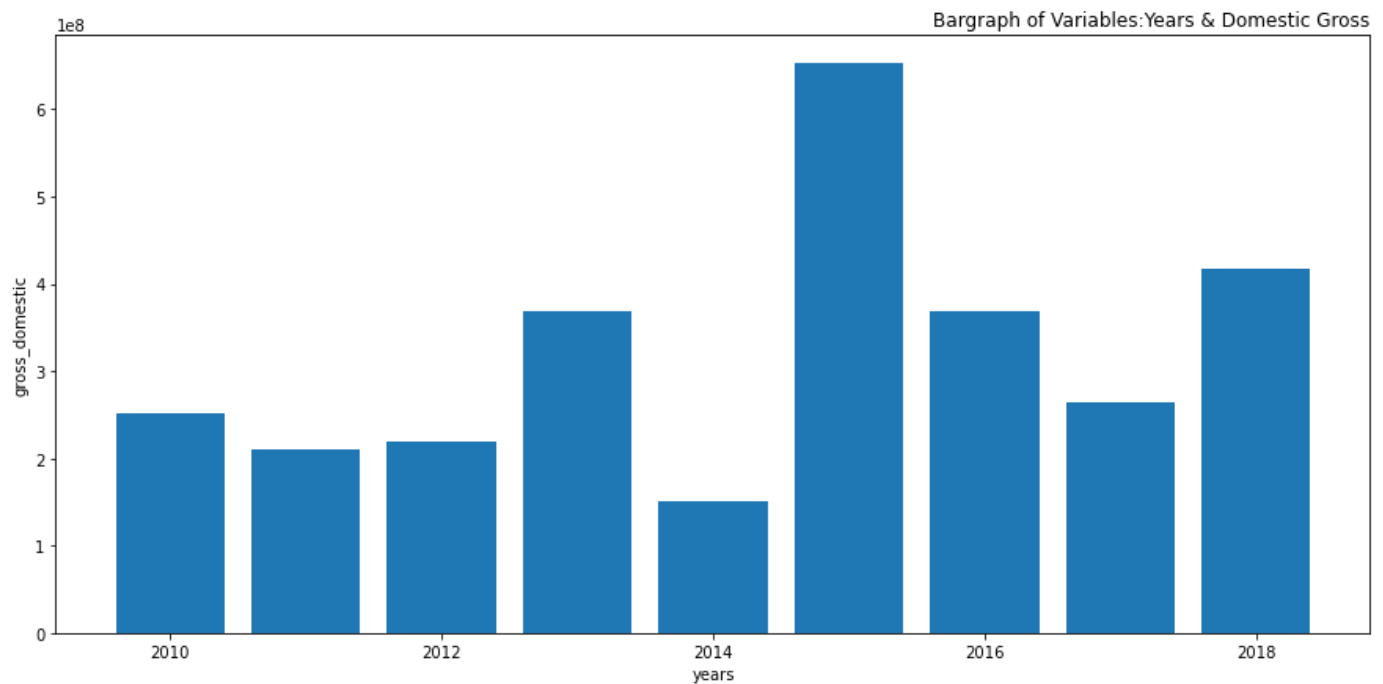
```
In [30]: studio2 = pd.DataFrame(studio2)

years = studio2['year']
grossd = studio2['domestic_gross']
grossf = studio2['foreign_gross']
```

```
fig = plt.figure(figsize= (15,7))
```

```
plt.bar(years, grossd)
plt.xlabel ('years')
plt.ylabel('gross_domestic')
plt.title('Bargraph of Variables:Years & Domestic Gross', loc = "right")
```

```
plt.show()
```



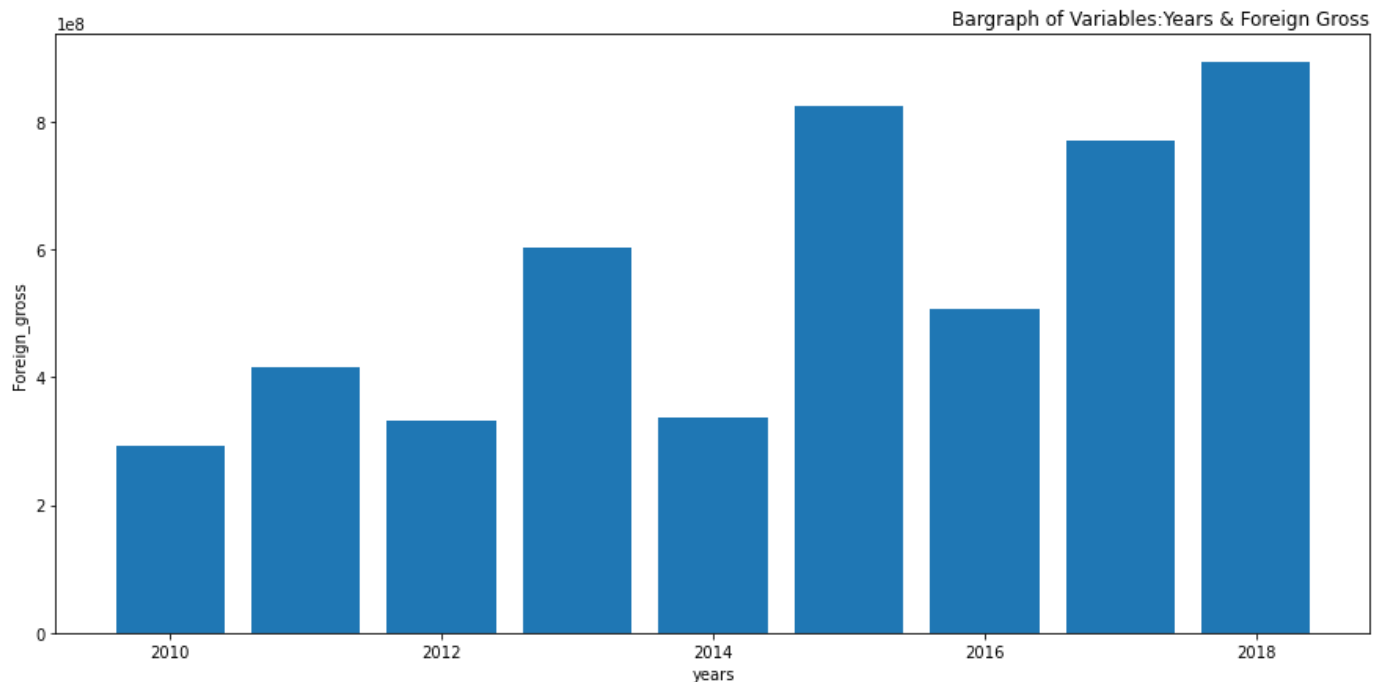
```
In [31]: studio2 = pd.DataFrame(studio2)

years = studio2['year']
grossd = studio2['domestic_gross']
grossf = studio2['foreign_gross']
```

```
fig = plt.figure(figsize= (15,7))
```

```
plt.bar(years, grossf)
plt.xlabel ('years')
plt.ylabel('Foreign_gross')
plt.title('Bargraph of Variables:Years & Foreign Gross', loc = "right")
```

```
plt.show()
```



```
In [32]: highest_gross = studio2['foreign_gross'].max()
print('the highest domestic grossing film by Universal studio is:', highest_gross)
```

the highest domestic grossing film by Universal studio is: 891800000.0

```
In [33]: target = 891800000
row = studio2[studio2['foreign_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

891800000 Jurassic World: Fallen Kingdom

```
In [34]: target = 'Jurassic World: Fallen Kingdom'
row = studio2.loc[studio2['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

Jurassic World: Fallen Kingdom 2018

3. WARNER BROS

WARNER BROTHERS

THE HIGHEST DOMESTIC GROSSING FILM BY THEM IS "DARK NIGHT RISES"

GROSSING AT: 448100000.0

PREMIERED IN 2012

ACCORDING TO THE BAR GRAPH 2012 WAS THE YEAR WITH THE MOST DOMESTIC GROSS INCOME

THE HIGHEST FOREIGN GROSSING FILM IS " Harry Potter and the Deathly Hallows Part 2"

GROSSING AT: 960500000.0

PREMIERED IN 2011

```
In [35]: studio3 = box_office.loc[box_office['studio']=='WB']
studio3
```

```
Out[35]:
```

	title	studio	domestic_gross	foreign_gross	year
2	Harry Potter and the Deathly Hallows Part 1	WB	296000000.0	664300000.0	2010
3	Inception	WB	292600000.0	535700000.0	2010
10	Clash of the Titans (2010)	WB	163200000.0	330000000.0	2010
35	Due Date	WB	100500000.0	111200000.0	2010
37	Yogi Bear	WB	100200000.0	101300000.0	2010
...
3161	12 Strong	WB	45800000.0	21600000.0	2018
3167	The 15:17 to Paris	WB	36300000.0	20800000.0	2018
3170	Teen Titans Go! To The Movies	WB	29800000.0	22300000.0	2018
3209	They Shall Not Grow Old	WB	18000000.0	19000000.0	2018
3264	2001: A Space Odyssey (2018 re-release)	WB	3200000.0	19000000.0	2018

140 rows × 5 columns

```
In [36]: highest_gross = studio3['domestic_gross'].max()
print('the highest domestic grossing film by WarnerBros is:', highest_gross)
```

the highest domestic grossing film by WarnerBros is: 448100000.0

```
In [37]: target = 448100000.0
row = studio3[studio3['domestic_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

448100000.0 The Dark Knight Rises

```
In [38]: target = 'The Dark Knight Rises'
row = studio3.loc[studio3['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

The Dark Knight Rises 2012

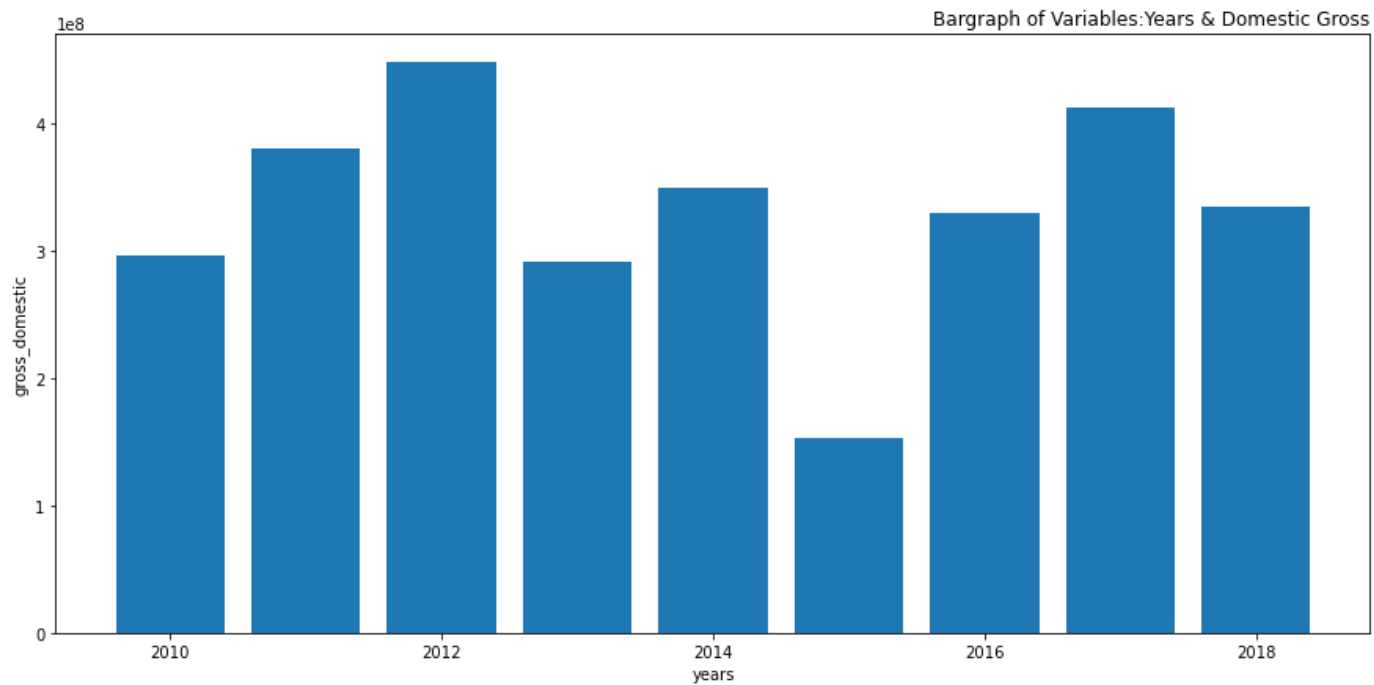
```
In [39]: studio3 = pd.DataFrame(studio3)

years = studio3['year']
grossd = studio3['domestic_gross']
grossf = studio3['foreign_gross']

fig = plt.figure(figsize= (15,7))

plt.bar(years, grossd)
plt.xlabel ('years')
plt.ylabel('gross_domestic')
plt.title('Bargraph of Variables:Years & Domestic Gross', loc = "right")

plt.show()
```



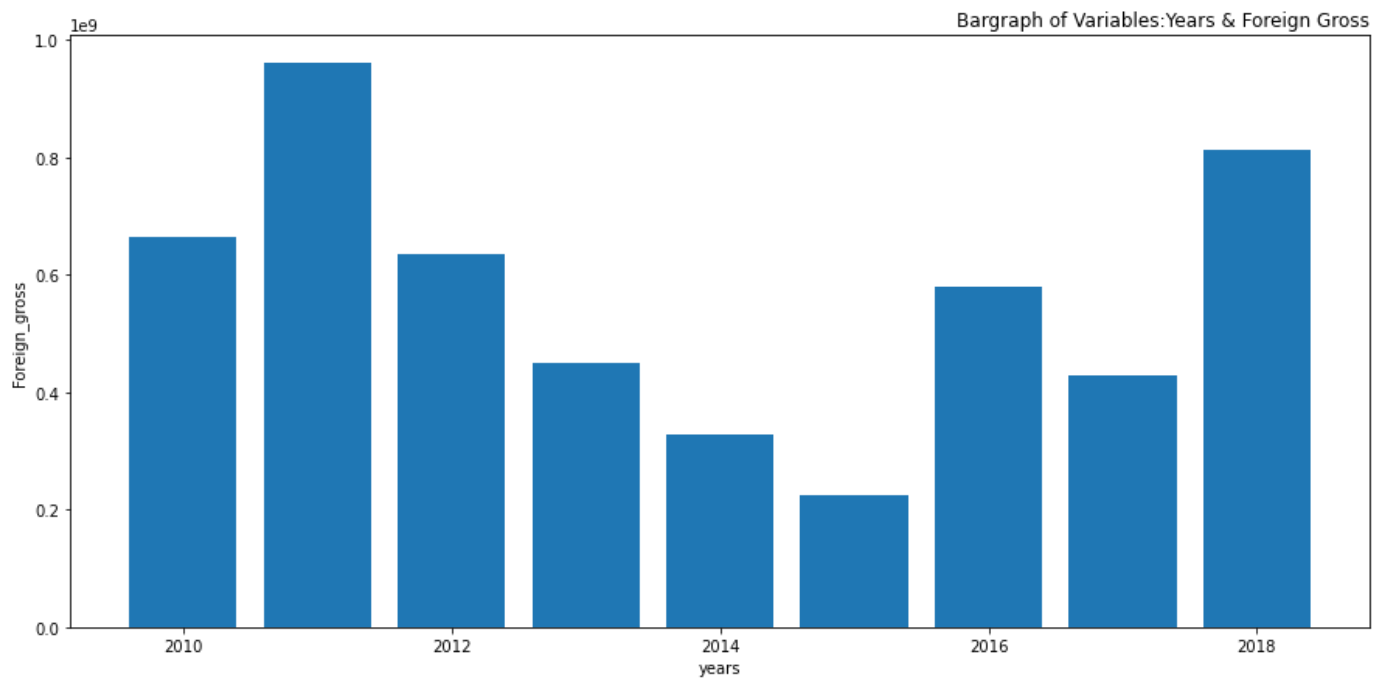
```
In [40]: studio3 = pd.DataFrame(studio3)

years = studio3['year']
grossd = studio3['domestic_gross']
grossf = studio3['foreign_gross']

fig = plt.figure(figsize= (15,7))

plt.bar(years, grossf)
plt.xlabel ('years')
plt.ylabel('Foreign_gross')
plt.title('Bargraph of Variables:Years & Foreign Gross', loc = "right")

plt.show()
```



```
In [41]: highest_gross = studio3['foreign_gross'].max()
print('the highest domestic grossing film by WarnerBros is:', highest_gross)

the highest domestic grossing film by WarnerBros is: 960500000.0
```

```
In [42]: target = 960500000.0
row = studio3[studio3['foreign_gross']== target]
if not row.empty:
    film = row['title'].values[0]
    print(target, film)
else:
    print('no',target)
```

960500000.0 Harry Potter and the Deathly Hallows Part 2

```
In [43]: target = 'Harry Potter and the Deathly Hallows Part 2'
row = studio3.loc[studio3['title']==target]
if not row.empty:
    year = row['year'].values[0]
    print(target, year)
else:
    print('no')
```

Harry Potter and the Deathly Hallows Part 2 2011

OVERVIEW PLOT OF BOX OFFICE DATA FRAME

```
In [44]: #plotting the domestic gross over years using plt.bar

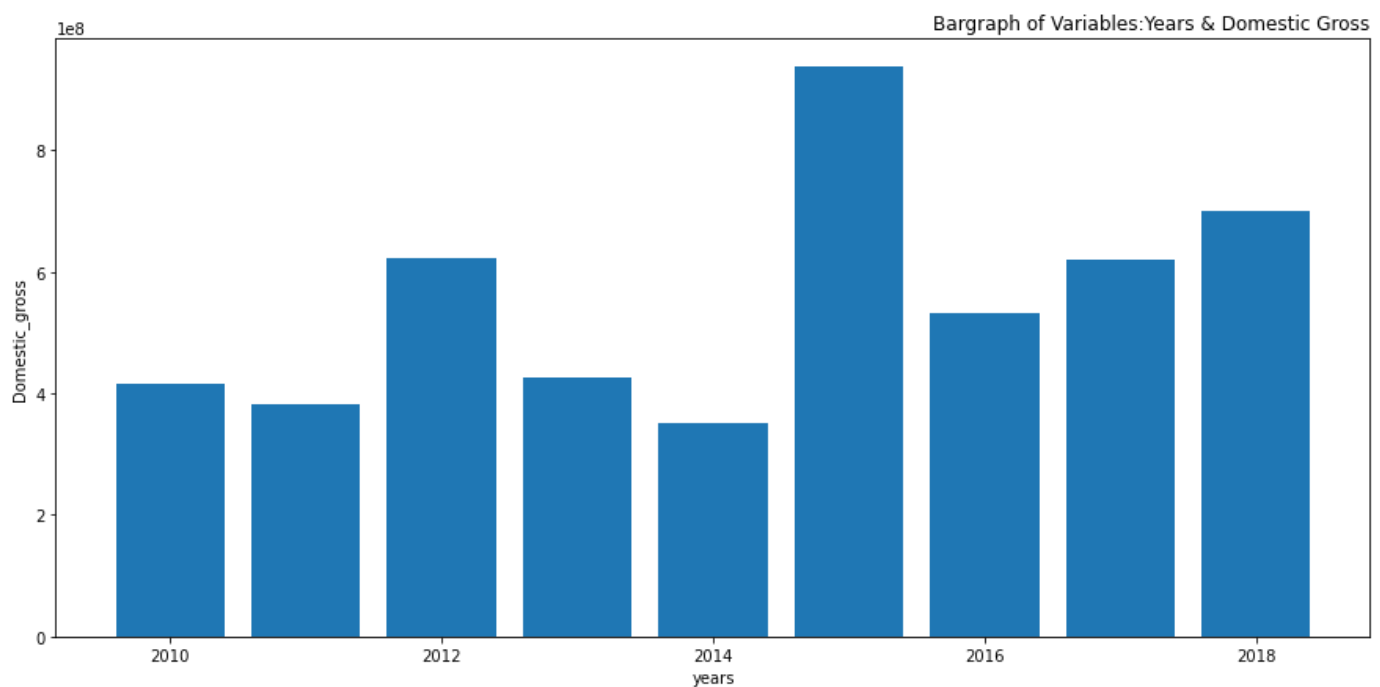
b = pd.DataFrame(box_office)

year= box_office['year']
dgross = box_office['domestic_gross']
fgross = box_office['foreign_gross']

fig = plt.figure(figsize=(15,7))

plt.bar(year, dgross)
plt.xlabel ('years')
plt.ylabel('Domestic_gross')
plt.title('Bargraph of Variables:Years & Domestic Gross', loc = "right")

plt.show()
```



overall in 2015 experienced the highest domestic grossing

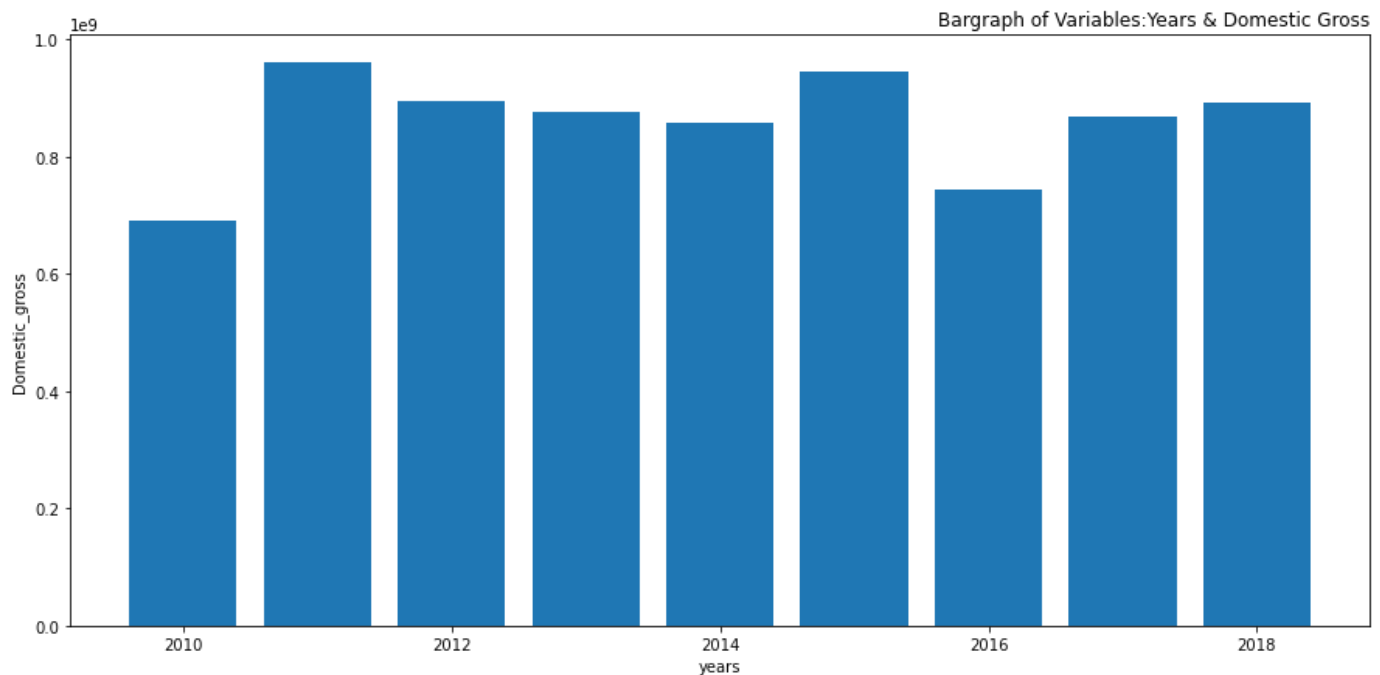

```
In [45]: c = pd.DataFrame(box_office)

year= box_office['year']
dgross = box_office['domestic_gross']
fgross = box_office['foreign_gross']

fig = plt.figure(figsize=(15,7))

plt.bar(year, fgross)
plt.xlabel ('years')
plt.ylabel('Domestic_gross')
plt.title('Bargraph of Variables:Years & Domestic Gross', loc = "right")

plt.show()
```



overall in 2011 and 2015 experienced the highest domestic grossing

CONCLUSION

domestic gross and foreign gross are positively correlated at 0.79

i analysed the top 3 studios: IFC, WARNER BROS AND UNIVERSAL STUDIOS

WARNER BROS: "The Dark Knight Rises" was the highest domestic grossing film in 2012, contributing to the highest domestic gross income that year. On the other hand, "Harry Potter and the Deathly Hallows Part 2" was the highest foreign grossing film in 2011, contributing to the highest foreign gross income in that year.

UNIVERSAL STUDIOS "Jurassic World" and "Jurassic World: Fallen Kingdom" were the highest domestic and foreign grossing films, respectively, for Universal Studios. These films contributed to the years 2015 and 2018 being the years with the highest domestic and foreign gross income, according to the provided bar graph.

IFC STUDIOS "Boyhood" and "Heartbreaker" were the highest domestic and foreign grossing films, respectively, for IFC Studio. These films contributed to 2014 being the year with the highest domestic gross income and 2010 being the year with the highest foreign gross income, according to the provided bar graph

if time permitted i would comment better and use the budgets data to come up with profit margins that these films have and use it to analyse our expected commercial success

In []: