

MOVIE PREDICT AND RECOMMENDATION

1. Statement and problem

Movies have been around us for long. First they come out with black and white movies, for many years they can make a color and sound effect like today. For today, we have so many film industries that want to bring their idea, and their time to help us entertainment but some movies have revenue less than others or review rating was less than they are expecting. So to help the company come up with ideas and plan for making a good movie. In this project I want to build a machine that can help predict the movie's success or fail with our data we have. And there are many people who are looking for movies to watch everyday, so it's hard for people to look for something they don't know, so I build the recommendation machine, so it can recommend the movie to the same type as the user who watched it.

2. Collection data

The data was in the kaggle composer, but data just has to the day it was posted as author collected from the API at the time he/she worked on it. So I want to have more data up to today, so I just generated a metadata movie for myself. First I got the last movie ID on the [website of TMDb](#). The Movie Database (TMDb) is a community built movie and TV database. Every piece of data has been added by our amazing community dating back to 2008, I generate over one by one movie ID. It may take a bit of time, but that can get me an idea how I can pull data from API.

3. Cleaning and transform data

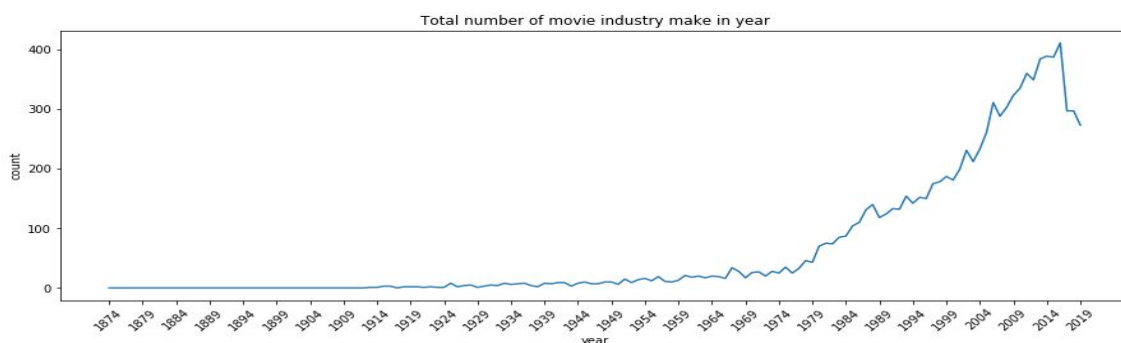
The data format was not usable yet, so we need to apply some methods to make usable data.

First, we needed to drop some columns that were not useful for us.

Second, you will see the data after the request from the API gives you a string of dictionaries. The information has the ID and name for those specified columns, but we just need the name only. So I just convert the string of dictionaries to dictionaries and get the list name from the key name. You can find the code I did in this [github](#).

4. Data Visualization

- a. Total the number of movie was increasing by year, as the graph show us that from 1874 is first movie was make, and they are likely make same number of movie each year to 1924, from 1924 to 1979 have more movie was make, the line begin increasing but not much. And after 1979 the number of movies was increasing a lot. As data we have, let's look at top 10 year revenue. This is just an idea to overlook the industry, because our data had a lot of missing revenue data.



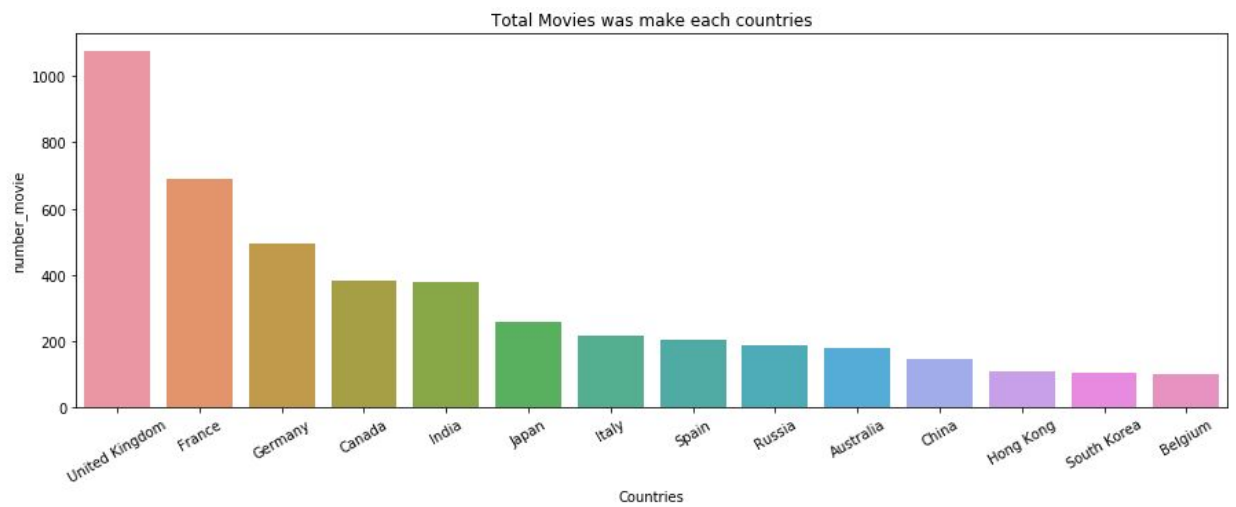
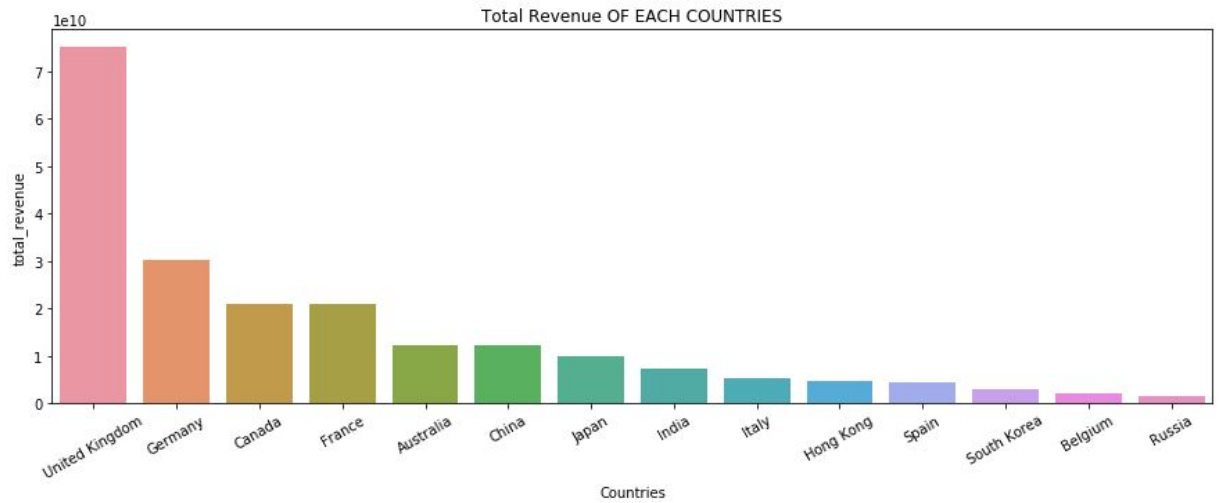
	year	title	runtime
34693	1874	Passage of Venus	1.0
34690	1878	Sallie Gardner at a Gallop	1.0
70946	1881	Athlete Swinging a Pick	1.0
41270	1883	Buffalo Running	1.0
135279	1885	L'homme machine	1.0

Recent movie, was have a time run longer, as you know standard most movie was runaround 90 to 120 min long.

	year	title	runtime
174446	2019	Sunday	13.0
174441	2019	Queen + Béjart - Ballet For Life	58.0
126731	2019	The Ocean Washed Open Your Grave	3.0
174467	2019	Entropia	28.0
210551	2019	Jorge	20.0

- d. Next, let us divide countries and see where they make more and have best success in the film industry on a data set. The chart below, you can see that the USA was the top 1 make movie, then after that was the UK and France, and Germany.

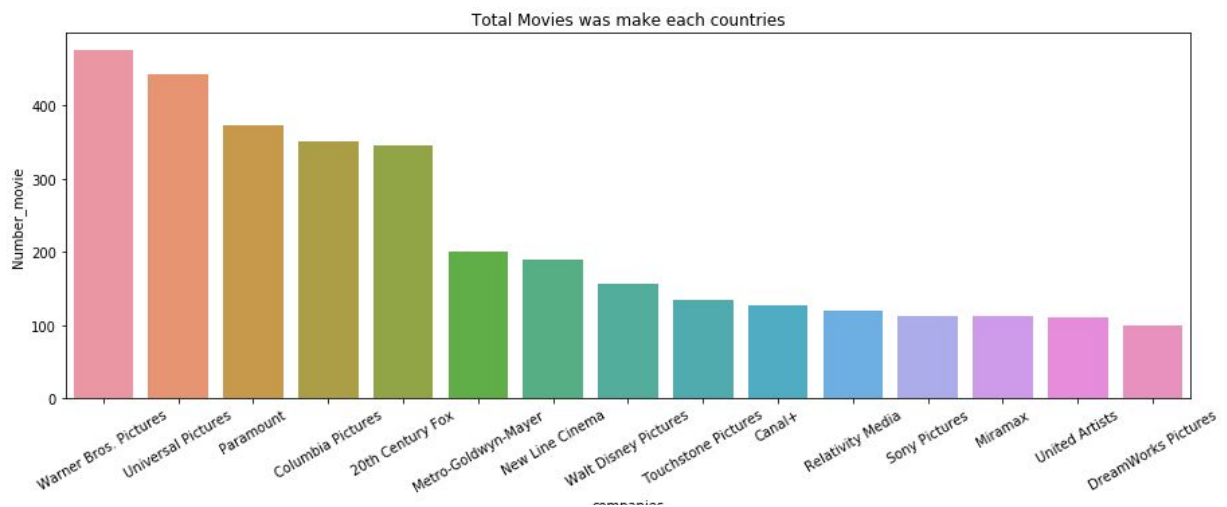
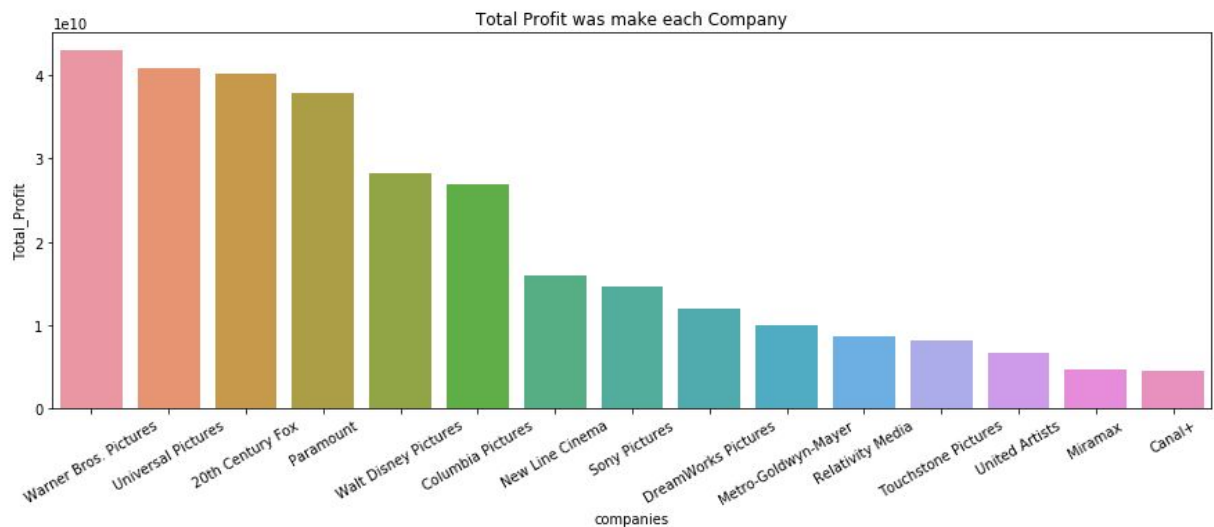
	average_revenue	total_revenue	number_movie
Countries			
United States of America	7.714978e+07	4.976932e+11	6451
United Kingdom	6.990096e+07	7.514353e+10	1075
France	3.021624e+07	2.084921e+10	690
Germany	6.078097e+07	3.014736e+10	496
Canada	5.485054e+07	2.089806e+10	381
India	1.911960e+07	7.246328e+09	379
Japan	3.894279e+07	9.969354e+09	256
Italy	2.382187e+07	5.169346e+09	217
Spain	2.157035e+07	4.378782e+09	203
Russia	8.155989e+06	1.517014e+09	186



France is rank 3 product movie, but their revenue was at rank 5. China is kind of the opposite, they were rank 12 of countries' product movies, but their revenue was rank 7.

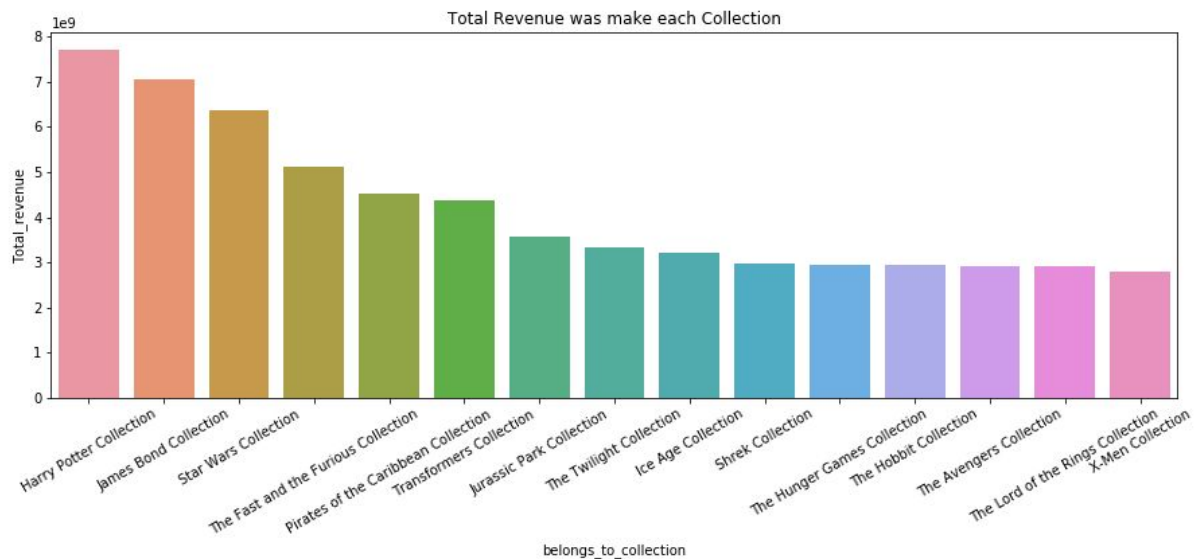
- e. Now we will get more detail about each company and which was best, so the chart below that you can see Warners Bros. Pictures, Universal Pictures, Paramount, Columbia Pictures, and 20th Century Fox are top 5 companies was make most profit and number movie they are make overall.

companies	Avergae_profit	Total_Profit	Number_movie
Warner Bros. Pictures	9.032809e+07	4.299617e+10	476
Universal Pictures	9.229253e+07	4.088559e+10	443
Paramount	1.014296e+08	3.793466e+10	374
Columbia Pictures	7.684656e+07	2.697314e+10	351
20th Century Fox	1.164543e+08	4.017672e+10	345
Metro-Goldwyn-Mayer	5.015126e+07	1.003025e+10	200
New Line Cinema	8.415754e+07	1.598993e+10	190
Walt Disney Pictures	1.797291e+08	2.821747e+10	157
Touchstone Pictures	6.081461e+07	8.209972e+09	135
Canal+	3.521476e+07	4.437060e+09	126

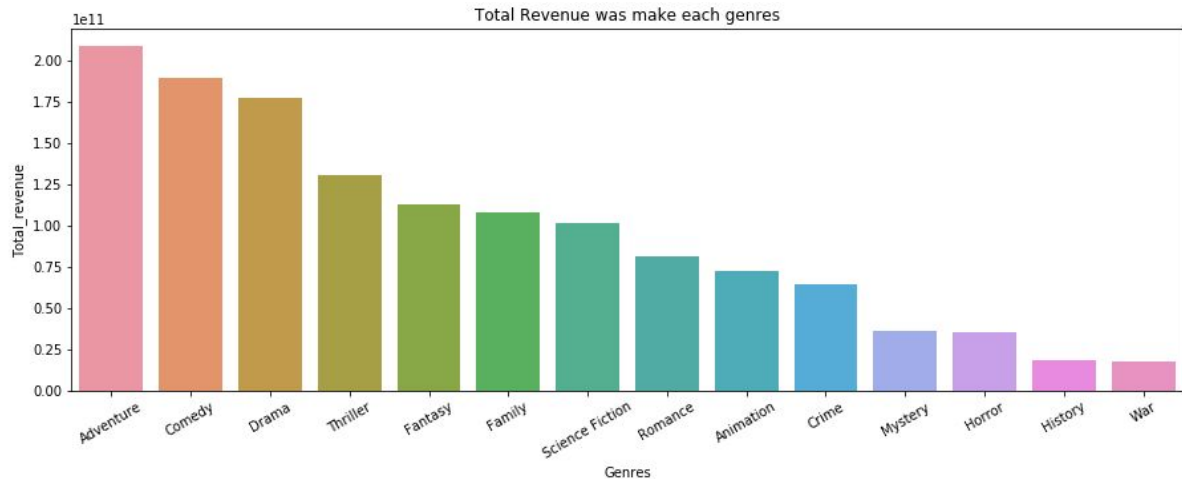


- f. Some movie they are so good, so they keep make a new seri after another. So you can see bellow that the top 1 have revenue is Harry Potter collection. And Jame Bond was number 2, after that Star War, the Fast and Furious. These movies you can see are so popular, so it made sense that it was in top 4 revenue of all the collection movies.

	Average_revenue	Total_revenue	Total_movie
belongs_to_collection			
Harry Potter Collection	9.633598e+08	7.706879e+09	8
James Bond Collection	2.827486e+08	7.068715e+09	25
Star Wars Collection	9.112054e+08	6.378438e+09	7
The Fast and the Furious Collection	6.406373e+08	5.125099e+09	8
Pirates of the Caribbean Collection	9.043154e+08	4.521577e+09	5
Transformers Collection	8.758574e+08	4.379287e+09	5
Jurassic Park Collection	8.948083e+08	3.579233e+09	4
The Twilight Collection	6.686215e+08	3.343107e+09	5
Ice Age Collection	6.433533e+08	3.216767e+09	5
Shrek Collection	7.411823e+08	2.964729e+09	4

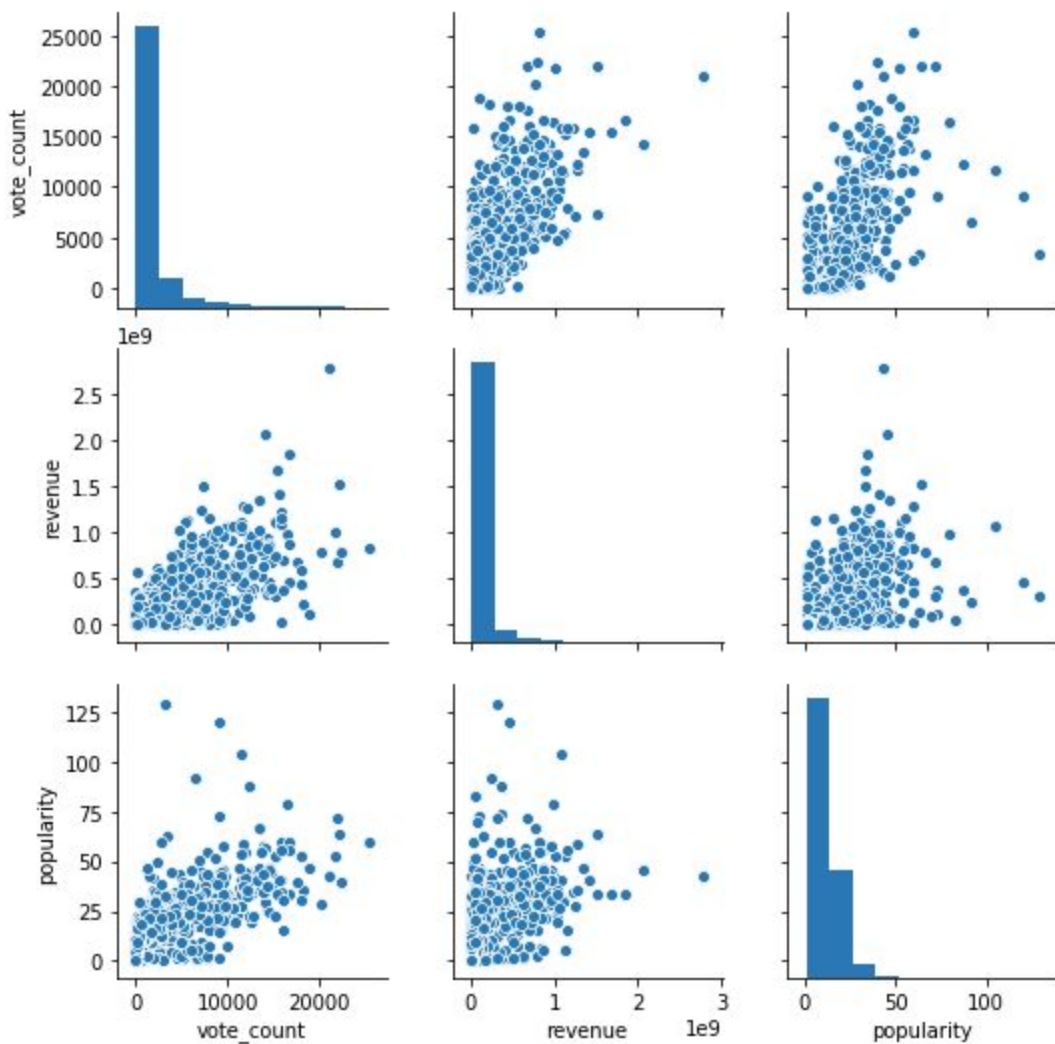
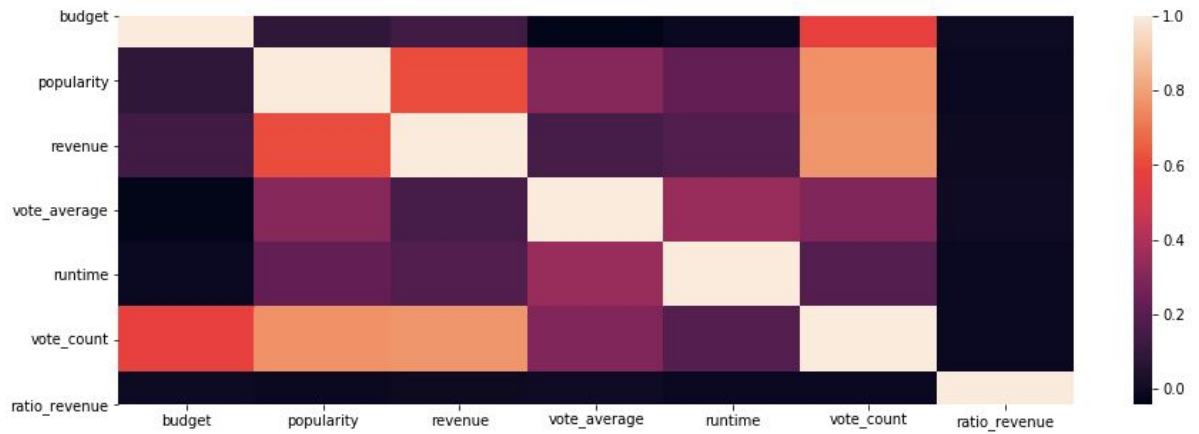


- g. Genres of action, adventures, and comedy are the top 3 movies that have most revenue. It can be said people are more likely interested in these kinds of movies.



	Avergae_revenue	Total_revenue	Total_movie
Genres			
Action	1.025900e+08	2.134897e+11	2081
Adventure	1.640720e+08	2.086996e+11	1272
Comedy	5.798689e+07	1.889213e+11	3258
Drama	3.835340e+07	1.773078e+11	4623
Thriller	6.120269e+07	1.302393e+11	2128
Fantasy	1.463286e+08	1.122341e+11	767
Family	1.299991e+08	1.080293e+11	831
Science Fiction	1.181741e+08	1.016297e+11	860
Romance	4.568787e+07	8.073046e+10	1767
Animation	1.467283e+08	7.189686e+10	490

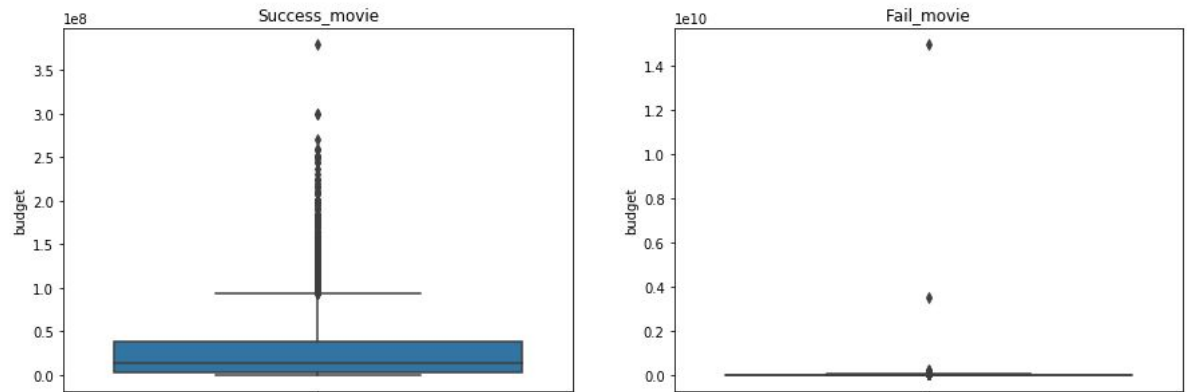
- h. Now, we can see the relationship between these numeric columns. As the heatmap plot below show us that they are not have much relation, but some of them are have most is around 0.6 or 0.7 correlation.



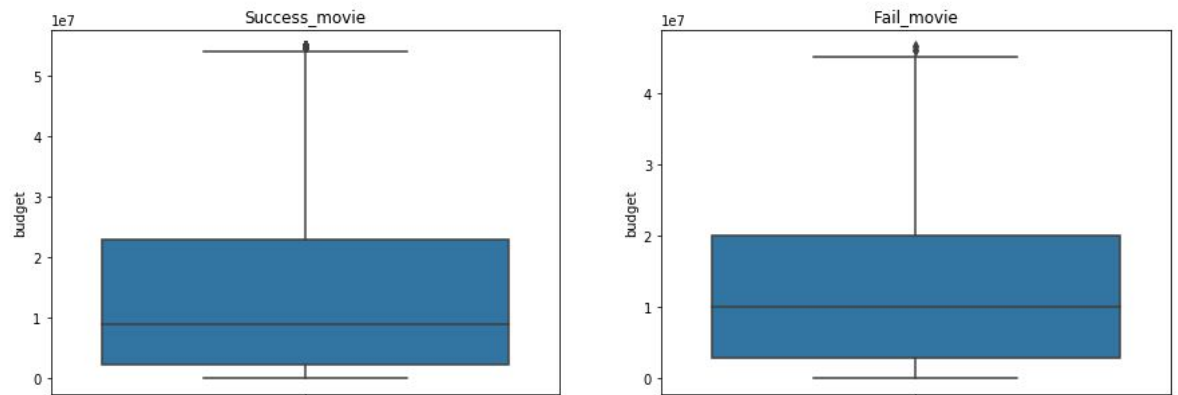
i. Inference Statistic

- We will take a look at what different the average budget of success movie and fail movie, but data has so many outliers. That we need to remove it.

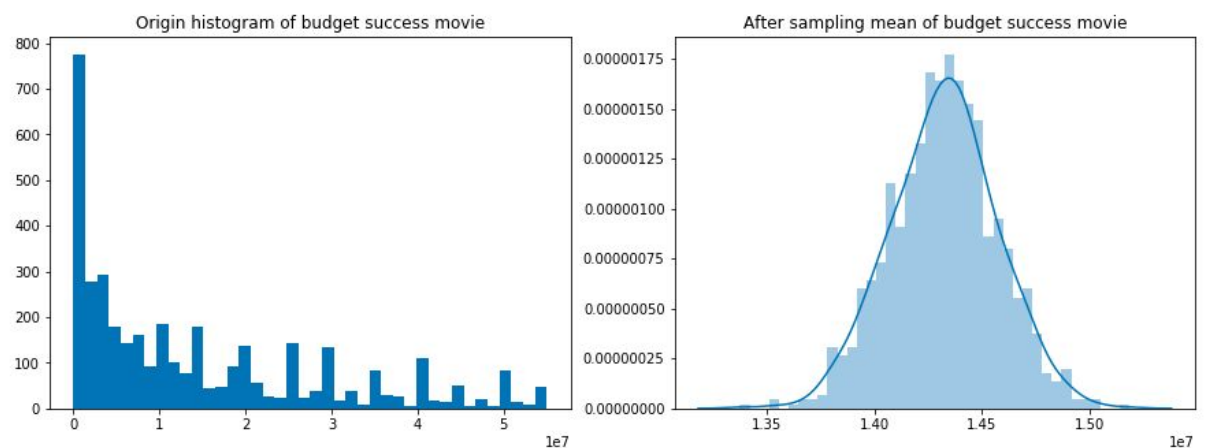
- + Plot box plot before we remove the outlier, so can notice that the failed movie has a largest gap of outliers. It is bad for our conclusion.

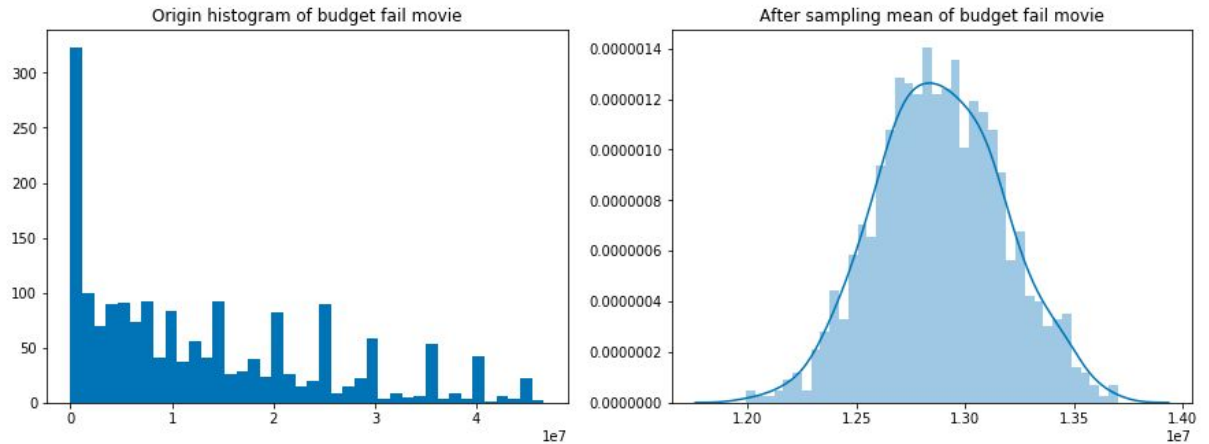


- + These plot below is when we remove the outlier, box plot show is more better



- We used the bootstrap resampling to tell that mean of success is 14M and fail movie 13M





- We perform the hypothesis test to test the different mean average of success and fail movies.
 - We identified the null hypothesis and alternative hypothesis.
 - H_0 : the different average of success and fail is ≥ 1438763
 - H_a : the different average is < 1438763
 - With significant is 5% or 0.05
 - After our test and we got the p-value = 0.501
 - Then we can't reject our null hypothesis.

