

# Rentomojo Assignment

Submitted by:-

Rishabh kapoor  
rishabhkapoor101@gmail.com

# Acknowledgement

First and foremost I would like to thank “RentMojo” for entrusting me with a project that gave me a chance to prove my value that I can provide to the company if I am chosen for the position. Also I would like to thank all my professors that taught me the tools necessary for such work. And last but not the least I would like to thank my family that has been pushing me forward since day one.

Without any further ado let's have a gander on the project.

(Also for a comfortable read please download full in pdf format from my github :- <https://github.com/rishabhkapoor101/MoviesAnalysis> )

# Dataset

So the data set was a csv formatted file that contained the data about the movies.  
It had 24 columns namely:-

1. index
2. budget
3. genres
4. homepage
5. id
6. keywords
7. original\_language
8. original\_title
9. overview
10. popularity
11. production\_companies
12. production\_countries
13. release\_date
14. revenue
15. runtime
16. spoken\_languages
17. status
18. tagline
19. title
20. vote\_average
21. vote\_count
22. cast
23. crew
24. director

With a row count of 4803.

The data set with a skim, seemed a piece of cake to work with, but after just digging a little, I realised that it was not the case.

I had to clean the data as there were “nan” values that essentially means “not a number” even when a number was expected.

This was not everything that was a challenge about the dataset, in fact we will see much more as we will progress through the dataset.

# Toolset

The tools I used for this assignment are as follows:-

- Python3.8
- Pandas
- Jupyter notebook
- Microsoft Excel
- Tableau
- Pdf viewer
- github

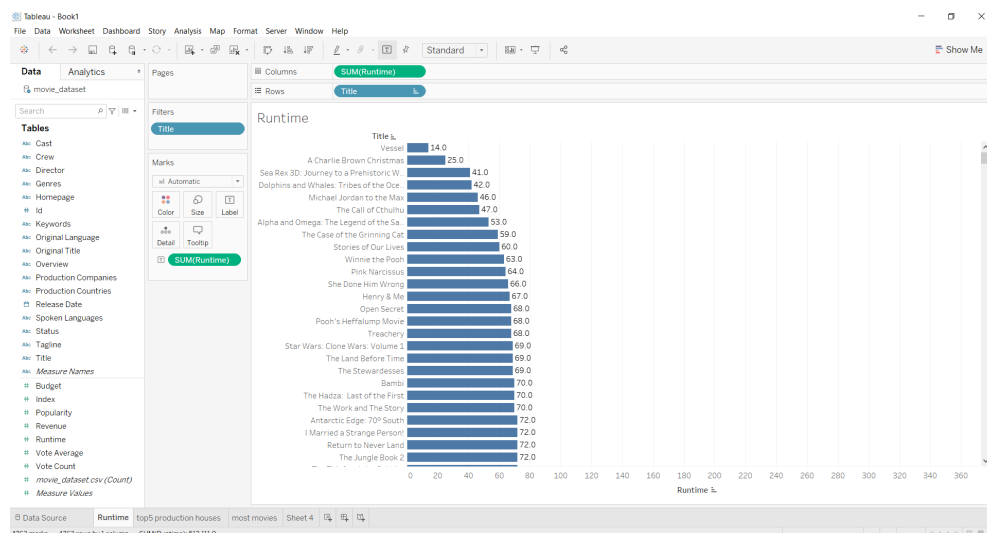
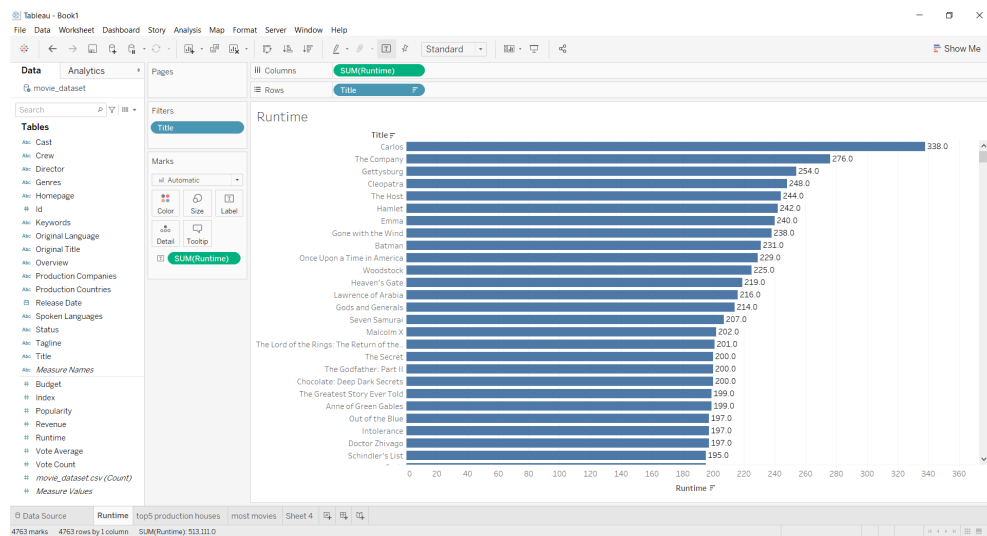
# Task 1

For task 1, I had to find the movies at contrasting ends of the runtime spectrum, i.e. I had to find the movies which had the highest and least runtime.

I did this without any difficulty in Tableau.

**Result :-** Movie with least runtime was **Vessel** with a runtime of **14 min** and Movie with highest runtime was **Carlos** with a runtime of **338 min**

I have also attached the screenshot of my findings as follows...



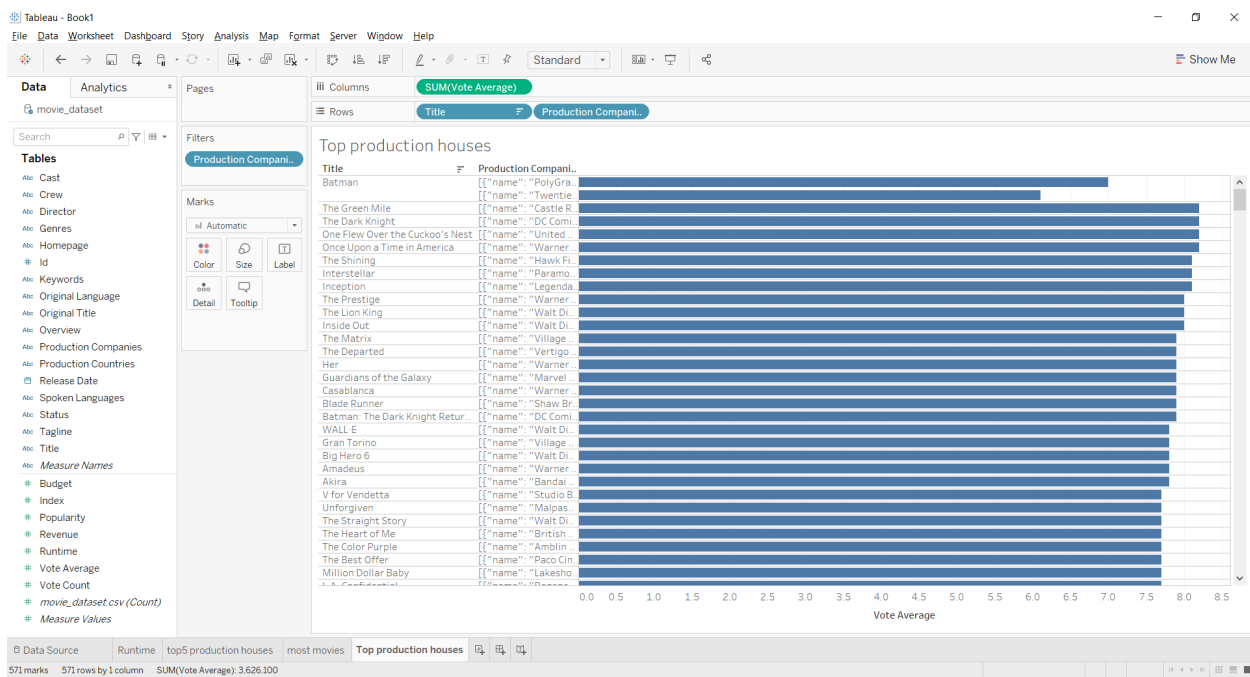
# Task 2

For the second task I had to find the top five production houses with highest vote average, revenue and then find their 5 movies each.

What I found was that the biggest corporations according to the budget were:

- Walt disney
- Warner bros
- Dc comics
- Marvel
- BBc

And their respective movies are attached as follows:-

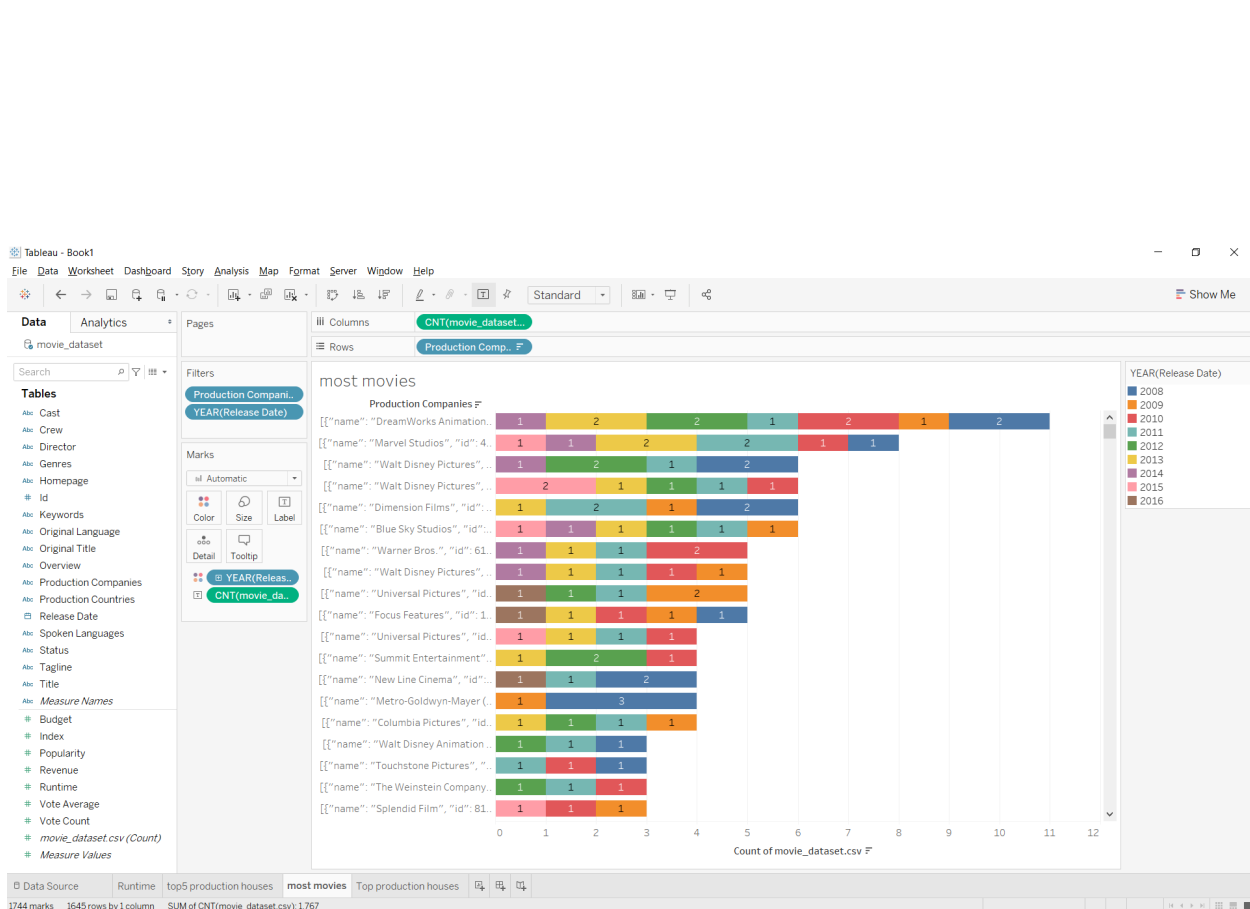


(for a cleaner look, please refer my book1.twb file)

# Task 3

For task three, I had to list the production houses that released most movies in the period of 2006-2016

That was fairly easy and the sheet below will be self explanatory.



(I may sound a bit redundant, but for a closer and cleaner look, have a look at my book1.twb)

# Taks 4

For this task I had to choose a company between DC comics and Marvel studios where I would invest all my money had I been given the option.

This was the question where I felt the most heat.

My approach was to compare the growth rate of the two companies rather than the raw revenue and that was a challenge in itself, I had to separate years from full dates, I had to classify them in production\_house categories, which essentially were Dc, Marvel or Neither. And then I had to plot new info in the tableau. For which I had to make a new dataframe from modified values for plotting which I did make with the name of **“modified.csv”**

Everything except for the plotting was done in python3.8 and pandas in jupyter notebook.

I have attached some screenshots of the same...

The screenshot shows a Jupyter Notebook window titled 'pybook2 - Jupyter Notebook'. The browser address bar shows 'localhost:8888/notebooks/pybook2.ipynb'. The notebook interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Help) and a toolbar with icons for file operations, running cells, and viewing code. The code in the notebook is as follows:

```
In [8]: print(len(li))
4803

In [9]: li[-1]
Out[9]: '[{"name": "rusty bear entertainment", "id": 87986}, {"name": "lucky crow films", "id": 87987}]'
```

```
In [10]: phouse = list()
for i in li:
    if "marvel" in i.lower():
        phouse.append("Marvel")
    elif "dc" in i.lower():
        phouse.append("Dc")
    else:
        phouse.append("Neither")
print(phouse)
```

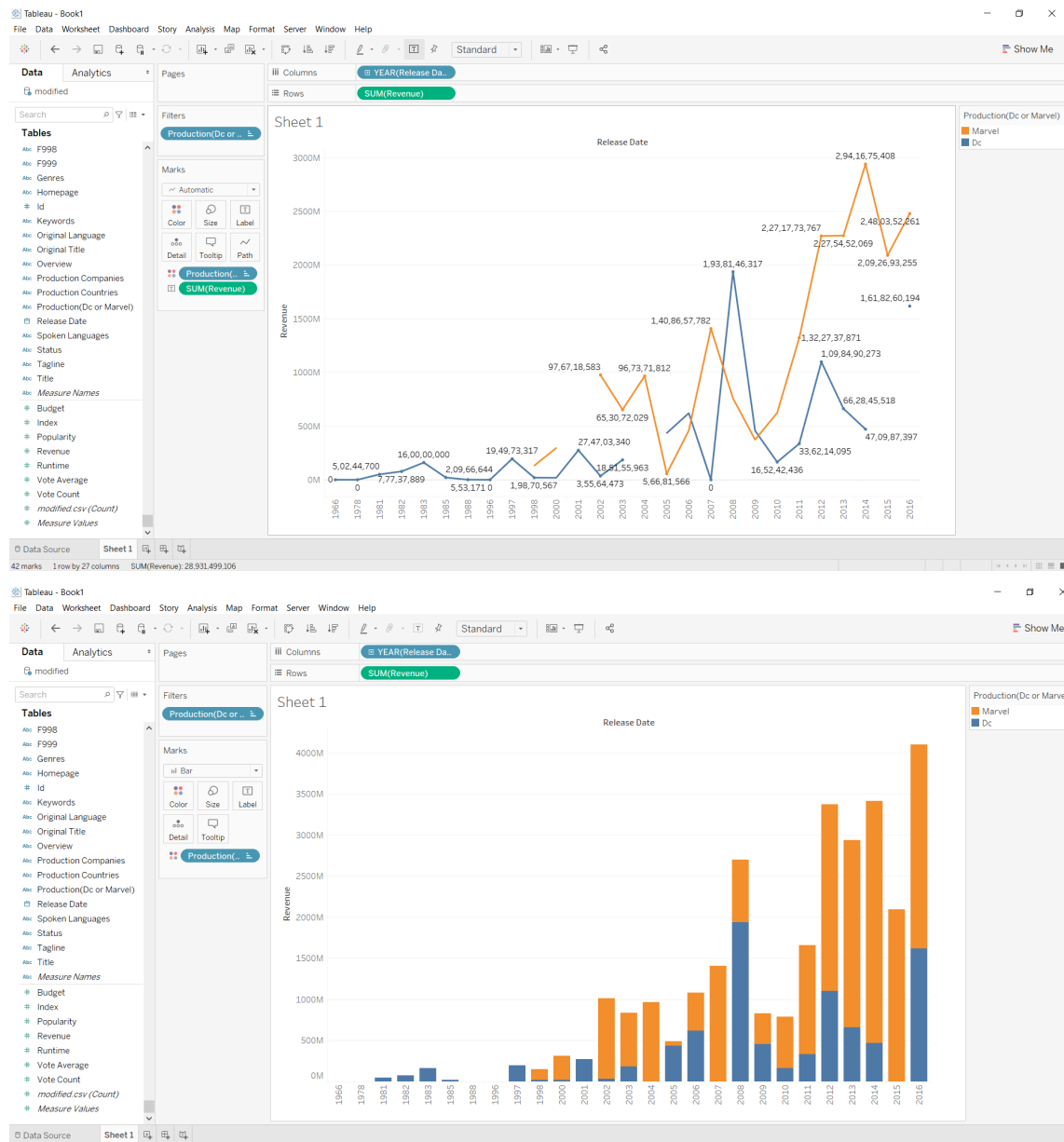
The output of the code is a long list of strings, each representing a production house classification. The list starts with several 'Neither' entries, followed by 'Dc', 'Marvel', and 'Neither' entries, and continues with many more 'Neither' entries. The list is truncated in the screenshot, showing only the first few elements of the output.





# CHARTS

Below, the charts compare the revenue year wise ie growth and needless to say Marvel is the winner even though Dc started prior to Marve, so if I am given a chance to invest my money, without a shred of doubt it is going to be marvel.



(for a clearer look, refer to book2.twb)

# Workflow...

As I already, have mentioned the tools I used, lets see how were they used,

- Python3.8 was very helpful in data cleaning and other manipulations.  
Also it was a massive help in getting to know the data on a very intimate level.
- Pandas was used via python to read and manipulate data frames easily.
- Jupyter notebooks were used for better formation visual representation of dataframes.
- And Github was used to make a public repo for this very project.

# Thank you.

I hope you enjoyed reading this report.

A report by:-

Rishabh kapoor.

[rishabhkapoor101@gmail.com](mailto:rishabhkapoor101@gmail.com)