# PROJECT REPORT

# ECONOMETRIC ANALYSIS OF MOVIE DATASET

## AKSHAY ANANDBABU

## RITVIK REDDY

## DA-5020 COLLECT, STORE, RETRIEVE DATA

## SUMMER-FULL 2017

# TABLE OF CONTENTS

## 1. Project Description

As movie enthusiasts, we were immediately drawn to the idea of using movies as our topic. The main motive of a movie producer is to earn more money than was invested in creating the movie. Through our analysis, we tried to gain insight on the thoughts of a producer and what are the keys to making a successful movie. We decided to work on a movie dataset to address the factors that affect the success of a movie worldwide. In our project, we base the success of a movie by how much it earns in the domestic and worldwide box office. Thus, the theme of our project is based around the worldwide box office collection and the various components that affect it.  We perform this econometric analysis using plots in R and creating a linear regression model to see the correlation between the different variables and the worldwide box office collection.

## 2. Selection Rationale:

To perform our analysis, we needed data on the production budget and worldwide box office collection and ratings for a variety of movies released internationally. Moreover, we needed a mix of movies with different genres and content ratings so as to widen our range of analysis. The *the-numbers.com* website provided a comprised list of the highest grossing movies of all time released globally. The following were the major reasons for choosing the website for our project:

- The website is open-source and does not require access authorization.
- The entire data was available in a single page as a search result
- It provided a comprehensive profile for each movie that included the domestic, international and worldwide income for the movie

- It provided both the domestic and international release dates separately so if the international release date for a movie was missing, we could make the assumption that the movie was only available in the country of release.

## 3. Process Implemented

1. Scrape the movie data from "the-numbers.com" using WebScraper extension in Google Chrome.

2. Download it as a CSV file and save it.

3. Import the CSV into the R environment.

4. Clean the data-set for abnormalities and store it in a new table.

5. Store the tidy dataset into a SQLLite database.

6. Retrieve the data using queries and perform analysis (using ggplot).

7. Create a linear regression model to find the linear relation between the variables.

## 3.1 Scraping the Dataset

We used the WebScraper chrome extension to perform scraping. The WebScraper provides a cohesive node selection process and allows the user to explore a link on a page to get data inside the link. This feature was not available in other tools like import.io, grepsr, dexi.io and we could not use the *rvest* package as we had to traverse through multiple links and it would be uneconomical to store each link and perform extraction. WebScraper provides an interesting option to browse the data while the scraper is running. This allowed us to pause the scraping process if there were discrepancies and check the current movie that is scraping. It also provides the option of downloading the scraped data as a structured CSV file.

**Fig 3.1.1: Selecting nodes to be scraped**



**Fig 3.1.2: Elements created after selecting nodes for each link.**



**Fig 3.1.3: Selector graph displaying the schema of the scraping process**

**Fig 3.1.4: Preview of the final data**

## 3.2 Importing the CSV file into the R environment

Once the scraping process was completed and downloaded as a CSV file, we imported the file into the R environment using the read_csv() function and stored it in a suitable object as a data frame.



**Fig 3.2: Final Scraped CSV file**

## 3.3 Cleaning the Data

- The movie title included the year of release and since we have release dates in separate columns, it was unnecessary to have the year of release along with the movie name. We used the gsub() function and general expressions to remove the movie year.

- The international and domestic release dates were not in the default R date-format so we converted them using the mdy() function in the lubridate package. The release dates also contained the place of release inside quotes that had to be removed, so we used the str_extract() function with general expressions to extract the date. We stored the clean variables in the same table.

- Cleaning the budget and box office columns were fairly uncomplicated as the values for production budget and box office collections for movies that were released outside the united states were already converted into USD. We cleaned the columns by removing the dollar sign "$" and comma "," from each row and stored it as a numeric to perform analysis.

- We converted the duration column into a numeric and extracted the number of minutes, removing unnecessary characters using general expressions.

## 3.4 Storing the clean data into SQLLite

We used **dbConnect** to create a connection R to SQLLite and created a database named "AR_project.sqlite". The clean dataset was written inside this database using the **dbWriteTable** *query.*

## 3.5 Retrieving the data and performing analysis

The clean data was retrieved using ***dbGetQuery*** and the required variables were selected using the select query in SQLLite. Our analysis consists of three parts distribution, econometric analysis and linear regression modelling.

- In our distribution analysis, we check for consistencies in the data and check if there are certain variables that skew the results. Using ggplot to create a histogram, we found that the distribution of IMDb ratings follow a near normal distribution with the mean lying around the 6.5-7.5 range. We also created a distribution plot for content ratings. From the plot, we see that movies that are rated R occur the highest in the data-set. Since the data set is a list of movies with highest box office collection, we can make a possible conclusion that R rated movies fare well in box office collections.



**Fig 3.5.1: Normal distribution of IMDb ratings**

**Fig 3.5.2: Scatterplot showing the number of movies released in each country**



**Fig 3.5.3 Barplot showing the number of movies based on content ratings**

- Our econometric analysis involves finding the effect of IMDb rating, genre on the worldwide box office. We used a scatterplot and a smoothening line to find a reliable trend. A boxplot was created to see the effects of genre on the worldwide box office. We decided a boxplot would provide an ideal visualisation as it would provide the range of income of each genre and the median line to assist in finding trends.



**Fig 3.5.4: Scatterplot showing the trend between IMDb rating and worldwide box office collection.**



**Fig 3.5.5 Boxplot showing the trend between Genre and worldwide box office**

- Our linear model finds a linear relation between the dependent variable(worldwide box office)and independent variables ( IMDb rating, production budget, domestic box office, international box office). Through our model we  find the correlation between these variables and the coefficient of determination which is given by the r-squared value.

```
    train_positions: function
    train_scales: function
    vars: function
    super:  <ggproto object: Class FacetNull, Facet>
--------------------------------
mapping: colour = Genre
geom_boxplot: outlier.colour = NULL, outlier.fill = NULL, outlier.shape = 19, outlier.size = 1.5, outlier.stroke =
0.5, outlier.alpha = NULL, notch = FALSE, notchwidth = 0.5, varwidth = FALSE, na.rm = FALSE
stat_boxplot: na.rm = FALSE
position_dodge


Call:
lm(formula = worldwide_box_office ~ production_budget + duration +
    imdb_score, data = data3)

Residuals:
      Min        1Q     Median        3Q        Max
-162032317  -49245715  -24475721   19612010  734865607

Coefficients:
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)        -2.142e+07  1.882e+07   -1.138    0.255
production_budget   2.331e+00  1.497e-01   15.570  < 2e-16 ***
duration            6.731e+05  1.324e+05    5.086 4.07e-07 ***
imdb_score         -3.143e+06  1.968e+06   -1.597    0.111
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 92410000 on 1684 degrees of freedom
  (2775 observations deleted due to missingness)
Multiple R-squared:  0.1632,    Adjusted R-squared:  0.1617
F-statistic: 109.4 on 3 and 1684 DF,  p-value: < 2.2e-16
```
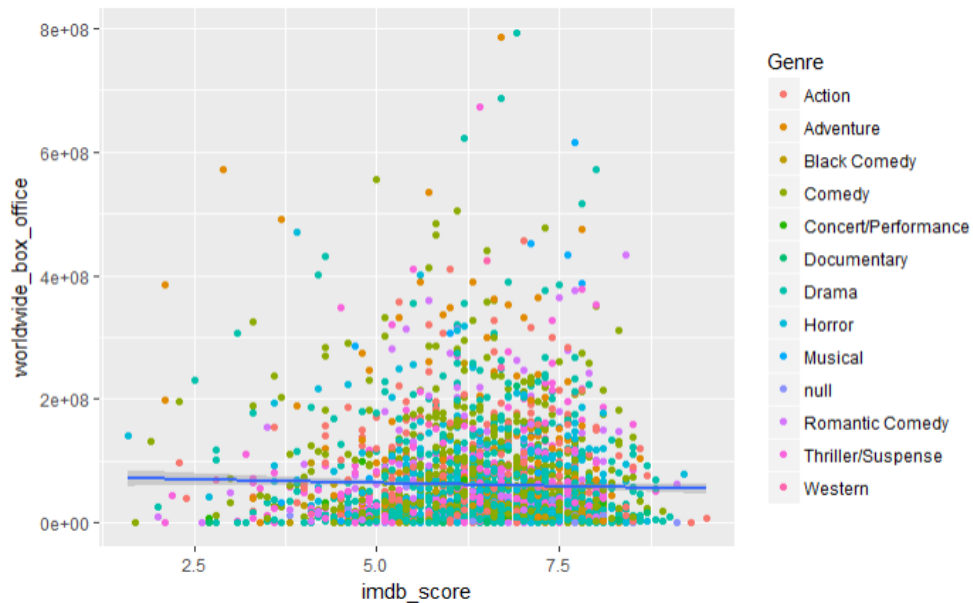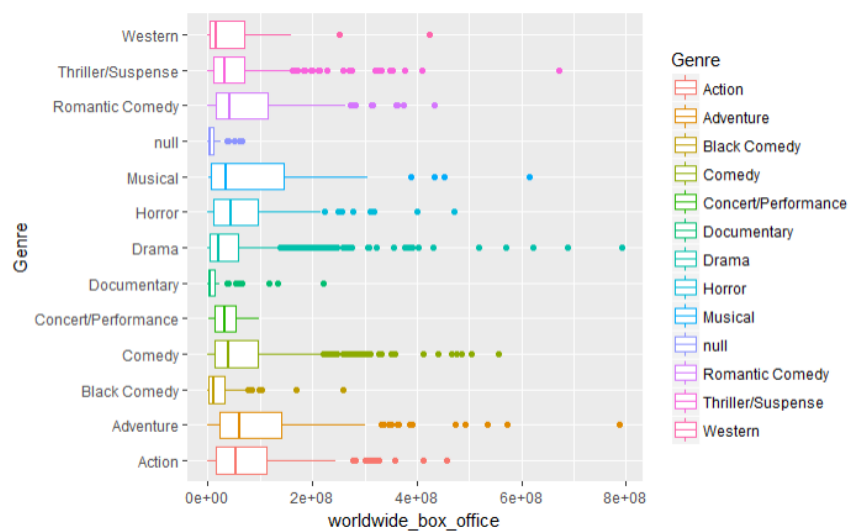
**Fig 3.5.6  R output displaying the results of the linear model**

## 4. Complications faced during the process

- **Long run times of WebScraper:**

Web-Scraper requires a long time to scrape data because it prints the results of each movie before moving on to the next movie. The scraper has to go through each link, scrape the data and then go back to the start link and repeat the process for each movie. The entire scraping process ran for upto 6 hours, during which we had to leave our computers on or it ended up in pausing the scrape.

- **Multiple international dates:**

  The data scraped has multiple international release dates for each country that was released globally though the release dates for each country were about equal. We solved this problem by extracting the first release date by using *str_extract().*

- **Cleaning Process:**

  There were a few intricate problems we faced while cleaning the scraped data set. The majority of the data types that we encountered were factors so we could not perform any string operations on it. We used *lapply[]* to change the variables into characters and then performed string operations on the variables to clean them. We also had to convert the world-wide box office collection given in dollars to numeric befor using it  for analysis.

## 5. Learning Outcomes and Future Work:

- Web Scraping using web-tools.

- Cleaning and formatting data using general expressions.

- Storing data in a database using SQLLite.

- Retrieving data from the database using queries.

- Creating plots and learning ways to make the plot look more informative using *stat_function().*

- Creating a linear regression model to find correlations between variables.

Future work for this project would be to create and implement an algorithm which would predict the overall box office collection of a movie by knowing its genre, production budget and the domestic box office collection. The project can be expanded

by performing a sentiment analysis on movie reviews and see the effect of movie

reviews on the rating of a movie which would in turn affect the box office collection.

## 6. <u>References</u>

- http://www.the-numbers.com/movie/budgets/all

- https://stackoverflow.com/questions/

- https://www.google.com

- http://webscraper.io/