

# CSE0448 Movies Revenue

Ahmad Nawar Droubi and Aydın Davutoğlu

**Abstract**—In this study, we present a statistical analysis about a movies data set composed of around 8,000 movies meta-data, which includes: the title of the movie; the budget, the revenue, run-time, the Motion Picture Association of America rating; genre of the movie; number of top 100 actors according to IMDb that were present in the Movie; languages that the movie was released in; countries in which the movie was released, Meta-score of the movie; IMDb rating, number of IMDb votes on the movie; production company and year of release. The data-set contains some missing values, which we had to clean, then we did some feature engineering to prepare the data-set for analysis and for creating the machine learning model. The statistical analysis of the data-set was for the aim of testing 9 hypothesis with the main aim of the project being finding what effects the revenue of a movie.

## I. INTRODUCTION

This study is based on the goal of predicting if a movie is going to be successful or unsuccessful based on revenue. First, we did basic feature engineering techniques on some of the features of the data-set, namely: the budget of the movie; famous actors that participated in the movie; genre of the movie; the rating of the movie; the IMDb votes on the movie; the amount of awards the movie won and which production company has produced said movie. Second, we used well known statistical tests in order to be able to predict the amount of the revenue the movie is going to generate.

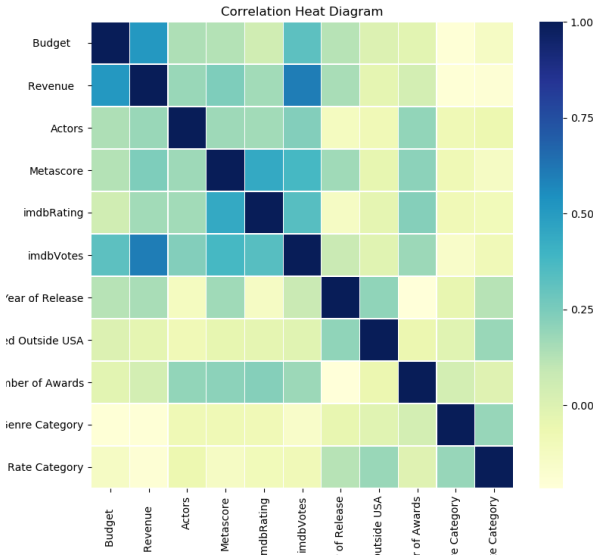


Fig. 1: Correlation Heat Diagram showing the relations between features in the data-set

## II. METHODOLOGY

The methodology used in creating/collecting the data-set was random sampling from a collection of top 10,00 movies by revenue as of the 5<sup>th</sup> of May, 2019. To test our main research hypothesis we created and tested 9 hypotheses in the aim of finding out which features of the data-set effected the movie revenue in order to include them in the ML model.

### A. Movie Dataset

In order to create the movie data-set we used The Movie Database API [1] to get a random list of 10,000 movies from the top 100,000 by revenue, then we used OMDb API [2] to collect the meta-data for the movies. The meta-data we used in this study are: 1) Title 2) Budget 3) Revenue 4) Run-time 5) Motion Picture Association of America rating 6) Genre 7) Top 100 Actors in Movie 8) Language 9) Country 10) Meta-score 11) IMDb Rating 12) IMDb Votes 13) Production Company 14) Year of Release

### B. Data Processing

The original data-set contained total of 10,000 movies, after filtering the data-set it reduced to 8,000 movie. We found a number of problems in the following features:

- 1) The Year column : was missing a lot of values from the OMDb API, where the one got from The Movie Database was complete, so we dropped the original column and extracted only the year from the date column which we got from TMD, then renamed it into “Year of Release” to better reflect the actual content of the column.
- 2) The Language column: was missing data for some movies. All the data we gathered was from Hollywood, so logically we changed the the missing values to English.
- 3) The Country column: was treated in a similar way to Language column. So, for movies with this value missing we assigned the value USA
- 4) The IMDB Rating and Votes columns: Zero has been assigned to the missing values, in purpose of protect the data-set from the missing value in this column.
- 5) The Actors column: This column had the name of the actors in that film. The main goal of this column is to see how many famous actor a movie include. To achieve this goal we gathered the top 100 famous actor and we checked each movie how many famous does it have.
- 6) The Production column: The missing values in this columns has been changed from null to others.
- 7) The Run-time column: The mean of this column has been assigned to the missing value.
- 8) The Genre column: To be able to use label encoding on

this column we choose the first value in this column which originally had multiple values.

### C. Statistical Analysis

First we found the statistical data for the numerical features as follows:

- 1) For IMDb Votes the mean is: 79283.1397 vote
- 2) For Production Company the mode is "Others" which indicates that no one production company had their movie be in top 100,000 movie by revenue.
- 3) For Revenue the mean is: 7.099426e+07.

Using the correlation heat diagram Figure 1 we noticed the features that had a likely significant relation to revenue. Those features were: the budget and IMDb votes. Furthermore, we tested 9 hypotheses which are (each with its NULL hypothesis, method of testing and result):

1) H0: The lower the number of famous actors in a movie, the lower the revenue is.

H1: The higher the number of famous actors in a movie, the higher the revenue is.

Method of testing: using bar plot and checking the Peterson's Correlation Coefficient which was 0.187.

Result: We rejected the alternative hypothesis.

2) H0: The budget of a movie doesn't effect the number of famous actors in it.

H1: The budget of a movie does effect the number of actors in it.

Method of testing: the Peterson's Correlation Coefficient which was 0.141.

Result: We rejected the alternative hypothesis.

3) H0: If a movie's genre is Action, it has a lower average revenue, than if it was of another genre.

H1: If a movie's genre is Action, it has a higher average revenue, than if it was of another genre.

Method of testing: calculating the overall mean, and comparing it with the mean of action genre movies.

Result: The alternative hypothesis passed and we reject the null hypothesis

4) H0: If a movie is released outside of USA, it's revenue will be lower than revenue if released in USA only.

H1: If a movie is released outside USA, it's revenue will be higher than revenue if released in USA only.

Method of testing: the Mann-Whitney U test with p-value = 2.44e-10

Result: We reject the alternative hypothesis

5) H0: The Higher the budget of a movie the Lower the revenue is.

H1: The Higher the budget of a movie the Higher the revenue is.

Method of testing: the Peterson's Correlation Coefficient

which was 0.51.

Result: The alternative hypothesis passed and we reject the null hypothesis

6) H0: The Higher number of famous actors in a movie, the Lower its rating is.

H1: The Higher the number of famous actors in a movie, the Higher its rating is.

Method of testing: bar plot

Result: The alternative hypothesis passed and the null hypothesis failed.

7) H0: The Higher the number of people who voted for a movie, the higher its rating is.

H1: The Higher the number of people who voted for a movie, the lower its rating is.

Method of testing: the Peterson's Correlation Coefficient which was 0.34.

Result: We reject the alternative hypothesis.

8) H0: If a movie is produced by others, then it will have higher than average revenue.

H1: If a movie is produced by others, then it will have lower than average revenue.

Method of testing: the Mann-Whitney U test with p-value = 2.44e-10

Result: we reject the alternative hypothesis

9) H0: If a movie has more than 1 award, it has a lower revenue.

H1: If a movie has more than 1 award, it will have a higher revenue.

Method of testing: using bar plot and checking the Peterson's Correlation Coefficient which was 0.04.

Result: We reject the alternative hypothesis

### D. Machine Learning Model

We tried in this part to create a ML model that could classify - with acceptable accuracy - movies into successful and successful movies, and for it to predict if the movie will be so. The criterion chosen to whether a movie is successful or not if it's value is revenue is double it's budget. We first split the data-set into training and testing data-sets (.7, .3 respectively). We used Logistic Regression to achieve our goal. The label (y) was the column "IsSuccessful" which indicate, as the name states, whether a movie is successful or not. The features used for final model (after multiple tries) were budget and imdbVotes.

Confusion Matrix		Target	
Model		True	False
	True	497	366
	False	173	1318

## III. CONCLUSION

From the statistical analysis explained before we can conclude that: the higher the budget is, the more famous actors

	Precision	Recall	F1-score	Support
0.0	0.74	0.58	0.65	863
1.0	0.78	0.88	0.83	1491
Accuracy			0.77	2354
Macro avg	0.76	0.73	0.74	2354
Weighted avg	0.77	0.77	0.76	2354

are in a movie, and the more famous actors in a movie the higher a movie's rating is. The action genre generates higher revenues overall than other genres. Furthermore, the higher the budget is the higher the revenue.

#### REFERENCES

- [1] T. M. Database, "The movie database api," <https://www.themoviedb.org/>, May 2019.
- [2] B. Fritz, "Omdb api," <https://www.omdbapi.com/>, May 2019.