# Predicting Movie Success

DATS6103 Project Presentation
Amna Gul, Hemanth Koganti , Madhuri Yadav
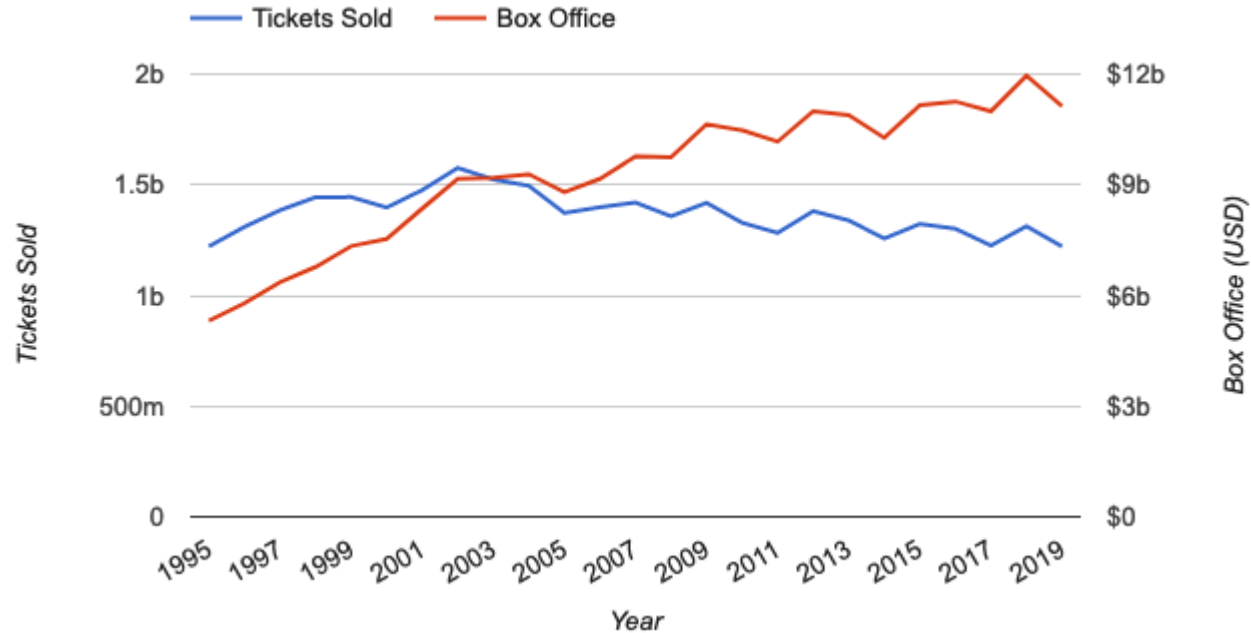
# Agenda

- ❏ Introduction - Problem Statement
- ❏ Data
- ❏ Preprocessing
- ❏ EDA
- ❏ Modeling
- ❏ Conclusion
- ❏ References

❏ Introduction

● Film industry is one of the top grossing industries in the world

❏ Problem Statement

- Billions of dollars are invested each year, expecting high margin of profit
- Whether is it possible to use machine learning algorithms to predict if movie will be a success?
  - Criteria for success
    - Movie is able to generate amount of revenue that is "greater" than the budget of the film

# Data Source

- Primary source of our data set is [Kaggle](#) (TMDB) till July 2017
- Columns were also added by parsing data made publicly available on [IMDb's website](#)
  - Recent data 2019

# Data Preprocessing/Cleaning

- Raw data: 45,000 rows and 24 columns
  - ➤ Remove irrelevant columns e.g. "home page", "poster path"

| poster_path |
| --- |
| /rhlRbceoE9lR4veEXuwCC2wARtG.jpg |
| /vzmL6fP7aPKNKPRTFnZmiUfciyV.jpg |
| /6ksm1sjKMFLbO7UY2i6G1ju9SML.jpg |
| /16XOMpEaLWkrcPqSQqhTmeJuqQl.jpg |
| /e64sOI48hQXyru7naBFyssKFxVd.jpg |
| /zMyfPUelumio3tiDKPffaUpsQTD.jpg |

  - ➤ Searched for corrupt values in remaining columns e.g. "budget" containing alpha-numeric values

| budget |
| --- |
| 300000000 |
| 260000000 |
| 260000000 |
| /zV8bHuSL6WXoD6FWogP9j4x |
| /zaSf5OG7V8X8gqFvly88zDdRr |
| 260000000 |

# Data Preprocessing/Cleaning

➢ Columns "Genre" and "Production_Companies" were in JSON format



➢ Converted columns to their proper data type e.g. "release_date" was converted date-time format instead of string. Month was extracted to create a separate column.

➢ Merged "average_rating" and "vote_count" from IMDb's website
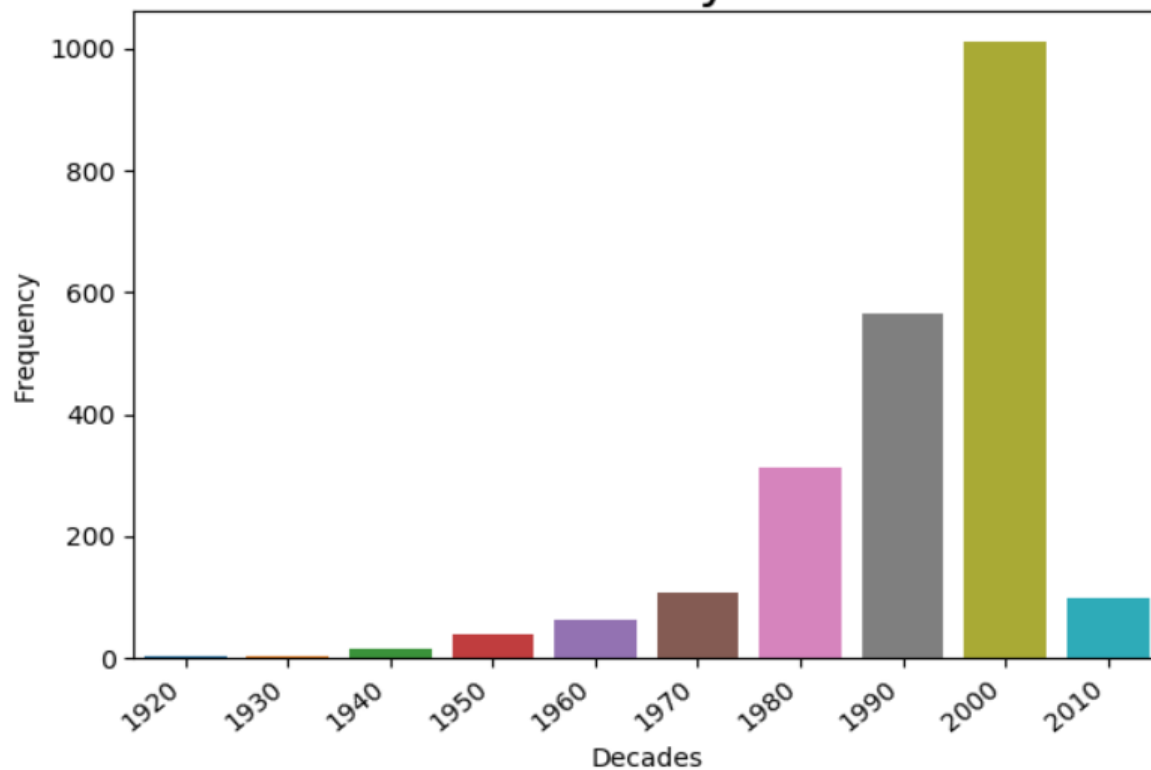
# Data Preprocessing/Cleaning

- ➢ Outliers: excluded all those rows for which budget or revenue value was unrealistically low
- ➢ Created our target column by dividing revenue by budget
- ➢ Removed duplicates
- ➢ End result: 2,222 rows & 17 columns

# EDA
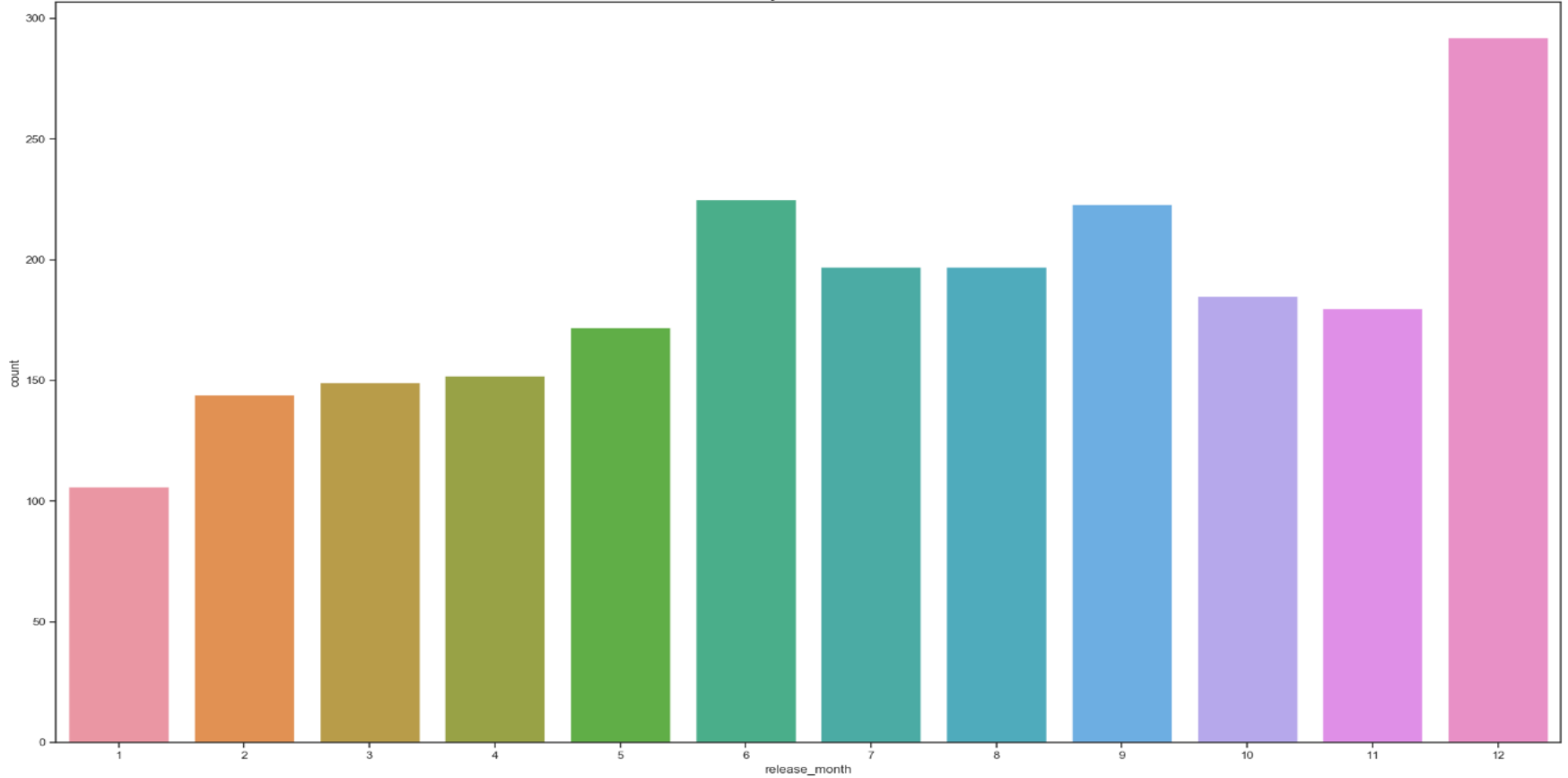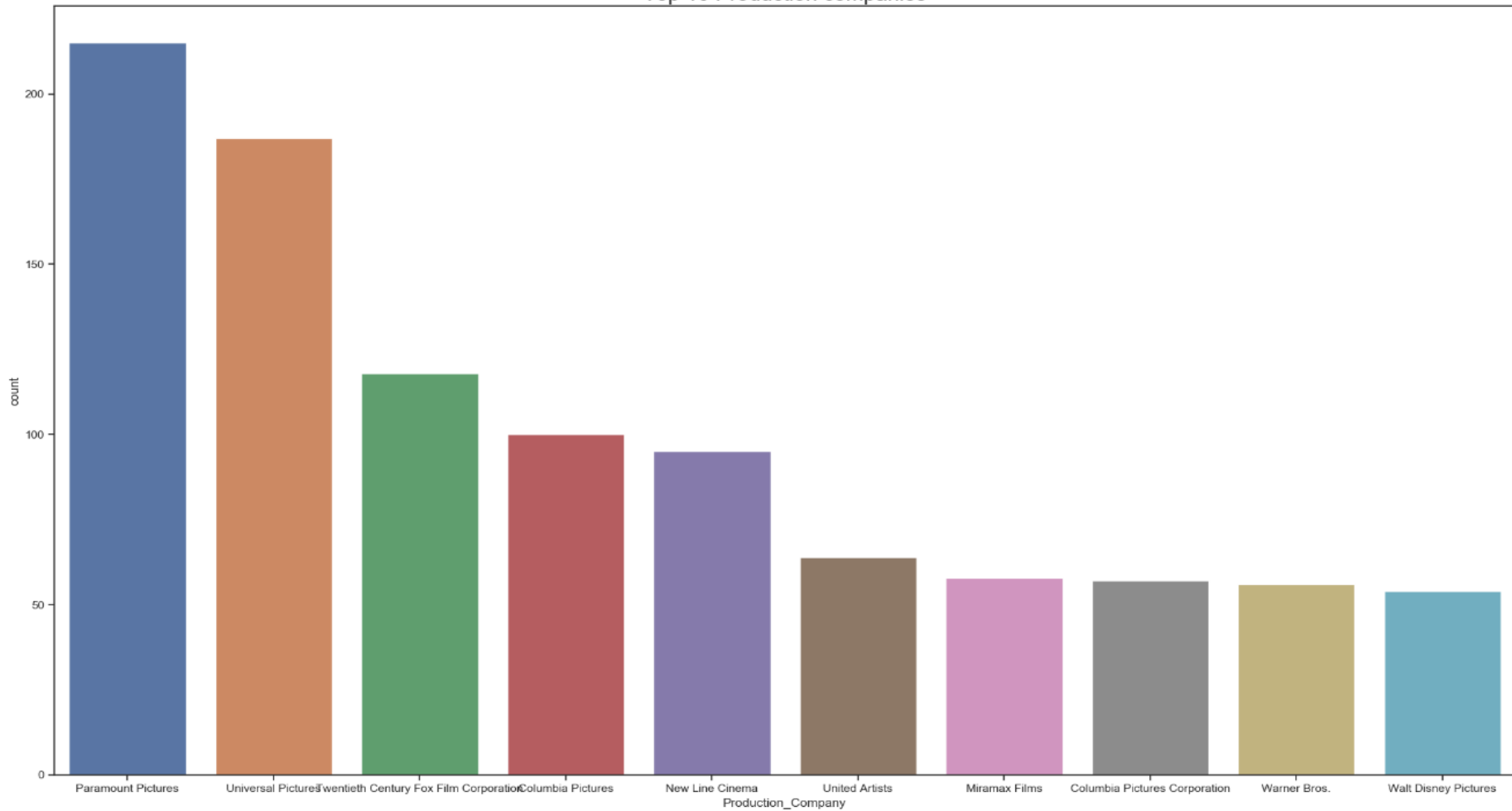


Movie Count by Decades

Movie Count by Genre
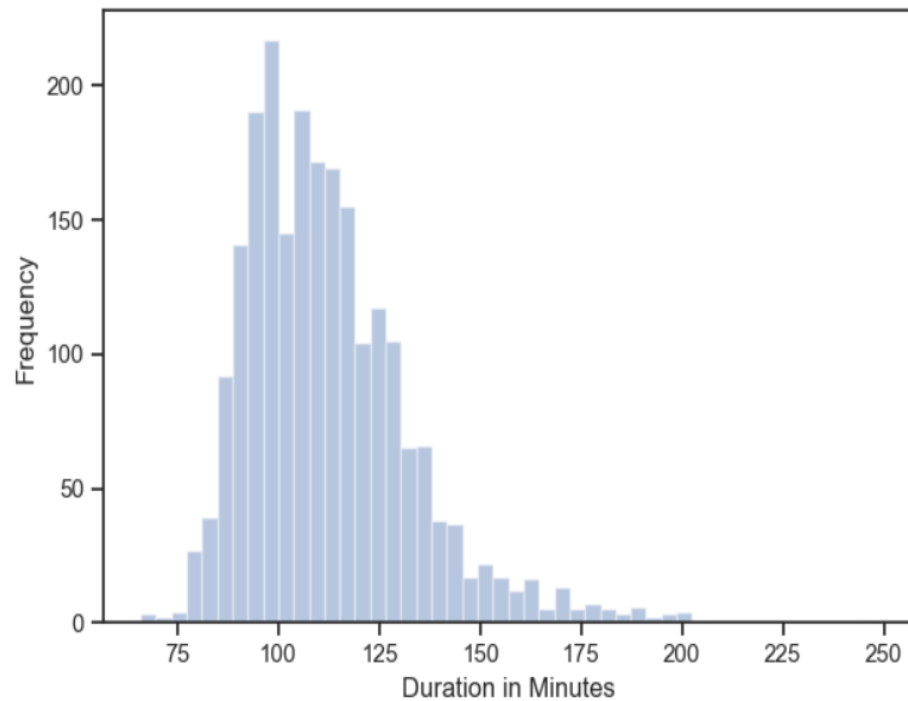
Movies by Release month

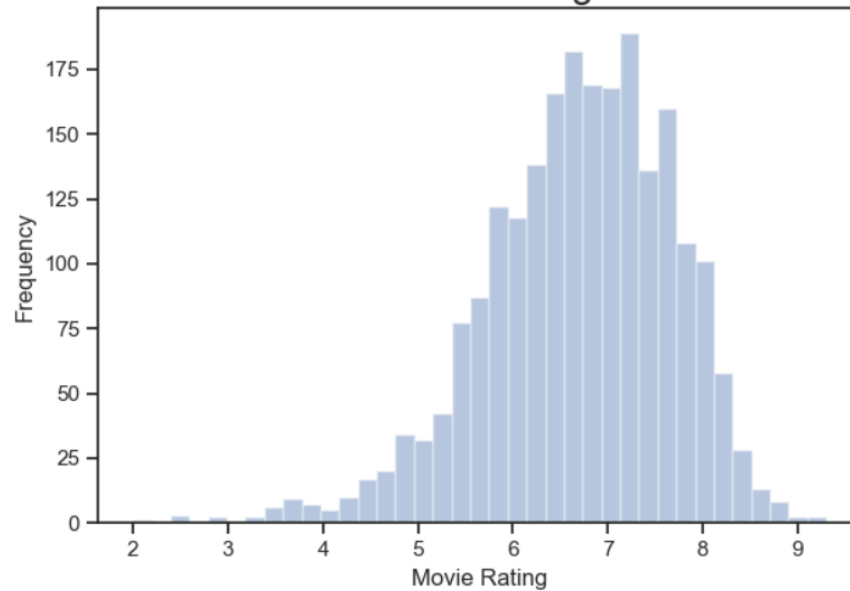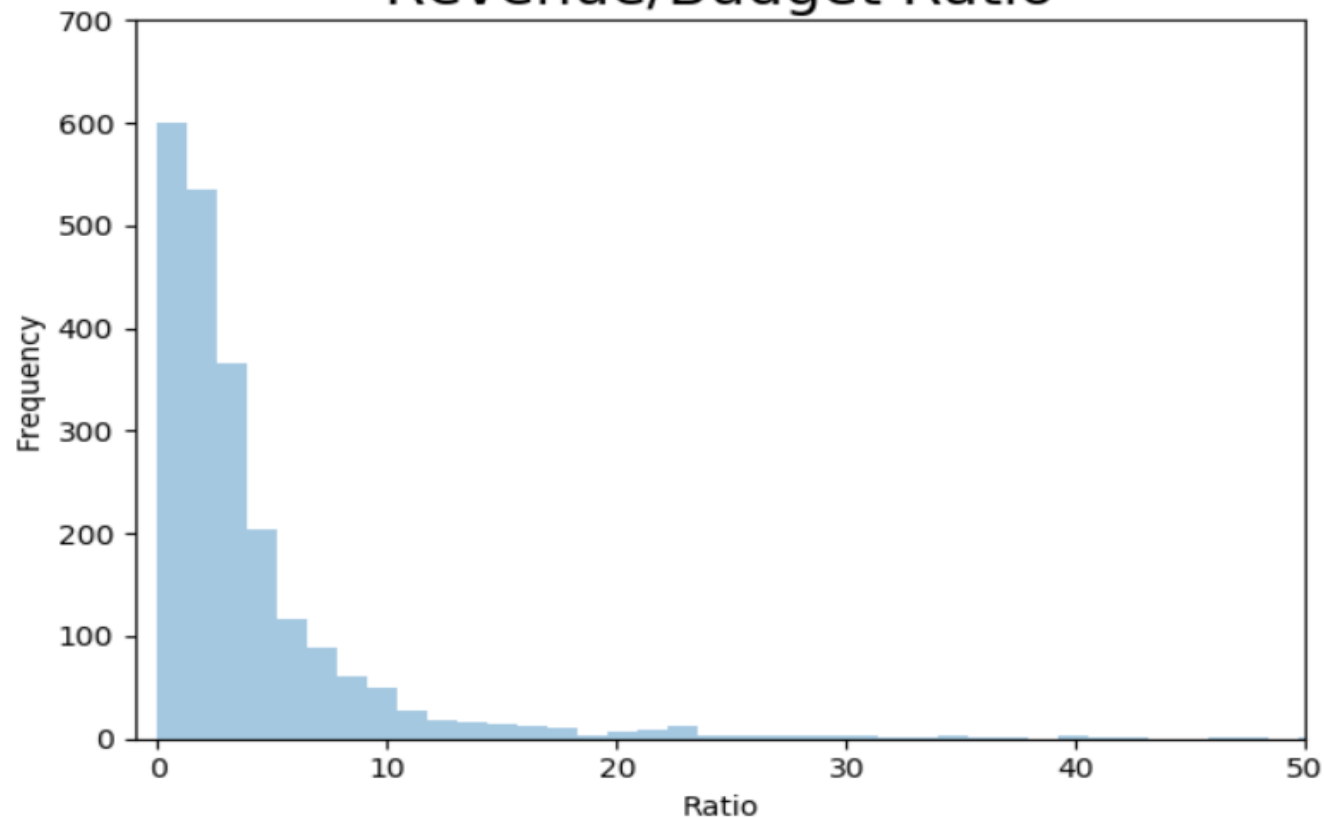Top 10 Production companies

Revenue/Budget Ratio
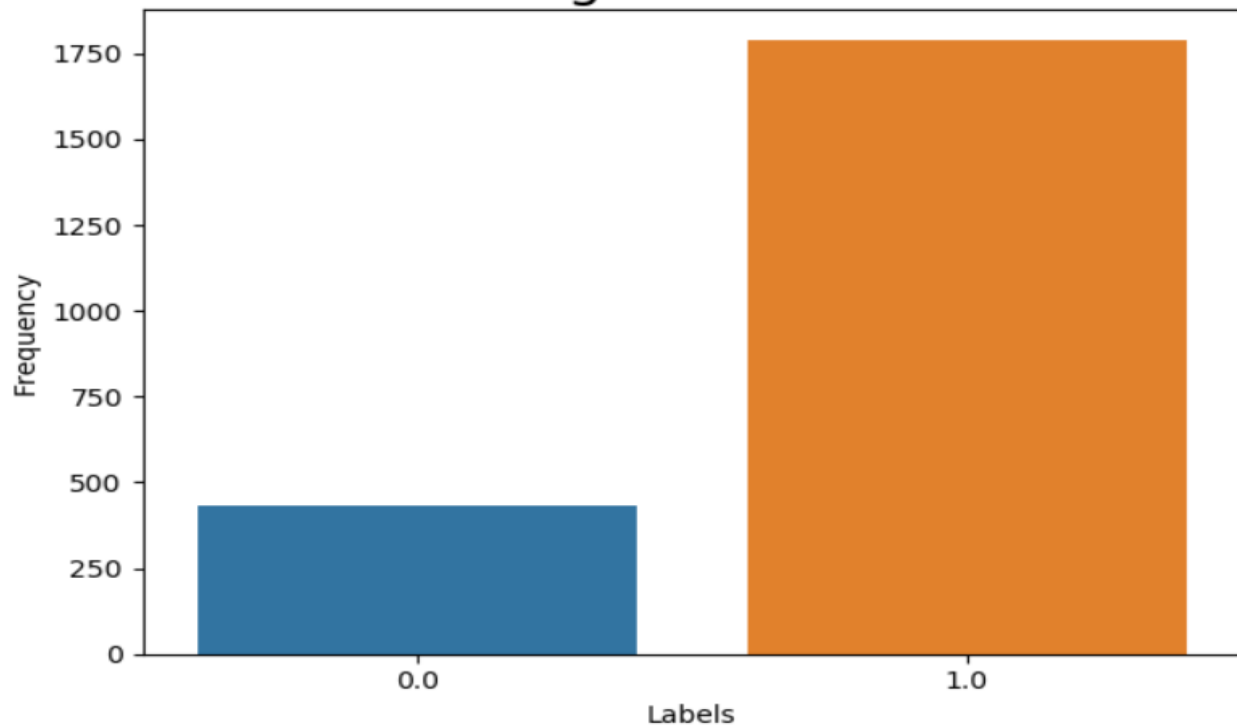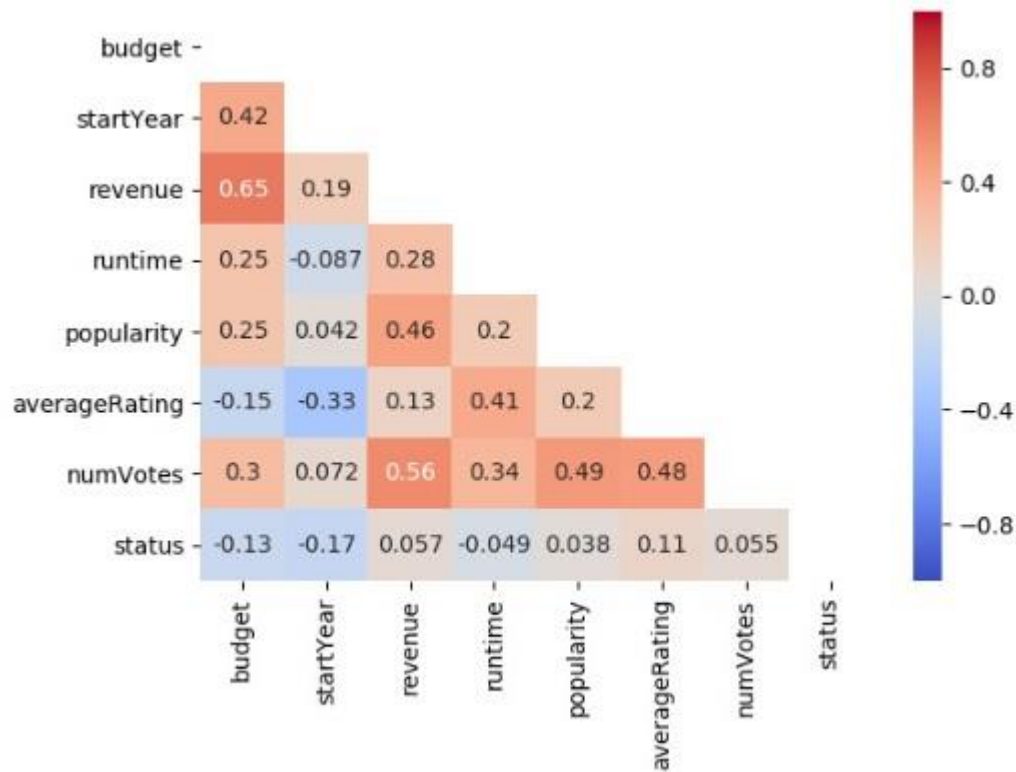
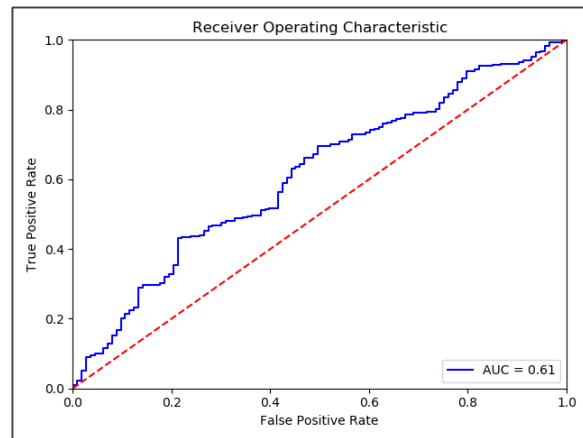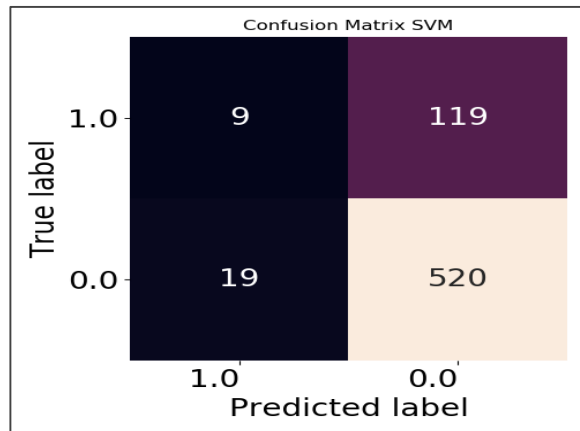Target Variable

# Correlation Heatmap

# Modeling

1. Data Splitting
2. Stratified Sampling
3. Oversampling
4. Label Encoding
5. Scaling
6. Feature selection
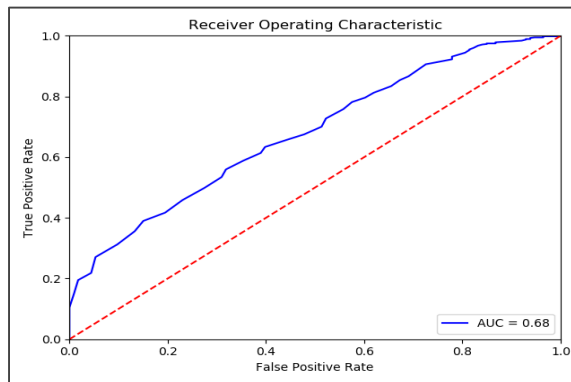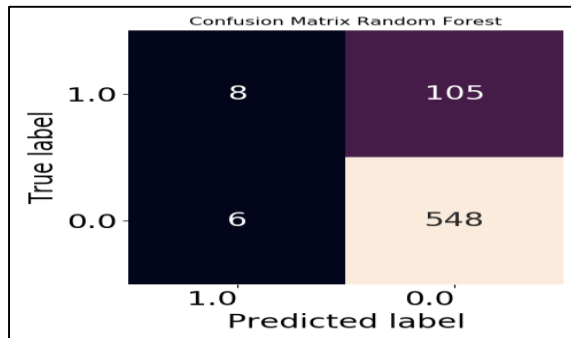
# Continued...

Results with four features:

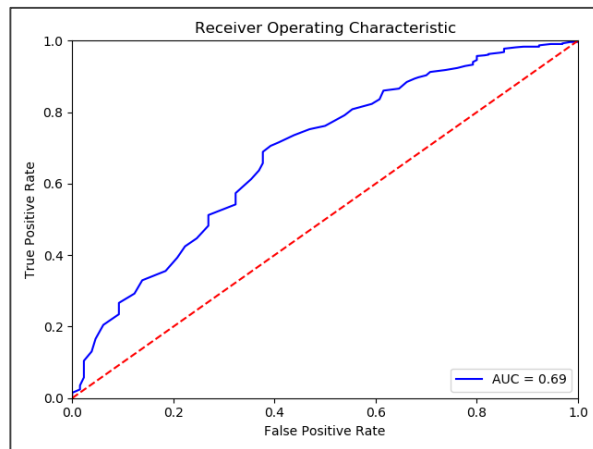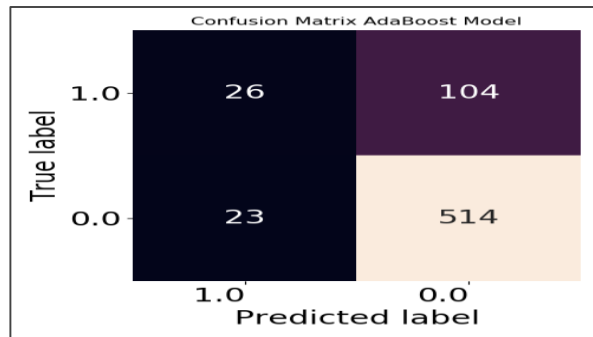| Model | Accuracy |
|---|---|
| Accuracy DT Entropy | 69.26% |
| Accuracy SVM | 81.40% |
| Accuracy RF | 80.35% |
| Accuracy KNN | 71.81% |
| Accuracy NB | 22.93% |





Area Under the Curve = 61%

# Continued…



| Model | Accuracy |
|---|---|
| Accuracy DT Gini | 73.91% |
| Accuracy SVM | 81.41% |
| Accuracy RF | 83.35% |



Area Under the Curve = 68%

# Continued…

| Model | Accuracy |
|---|---|
| Decision Tree (Entropy) | 72.56% |
| Support Vector Machine | 69.12% |
| Random Forest | 79.91% |
| Bagging(Mode) | 78.41% |
| Adaptive Bootstrap | 80.96% |



Confusion Matrix AdaBoost Model



Receiver Operating Characteristic

Area Under the Curve = 69%
K = 20.57%

# Limitation

- Biased Label
- Missing Values
- Invalid Data
- Additional Features

# References

https://www.the-numbers.com/market/

http://www.diva-portal.org/smash/get/diva2:1106715/FULLTEXT01.pdf

https://io9.gizmodo.com/how-much-money-does-a-movie-need-to-make-to-be-profitab-5747305

Maklin, C. (2019). *AdaBoost Classifier Example In Python*. [online] Medium. Available at: https://towardsdatascience.com/machine-learning-part-17-boosting-algorithms-adaboost-in-python-d00faac6c464