Thomas Preston
DSC680


**Predicting Movie Success**

Introduction:

Predicting whether a movie will succeed before it hits theaters has always been a tricky task for the film industry. Historically, studios and producers mainly relied on star power and experience, but these approaches don't always lead to perfect results. However, with advancements in data science and machine learning, we now have the chance to learn from historical data and statistical models to make better predictions about a film's potential success.

In this project, I used a dataset containing over 5,000 films, using machine learning techniques to predict a movie's performance based on factors like budget, genre, runtime, and release date. I primarily measured success through box office revenue, and I also categorized movies as "hits" or "flops" based on specific revenue data. By understanding these trends, studios will be able to refine their production and marketing strategies and possibly save millions in budgets and improve their overall profits.
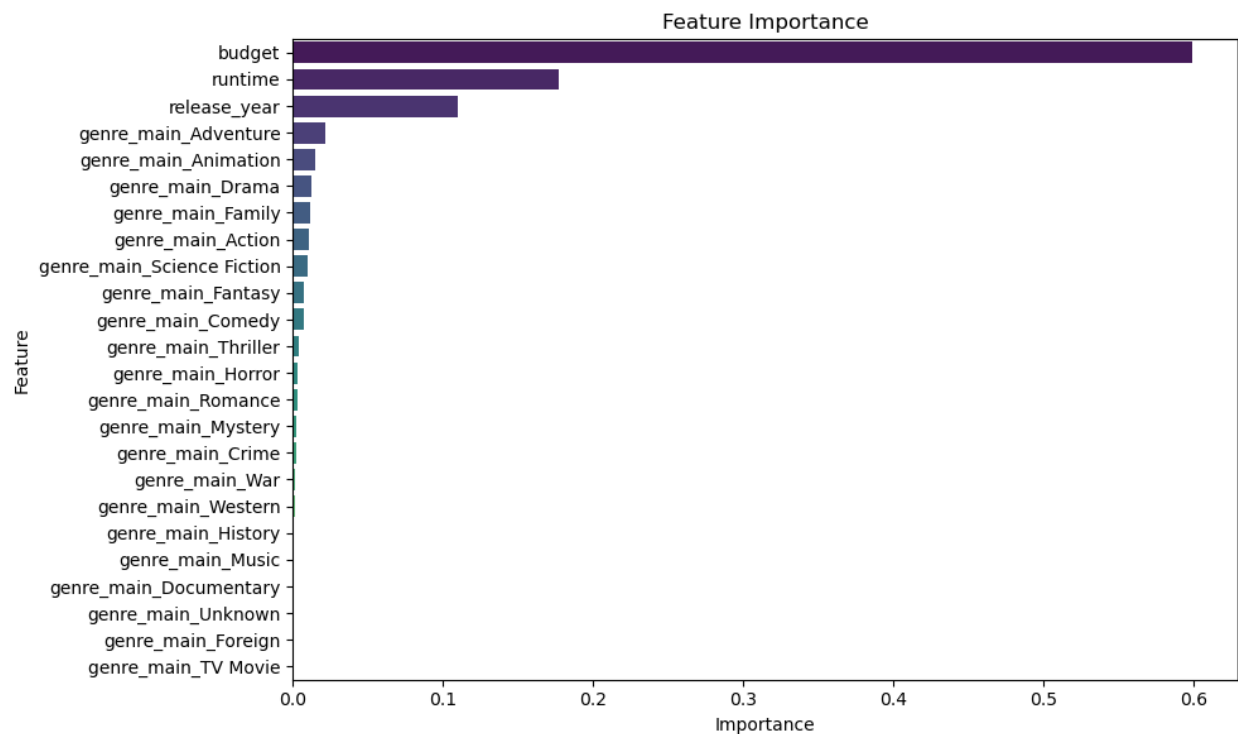
Data & Methods:

The dataset used for this project is the (TMDB 5000 Movie Dataset, 2017), which includes detailed information on thousands of movies spanning multiple decades. The data contains important categories such as budget, revenue, genres, release dates, and runtime. Initial data cleaning involved removing entries with missing or zero values in critical columns like budget and revenue, as these would most likely skew the models training and evaluation. The genre information, originally stored as a list of multiple genres per movie, was simplified to the primary genre to reduce complexity. I changed the release dates from strings to datetime objects, which made it easier to pull out the release years as a numerical feature. For modeling, I used two approaches. A random Forest Regressor was trained to predict continuous revenue values and a Random Forest Classifier was used to categorize movies as hits or flops. Hits are defined as movies earning over 100 million dollars in box office revenue. I chose these models because they're really reliable and can handle nonlinear relationships well, without needing a lot of complicated data preparation.
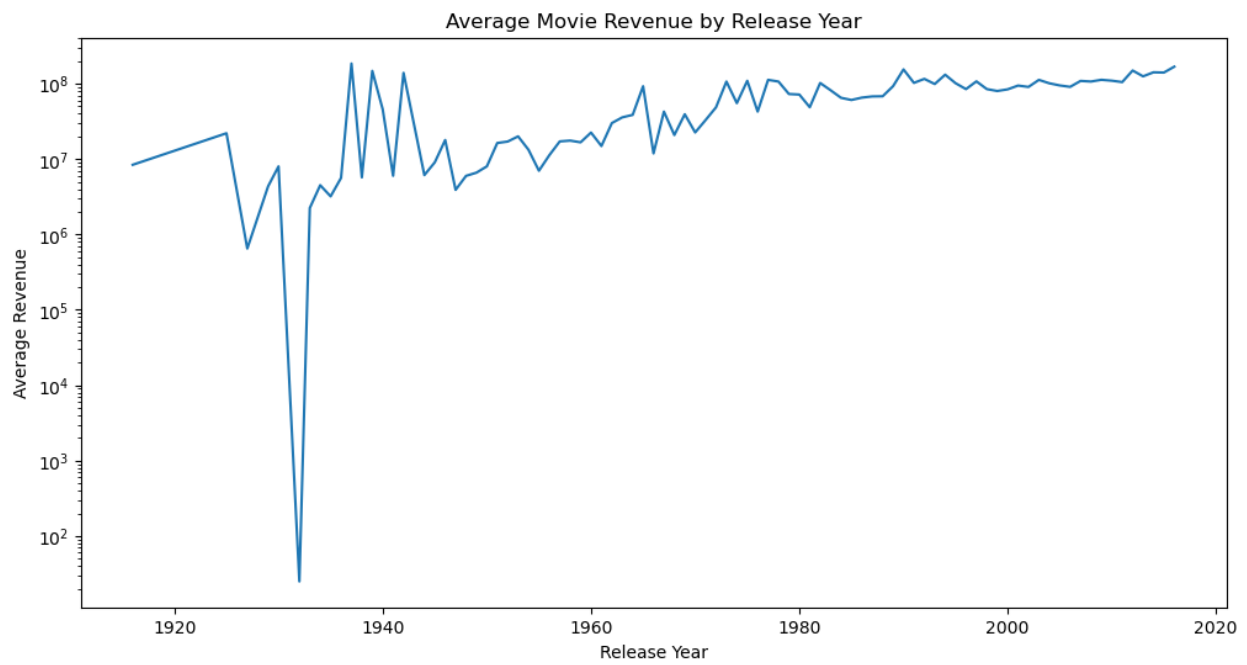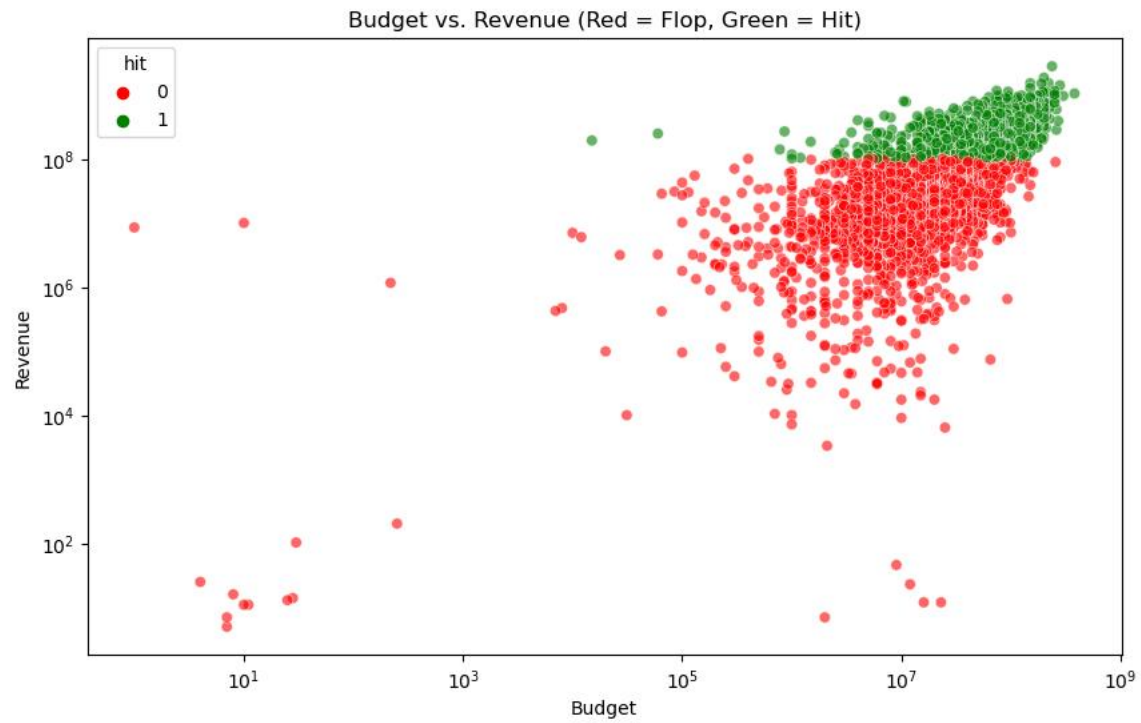
Results:

The regression model turned out very nicely, with an RMSE of around 139 million. This makes sense because movie revenues can vary a lot from a few thousand dollars to millions. But overall this means it can predict movie revenues pretty accurately, which is impressive given how tricky that can be. On another note, my classification model did a solid job of distinguishing hits from flops, achieving precision and recall scores of 0.75 and 0.59. This means my model did a good job of spotting which movies are likely to be hits, while also avoiding too many false positives. When I looked at what factors mattered most, budget stood out as the most

Thomas Preston
DSC680

important, which makes total sense since we all know that bigger budgets can lead to bigger box office returns. But it's not just about the budget, things like runtime and genre also made a big difference. So, it looks like the length of a movie and its category really do matter when it comes to finding success at the box office.

Visualizations further supported these findings, the scatter plot of budget versus revenue showed a positive correlation, while the hits versus flops count plot illustrated the difference between successful and unsuccessful movies in the dataset. Additional plots like revenue distribution and average revenue over release years provided perspective about trends and data skewness.

Thomas Preston
DSC680

## Budget vs. Revenue (Red = Flop, Green = Hit)



## Average Movie Revenue by Release Year

Thomas Preston
DSC680

**Runtime vs Revenue**



Discussion:

These results suggest that while budget remains the main predictor of success, other factors contribute majorly to the outcome, emphasizing the complicated nature of the movie industry. The model's limitations should also be acknowledged. Important variables such as marketing budget, star and director power, critical reviews, and audience response were not included due to data availability limitations, yet these could greatly influence a movie's performance. The dataset itself may also have biases. For example, with older movies, there could be issues like missing financial data or different market conditions that can throw off our model's accuracy. Also, deciding what revenue level counts as a "hit" is a bit subjective and could definitely be adjusted based on the genre or adjusted for inflation. Looking ahead, I could enhance the approach by incorporating more data sources, trying out more machine learning methods like gradient boosting or neural networks or taking a closer look at feature engineering to really understand the complex relationships involved. Also, by using cross-validation and tweaking the model's settings, I could make it even more dependable.

Conclusion:

This project really shows how powerful data can be in predicting whether a movie will be successful. It offers some great insights for filmmakers and studios looking to make smarter choices. By focusing on important features and using machine learning models, it shows that we can get a pretty good idea of how a movie will perform financially and can even sort them into hits or flops. Of course, there are some limitations and there is plenty of room to grow, but these models can be great tools to support the traditional methods in the entertainment world.

Thomas Preston
DSC680

As we continue to develop and bring in a wider variety of data, we can make these predictions even more accurate and useful. Ultimately, this can help the film industry minimize risks and boost profits.

References:

TMDB 5000 Movie Dataset. (2017). Kaggle. Retrieved from
www.kaggle.com/datasets/tmdb/tmdb-movie-metadata

Questions:
1. How accurate is your model when it comes to predicting movie revenue and figuring out which ones are hits or flops?
2. Why did you decide on $100 million as the cutoff for what counts as a "hit"?
3. What role does genre play in how successful a movie is, according to your model?
4. How does your model deal with changes in the movie industry over the years?
5. Did you think about including other factors like marketing budgets, star power, or critical reviews?
6. What kind of challenges did you encounter while cleaning and prepping the dataset?
7. Looking ahead, how do you think you could boost the model's performance in the future?
8. How do you think studios or investors could actually use this model in real life?
9. What do you see as the limitations of using past data to predict whether a movie will be successful?
10. Did you check out any other machine learning algorithms besides Random Forests?

Answers:
1. The Random Forest Regression model had an RMSE of about 139 million. This is reasonable given how much a movie can range so the model did well in identifying hits versus flops. It was more accurate at predicting flops though because there were more in the dataset.
2. I looked it up and it is a commonly used industry benchmark defining film success. It made it easy to split the classification as well.
3. It had an impact but it wasn't the most important factor. Budget, runtime and release year were much more important and better predictors.
4. Using release year helps account for shifts in trends over time. But it's hard to truly account for some new variables.
5. Yes, those are very important predictors but the dataset did not include them. They would be great for future research.
6. Missing or zero values was the main one. A budget listed as zero is not usable. Also the genre column was stored as strings and needed to be extracted into binary. Dates were in various formats as well.
7. Using richer or more data sources. Other model types like neural networks could help on the hit/flop side.

8. It could be a good support tool. You could get rough forecasts using budget, genre and release window. But it is not a replacement for expert judgement.
9. The movie industry evolves very quickly. What worked a decade ago might not work now.
10. Yes, I also tested out Linear Regression for revenue prediction. It performed worse than Random Forest.