

# Using Data Mining to Predict Popularity and Market Success of Movies

Harpreet Singh, Jian Wu, Xinhai Zhou, Zelan Xiang

Department of Computer Science, University of Victoria

Victoria, BC, Canada V8P 5C2

**Abstract** - This document is our proposal to start a research project on applying data mining techniques and methods to predict popularity and market success of a specific movie with some given features such as production company, genre, cast, vote count etc.

## Introduction

Even though movies are an artistic endeavour and not purely a way to make money, people involved in production of movies dream of a solution that can help them predict whether a movie is going to be successful before they spend a lot of effort and money making it.

We propose a project using algorithms and techniques of data mining to predict different measures of a movie's success before its production. As part of this project we would analyze performance of several models on our dataset, and then write the best performing model ourselves from scratch. We would also build a website with a UI to pick different characteristics of a movie like Director, Budget, Genre, etc and then provide a prediction on the different measures of success i.e. Imdb rating, popularity and gross profit.

We all are students from the Department of Computer Science at the University of Victoria. There are four members in our team:

- Harpreet Singh (Grad., Email: [hsingh91@uvic.ca](mailto:hsingh91@uvic.ca)),
- Jian Wu (Grad., Email: [wujian@uvic.ca](mailto:wujian@uvic.ca)),
- Xinhai Zhou (Undergrad., Email: [zhouxinh@uvic.ca](mailto:zhouxinh@uvic.ca)),
- Zelan Xiang ( Undergrad., Email: [zelan@uvic.ca](mailto:zelan@uvic.ca))

## Related Work

There has been a substantial amount of research on this topic. MIAS [1] provides early predictions of movie profitability. Ramesh Sharda uses neural networks to predict financial performance of a movie but considers the problem to be a classification problem rather than forecasting a point estimate [2]. Márton Mestyán and Taha Yasseri measured and analyzed the activity level of editors and viewers of the corresponding entry to the movie in Wikipedia and built a multivariate linear regression model for predicting box office revenue [3]. Kyuhan Lee and Jinsoo Park use a non conventional approach in which they add a new feature that incorporates transmedia storytelling and apply an ensemble approach. As a result of this they propose a new model which outperforms models from previous studies [4].

### **Data Source**

Instead of scraping the TMDB website to obtain the dataset, we will use the existing kaggle TMDB 5000 movies dataset [5]. This dataset was generated from The Movie Database API by kaggle, but is not endorsed or certified by TMDB. The dataset include two data, movie metadata and movie credits data. The movie metadata attributes include: budget, genres, homepage, id, keywords, original-language, original-title, overview, popularity, production\_companies, production\_countries, release\_date, revenue, runtime, spoken\_languages, status, tagline, title, vote\_average, and vote\_count; the movie credits data attributes include: movie\_id, title, cast, and crew.

### **Methodology**

We will first collect movie data as our dataset from kaggle, and we will separate the set into training set and testing set according to the knowledge learned from the data mining course. Then we apply different existing classification and regression models in Scikit-learn [6] with our training data set. After comparing the results from different models in Scikit-learn, we will choose the one with the best performance, and then try to implement the algorithm in Python from scratch. Finally, we will build a UI that allow the user to input features of a movie and get the predicted result.

Overall, in this project we will use tools such as Python, Scikit-learn, Flask for backend and Angular [7] for the webapp.

### **Estimated Time Line**

Our project starts on Oct 9th, 2017 and is estimated to finish by Nov 28th, 2017. Here is the estimated time line of our project:

Oct 13th, 2017: Decide the group members and the topic of the project

Oct 18th, 2017: Detailed tasks assigned to each team member

Oct 24th, 2017: Finish data collection and pre-processing

Nov 7th, 2017: Come up a classifier prototype

Nov 14th, 2017: Finish preliminary analysis of the data and improve the classifier

Nov 20th, 2017: Finish the improvement of the classifier

Nov 23th, 2017: Done with data analysis

Nov 26th, 2017: Finish the final report

Nov 28th, 2017: Final presentation

## References

- [1] M. Lash and K. Zhao, "Early Predictions of Movie Success: The Who, What, and When of Profitability", *Journal of Management Information Systems*, vol. 33, no. 3, pp. 874-903, 2016.
- [2] R. SHARDA and D. DELEN, "Predicting box-office success of motion pictures with neural networks", *Expert Systems with Applications*, vol. 30, no. 2, pp. 243-254, 2006.
- [3] M. Mestyán, T. Yasseri and J. Kertész, "Early Prediction of Movie Box Office Success Based on Wikipedia Activity Big Data", *PLoS ONE*, vol. 8, no. 8, p. e71226, 2013.
- [4] K. Lee, J. Park, I. Kim and Y. Choi, "Predicting movie success with machine learning techniques: ways to improve accuracy", *Information Systems Frontiers*, 2016.
- [5] "TMDB 5000 Movie Dataset | Kaggle", *Kaggle.com*, 2017. [Online]. Available: <https://www.kaggle.com/tmdb/tmdb-movie-metadata>. [Accessed: 14- Oct- 2017].
- [6] P. Fabian, et al. "Scikit-learn: Machine learning in Python." *The Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, 2011.
- [7] "Angular Docs", *Angular.io*, 2017. [Online]. Available: <https://angular.io/about?group=Angular>. [Accessed: 14- Oct- 2017].