
Selecting Movie Features to Predict Monetary Success and Popularity

Kamarion Porter
Harvey Mudd College
kporter@g.hmc.edu

Abstract

1 In this paper, I seek to discover the similarities and differences between monetary
2 success and critical acclaim in the film industry. To accomplish this, I first find
3 two columns in the IMDB Movies Dataset to represent monetary success and
4 popularity, taking these to be possible output variables. Next, I extract features
5 from the remaining columns; the column that is not being treated as the output
6 variable is included in this group. Then, I use the chi-square statistic in order to
7 test how closely related the features are to the output variable. I also make logistic
8 regression and Naive Bayes models using these features. Finally, I perform a
9 comprehensive analysis on the features and models for each output.

10 **1 Introduction**

11 From their humble beginnings as reels of sewn-together film, motion pictures have grown to become
12 a staple of entertainment in the modern day. Its crafting has become an art form that is cherished
13 by many. However, these films have also cultivated a billion-dollar industry and the goals of the
14 companies at its helm center around making a profit. Sometimes this goal pairs well with providing
15 the audience with a memorable experience, and other times one is sacrificed for the sake of the other.
16 Using the IMDB Movies Dataset, which boasts over 2.5 million movies and TV series, I seek to
17 discover the similarities and differences between monetary success and critical acclaim in the film
18 industry.

19 **2 Related Work**

20 Much work has got in predicting the attributes of motion pictures. Latif and Afzal (2016) [4]
21 applied feature selection and various classifiers, including logistic regression, Naive Baynes, and
22 PART, upon data that they extracted from IMDb in order to predict movie popularity. I also use IMDb
23 data to perform feature selection and classification, however, I employ multiple classifiers to compare
24 multiple outputs. Abidi et al. (2020) [2] extracted a small sample of IMDb data and performed
25 feature selection to find the best inputs and outputs features within a movie for predicting the movie's
26 popularity. I also attempt to predict a movie's popularity using through its features, however, I work
27 with a large dataset featuring movies and TV series. Dhir and Raj (2018) [3] utilized correlation
28 analysis and several machine learning classifiers ob the IMDb 5000 Movie Dataset in order to predict
29 a movie's success, which they represent by the IMDb rating (or score as they reference it). Along
30 with the different strategies that we employ and outputs that we focus on, we differ in that they rely
31 on the IMDb rating as a reliable measure of success, which I disagree with.

32 **3 Procedure**

33 **3.1 Preprocessing**

34 **3.1.1 Choosing Output Variables**

35 My initial intention was to use two of the set's columns, 'gross_income' and 'rating', as my
36 measures of monetary success and popularity respectively. IMDb records gross income as the profit
37 (the total revenue - the total cost) earned by a motion picture. As any net loss results in no profit,
38 net losses are recorded as zero. IMDb's rating system is also simple. Any individual who visits the
39 site can rate a choice film on an integer scale of 1-10. The average of all of the individual ratings
40 becomes the movie's overall rating. In this dataset, the number of people who rated that movie is
41 stored in the column 'votes'. Because two films can have the same rating but a drastically different
42 number of votes, a film's rating, on its own, cannot be used as a measure of popularity.

43 However, I was able to create such a statistic using the 'rating' and 'votes' columns called the
44 adjusted rating. The adjusted rating is the product of the film's rating - 5.5 and the natural log of the
45 number of votes. The rating is initially subtracted by 5.5 as 5.5 lies at the center of the rating scale,
46 thus any motion pictures rated higher than 5.5 have an overall positive rating, those lower than 5.5
47 have an overall negative rating, and those at 0 have an overall neutral rating. Votes are a factor within
48 this derivation to account for the popularity/infamy of the film. We then multiply this value by the
49 natural log of the number of votes for the difference between 10 and 100 votes to be more notable
50 than the difference between 10,000 and 10,090 as the former displays a much bigger difference in
51 how well-known the corresponding films are.

52 **3.1.2 Data Filtering**

53 According to Kaggle.com [1], the IMDB dataset, developed by Ashish Jangra, contains over 2.5
54 million movies and TV series as well as a usability rating of 8.2. However, as I delved further into
55 the dataset, I discovered that the data did not exist as advertised. To start, the year column contained
56 not just the motion picture's start and end year, but also the type of media it was. From this, I learned
57 that in addition to motion pictures (movies and TV series), this dataset contained podcast series,
58 video games, and music videos. Furthermore, within the remaining group of true motion pictures,
59 the vast majority of them had both no gross income and no votes, meaning that the film made no
60 money (potentially resulting in a net loss) and had no ratings on IMDB. After filtering the data
61 once more, I was left with only around 550 thousand shows (many of which had no gross income),
62 however, they all had at least 5 votes. This is reassuring as for a motion picture to have any sort of
63 popularity or monetary success, it has to be known, so the fact that within our dataset, films only
64 have a positive gross income (our measure of monetary success) if they have votes (a part of our
65 measure of popularity) reaffirms the usability of the dataset for this project. The remaining dataset
66 still contained some invalid entries in the form of invalid years (a placeholder of 2100 was used) and
67 invalid ratings (a joke rating of 11.0 was applied to certain shows), however, these made up a very
68 small portion of the data. With this, I concluded the data filtering process.

69 **3.1.3 Normalization and Discretization**

70 One final preprocessing step that I had to take was normalizing and discretizing the columns. A
71 good amount of the columns had unique numerical formats (some with mixed types) and others
72 had multiple categories ('genre', 'director_id', 'stars_id'). Because I had a mix of categorical and
73 numerical data, I discretized the numerical columns by making ordered categories that corresponded
74 to numerical ranges. For each column, the ranges were selected such that there was a proportionate
75 amount of data in each category, and the categories were semantically meaningful. For example,
76 the duration column contains the runtime for a movie, the number of minutes in an episode for an
77 ongoing TV series, and the total number of minutes for a completed TV series. After discretization,
78 the column's categories became hour-long windows, with anything 4 hours and above grouped
79 together. Since the runtime for most TV episodes is under an hour, movies are typically 1-4 hours

80 long, and most TV series contain at least 4 hours of content, these categories fit nicely to these media
81 types (although not one-to-one). This also leaves room for separation with regards to shorter and
82 longer movies.

83 **3.1.4 Issues with Multi-Labeled Data**

84 For the multi-labeled data, I first tried to make a column that corresponded to each label. Unfor-
85 tunately, this was only plausible for the ‘genre’ column as the other columns (‘directors_id’ and
86 ‘stars_id’) had too many unique values for this to be feasible. I then tried to use the first label within
87 each entry. Since the labels were sorted by the corresponding director or star’s placement in the
88 credits, this would’ve been an analytically useful feature. However, the caveat to this was that because
89 these columns contain unordered categorical data, I had to encode these columns using a one hot
90 encoding. A one-hot encoding makes a one-to-one mapping between a label and a vector that is the
91 length of all of the unique labels in the set. This resulted in hundreds of thousands of columns being
92 added to the encoded training set. As my laptop could not process this in a reasonable amount of
93 time, I was forced to drop these columns.

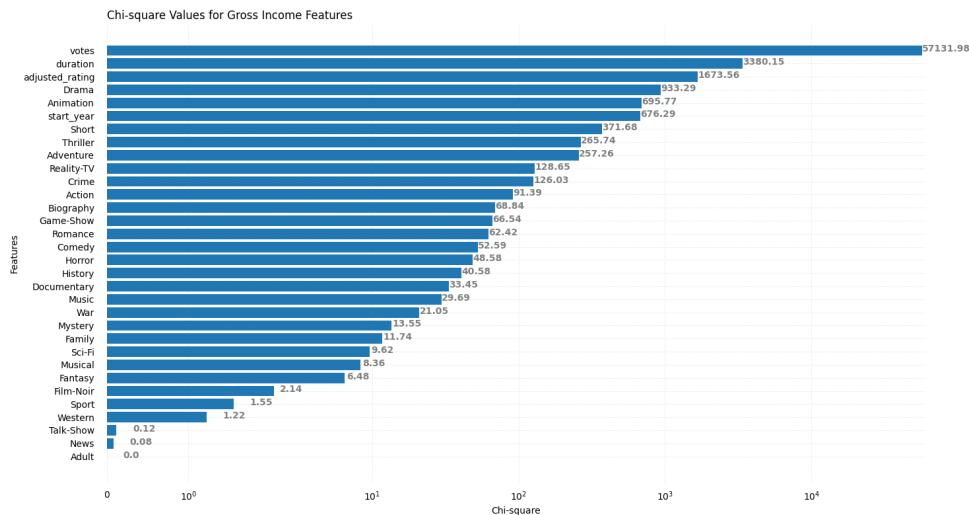
94 **4 Creating Metrics and Models**

95 All in all, the features that I was able to extract from the dataset were ‘gross_income’, ‘ad-
96 justed_rating’, ‘duration’, ‘votes’, ‘start_year’ (the year in which the show debuted), and all of the
97 genre categories, which (besides the genres) are all ordered variables. Encoding the ordered variables
98 was much easier and less time-consuming than the unordered variables. Using two-thirds of the
99 data as the training set and one-third as the test set, I was able to encode both sets using an ordinal
100 encoding, which mapped each label to a single integer. With the features prepared, I measured the
101 chi-square statistic for each feature. For the 3 features with the best chi-square values, I made a
102 multi-bar chart of the percent frequency of the output variable categories (gross income or adjusted
103 rating) over the feature variable categories. For example, the percentage of entries in the dataset
104 in which the gross income is less than \$1,000 given that the duration is less than an hour would
105 be a single bar in a plot if duration was a top feature. And finally, using either all of the encoded
106 features or the three with the best chi-square scores, we fit either a logistic regression or Naive
107 Bayes classifier to the training data, predicted the output values using the testing data, and then
108 evaluated the strength of this prediction by calculating the model’s accuracy, balanced accuracy,
109 recall, precision, and F1-score. Besides accuracy and balanced accuracy, these scores are calculated
110 using a macro-average, which means that the given statistic is calculated for each output category,
111 and the average of their percentages is taken as the value. This provides more accurate data regarding
112 each model’s effectiveness on imbalanced datasets, such as this one. The balanced accuracy is the
113 average of the micro-averaged recall (the more standard definition of recall) and the true negative rate
114 (the proportion of true negative values over all of the negative values).

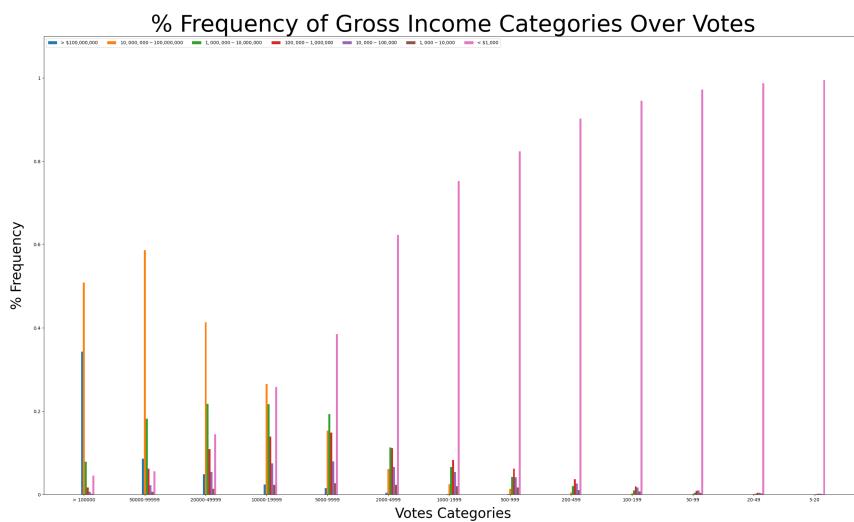
115 **5 Results**

116 **5.1 Chi-Square-selected Features**

117 As we can see from horizontal bar plot below, the top 5 features for gross income are votes,
118 duration, adjusted rating, the Drama genre, and the Animation genre, all with large chi-square scores,
119 although votes has the largest by a wide margin. The worst 5 features for gross income are the Adult
120 genre, the News genre, the Talk-show genre, the Western genre, and the Sports genre, with a large
121 jump between the worst 3 features and the rest. Underneath the chi-square plot lies a multi-bar plot
122 of the percent frequency of the gross income categories over the categories of its top feature, which
123 in this case is votes. For example, the chance that a movie in this dataset has a gross income under
124 \$1,000 given that it has 100-199 votes is the value of a bar in this graph.

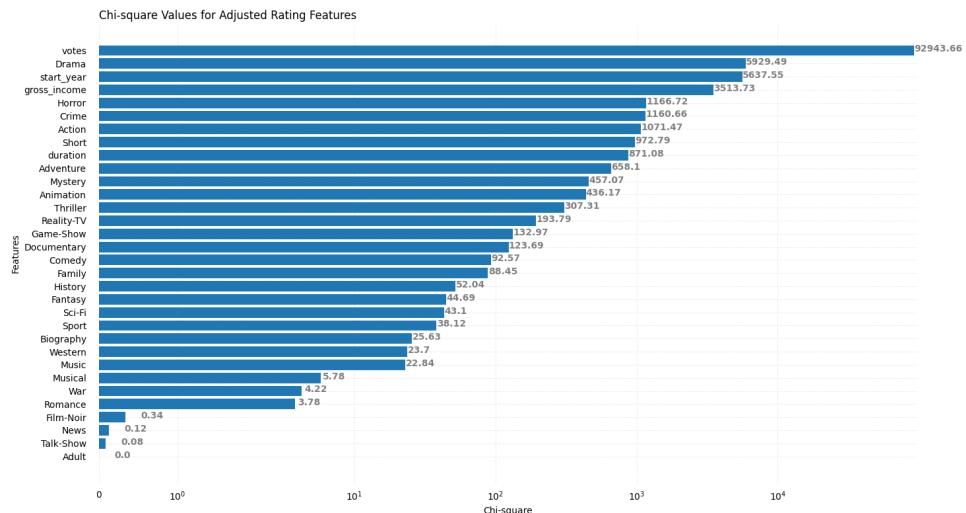


125

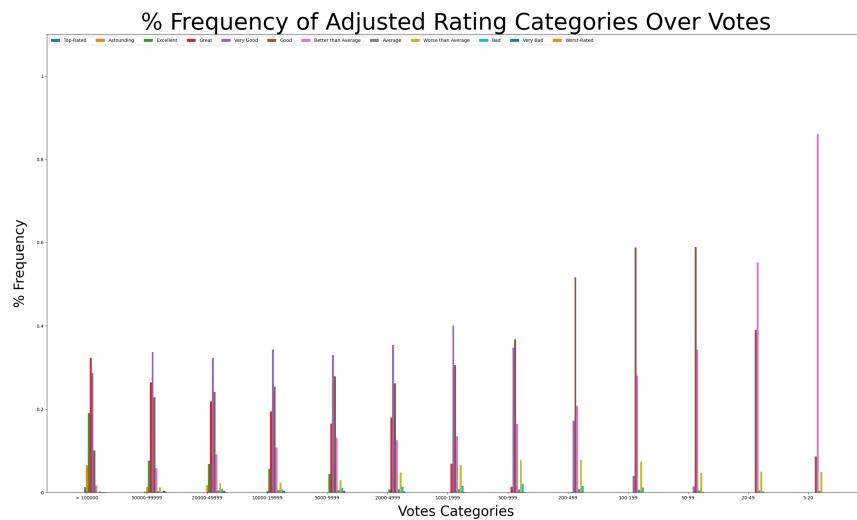


126

127 Shifting to the next plot, we can see that the top 5 features for the adjusted rating are votes, the
128 Drama genre, the start year, gross income, and the Horror genre, all with large chi-square scores,
129 although votes, once again, has the largest by a wide margin. The worst 5 features for the adjusted
130 rating are the Adult genre, the Talk-show genre, the News genre, the Film Noir genre, and the
131 Romance genre, with a large jump between the Romance genre and the rest. Underneath the chi-
132 square plot lies a multi-bar plot of the percent frequency of the adjusted income categories over the
133 categories of its top feature, which in this case is also votes.



134



142

Statistic (%)	Adj. Rating (All)	Adj. Rating (Top 3)
Accuracy	66.61	62.44
Balanced Accuracy	18.32	13.81
Precision	28.56	15.65
Recall	18.32	13.81
F1	20.02	14.02

143

144 5.3 Naive Bayes Models

145 Now, we'll turn our attention to the Naive Bayes models. The table below uses the same format as
 146 the previous one.

147

Statistic (%)	Gross Income (All)	Gross Income (Top 3)
Accuracy	92.50	93.24
Balanced Accuracy	32.63	27.41
Precision	32.30	28.61
Recall	32.63	27.41
F1	30.84	26.63

148

149

Statistic (%)	Adj. Rating (All)	Adj. Rating (Top 3)
Accuracy	64.64	65.26
Balanced Accuracy	21.13	16.28
Precision	22.01	16.96
Recall	21.13	16.28
F1	21.18	15.69

150

151 6 Discussion

152 The most apparent discovery from this dataset is that the number of votes is the most powerful
 153 feature for both gross income and the adjusted rating. Now, with regards to the adjusted rating, while
 154 votes is one of factors used to calculate it, the adjusted rating can be positive or negative based on
 155 how high or low the rating is. As such, it is only guaranteed that the magnitude of the adjusted rating
 156 is proportional to that of the number of votes. Furthermore, from the percent frequency plot for
 157 adjusted income and votes, as the number of votes increase, the percentage of larger adjusted rating
 158 scores increases. This implies that there is a positive correlation between the two.

159 It is also interesting that gross income and adjusted rating both had high chi-square values with
 160 respect to the other variable, earning third and fourth on the chart respectively, which implies a
 161 reflexive relationship between the two. Many shows become profitable because they're good, and
 162 many shows become good because they're profitable. While the former idea is pretty easy to
 163 comprehend, the latter could use more of an example. Take the somewhat recent Sonic the Hedgehog
 164 movie. An early trailer used to promote the movie became infamous because of how bad the live-
 165 action Sonic looked. Because of the public backlash, the movie's release was delayed in order to
 166 improve the film, and it became a well-received movie, which for a video-game adaption was a rare

167 feat. These improvements only occurred because of the profitability of the Sonic the Hedgehog IP.
168 The movie's high earning-potential made it better.

169 One surprising difference between the two outputs is that while the chi-square values for the
170 adjusted rating are generally higher, the logistic regression and Naive Bayes models for gross income
171 performed much better in every statistic. This implies that while the features cover more of the
172 variation with regard to the adjusted income, the features are better suited for predicting gross income
173 (at least within these models), which is strange. Despite these differences, the models for both outputs
174 also share quite a few similarities. With the exception of accuracy, their statistics are all very low
175 with precision being the next highest and balanced accuracy approximately equal to recall. Note
176 that the difference between precision and recall is much larger across the logistic regression models
177 compared to the Naive Bayes models. Another similarity across all of the models is that those that
178 only used the top 3 features had lower statistics than their counterparts that used all of the features
179 but not by a wide margin (aside from the precision metric for the adjusted rating). Another thing to
180 note is the accuracy being so much larger than any other statistic for all of the models reveals how
181 imbalanced this dataset is. While the performance of these models is lack-luster at best, they do
182 have a saving grace in that precision is its highest, accurate metric. When attempting to find motion
183 pictures of a certain quality (high income, low performance, etc.), it is more important for the results
184 to be accurate rather than comprehensive. Therefore, I conclude that the logistic regression model
185 performed best in all cases due to its higher precision. Also, the lack of a significant difference (with
186 one exception) between the full models and the 3-feature models displays how strong the relationship
187 between the top 3 features and their output is.

188 **7 Conclusion**

189 My goal for this project was to discover what features go into making a monetarily successful
190 motion picture and a popular one then compare their similarities and differences, and I believe that,
191 in a way, I succeeded. I discovered that these two pursuits had a lot more in common than I initially
192 thought, and that surface-level information about a show is ultimately not enough to determine
193 whether it will be good or profitable.

194 **8 References**

- 195 [1] <https://www.kaggle.com/datasets/ashishjangra27/imdb-movies-dataset>
- 196 [2] Abidi, Syed Muhammad Raza, et al. "Popularity Prediction of Movies: From Statistical Modeling
197 to Machine Learning Techniques - Multimedia Tools and Applications." SpringerLink, Springer US,
198 6 Jan. 2020, <https://link.springer.com/article/10.1007/s11042-019-08546-5>.
- 199 [3] Dhir, Rijul, and Anand Raj. Movie Success Prediction Using Machine Learning ... - IEEE Xplore.
200 2018, <https://ieeexplore.ieee.org/abstract/document/8703320>.
- 201 [4] Latif, Muhammad Hassan, and Hammad Afzal. Prediction of Movies Popularity Using
202 Machine Learning Techniques. Aug. 2016, https://www.researchgate.net/profile/Hammad-Afzal/publication/311913687_Prediction_of_Movies_popularity_Using_Machine_Learning_Techniques/links/586253ce08ae6eb0f-Movies-popularity-Using-Machine-Learning-Techniques.pdf.