

COMP30027 Report

Anonymous

1 Introduction

2 Methodology

2.1 Understanding the Data

2.1.1 Class Distribution

By counting the frequency of each class and understanding the distribution of the classes in the training data, we can gain valuable insights when designing and evaluating models.

2.2 Data Preprocessing

2.2.1 Normalisation

For most of the numeric values in the dataset, the important information is how it compares to other values, rather than the actual value itself. These values should be normalised. Since some models prefer 0 mean data, such as logistic regression and support vector machines, standardisation will be the normalisation method of choice. However, attributes like `title_year` or `average_degree centrality` are more valuable as the value itself, and so these will not be standardised.

2.2.2 Removing Highly Correlated Attributes

There are a few reasons to remove all but one attribute from a set of highly correlated attributes. Not only does removing the attributes decrease dimensionality and makes the models faster, but having multicollinearity among independent variables can result in less reliable statistical inferences.

2.2.3 Dropping Unuseful Values

Considering that the data is likely to be biased towards certain regions of the world, there will likely be certain attributes the exhibit an overwhelming majority of one variation, with many other unique but less prevalent values. These attributes would not be very useful to the data, and the best way to approach them is dropping them from the dataset. This also has the additional advantage of decreasing the dimensional-

ity of the data, and hence increasing the speed of the models.

2.2.4 One Hot Encoding

2.3 Baseline Models

To provide a baseline for benchmarking, we can develop a few simple models. The simplest that can be developed is just a model that guesses the most common class. Following this, a One Rule model could be developed on the categorical attributes, such as film categories.

2.4 Model Design

2.4.1 Models

Three models were picked, these were:

- Support Vector Machine
- Logistic Regression
- Random Forest

For each of the models, the below steps were followed to find the best model:

1. The training data was split into a training set and validation set.
2. For each hyperparameter model, a few reasonable options were selected.
3. Iterating over all combinations of hyperparameters, the combination that resulted in the highest accuracy over the training set was picked.
4. The model is evaluated over the validation set.
5. Using sequential forward selection and the best hyperparameters for the model, the best combination of features is selected for the training of the model.
6. Using the best hyperparameters and best features, the model is evaluated against

3 Results

3.1 Class Distribution

Figure 1 shows the distribution of the classes in the training dataset. It's clear that a rating of 2 is by far the most common, followed by 3. However, it's only a 61.2% majority, which means a model that always guesses 2 won't have the best accuracy.

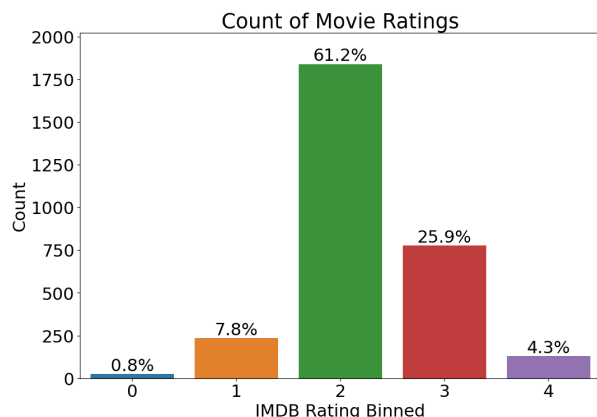


Figure 1: Count plot showing the distribution of classes in the training data

3.1.1 Category One Rule

3.2 Preprocessing

3.2.1 Unhelpful Features

Since most films are in english and there are many other unique values with very low frequency, the language feature is not very useful and should be dropped. A similar occurrence also exists with the country feature, where the USA is an overwhelming majority.

Similarly, a large majority of movies are either R or PG-13 rated, and so the rating of a movie would not be very helpful and can lead to overfitting. As such `content_rating` was also dropped from the dataset.

3.2.2 Names

Through filtering out all of the text attributes, we are also filtering out the names of actors and the movie director. Whilst the quality of film can be correlated to its director, and people often enjoy certain movies with certain actors more, the names alone are not enough to generalise the IMDB rating of a movie. This could instead result in the model memorising which directors and actors are better rated, and since there's a limited number of names, would result in overfitting to the training data.

3.3 Model 1: Support Vector Machine

The best hyperparameters were found to be:

C: 10

Gamma: 0.01

Kernel: RBF

These hyperparameters resulted in an accuracy of 70.0% and the below:

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.00	0.00	0.00	48
2	0.70	0.94	0.81	377
3	0.67	0.37	0.48	152
4	0.69	0.47	0.56	19

Table 1: *Classification Report of the Best SVM Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.41	0.36	0.37
Weighted Avg	0.68	0.70	0.64

Table 2: *Summary of Classification Report of the Best SVM Model without Feature Selection*

The best feature set found was:

- Number of Critics for Reviews
- Duration
- Director Facebook Likes
- Gross
- Number of Voted Users
- Cast Total Facebook Likes
- Face Number in Poster
- Title Year
- Movie Facebook Likes

These features resulted in an accuracy of 70.0% and the below:

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.00	0.00	0.00	48
2	0.70	0.94	0.81	377
3	0.67	0.37	0.48	152
4	0.69	0.47	0.56	19

Table 3: *Classification Report of the Best SVM Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.41	0.36	0.37
Weighted Avg	0.68	0.70	0.64

Table 4: *Summary of Classification Report of the Best SVM Model without Feature Selection*

3.4 Model 2: Logistic Regression

The best hyperparameters were found to be:

C: 0.1

Max Iterations: 1000

Penalty: L2

Solver: Saga

These hyperparameters resulted in an accuracy of 65.2% and the below:

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.00	0.00	0.00	48
2	0.65	0.99	0.79	377
3	0.69	0.12	0.20	152
4	0.00	0.00	0.00	19

Table 5: *Classification Report of the Best Logistic Regression Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.27	0.22	0.20
Weighted Avg	0.58	0.65	0.54

Table 6: *Summary of Classification Report of the Best Logistic Regression Model without Feature Selection*

The best feature set found was:

- Number of Critics for Reviews
- Duration
- Director Facebook Likes
- Gross
- Number of Voted Users
- Cast Total Facebook Likes

These features resulted in an accuracy of 66.8% and the below:

3.5 Model 3: Random Forest

The best hyperparameters were found to be:

Max Depth: 30

N Estimators: 200

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.00	0.00	0.00	48
2	0.68	0.94	0.79	377
3	0.56	0.27	0.36	152
4	0.89	0.42	0.57	19

Table 7: *Classification Report of the Best Logistic Regression Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.43	0.33	0.34
Weighted Avg	0.60	0.67	0.60

Table 8: *Summary of Classification Report of the Best Logistics Regression Model without Feature Selection*

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.50	0.04	0.08	48
2	0.73	0.91	0.81	377
3	0.64	0.50	0.56	152
4	0.67	0.42	0.52	19

Table 9: *Classification Report of the Best Random Forest Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.51	0.37	0.39
Weighted Avg	0.68	0.71	0.67

Table 10: *Summary of Classification Report of the Best Random Forest Model without Feature Selection*

These hyperparameters resulted in an accuracy of 71.2% and the below:

The best feature set found was:

- Number of Critics for Reviews
- Number of Voted Users

These features resulted in an accuracy of 60.2% and the below:

4 Discussion & Critical Analysis

5 Conclusion

References

Christopher M Bishop and Nasser M Nasrabadi. 2006. *Pattern recognition and machine learning*, volume 4. Springer.

Ratin	Prec.	Recall	F1-Score	Support
0	0.00	0.00	0.00	5
1	0.09	0.04	0.06	48
2	0.68	0.79	0.73	377
3	0.43	0.36	0.39	152
4	0.62	0.42	0.50	19

Table 11: *Classification Report of the Best Random Forest Model without Feature Selection*

	Prec.	Recall	F1-Score
Macro Avg	0.36	0.32	0.33
Weighted Avg	0.56	0.60	0.58

Table 12: Summary of Classification Report of the Best Random Forest Model without Feature Selection

This is a report template, suitable for L^AT_EX. Don't use fonts smaller than this one (11pt). Don't include a title page, table of contents, abstract, or other similar front matter.

Please don't include your name and/or student ID in the title or header; your report should be anonymised for the reviewing process..

Use **bold** for **emphasis**, but use sparingly.

Short quotations “*are included in the main text, in normal paragraph style, between double quotes and italicized*” All quotes should be properly referenced. Note that the citation style is defined in the accompanying style file; it is similar to AAAI house style. But you may use other (formal) citation styles if you prefer.

You can cite related papers or books like this (Bishop and Nasrabadi, 2006). Text¹ with footnotes at bottom of page.

5.0.1 Subsubsection

Figures should be placed in the text, not at the end. Figures must be captioned and explicitly mentioned in the text .

¹Footnote text