
CS 506 Midterm Report

Jae Hong Lee
Boston University
jhonglee@bu.edu

Abstract

1 For the CS 506 Midterm competition I'm using Principal-Component Analysis
2 (PCA) and XGBoost (Extreme Gradient Boosting) to predict the star rating associ-
3 ated with user reviews from Amazon Movie Reviews using the available features
4 dataset provided by Kaggle class midterm. This is my GitHub repository address.

5 1 Introduction

6 1.1 Problem

7 The goal of this competition is to predict the star rating associated with user reviews from Amazon
8 Movie Reviews using the available features.

9 1.2 Dataset

10 The movie rating dataset used for this midterm has the following attributes:

- 11 • Id: a unique identifier associated with a review
- 12 • Product Id: unique identifier for the product
- 13 • User Id: unique identifier for the user
- 14 • Helpfulness Numerator: number of users who found the review helpful
- 15 • Helpfulness Denominator: number of users who indicated whether they found the review
16 helpful
- 17 • Time: timestamp for the review
- 18 • Summary: brief summary of the review
- 19 • Text: text of the review
- 20 • Score: rating between 1 and 5

21 The dataset for movie ratings, which will be used in this midterm, consists of 1,697,533 rows. I need
22 to split it based on the score column. Out of the total data, 1,485,341 rows have a score, forming the
23 Training Set, while 212,192 rows have missing scores, which make up the Testing Set. After splitting,
24 the Training Set has 54 missing text columns and 28 missing summary columns, while the Testing
25 Set has 8 missing text columns and 4 missing summary columns.

26 Due to missing values in the dataset, I believe XGBoost can be effectively applied, as it is capable of
27 handling missing attributes in the data. The method can still make use of the available data, even when
28 some features are missing. Certain attributes in the Training Set, such as Id, Product Id, and User Id,
29 are simply unique identifiers and do not provide valuable information for the model. These fields are
30 more administrative and don't predict outcomes. Therefore, excluding them from the analysis could
31 streamline the process. With these attributes removed, a decision tree might be a more appropriate
32 method, as it can focus on the remaining meaningful features to make better predictions.

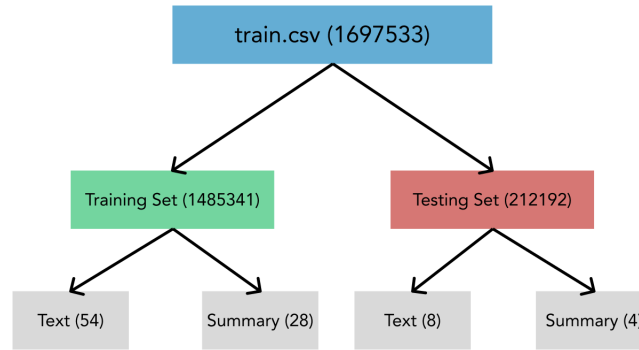


Figure 1: Data Split of train.csv

1.3 Model Overview

Decision tree models offer several advantages. They can handle both numerical and categorical data and perform automatic feature selection, making them well-suited to complex problems where interactions between variables are critical. Decision trees also capture non-linear relationships, enabling them flexible in various scenarios. Once trained, they make predictions quickly and provide clear, interpretable decision paths. However, they can be prone to overfitting, particularly on small datasets, as they may create overly complex trees that capture noise in the data. Additionally, decision trees tend to struggle with small variations or noisy data, though techniques like pruning can reduce this risk.

On the other hand, K-Nearest Neighbors (KNN) is a simple and intuitive algorithm that doesn't require a training phase, making it computationally light at the start. It is non-parametric, meaning it makes no assumptions about the underlying data distribution, which can be useful when the true relationship between features is unknown. KNN has its disadvantages. It can be computationally expensive during the prediction phase, as it requires storing and searching through the entire training set for each new prediction. KNN also struggles when irrelevant features are present or when the data is not scaled properly, as distance calculations become less meaningful. Additionally, KNN's performance declines when dealing with high-dimensional data (many attributes), as the algorithm becomes less effective at measuring meaningful distances.

2 Preprocessing Features

In this project, the preprocessing steps were designed to capture significant features and ensure that both text and numerical attributes are adequately represented. The following steps outline the approach taken:

2.1 Helpfulness Feature Creation

From the starter code, I used this feature. To assess user feedback reliability, an additional feature called Helpfulness was created. This feature is the ratio of 'HelpfulnessNumerator' to 'HelpfulnessDenominator', providing a proportionate measure of helpfulness. Missing values were replaced with 0 to avoid errors in the model. After creating this feature, the original 'HelpfulnessNumerator' and 'HelpfulnessDenominator' columns were dropped.

2.2 Text Vectorization with GloVe and Custom Embeddings

To capture semantic information from the Text column, GloVe embeddings and a custom Word2Vec model were utilized. GloVe embeddings provide pre-trained, context-independent word vectors, while the custom Word2Vec model was trained specifically on the project's dataset to capture its own patterns and nuances.

The word vectors from GloVe and Word2Vec were combined using different methods ('concat', 'average', or 'weighted'), allowing flexibility in representing each word based on its contextual

68 meaning. These word vectors were then averaged to create a single vector representation for each
69 sentence, effectively transforming text data into meaningful vectorized form.

70 2.3 Feature Extraction: Text & Summary Length and Sentiment Scores

71 For both the Text and Summary columns, additional attributes were extracted, including text length
72 (character count) and sentiment scores. VADER sentiment analysis was applied to determine the
73 sentiment polarity of both 'Text' and 'Summary', capturing positive, negative, or neutral sentiments.
74 This added feature helps the model learn the emotional tone of the reviews. Both text length and
75 summary length were also computed to assess verbosity and detail in reviews, providing additional
76 context for model learning.

77 2.4 TF-IDF and Count Vectorization in Summary

78 For the Summary column, TF-IDF and CountVectorizer with N-grams were applied to analyze
79 word and phrase patterns. TF-IDF captures term importance by weighting words that are frequent
80 in the document but rare in the overall dataset, helping the model focus on meaningful keywords.
81 CountVectorizer with N-grams captures common word pairs, such as "not good," providing the model
82 with more detailed phrase-level context. This combination of TF-IDF and CountVectorizer features
83 allows the model to learn both the importance of individual words and the structure of phrases in
84 summaries.

85 2.5 Final Data Preparation and Combination

86 After transforming all features, the text vectors, TF-IDF, CountVectorizer, and other numeric features
87 (like sentiment scores and length features) were combined into a single sparse matrix for efficient
88 storage and processing. The final combined dataset includes both the original numerical features
89 and the vectorized text attributes, ready for model training. This preprocessing pipeline effectively
90 transforms raw text and numerical data into a structured format that captures both text semantic
91 information and the quantitative insights from numerical attributes. The combined features are
92 expected to enhance the model's ability to accurately interpret and predict based on the data.

93 3 Experiments

94 3.1 XGBoost Model

95 The experiment dataset underwent Preprocessing, and the eight columns after the features processing,
96 the highest variance were selected. These datasets were then used to split the data into a training
97 set and a test set in a 3:1 ratio. This ensured that we could accurately train our data, but wouldn't
98 overfit so testing would still be accurate. The models were trained on this split data, and accuracy
99 was evaluated.

100 I conducted experiments using xgboost models: learning rate, max_depth, min_child_weight, and etc
101 parameters were used before calculating final accuracy

102 The initial parameter I used was the followings:

Initial Parameter Results	
Parameters	Values
N_estimator	100
Learning Rate	0.1
Max_depth	5
Min_child_weight	5
Gamma	0
subsample	0.8
colsample_bytree	0.8

104 The Accuracy that I got from this was **0.602045**

105 3.1.1 Hyper Parameter Tuning

106 XGBoost models, contain several parameters, but I will only mention the parameters I touched.

- N_estimator : { '100', '200', '300', '400', '500' } Number of gradient boosted trees
- Learning Rate : { 0.01, 0.05, 0.1, 0.2, 0.3 } Boosting learning rate.
- Max_depth : { 3, 4, 5, 6, 7, 8, 9, 10 } Maximum tree depth for base learners.
- Min_child_weight : { 1, 3, 5, 7 } Minimum sum of instance weight(hessian) needed in a child.
- Gamma : { 0, 0.1, 0.2, 0.3, 0.4 } Minimum loss reduction required to make a further partition on a leaf node of the tree.
- subsample : { 0.6, 0.7, 0.8, 0.9, 1.0 } Subsample ratio of the training instance.
- colsample_bytree : { 0.6, 0.7, 0.8, 0.9, 1.0 } Subsample ratio of columns when constructing each tree.

3.2 Result

I applied GridSearchCV through the parameters of XGBoost. Based on the results of the experiments, the final accuracy scores are as follows:

Initial Parameter Results	
Parameters	Values
N_estimator	500
Learning Rate	0.1
Max_depth	4
Min_child_weight	3
Gamma	0.1
subsample	0.8
colsample_bytree	1.0

XGBoost Fine Tuning score was: 0.6409585

As observed, this was my best result. I used this model to train the X_submission, and the result of the public leader board was 0.64315

124 **References**

- 125 [1] xgboost developers. *Frequently Asked Questions*. 2022. URL: <https://xgboost.readthedocs.io/en/stable/faq.html>.
126
- 127 [2] Jeffrey Pennington. *GloVe: Global Vectors for Word Representation*. 2014. URL: <https://nlp.stanford.edu/projects/glove/>.
128
- 129 [3] *SciKitLearn Documentation*. Accessed 25 Oct. 2024. URL: <https://scikit-learn.org/stable/index.html>.
130
- 131 [4] *SciKitLearn TfidfVectorizer Documentation*. Accessed 25 Oct. 2024. URL: https://scikit-learn.org/1.5/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html.
132
133