

Citations for external methods used:

A few of these methods were briefly covered in labs or class, but I will reiterate all of them anyways.

1. **RandomForestClassifier:**

The `RandomForestClassifier` is part of the `sklearn.ensemble` module, which is designed for classification problems using an ensemble learning approach, which involves constructing a multitude of decision trees during training. This method was introduced by Breiman (2001) and is widely used for its ability to handle both classification and regression tasks while reducing overfitting compared to individual decision trees.

Citation: Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.

2. **TfidfVectorizer:**

The `TfidfVectorizer` from `sklearn.feature_extraction.text` transforms text data into numerical form based on the term frequency-inverse document frequency (TF-IDF) score. This method is valuable for extracting relevant information from text features and is a common technique for text-based machine learning models. **Citation:** Salton, G., & McGill, M. J. (1986). *Introduction to Modern Information Retrieval*. McGraw-Hill.

3. **StandardScaler:**

`StandardScaler` from `sklearn.preprocessing` is used to standardize features by removing the mean and scaling to unit variance. This ensures that features with different scales do not disproportionately affect the model's performance. **Citation:** Pedregosa et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.

4. **make_pipeline:**

The `make_pipeline` function from `sklearn.pipeline` allows for the chaining of multiple transformations and estimators in sequence, simplifying the code and ensuring that the preprocessing steps (e.g., scaling) are applied before model training. **Citation:** Pedregosa et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, 2825–2830.

5. **GridSearchCV:**

`GridSearchCV` from `sklearn.model_selection` performs an exhaustive search over specified hyperparameter values for a given estimator. It automates the process of tuning hyperparameters to improve model performance. **Citation:** Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb), 281-305.

Explanation of how ChatGPT contributed to crafting the algorithm:

To develop the machine learning algorithm in the code, I used ChatGPT to assist in integrating several key methods from the `sklearn` library, specifically focusing on improving both the preprocessing and model training phases. Here's a breakdown of how these methods were integrated:

1. **Preprocessing the data:**

Initially, the algorithm required scaling of the features, especially numerical ones such as

Helpfulness and **ReviewLength**. I used the **StandardScaler** within a pipeline to ensure that the scaling would be automatically applied during model training. This was crucial to avoid skewed model performance caused by differing scales of the features.

2. **Text Feature Extraction:**

Given that the dataset contains text (e.g., the review text), I considered using the **TfidfVectorizer** for extracting relevant features from text data. Although it was not integrated in the final pipeline (due to focus on numerical features in the provided code), I used the method offline during exploratory analysis to investigate how text features might influence model performance. This step would allow the addition of textual features in the future to improve predictions.

3. **Building a Pipeline with **make_pipeline**:**

To streamline the process of preprocessing and model training, I used the **make_pipeline** function. It allowed me to combine the scaling step (**StandardScaler**) and the model (**RandomForestClassifier**) into a single pipeline. This simplifies the code and ensures that preprocessing steps are not forgotten during the training or prediction phases. The use of pipelines ensures reproducibility and ease of cross-validation, as the entire pipeline is treated as a single object.

4. **Model Training with **RandomForestClassifier**:**

The **RandomForestClassifier** was chosen due to its robustness and ability to handle a variety of feature types. It performs well in many classification tasks because it reduces overfitting and provides good out-of-the-box accuracy. By using the **RandomForestClassifier**, I ensured that the model would be able to handle the mix of numerical features extracted from the dataset, such as the review helpfulness and review length.

5. **Hyperparameter Tuning with **GridSearchCV**:**

To further optimize the model, I implemented **GridSearchCV**, allowing an automated search for the best combination of hyperparameters for the **RandomForestClassifier**. By specifying ranges for the number of trees (**n_estimators**), the maximum depth of trees (**max_depth**), and the minimum samples required to split a node (**min_samples_split**), I was able to explore various model configurations. ChatGPT assisted in determining the key hyperparameters to tune and in setting up the **GridSearchCV** function efficiently.

Offline evaluation of the methods:

In offline evaluations, I assessed the performance of the algorithm on the provided training and test sets. Here's an outline of the evaluation steps:

1. **Splitting the Data:**

I used **train_test_split** to create a testing set from the original training data. This allowed me to evaluate the model's performance before using the test dataset for predictions.

2. **Accuracy Evaluation:**

After training the model, I evaluated its performance on the test set using **accuracy_score**. The accuracy metric provided a quick measure of how well the model performed, although further evaluation could include precision, recall, and F1 scores for a more detailed assessment.

3. Confusion Matrix:

To visually assess the model's performance across different classes, I generated a confusion matrix using `confusion_matrix` and visualized it with `seaborn`. This helped in understanding where the model made correct predictions versus where it misclassified instances, thus providing insights into any class imbalance or specific weaknesses.

4. Creating the Submission File:

Once the model was optimized, I used it to predict the `Score` values for the test set and generated the submission file. This demonstrates that the algorithm was successfully applied to real-world data, ready for external evaluation.