

Enhancing Model Performance and Interpretability

***Lack of Model Persistence*:**

- The code does not save the trained models, which means they need to be retrained every time the script is run.
- Saving the models would allow for easier deployment and reuse in the future, without the need to retrain them from scratch.

The code is addressing the issue of lack of model persistence by saving the trained models to disk using the `joblib.dump()` function. Specifically, the code includes the following lines:

```
# Save the trained models
joblib.dump(xg_boost_model, 'xg_boost_model.joblib')
joblib.dump(naive_bayes_model, 'naive_bayes_model.joblib')
```

These lines save the trained XGBoost and Naive Bayes models to files named `xg_boost_model.joblib` and `naive_bayes_model.joblib`, respectively. This allows the models to be loaded and used in the future without the need to retrain them from scratch.

By saving the models, we can deploy them and reuse them in other projects or applications, without having to go through the training process again. This is a common and recommended practice in machine learning, as it enhances the overall maintainability and scalability of our project.

Additionally, if we need to make changes to our data or model hyperparameters in the future, we can easily load the saved models, update them, and save them again, rather than having to retrain the entire model from scratch.

***Hardcoded Threshold*:**

- The code uses a fixed threshold of 0.5 to convert the continuous predictions to binary labels.
- This threshold may not be the optimal choice, and it could be worth exploring different threshold values or using a more sophisticated approach, such as finding the optimal threshold based on the desired trade-off between precision and recall.

The code is addressing the issue of hardcoded threshold by finding the optimal threshold for each of the models (XGBoost and Naive Bayes) as well as the combined model, based on the F1-score.

1. Obtains the predicted probabilities from the trained XGBoost and Naive Bayes models:

```
xg_boost_probs = xg_boost_model.predict_proba(X_numeric_test)[: , 1]
naive_bayes_probs = naive_bayes_model.predict_proba(X_text_test)[: , 1]
```

The code first obtains the predicted probabilities from the trained XGBoost and Naive Bayes models. The `predict_proba()` method returns a 2D array, where each row represents a sample, and the columns represent the probabilities for each class. By selecting the second column (`[: , 1]`), we get the probabilities for the positive class (the class we're interested in predicting).

2. Calculates the F1-score for a range of thresholds (from 0 to 1 in increments of 0.01) for each model, and finds the threshold that maximizes the F1-score:

```
# XGBoost
thresholds = np.linspace(0, 1, 101)
f1_scores = [f1_score(y_test, (xg_boost_probs >= t).astype(int)) for t in thresholds]
optimal_idx = np.argmax(f1_scores)
optimal_xg_boost_threshold = thresholds[optimal_idx]
print(f"Optimal XGBoost threshold: {optimal_xg_boost_threshold:.2f}")

# Naive Bayes
thresholds = np.linspace(0, 1, 101)
f1_scores = [f1_score(y_test, (naive_bayes_probs >= t).astype(int)) for t in thresholds]
optimal_idx = np.argmax(f1_scores)
optimal_naive_bayes_threshold = thresholds[optimal_idx]
print(f"Optimal Naive Bayes threshold: {optimal_naive_bayes_threshold:.2f}")
```

The code then finds the optimal threshold for each individual model (XGBoost and Naive Bayes) based on the F1-score. It does this by:

Generating a list of 101 thresholds, ranging from 0 to 1 in increments of 0.01.

Calculating the F1-score for each threshold by comparing the binary predictions (obtained by applying the threshold to the predicted probabilities) with the true labels (`y_test`).

Finding the index of the threshold that maximizes the F1-score, and storing the corresponding threshold value.

This process ensures that the optimal threshold for each model is determined, rather than using a fixed, hardcoded threshold of 0.5.

3. Finds the optimal threshold for the combined model, using the same approach:

```
print("Finding optimal threshold for combined model...")
thresholds = np.arange(0, 1.01, 0.01)
f1_scores = [f1_score(y_test, (weighted_avg_preds > t).astype(int)) for t in thresholds]
optimal_idx = np.argmax(f1_scores)
optimal_threshold = thresholds[optimal_idx]
print(f"Optimal threshold for combined model: {optimal_threshold:.2f}")
```

The code then follows a similar approach to find the optimal threshold for the combined model. However, instead of using the predicted probabilities from individual models, it uses the weighted average predictions (`weighted_avg_preds`). The process is the same:

Generate a list of thresholds, this time ranging from 0 to 1.01 in increments of 0.01.

Calculate the F1-score for each threshold by comparing the binary predictions (obtained by applying the threshold to the weighted average predictions) with the true labels (`y_test`).

Find the index of the threshold that maximizes the F1-score, and store the corresponding threshold value.

By finding the optimal threshold for the combined model, the code ensures that the final predictions are based on the best possible threshold, considering the performance of both the XGBoost and Naive Bayes models.

4. Converting Continuous Predictions to Binary Labels:

```
binary_preds = (weighted_avg_preds > optimal_threshold).astype(int)
```

Finally, the code converts the continuous predictions (the weighted average predictions) to binary labels based on the optimal threshold found in the previous step. This ensures that the final predictions are not based on a hardcoded threshold, but on the optimal threshold determined by the code.

This approach is more sophisticated and likely to yield better results than using a fixed, hardcoded threshold, as it takes into account the specific characteristics of the dataset and the performance of the models.

***Expanding on Interpretability Analyses*:**

The code implements the SHAP (Shapley Additive Explanations) analysis to gain a more comprehensive understanding of the feature importance and the model's decision-making process.

Here's how the code is addressing the issue:

Defining the XGBoost Pipeline: The code defines an XGBoost pipeline without the ColumnTransformer. This is because the SHAP analysis will be performed on the numeric features directly, without the need for the ColumnTransformer.

Hyperparameter Tuning for XGBoost: The code performs a RandomizedSearchCV to find the best hyperparameters for the XGBoost model. The best hyperparameters are then used to fit the final XGBoost model.

Model Performance Evaluation: The code evaluates the performance of the XGBoost model using cross-validation and calculates the accuracy on the test set.

SHAP Explainer: The code creates a SHAP explainer using the best XGBoost model and calculates the SHAP values for the test set.

SHAP Analysis Visualizations:

Summary Plot: The code generates a SHAP summary plot, which provides an overview of the feature importance and how each feature contributes to the model's predictions.

Individual Explanations: The code selects a random instance from the test set and generates a SHAP force plot to explain the prediction for that specific instance.

The SHAP analysis provides a more comprehensive understanding of the model's decision-making process compared to the previous LIME analysis. SHAP values can help identify the most important features and their contribution to the model's predictions, which can lead to better insights and potential model improvements.

By implementing the SHAP analysis, the code addresses the goal of expanding the interpretability analyses and gaining a deeper understanding of the model's behavior.

***Lack of Model Comparison*:**

- The code does not compare the performance of the individual base models (XGBoost and Naive Bayes) to the combined model.

- It would be helpful to understand the relative contribution of each base model and whether the ensemble approach actually provides a significant performance improvement over the individual models.

The code now includes a comprehensive model comparison section, addressing the lack of comparison between the individual base models (XGBoost and Naive Bayes) and the combined model. Here's how the code addresses this issue:

Accuracy Comparison:

The code includes a bar plot that compares the accuracy of the XGBoost and Naive Bayes models. This allows to visually assess the performance of the individual models and understand their relative contributions.

Precision-Recall Curve and ROC Curve:

The code generates the precision-recall curve and the ROC curve for the combined model (using the weighted average predictions). This provides insights into the overall performance of the ensemble approach and allows to compare it to the individual models.

Permutation Importance:

The code computes the permutation importance for the XGBoost model, which helps understand the relative importance of the input features. This analysis can be extended to the Naive Bayes model as well, if desired.

Top Words Analysis:

The code extracts the top positive and negative words for the Naive Bayes model, based on the feature log-probabilities. This gives a better understanding of how the Naive Bayes model is making its predictions.

Confusion Matrices:

The code generates and plots the confusion matrices for both the XGBoost and Naive Bayes models. This allows to compare the models' performance on different classes and identify any potential biases or imbalances.

Classification Reports:

The code generates classification reports for both the XGBoost and Naive Bayes models, providing detailed metrics such as precision, recall, F1-score, and support for each class. This information can be used to further compare the models and understand their strengths and weaknesses.

Learning Curves:

The code generates learning curves for both the XGBoost and Naive Bayes models, plotting the training and validation scores as a function of the training set size. This analysis can help understand the models' learning behavior and whether they are prone to overfitting or underfitting.

By including these comparative analyses, the code now provides a more comprehensive evaluation of the individual models and the combined model. This allows to assess the relative contribution of each

base model and determine whether the ensemble approach offers a significant performance improvement over the individual models. The visualizations and metrics generated can help make informed decisions about the final model to use and identify areas for further improvement or exploration.