

## Model Tuning

### 1. Overfitting Potential:

- The code does not use any techniques to prevent overfitting, such as cross-validation or regularization.
- Cross-validation would help assess the model's performance on unseen data and identify potential overfitting.
- Regularization techniques, such as L1 or L2 regularization, could be applied to the XGBoost and Naive Bayes models to discourage model complexity and improve generalization.

### 2. Lack of Hyperparameter Tuning:

- The code uses the default hyperparameters for the XGBoost and Naive Bayes models, which may not be the optimal settings for the given problem.
  - Hyperparameter tuning, using techniques like grid search or random search, could help find the best combination of hyperparameters for each model, potentially leading to improved performance.
- 

## Hyperparameter Tuning and Cross-Validation

The updated code addresses the issues of overfitting potential and lack of hyperparameter tuning as follows:

### 1. Overfitting Potential:

- The code now uses cross-validation to assess the model's performance on unseen data and identify potential overfitting.
- The use of 6-fold cross-validation provides a more reliable estimate of the models' generalization ability, as they are evaluated on a larger portion of the dataset.
- The hyperparameter tuning process, which includes regularization parameters like ``reg_alpha`` and ``reg_lambda`` for the XGBoost model, helps discourage model complexity and improve generalization.

### 2. Lack of Hyperparameter Tuning:

- The code now performs hyperparameter tuning for both the XGBoost and Naive Bayes models.
- For the XGBoost model, the code uses `RandomizedSearchCV` to efficiently explore the large hyperparameter space and find the best combination of hyperparameters.
- For the Naive Bayes model, the code uses `GridSearchCV` to exhaustively evaluate all possible combinations of the smaller hyperparameter space.
- The tuned hyperparameters, such as ``max_depth``, ``learning_rate``, ``n_estimators``, ``alpha``, and ``fit_prior``, are then used to fit the final models, which can lead to improved performance compared to using the default hyperparameters.

**The key points are:**

**1. Cross-Validation:**

- The use of 6-fold cross-validation helps assess the models' performance on unseen data and identify potential overfitting.
- Cross-validation provides a more reliable estimate of the models' generalization ability.

**2. Hyperparameter Tuning:**

- The code uses RandomizedSearchCV for the XGBoost model and GridSearchCV for the Naive Bayes model to find the best combination of hyperparameters.
- The tuned hyperparameters, including regularization parameters, help improve the models' performance and prevent overfitting.

By incorporating these changes, the updated code addresses the issues of overfitting potential and lack of hyperparameter tuning. The use of cross-validation and regularization helps mitigate overfitting, while the hyperparameter tuning process aims to find the optimal settings for each model, potentially leading to improved performance.

---

## **Approaches used to tackle Overfitting Potential and Hyperparameter Tuning**

**1. Grid Search:**

- Grid search is a technique for hyperparameter optimization that exhaustively searches through a manually specified subset of the hyperparameter space.
- It creates a grid of all possible combinations of the specified hyperparameters and evaluates the model's performance for each combination.
- Grid search is a comprehensive approach, as it evaluates all possible combinations, but it can be computationally expensive, especially when the hyperparameter space is large.

**2. Randomized Search:**

- Randomized search is an alternative approach to hyperparameter optimization that samples the hyperparameter space randomly.
- Instead of evaluating all possible combinations, it randomly selects a fixed number of hyperparameter configurations to evaluate.
- Randomized search is more efficient than grid search, as it can explore a larger hyperparameter space in a shorter amount of time.

**Reason for using Randomized Search:**

In the provided code, the reason for using randomized search instead of grid search is that the grid search was taking a long time (30 minutes) to complete and did not provide any results.

The key differences between grid search and randomized search that led to this decision are:

### **1. Computational Efficiency:**

Randomized search is generally more computationally efficient than grid search, especially when the hyperparameter space is large. Grid search evaluates all possible combinations, which can be extremely time-consuming, while randomized search samples the hyperparameter space randomly, reducing the number of evaluations required.

### **2. Memory Efficiency:**

Grid search can be memory-intensive, as it needs to store all the possible hyperparameter combinations and their corresponding evaluation results. In contrast, randomized search only needs to store the current set of sampled hyperparameter configurations and their results, making it more memory-efficient.

In the case where the grid search was taking a long time to complete and not providing any results, it's likely that the hyperparameter space was too large, and the grid search algorithm was unable to explore it efficiently within the available time and memory constraints. By using randomized search, the code can explore a larger hyperparameter space in a more efficient manner, potentially finding better hyperparameter configurations in a shorter amount of time.

The use of randomized search in this case is a memory-efficient approach that can provide a good approximation of the optimal hyperparameters, without the need to exhaustively evaluate all possible combinations as in the grid search method.

---

## **Final Approach Used**

### **1. Hyperparameter Tuning for XGBoost:**

- The code defines a set of hyperparameters to try out for the XGBoost model. These include things like the maximum depth of the decision trees, the learning rate, the number of trees to use, and various regularization parameters.
- It then uses a technique called "RandomizedSearchCV" to find the best combination of these hyperparameters. This means it randomly tries out different combinations and evaluates how well the model performs on the training data.
- The `RandomizedSearchCV` parameters tell the algorithm how many different combinations to try (20 in this case) and how many times to split the data for cross-validation (6 folds).

### **2. Hyperparameter Tuning for Naive Bayes:**

- The code defines a smaller set of hyperparameters to tune for the Naive Bayes model, including the smoothing parameter `alpha` and whether to learn class prior probabilities or not.
- Instead of using RandomizedSearchCV, the code uses "GridSearchCV" to tune the Naive Bayes hyperparameters. This means it will try out all possible combinations of the specified hyperparameters.

- The `GridSearchCV` parameters are similar to the ones used for XGBoost, including the number of cross-validation folds (6).

### **3. Model Fitting with Best Hyperparameters:**

- After finding the best hyperparameters for each model, the code sets those hyperparameters and then trains the final models on the training data.

### **4. Cross-Validation:**

- The code performs 6-fold cross-validation on both the XGBoost and Naive Bayes models to get a better estimate of how well they will perform on new, unseen data.

- Cross-validation means splitting the data into 6 parts, training the model on 5 of those parts, and then testing it on the remaining part. This is done 6 times, and the results are averaged to get a more reliable performance estimate.

#### **The key point:**

#### **Naive Bayes Hyperparameter Tuning:**

- The code uses GridSearchCV for the Naive Bayes model, as it has fewer hyperparameters to tune compared to XGBoost. GridSearchCV is better suited for smaller hyperparameter spaces.

The goal of these changes is to find the best possible hyperparameters for each model and get a more reliable assessment of their generalization ability, which helps prevent overfitting.

---

## **Explaining the Code in Depth**

### **1. Hyperparameter Tuning for XGBoost:**

- The code defines a dictionary `xg\_boost\_params` that contains the hyperparameters to be tuned for the XGBoost model.

- The hyperparameters include:

- `classifier\_\_max\_depth`: The maximum depth of the tree. Larger values can lead to more complex models that are more likely to overfit.

- `classifier\_\_learning\_rate`: The step size at each iteration while moving toward a minimum of a loss function. Smaller values can lead to a more robust model, but may require more iterations to converge.

- `classifier\_\_n\_estimators`: The number of trees to be generated. Increasing the number of trees can improve the model's performance, but too many trees can lead to overfitting.

- `classifier\_\_min\_child\_weight`: The minimum sum of instance weight (hessian) needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than `min\_child\_weight`, the building process will give up further partitioning. This parameter helps control the complexity of the model and prevent overfitting.

- `'classifier__gamma'`: The minimum loss reduction required to make a further partition on a leaf node of the tree. Larger values of `'gamma'` can lead to a more conservative model that is less likely to overfit.
- `'classifier__subsample'`: The fraction of samples to be used for fitting the individual base learners. Reducing the subsample size can help prevent overfitting by introducing more randomness into the model.
- `'classifier__colsample_bytree'`: The fraction of columns to be used when constructing each tree. Reducing the column sample size can also help prevent overfitting.
- `'classifier__reg_alpha'`: The L1 regularization term on the weights. Larger values lead to more regularization, which can help prevent overfitting.
- `'classifier__reg_lambda'`: The L2 regularization term on the weights. Larger values lead to more regularization, which can help prevent overfitting.
- `'classifier__scale_pos_weight'`: A parameter to address class imbalance. It sets the scale of positive weights, which can help improve the performance on the minority class and prevent overfitting.
- The code then creates a `'RandomizedSearchCV'` object `'xg_boost_random_search'` to perform the hyperparameter tuning.
- The `'RandomizedSearchCV'` parameters:
  - `'xg_boost_model'`: The model pipeline to be tuned.
  - `'xg_boost_params'`: The hyperparameter space to be explored.
  - `'cv=6'`: The number of cross-validation folds to use (6 in this case).
  - `'scoring='accuracy'`: The metric to be used for evaluating the model performance (accuracy in this case).
  - `'n_jobs=1'`: The number of parallel jobs to run (1 in this case).
  - `'n_iter=20'`: The number of hyperparameter configurations to sample (20 in this case).
  - `'random_state=42'`: The random state for reproducibility.
- The `'xg_boost_random_search.fit(X_numeric_train, y_train)'` line fits the XGBoost model with the best hyperparameters found during the randomized search.

## 2. Hyperparameter Tuning for Naive Bayes:

- The code defines a dictionary `'naive_bayes_params'` that contains the hyperparameters to be tuned for the Naive Bayes model.
- The hyperparameters include:
  - `'classifier__alpha'`: A smoothing parameter. Larger values can lead to a more robust model, but may underfit the data.
  - `'classifier__fit_prior'`: A boolean indicating whether to learn class prior probabilities or not.
- The code then creates a `'GridSearchCV'` object `'naive_bayes_grid_search'` to perform the hyperparameter tuning.
- The `'GridSearchCV'` parameters are the same as for the XGBoost model, except for the use of grid search instead of randomized search.
- The `'naive_bayes_grid_search.fit(X_text_train, y_train)'` line fits the Naive Bayes model with the best hyperparameters found during the grid search.

## 3. Model Fitting with Best Hyperparameters:

- The code sets the best hyperparameters found during the hyperparameter tuning for both the XGBoost and Naive Bayes models using the `set_params()` method.
- The `xg_boost_model.set_params(xg_boost_random_search.best_params_)` line sets the best hyperparameters for the XGBoost model.
- The `naive_bayes_model.set_params(naive_bayes_grid_search.best_params_)` line sets the best hyperparameters for the Naive Bayes model.
- The `xg_boost_model.fit(X_numeric_train, y_train)` and `naive_bayes_model.fit(X_text_train, y_train)` lines fit the final models using the training data.

#### **4. Cross-Validation:**

- The code performs 6-fold cross-validation on the XGBoost and Naive Bayes models to assess their performance.
- The `cross_val_score()` function is used, with the following parameters:
  - `xg_boost_model` and `naive_bayes_model`: The models to be evaluated.
  - `X_numeric` and `X_text`: The input data for the respective models.
  - `y`: The target variable.
  - `cv=6`: The number of cross-validation folds (6 in this case).
  - `scoring='accuracy'`: The metric to be used for evaluating the model performance (accuracy in this case).
  - `n_jobs=1`: The number of parallel jobs to run (1 in this case).
- The `xg_boost_cv_scores` and `naive_bayes_cv_scores` variables store the cross-validation scores for each model.
- The mean accuracy for each model is then printed using `xg_boost_cv_scores.mean()` and `naive_bayes_cv_scores.mean()`.