

2. a. Randomly select 80% data instances as training, and the remaining 20% data instances as validation. Change the parameter setting on criterion ("gini", "entropy"). Draw a figure showing the training error and validation error w.r.t. different parameter values

The code reads in a CSV file "news-train.csv" using pandas and uses sklearn's CountVectorizer to preprocess the "Text" column of the DataFrame. Then, it splits the data into training and validation sets using train_test_split, and trains a decision tree classifier with different criterion values (gini and entropy). It then computes and saves the training and validation errors for each criterion value. Finally, it plots the training and validation errors for each criterion value on a graph using matplotlib.

The resulting plot shows the training and validation error rates for the decision tree classifier with gini and entropy criterion values. The plot helps to visualize the model performance and compare the training and validation errors for each criterion value. The goal is to choose the criterion value that yields the lowest validation error, and use that to train the final model.

b. Evaluate the decision tree using 5-fold cross-validation. Draw a figure showing the average training error and average validation error w.r.t. different parameter values. The parameters you should tune include: criterion, min_samples_leaf, and max_features (as int). We recommend Google Colab and scikit-learn package

The code reads in a CSV file "news-train.csv" using pandas and uses sklearn's CountVectorizer to preprocess the "Text" column of the DataFrame. It sets a decision tree classifier model and sets different hyperparameters, such as criterion, min_samples_leaf, and max_features, to be evaluated. It then initializes a KFold object with 5 splits to perform 5-fold cross-validation. The model is trained and evaluated on each fold for each combination of hyperparameters, and the average training and validation errors are computed for each model. Finally, the average training and validation errors are plotted on a graph using matplotlib.

The resulting plot shows the average training and validation errors for the decision tree classifier with different hyperparameter combinations. The plot helps to visualize the model performance and compare the average training and validation errors for each combination of hyperparameters. The goal is to choose the hyperparameters that yield the lowest validation error, and use those to train the final model.

Note that this code performs an exhaustive search over all possible combinations of hyperparameters, which can be computationally expensive. Therefore, it's important to be mindful of the computation time and use techniques like random search or Bayesian optimization to search for the optimal hyperparameters more efficiently.

3. Evaluate random forests model on pre-processed training data.

a. Describe your parameter setting.

b. Use 5-fold cross-validation to evaluate the performance w.r.t. the number of trees (n_estimators).

c. Use 5-fold cross-validation to evaluate the performance w.r.t. the minimum number of samples required to be at a leaf node (min_samples_leaf).

d. Report the average accuracy and standard deviation for different parameter values (organize the results in a table).

e. Draw a figure showing the result.

This code uses the RandomForestClassifier algorithm from the scikit-learn library to classify text documents into different categories. It loads the training data from a CSV file, converts the text into numerical features using bag-of-words representation, and performs 5-fold cross-validation for different values of the hyperparameters `n_estimators` and `min_samples_leaf`. It reports the average accuracy and standard deviation of the cross-validation results for each hyperparameter setting and draws a figure to visualize the results. The first figure shows how the accuracy of the model varies with the number of trees in the forest, and the second figure shows how the accuracy varies with the minimum number of samples required to be at a leaf node.

4. Predict the labels for the testing data (using raw training data and raw testing data).
a. Describe how you pre-process the data to generate features. b. Describe how you choose the model and parameters. c. Describe the performance of your chosen model and parameter on the training data. d. The final classification models to be used in this question are limited to decision trees, random forests, and boosting trees (AdaBoost, or GradientBoostingTree). It is OK to use other models/methods to do feature engineering (e.g., using word embeddings).

The code performs a grid search using the GridSearchCV to find the best hyperparameters for a random forest classifier trained on text data from a CSV file. The hyperparameters being optimized are the number of estimators and the maximum depth of the decision trees in the forest. The best hyperparameters found are printed out, along with the training accuracy of the model with those hyperparameters. Finally, the best classifier is used to make predictions on the test data and the predictions are stored in the predictions variable. However, since the `test_data` variable is loaded with the same CSV file as `train_data`, the predictions variable contains the predictions made on the training data rather than the actual test data. To evaluate the performance of the model on the test data, a separate test set should be used.