

Hyperparameter Optimization and Evaluation of Gradient Boosting Classifier on Wine Dataset

Abstract

This document presents a detailed walk-through of a Python code implementation of a Gradient Boosting Classifier applied to the wine dataset from sklearn's datasets module. The focus is on hyperparameter optimization using RandomizedSearchCV and model evaluation using accuracy and confusion matrix.

1. Introduction

Gradient boosting is a powerful machine learning algorithm used for classification and regression tasks. This document delves into the application of a Gradient Boosting Classifier on the wine dataset, with a focus on optimizing the model's hyperparameters using RandomizedSearchCV.

2. Code Walk-through

2.1 Importing Necessary Libraries

The required libraries for data handling, machine learning modeling, and plotting are imported. This includes matplotlib, numpy, sklearn's datasets, train_test_split, GradientBoostingClassifier, accuracy_score, confusion_matrix, RandomizedSearchCV, and scikitplot's plot_confusion_matrix.

2.2 Function Definition

A function named `'optimize_parameters'` is defined which incorporates the whole process of loading data, splitting it, optimizing the model parameters, making predictions, and evaluating the model.

2.3 Loading and Splitting the Dataset

The wine dataset is loaded from sklearn's datasets, and the features and labels are stored in X and y respectively. The dataset is then split into a training set and a testing set using a split ratio of 75:25.

2.4 Model Definition and Parameter Optimization

A GradientBoostingClassifier model is defined without any initial parameters. For parameter optimization, a dictionary of possible parameters is created which includes 'learning_rate', 'subsample', 'n_estimators', and 'max_depth'. Using these parameters, a RandomizedSearchCV object is created and fit to the training data to find the best parameters.

2.5 Model Evaluation

The model's performance is evaluated on the test set using accuracy and confusion matrix. A confusion matrix plot is also generated using the `'plot_confusion_matrix'` function from scikitplot.

3. Code Execution

Upon calling the `'optimize_parameters'` function, the code carries out the model training, parameter optimization, and evaluation process, and provides the best parameters, accuracy, and a graphical confusion matrix as output.

4. Conclusion

This code provides an example of how to use Gradient Boosting for classification tasks, with a particular emphasis on hyperparameter optimization and model evaluation. By implementing this code, users can understand how to effectively optimize machine learning models and evaluate their performance.