Lab -3

PRML AY 2020-21 Trimester - III

March 21, 2021 Deadline: March 26, 2021, 11:59

Random Forest and Bagging

Dataset:- Consider the credit sample <u>dataset</u>, and predict whether a customer will repay their credit within 90 days. This is a binary classification problem; we will assign customers into good or bad categories based on our prediction.

Data Description:-

Features	Variable Type	Value Type	Description
Age	Input Feature	integer	Customer age
Debt Ratio	Input Feature	real	Total monthly loan payments (loan, alimony, etc.) / Total monthly income percentage.
Number_Of_Time_3 0-59_Days_Past_D ue	Input Feature	integer	The number of cases when a client has overdue 30-59 days (not worse) on other loans during the last 2 years.
Number_Of_Time_6 0-89_Days_Past_D ue	Input Feature	integer	A number of cases when the customer has 60-89dpd (not worse) during the last 2 years.
Number_Of_Times_ 90_Days_Late	Input Feature	integer	Number of cases when a customer had 90+dpd overdue on other credits
Dependents	Input Feature	integer	The number of customer dependents
Serious_Dlq_in_2yr s	Target Variable	Binary: 0 or 1	The customer hasn't paid the loan debt within 90 days

Perform the following tasks for this dataset:-

Question-1 (Random Forest): (Total 20 Marks)

- 1. Preprocessing the data. (5 Marks)
 - a. Plot the distribution of the target variable.
 - b. Handle the NaN values.
 - c. Visualize the distribution of data for every feature.
- 2. Train the Random Forest Classifier with the different parameters, for e.g.:- (5 Marks)
 - i. Max features = [1,2,4]
 - ii. $Max_depth = [2,3,4,5]$
- 3. Perform 5 fold cross-validation and look at the ROC AUC against different values of the parameters (you may use Stratified KFold function for this) and Perform the grid-search for the parameters to find the optimal value of the parameters. (you may use GridSearchCV for this) (5 Marks)
- 4. Get the best score from the grid search. (2 Marks)
- 5. Find the feature which has the weakest impact in the Random Forest Model. Briefly justify your answer. (3 Marks)

Question-2 (Bagging): (Total 20 Marks)

- 6. Perform bagging-based classification using Decision Tree as the base classifier. (15 Marks)
 - a. The number of trees to be considered is {2,3,4}.
 - b. Perform 5 fold cross-validation using ROC AUC metric to evaluate the models and collect the cross-validation scores (use function cross_val_score for this).
 - c. Summarize the performance by getting mean and standard deviation of scores
 - d. Plot the model performance for comparison using boxplot.
- 7. Compare the best performance of bagging with random forest by plotting using boxplot. (5 marks)

Here is the Colab notebook <u>attached</u> for your reference.

Instructions:-

Please Submit the necessary code(s) (Notebook) and a PDF explaining and analyzing the steps in both the questions along with necessary plots/figures.

Note:- No submission will be accepted after the final deadline.