**Lab Assignment 6**

**Pattern Recognition and Machine Learning**

**Part 1:-**

**Classification metrics**

Choosing right evaluation metrics for the problem is one of the most important aspect of machine learning. Choice of metrics allows us to compare performance of different models and helps in model selection.

In this task, we will explore following metrics:

- confusion matrix

- accuracy

- precision

- recall

- f1 score

**Dataset**

The training and Testing dataset is available at

https://github.com/mishravipul/data/raw/main/ozone_levels_train.csv

https://github.com/mishravipul/data/raw/main/ozone_levels_test.csv

The dataset is **modified version** of the dataset 'ozone level' on provided by UCI Machine Learning repository.

Original dataset: https://archive.ics.uci.edu/ml/datasets/Ozone+Level+Detection

**Objective**

To learn about classification metrics and compare logistic regression and decision tree on the same dataset

**Tasks**

- define X(input) and Y(output)

- train the decision tree model

- train the logistic model

- construct a confusion matrix

- calculate the classification accurace

- calculate the Precision

- calculate the Recall

- calculate the F1 score

- calculate Area Under ROC Curve

**Further fun**

- Calculate precission and recall

- find the area under the curve for Roc metrics

- impliment below metrics using inbuilt librarires confusion matrix accuracy precision recall f1 score

**Helpful links**

- Classification metrics with google developers: https://developers.google.com/machine-learning/crash-course/classification/true-false-positive-negative

- classification metrics: https://www.kdnuggets.com/2020/04/performance-evaluation-metrics-classification.html

- pd.get_dummies() and One Hot Encoding: https://queirozf.com/entries/one-hot-encoding-a-feature-on-a-pandas-dataframe-an-example

- Differences between Logistic Regression and a Decision Tree: https://www.geeksforgeeks.org/ml-logistic-regression-v-s-decision-tree-classification/

- Decision Tree Classifier by Sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

- Understanding classification metrics like Precision, Recall, F-Scores and Confusion matrices: https://nillsf.com/index.php/2020/05/23/confusion-matrix-accuracy-recall-precision-false-positive-rate-and-f-scores-explained/

- Understanding the ROC Curve: https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc

Part 2-

**Dataset**

The dataset is available at

Original Dataset : https://www.kaggle.com/itsmesunil/bank-loan-modelling

**Features (X)**

1.  Age - Customer's age in completed years. (Numeric)

2.  Experience - No. of years of professional experience (Numeric)

3.  Income - Annual income of the customer. (Numeric)

4.  ZIPCode - Home Address ZIP code. (Numeric)

5.  Family - Family size of the customer. (Numeric)

6.  CCAvg Avg. - Spending on credit cards per month (Numeric)

7.  Education -

    - Education Level (Categorical|Multiclass):

        - 1: Undergrad

        - 2: Graduate

        - 3: Advanced/Professional

8.  Mortgage - Value of house mortgage if any. (Numeric)

9.  Securities Account - Does the customer have a securities account with the bank? (Categorical | Binary)

10. CD Account - Does the customer have a certificate of deposit (CD) account with the bank? (Categorical | Binary)

11. Online - Does the customer use internet banking facilities? (Categorical | Binary)

12. CreditCard - Does the customer uses a credit card issued by UniversalBank? (Categorical | Binary)

**Target (y)**

- Personal Loan : Did this customer accept the personal loan offered in the last campaign? (Binary)

**Objective**

- To apply Logistic Regression and Decision Tree Algorithms on the given imbalanced dataset. and compare the algorithms used on our dataset on the

basis of appropriately used evaluated metrics while presenting a summarized analysis of what you find.

**Tasks**

- Download and load the data (csv file).

- Process the data according to guidelines given in the comments of the respective cells.

- Split the dataset into 80% for training and rest 20% for testing. (sklearn.model_selection.train_test_split function).

- Initialize Logistic Regression and Decision Tree Models (With parameters given in the cell).

- Train the models on the same dataset.

- Compute the confusion matrix for both models and compare.

- Compute a classification report (Precision, Recall and F-1 score) for both models and compare.

- Compute and plot the ROC Curve of both curves and simultaneously compute the ROC-AUC for both models and thereby compare.

- Summarize your findings and give reasoning for your results (comparing task_1 and task_2).

**Further Fun**

- Train model on different train-test splits such as 60-40, 50-50, 70-30, 80-20, 90-10, 95-5 etc. and observe accuracies on both X_train and X_test.

- Shuffle training samples with different random seed values in the train_test_split function. Check the model error for the testing data for each setup.

- Explore ways to deal with imbalanced dataset. Use different methods (such as eliminating outliers and such) to experiment with the given dataset.

**Helpful links**

- pd.get_dummies() and One Hot Encoding: https://queirozf.com/entries/one-hot-encoding-a-feature-on-a-pandas-dataframe-an-example

- Differences between Logistic Regression and a Decision Tree: https://www.geeksforgeeks.org/ml-logistic-regression-v-s-decision-tree-classification/

- When are Decision Trees better than Logistic Regression?: https://www.displayr.com/decision-trees-are-usually-better-than-logistic-regression

- How to choose between Logistic Regression and Decision Trees given a dataset: https://datascience.stackexchange.com/questions/6048/should-i-use-a-decision-tree-or-logistic-regression-for-classification

- Decision Tree Classifier by Sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

- Understanding classification metrics like Precision, Recall, F-Scores and Confusion matrices: https://nillsf.com/index.php/2020/05/23/confusion-matrix-accuracy-recall-precision-false-positive-rate-and-f-scores-explained/