

Programming Assignment 2: Classification Task and Performance Evaluation (10 points)

October 3, 2024

- In this assignment, you will be using the dataset assigned to you in Assignment 1.
- You will be assigned three classification methods from the following classification methods: **Naive Bayes Classifier**, **Support Vector Machine (SVM)**, **Decision Tree**, **Neural Network**, **Random Forest**, **Adaboost**
- Scikit-learn (https://scikit-learn.org/stable/user_guide.html) will be used in this assignment.

1. Naive Bayes Classifier: **GaussianNB** with default parameters.
2. Support Vector Machine (SVM): **LinearSVC** with default parameters.
3. Decision Tree: **DecisionTreeClassifier** with parameter **max_depth=10** and default values for the other parameters.
4. Neural Network: **MLPClassifier** with parameter *hidden_layer_sizes* = (10, 10, 10,) (i.e., 3 hidden layers with 10 nodes each) and default values for the other parameters.
5. Random Forest: **RandomForestClassifier** with default parameters.
6. Adaboost: **AdaBoostClassifier** with default parameters.

1. Use images from ALL FOUR classes.
2. Convert the images to edge histograms. (Assignment 1 - These will be the vector representations of the images). This will be your dataset for Part 3. (0.25 point)
3. Split the dataset into a training set and a test set: For each class, perform a training/test split of 80/20. (0.25 point)
4. Perform standardization on the training dataset. (see <https://scikit-learn.org/stable/modules/preprocessing.html>) (0.25 point)
5. Perform standardization on the test dataset using the means and variances you obtained from the training dataset.
6. (Performance Comparison) Perform stratified 5-fold cross-validation on the 4-class classification problem using the three classification methods (available on canvas) assigned to you. Plot the (3) confusion matrices for using three approaches (clearly label the classes) on the test set (See Figure 1). (If you use code from any website, please do proper referencing. You will get 0 point for this assignment without proper referencing) (3.75 points)

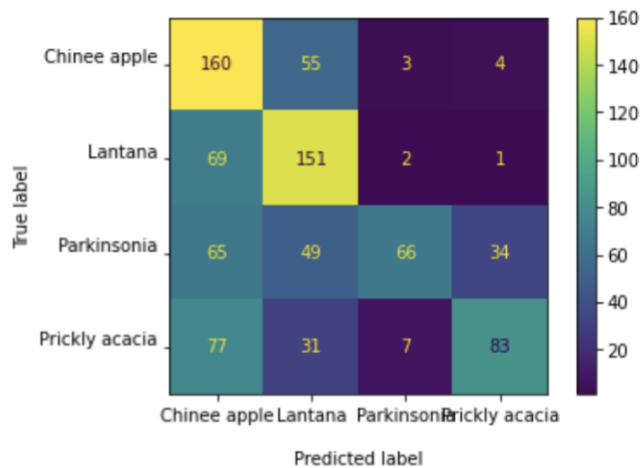


Figure 1:

- By visually comparing (e.g., looking at the color on the diagonal values, etc.) the three confusion matrices (on the test set), which do you think is the best method? Why? (0.50 point)
 - Based on the mean validation accuracies (from the 5-fold cross-validation) for the three methods. Which is the best method? (0.25 point)
 - Compute the accuracies for the three methods on the test set. Which is the best method? (0.25 point)
 - Compute the F-measure for the three methods on the test set. Which is the best method? (0.25 point)
7. (Model Selection) Use images from **TWO** classes. Perform a standard 5-fold cross-validation and a stratified 5-fold cross-validation on the **training set** (i.e., the standardized edge histogram dataset obtained from the training set) for Support Vector Classifiers using **LinearSVC** such that parameter $C = 0.1, 1, 10, 100$ and other parameters set as default. (2.5 points)
- Plot a graph (x-axis: C ; y-axis: mean validation/training error (%)) containing four error curves (2 validation error curves and 2 training error curves - label them clearly using a legend to define the curves). Which C has/have the lowest mean error for each curve? Comment about (1) the model complexity for SVM in relation to C , and (2) when/whether there is overfitting/underfitting. (1.5 points)
 - Use the C value with the lowest mean validation error for your SVM classifier from the stratified 5-fold cross-validation. What is the error for the test dataset (i.e., the standardized edge histogram dataset obtained from the test set)? (0.25 point)