

The provided dataset contains 569 data instances. Each data instance has 30 features that are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe the characteristics of the cell nuclei present in the image. (UCI Dataset). Each instance is classified into malignant or Benign.

Name your Jupyter notebook as **DATA240_Assignment2_Name_SJSU_ID**.

Split the dataset into 70% training and 30% test and provide the following experiments. Report the accuracy for the test set as a performance measurement for all the following tasks

1. Use “from sklearn.tree import DecisionTreeClassifier”.
 - a) Train a DT classifier with Entropy (C1) and GINI (C2) and compare the performance.
 - b) Visualize the C1 and C2 by using the “graphviz” library
 - c) Prune C1 and C2 by limiting the depth and compare their performance with the unpruned versions.
 - d) Use depth 1,...,20 and plot the performance for C1 and C2 separately.
 - e) Choose the best value for depth and visualize C1 and C2.
2. Use “from sklearn.ensemble import RandomForestClassifier”.
 - a) Train an RF classifier with 10 estimators and compare the performance for the test set with C1.
 - b) Change the number of estimators from 10,50,100,500, 1000, and plot the performance.
 - c) Perform 5 fold cross-validation and report the performance for RF classifier with 50 estimators
 - d) Plot the feature importance for RF with 200 estimators using the mean decrease in impurity and also feature permutation and explain the plots.
3. Use “from sklearn.ensemble import AdaBoostClassifier”.
 - a) Train a classifier with 10 estimators and compare the performance with C1 and RF in 2a.
 - b) Change the number of estimators from 10,50,100,500, 1000, and plot the performance.
 - c) Perform 5 fold cross-validation and report the performance for classifier with 50 estimators
4. Use “from sklearn.naive_bayes import GaussianNB”.
 - a) Train a classifier and compare the performance for the test set with C1 and 2a and 3a.
5. Use PCA and print the Cumulative proportion. Using Cumulative proportion, only keep the features that account for more than 95% (ratio of variance to keep) of the total variation associated with all the original variables.

- b) Train an RF classifier with 100 estimators using the dataset with reduced features and compare the performance with RF with 100 estimators using all the features.