

Analysis of the dataset:-

Data set is having imbalance. So, to avoid splitting with imbalance data is being shuffled then splitting of 80:20 ration for training set and test set is done.

As, data is shuffled and splitted , that might change our performance analysis by few figures.

Assumption:

We will consider precision and recall as performance metric for model selection because our data is imbalanced.

[A]

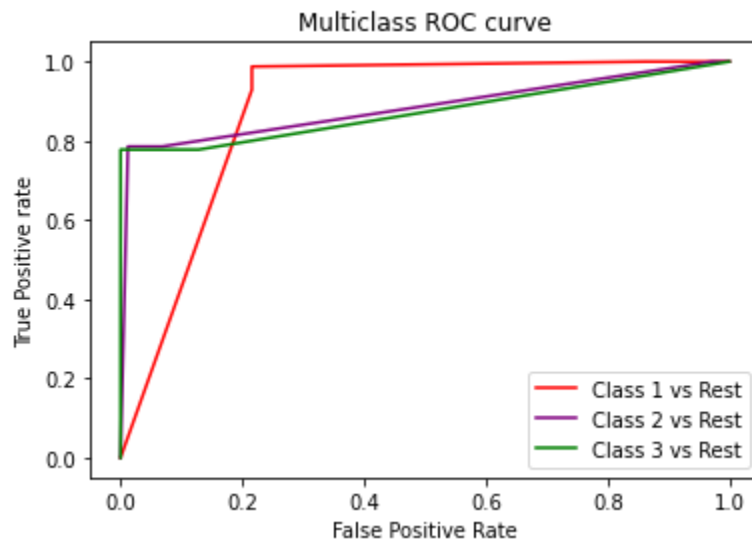
Question 1

Precision: 0.9310257385571193

Recall: 0.8699386774700582

Accuracy: 0.9517241379310345

Auc-Roc curve

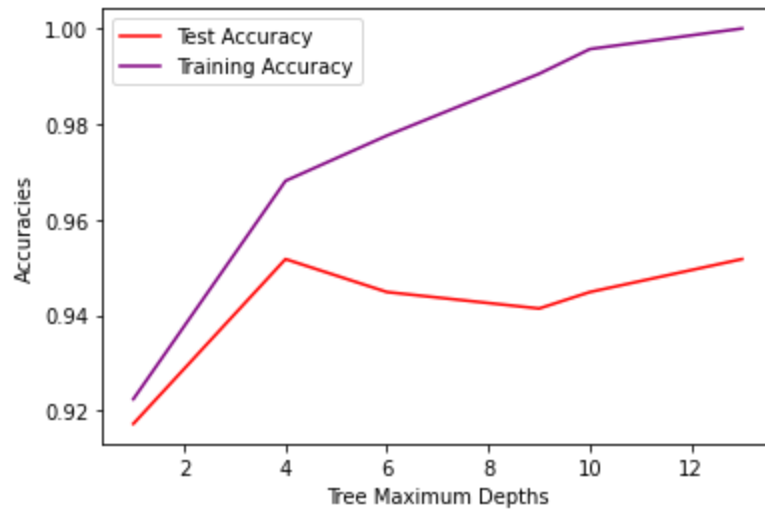


Tree image name:- DT_A_1.jpg.

Question 2

depth=[1, 4, 6, 9,10,13]

Accuracy vs Depth graph.



Here we observe as the depth of tree increases, graph of test accuracy and training accuracy diverges indicating that pruning is required to improve accuracy.

Question 3

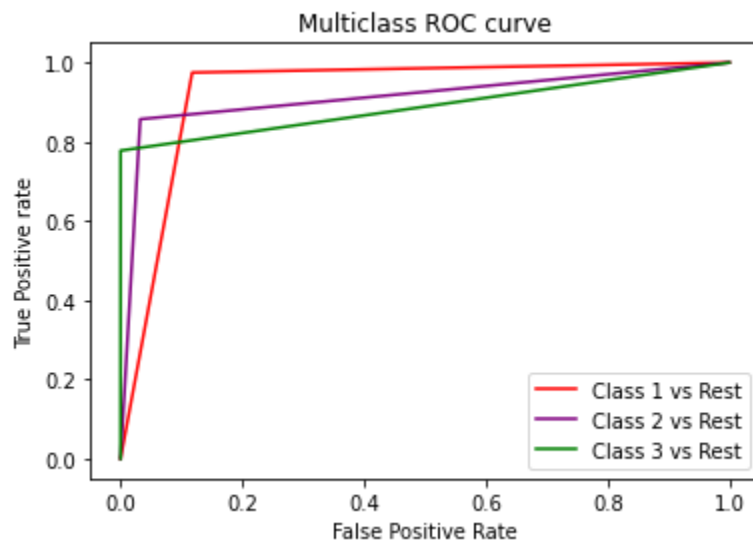
1.Hyper Parameter- Criterion:

Precision: 0.91659660921956

Recall: 0.8302561377875186

Accuracy: 0.9344827586206896

Auc-Roc curve:



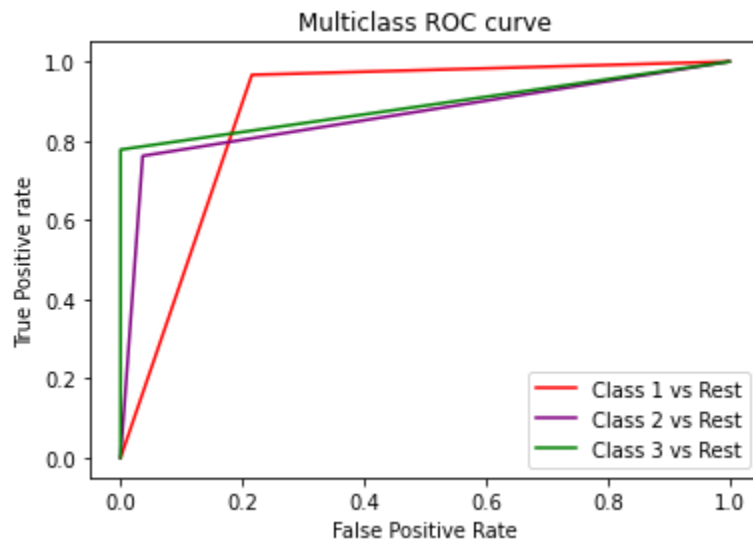
2.HyperParameter-Splitter

Precision: 0.9116777531411678

Recall: 0.8354032454450864

Accuracy: 0.9310344827586207

Auc-Roc curve:



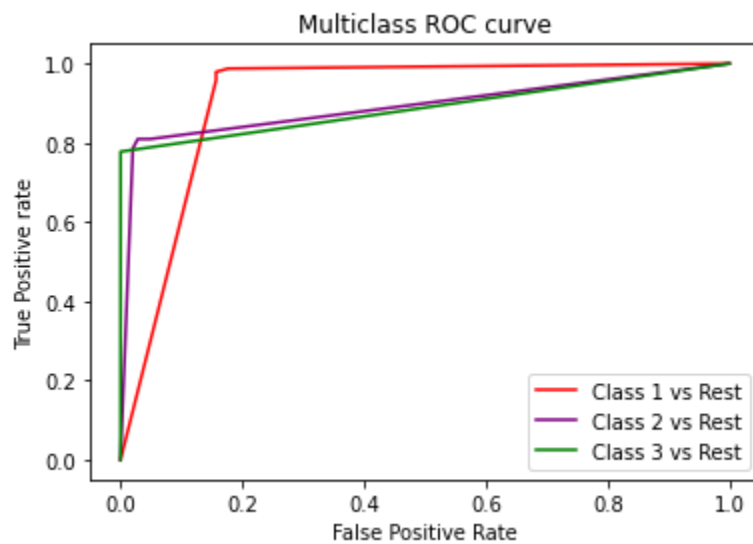
3.HyperParameter-min samples split:

Precision: 0.9438954529180092

Recall: 0.8503132540789444

Accuracy: 0.9517241379310345

Auc-Roc curve:



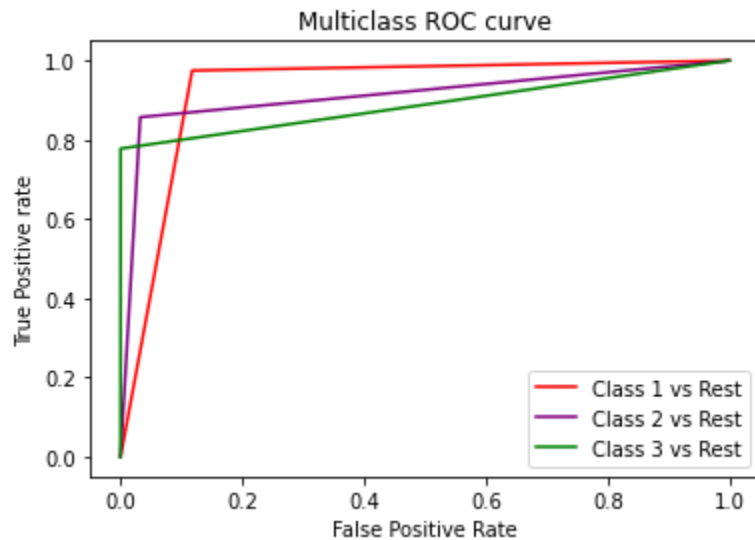
4.HyperParameter-max depth

Precision: 0.9310257385571193

Recall:0.8699386774700582

Accuracy: 0.9517241379310345

Auc-Roc curve:



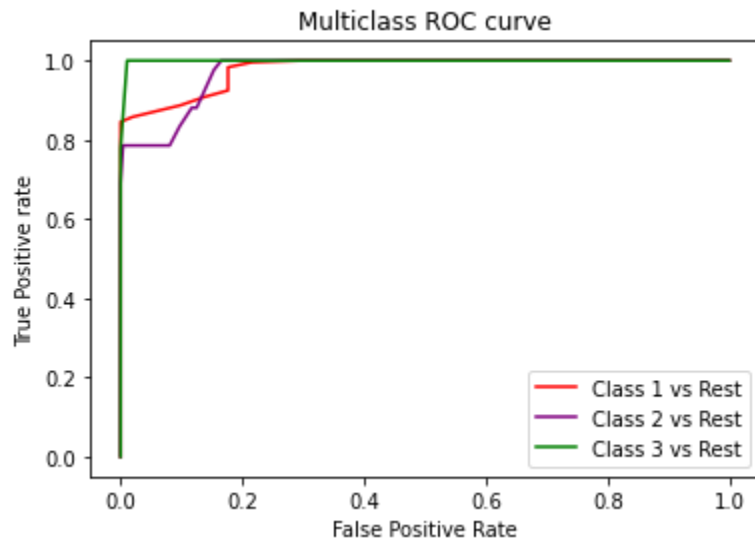
5.HyperParameter-min samples leaf

Precision: 0.9754705094889361

Recall:0.8531026543578845

Accuracy: 0.9586206896551724

Auc-Roc curve:



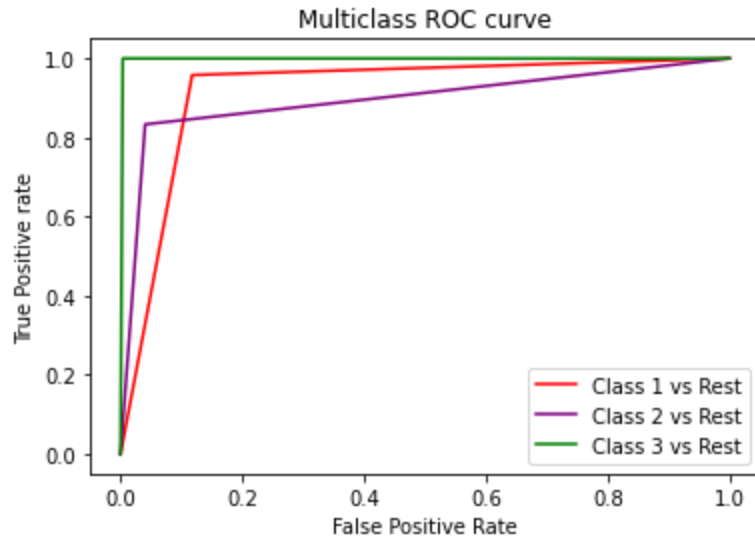
6. HyperParameter-max features

Precision: 0.8840819542947203

Recall:0.9304974430497444

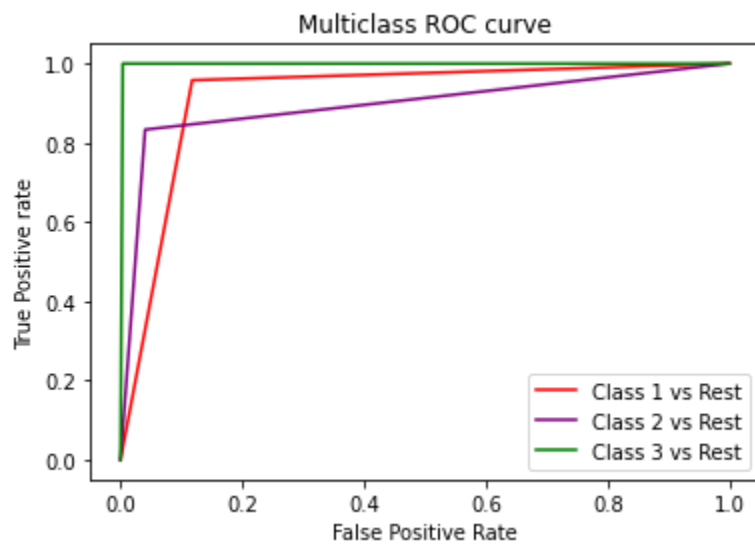
Accuracy:0.9413793103448276

Auc-Roc curve:

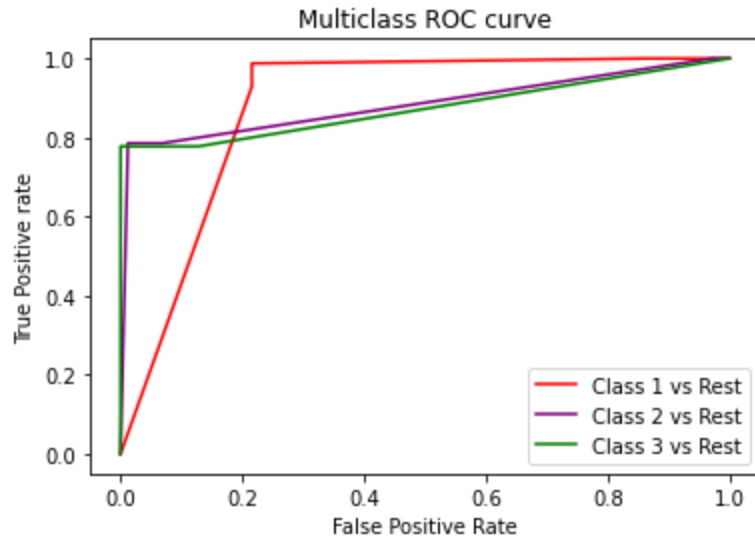


7. HyperParameter-class weight
Precision: 0.9389807162534435
Recall: 0.8924973987735494
Accuracy: 0.9517241379310345

Auc-Roc curve:



8. HyperParameter-max leaf node
Precision: 0.9573774179037337
Recall: 0.8503132540789444
Accuracy: 0.9517241379310345
Auc-Roc curve:



DT-A Selection:-

Hyper-parameters selected for DT-A are max_depth=16,max_features='log2' and min_samples_split= 8

These are selected because by applying them our precision and recall is improving

B. Pruning

Q1.

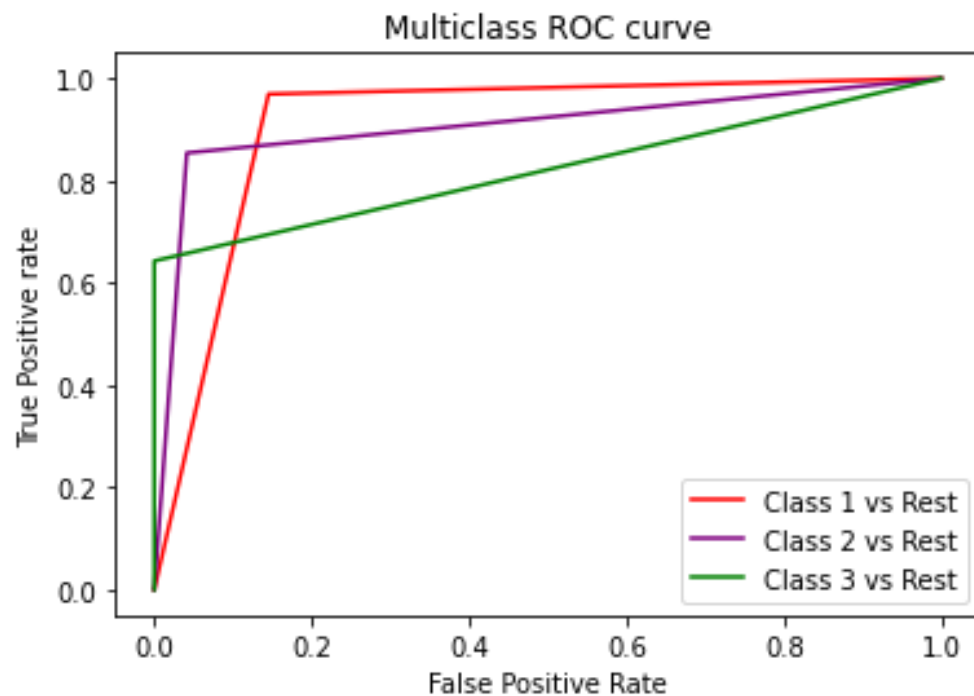
Before Deleting random node from tree

Precision- 0.9215970446149475

Recall – 0.8221073517126148

Accuracy - 0.9344827586206896

AUC-ROC Curve without deleting any random node:

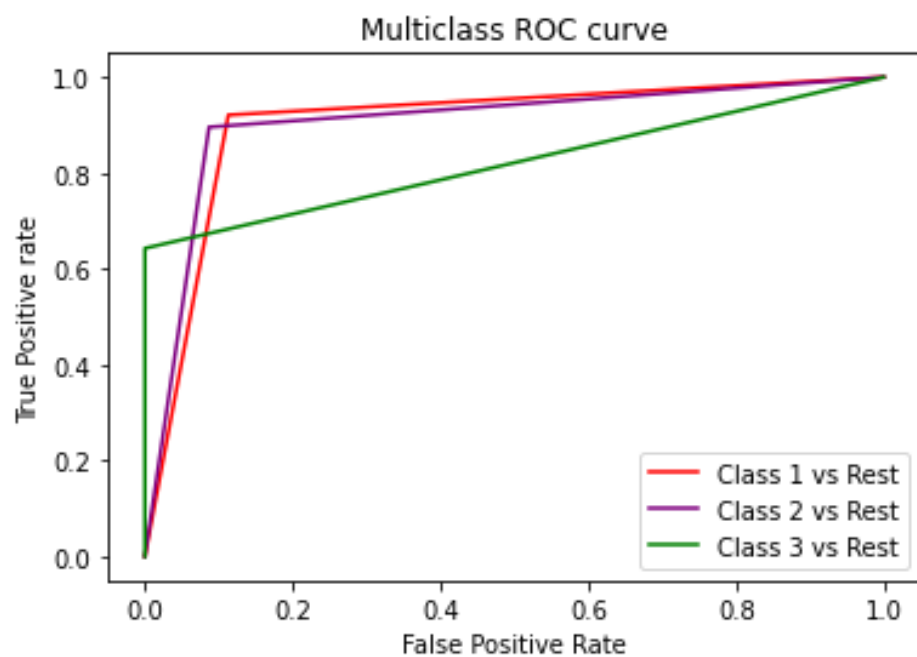


After Deleting random node from tree:-

Precision -0.879872311827957

Recall -0.8199143692564745

Accuracy - 0.903448275862069



After, deleting a random node from the tree precision, recall and accuracy decreased but this will change accordingly which random node is getting pick.
Sometimes there is no change in the performance after deletion of random node

Question 2:

Here,

Post Pruning

Cost complexity Prunning:

Alpha: 0.00090000000000000001

Precision: 0.9104859335038363

Recall: 0.8221757322175732

Accuracy: 0.9379310344827586

Pre prunning:

min_samples_split:2

Precision: 0.8840819542947203

Recall: 0.9304974430497444

Accuracy: 0.9413793103448276

DT-A

Precision: 0.9310257385571193

Recall: 0.8699386774700582

Accuracy: 0.9517241379310345

Observation:

Q3. Hybrid Sliq Pruning on DT-A and name of this tree after this pruning is DT-B-3
Comparison between trees on the behalf of precision,recall,accuracy

Decision Tree	Precision	Recall	Accuracy
DT-A	0.9310257385571193	0.8699386774700582	0.9517241379310345
DT-B-2-CC	0.9104859335038363	0.8221757322175732	0.9379310344827586
DT-B-2-min-samples-split	0.8840819542947203	0.9304974430497444	0.9413793103448276
DT-B-3	0.016091954022988506	0.3333333333333333	0.04827586206896552

Decisionont Tree pruned with hybrid sliq pruning is save as image DT_B_3.jpg

Learning's:

1. Learned how to change decision tree and how does it works in depth
2. A new technique was there hybrid sliq for the learning
3. Learned in understanding of pruning techniques

Reference:

https://devdocs.io/scikit_learn/modules/generated/sklearn.tree.decisiontreeclassifier#sklearn.tree.DecisionTreeClassifier