

TASK 2: Comparison of Classifiers on multiclass classification.

I used wandb platform to visualize plots.

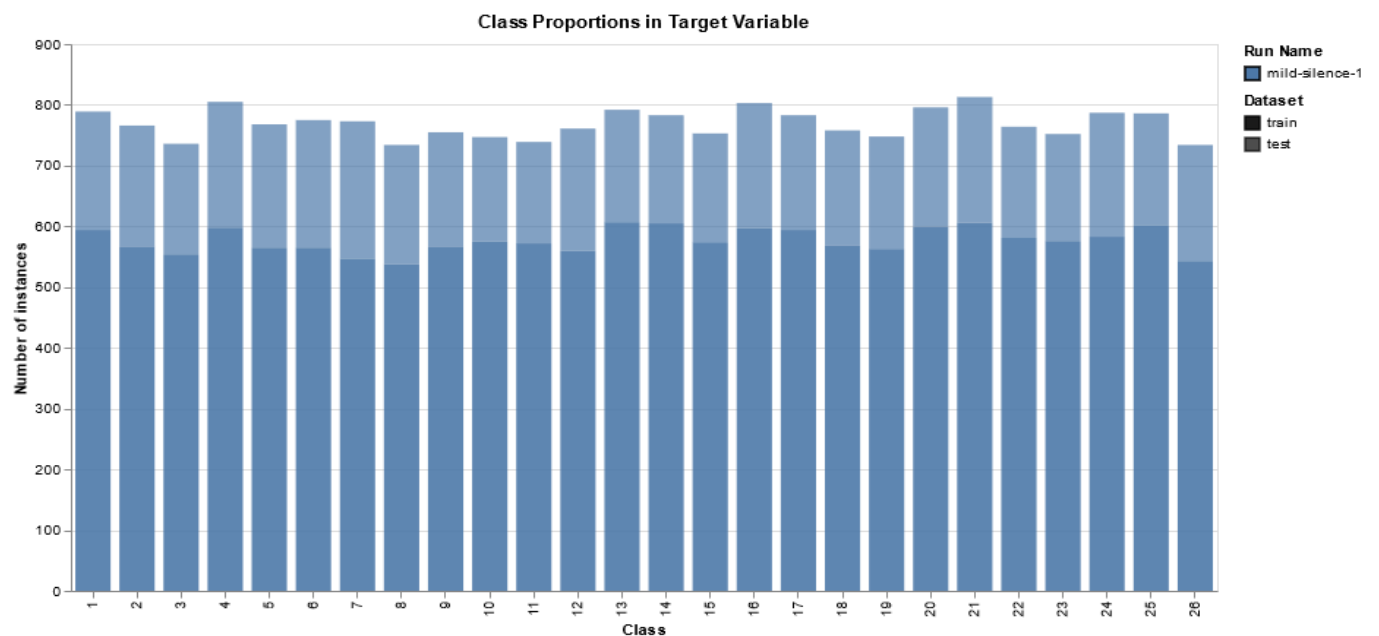
Disclaimer: I did not use seed so the results are variational.

The data statistics:

data	size	# of features
Training	15,000	16
testing	5,000	16

Class_proportion_plot:

Given plot shows the distribution of target classes in training and test sets. I believe modeling target won't be an issue, apparently, classes are equally distributed.



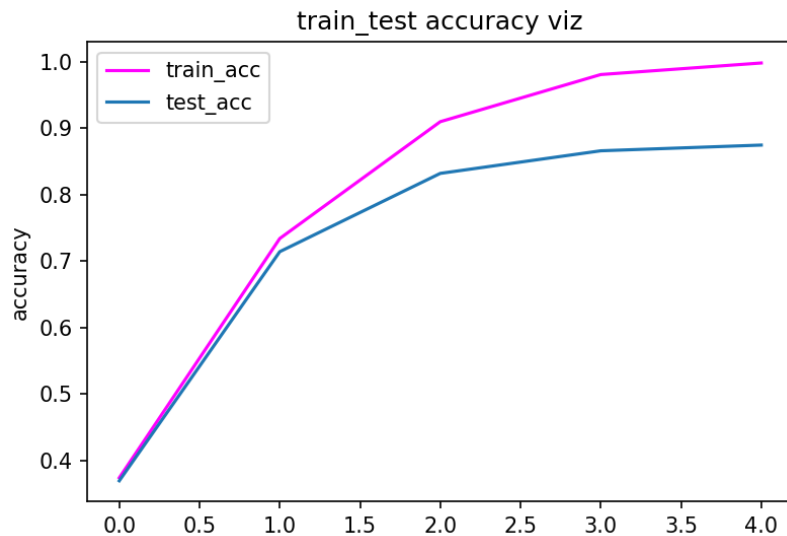
Preprocessing:

- Decision trees don't require scaling the data.
- No null-values were found, so we can go ahead
- with the split is 0.66/0.33

Decision Tree:

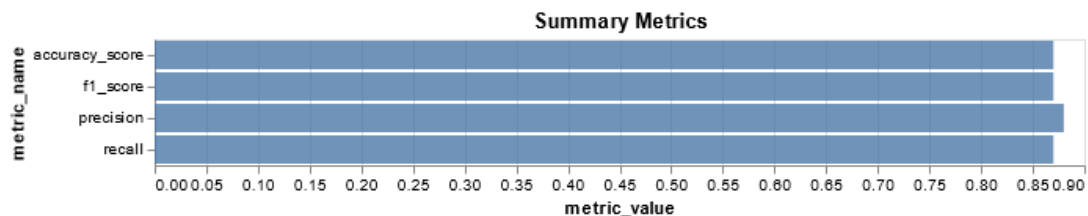
Classifier parameters: for different max_depth (height of the tree) the accuracies are shown below.

1. Visualizing model performance for different max_depth:

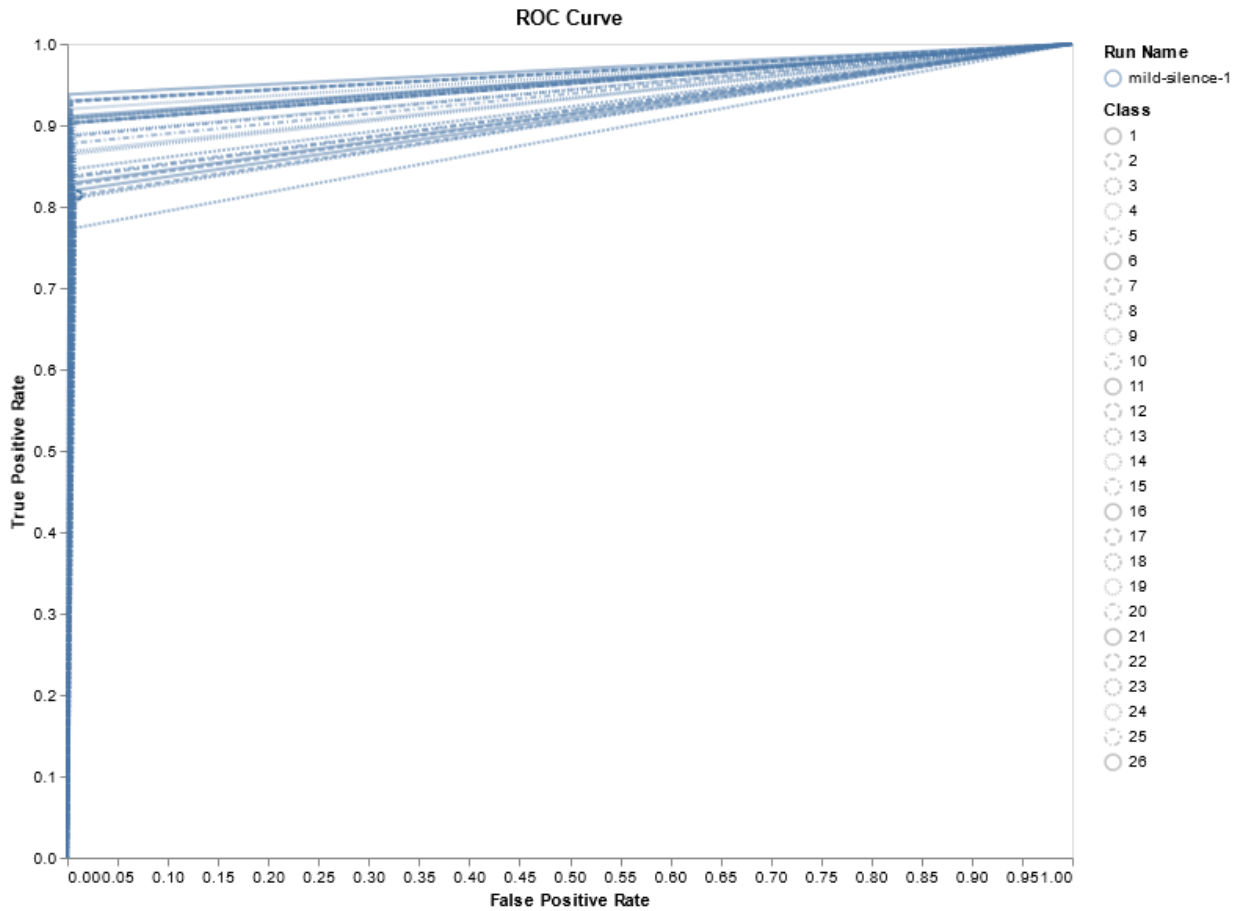


Best model(max_depth = 25, high accuracy):

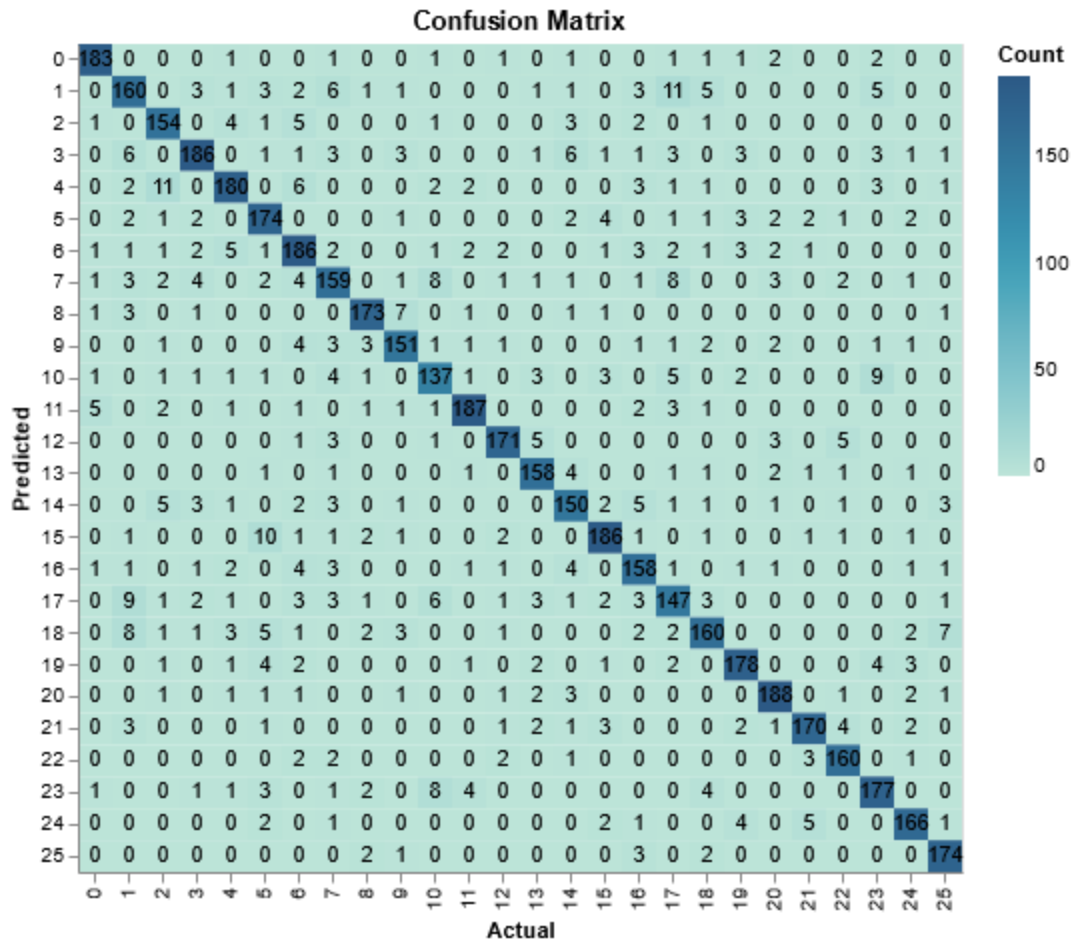
Summary metrics: here we can see accuracy, f1, precision, recall for the classification task. Precision is the highest.



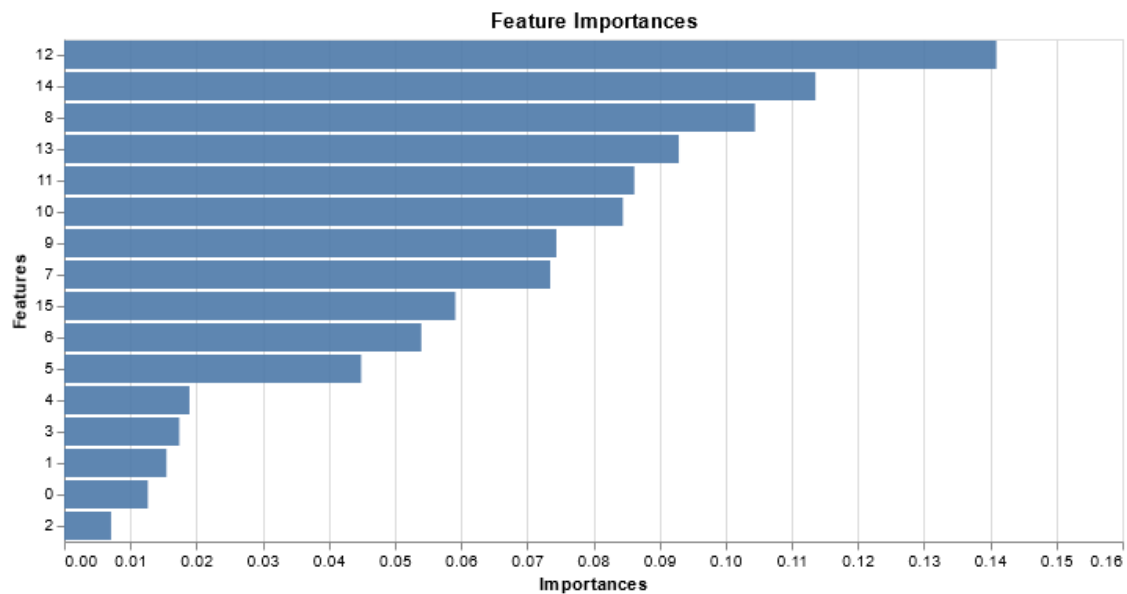
ROC curves plot true positive rate(y-axis) vs false positive rate(x-axis) for each class. The ideal score is a TPR=1 and FPR=0, which is the point on the top left. The greater the AUC-ROC the better.



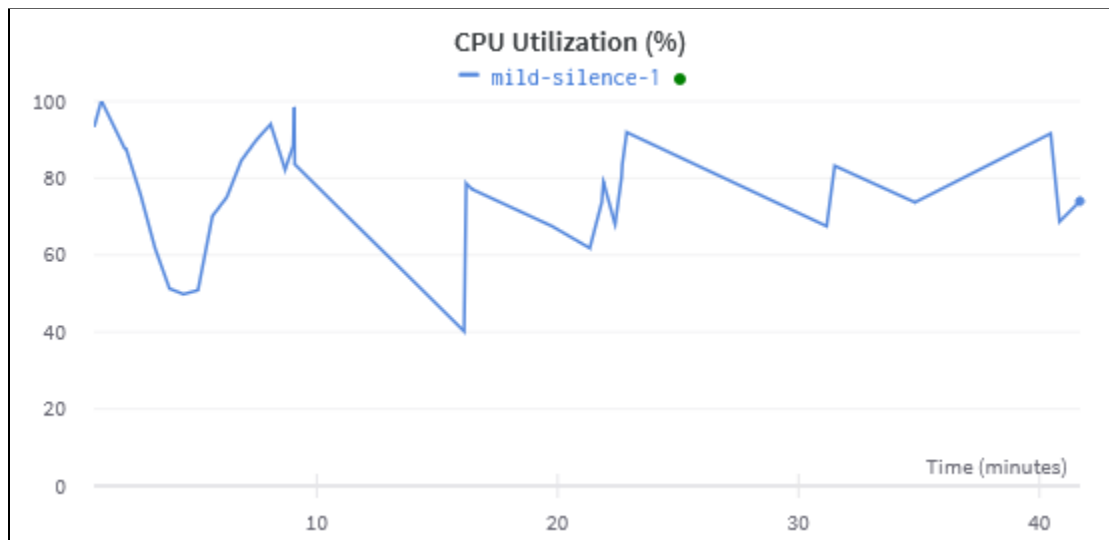
The following plot is the confusion matrix to evaluate the accuracy of classification. Intuitively, we can able to check the quality of model predictions and find patterns in the predictions the model gets wrong. The darker cell —> accurately classified.



Given plot, plots the importance of each feature for the classification task. Available for trees only. Seems target is influenced mostly by feature_12, feature_14,...

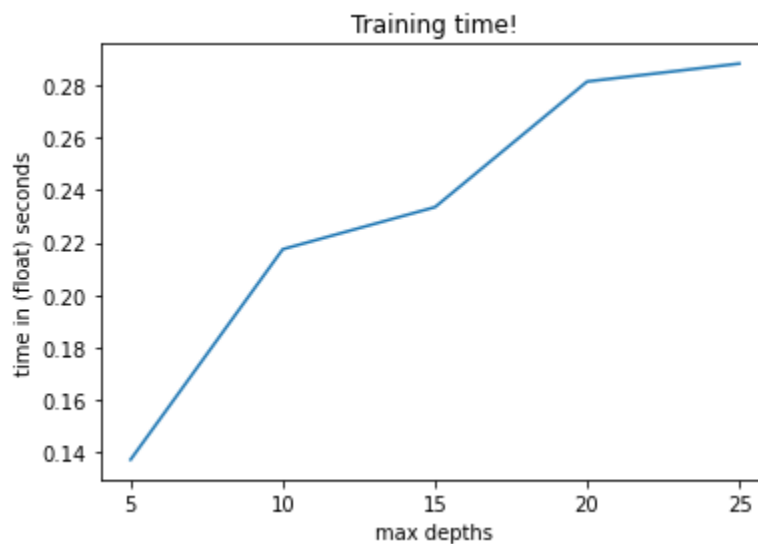


For the overall decision tree experiment, the CPU utilization is (including some running gaps):



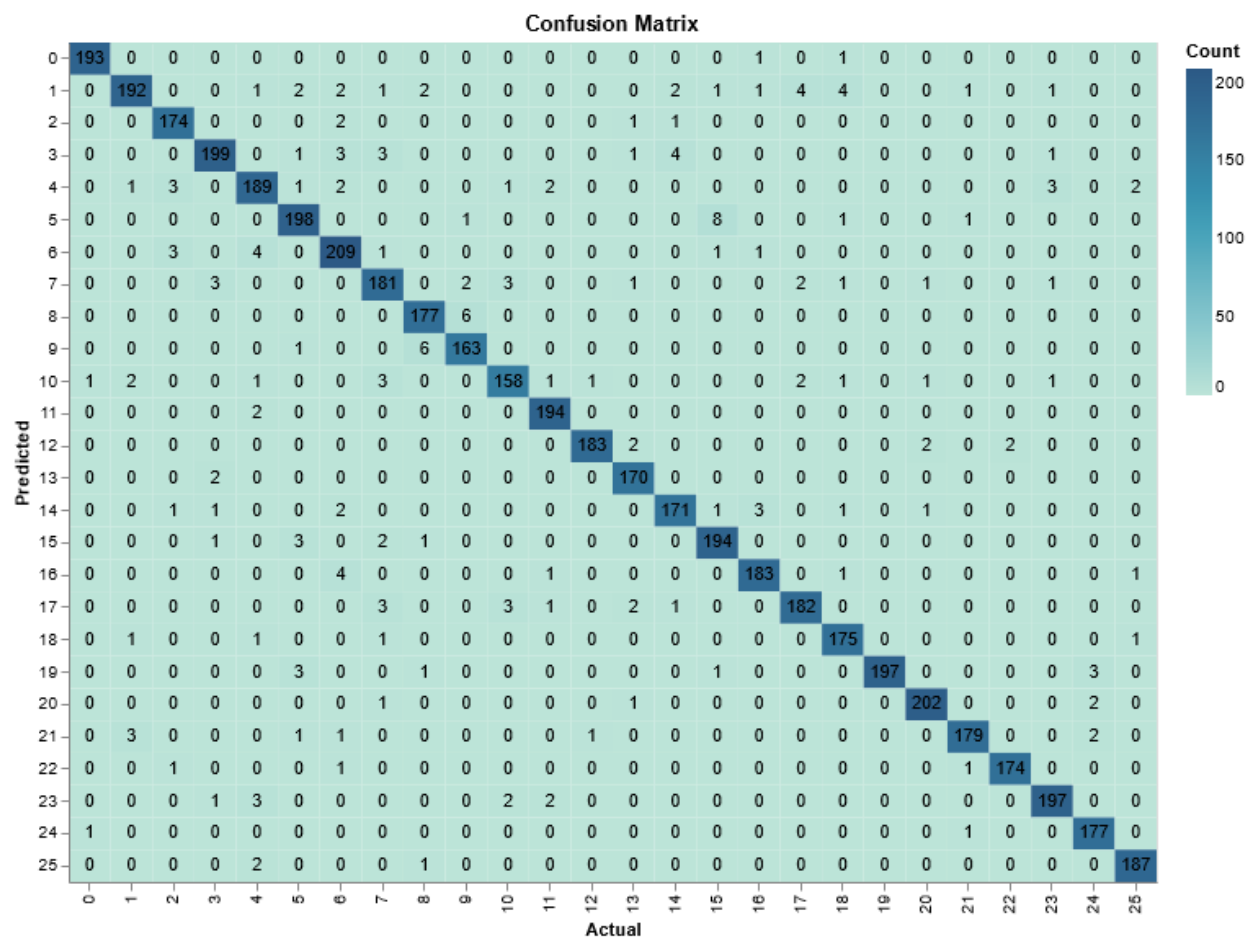
Training_time VS max_Depth:

It's obvious that the greater the max_depth the greater the training time because that many decision boundary gets added.

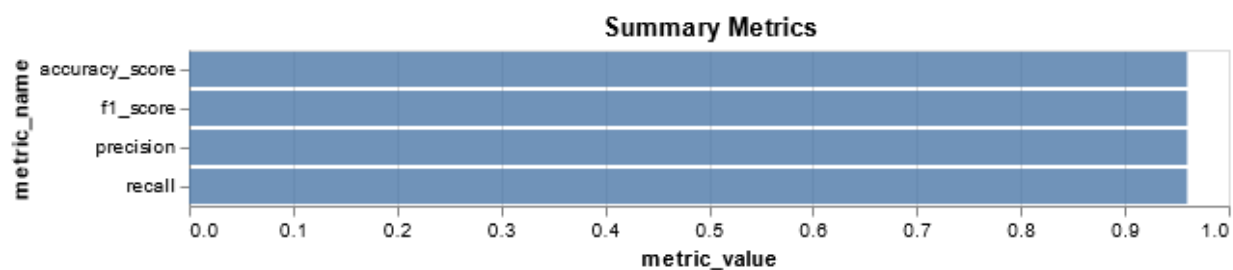


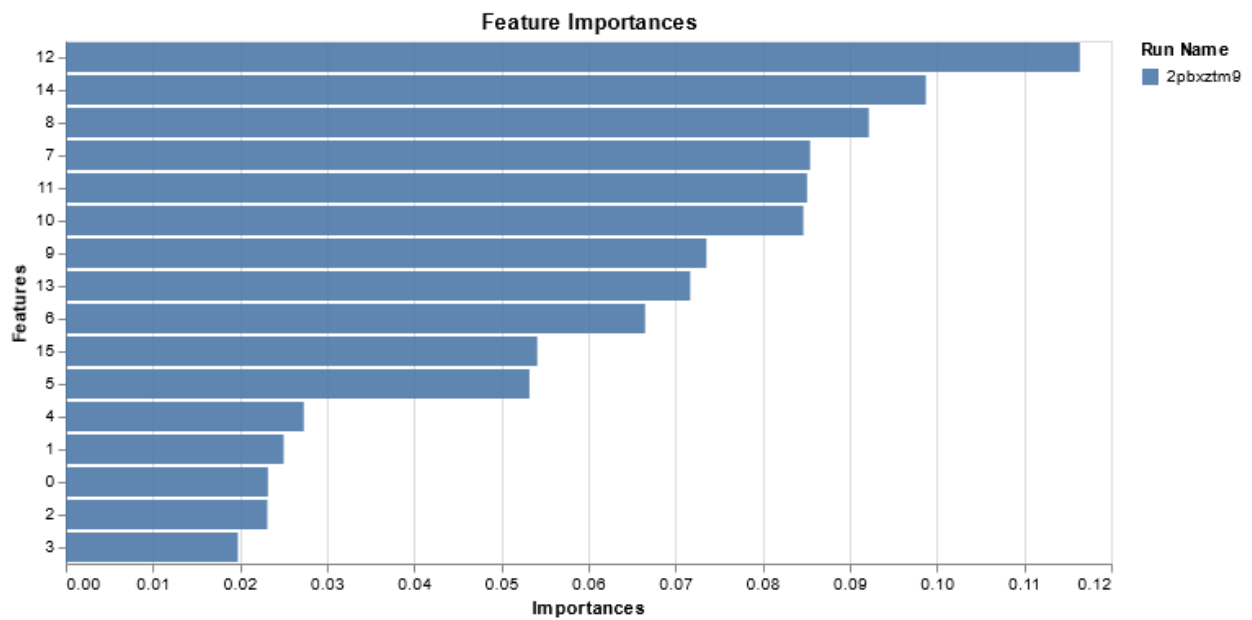
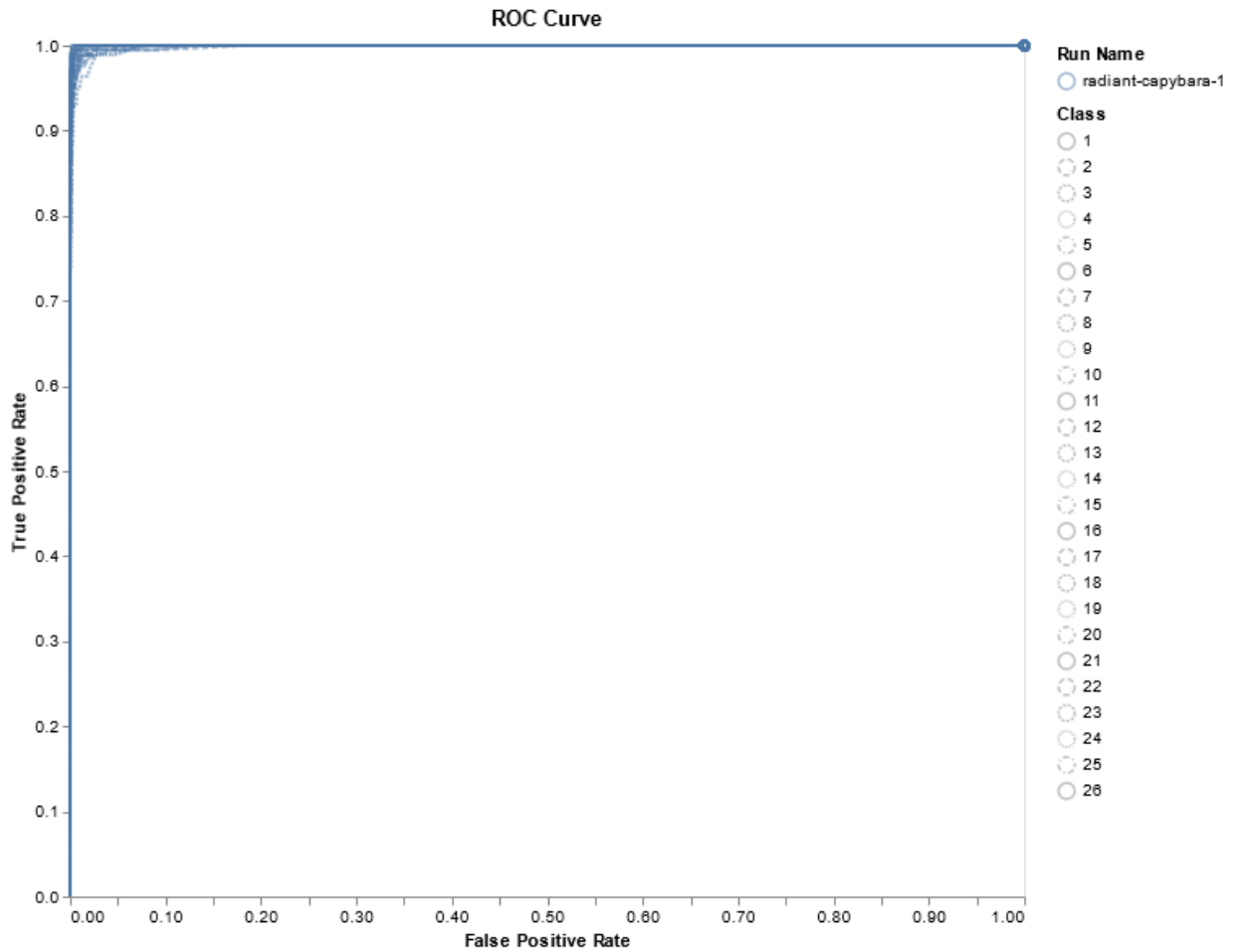
A decision tree with criterion entropy

Confusion matrix, we can see the good quality of predictions:

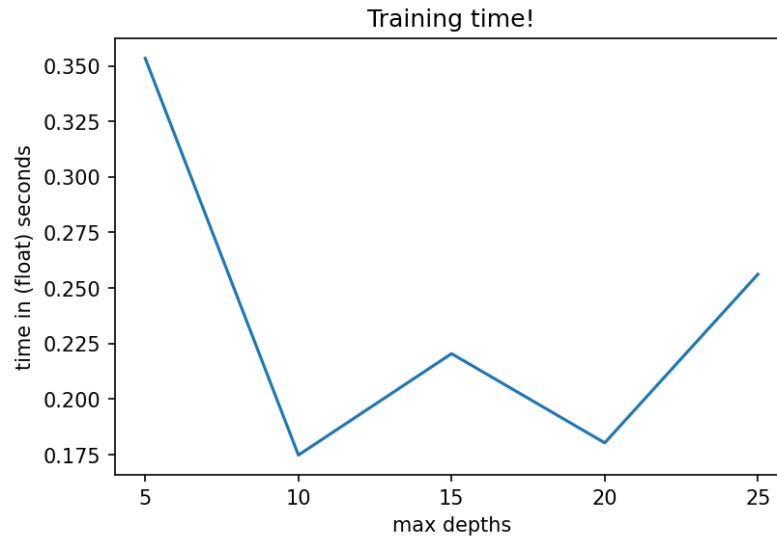


All metrics are better than gini criterion model:

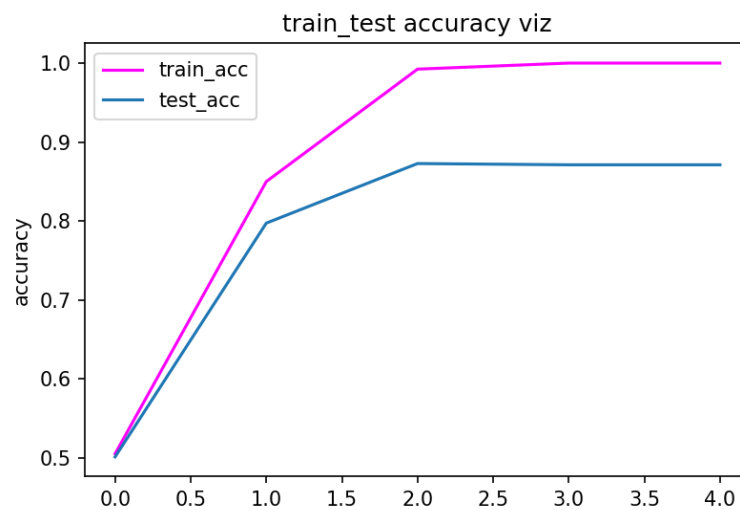




The training time is quietly fluctuating with max_depth, apparently, it's low for high max_depth.

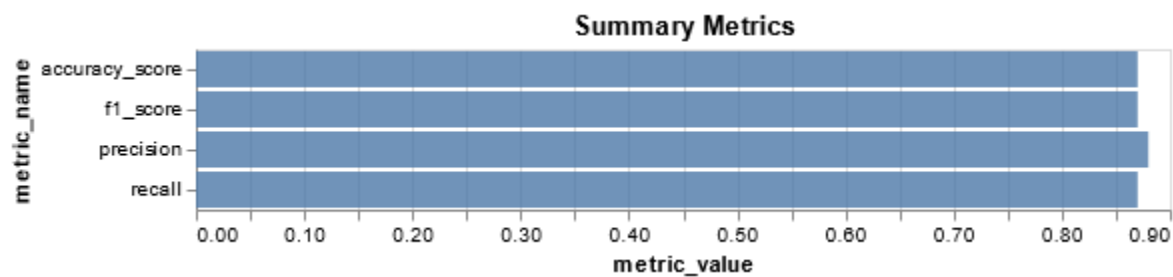
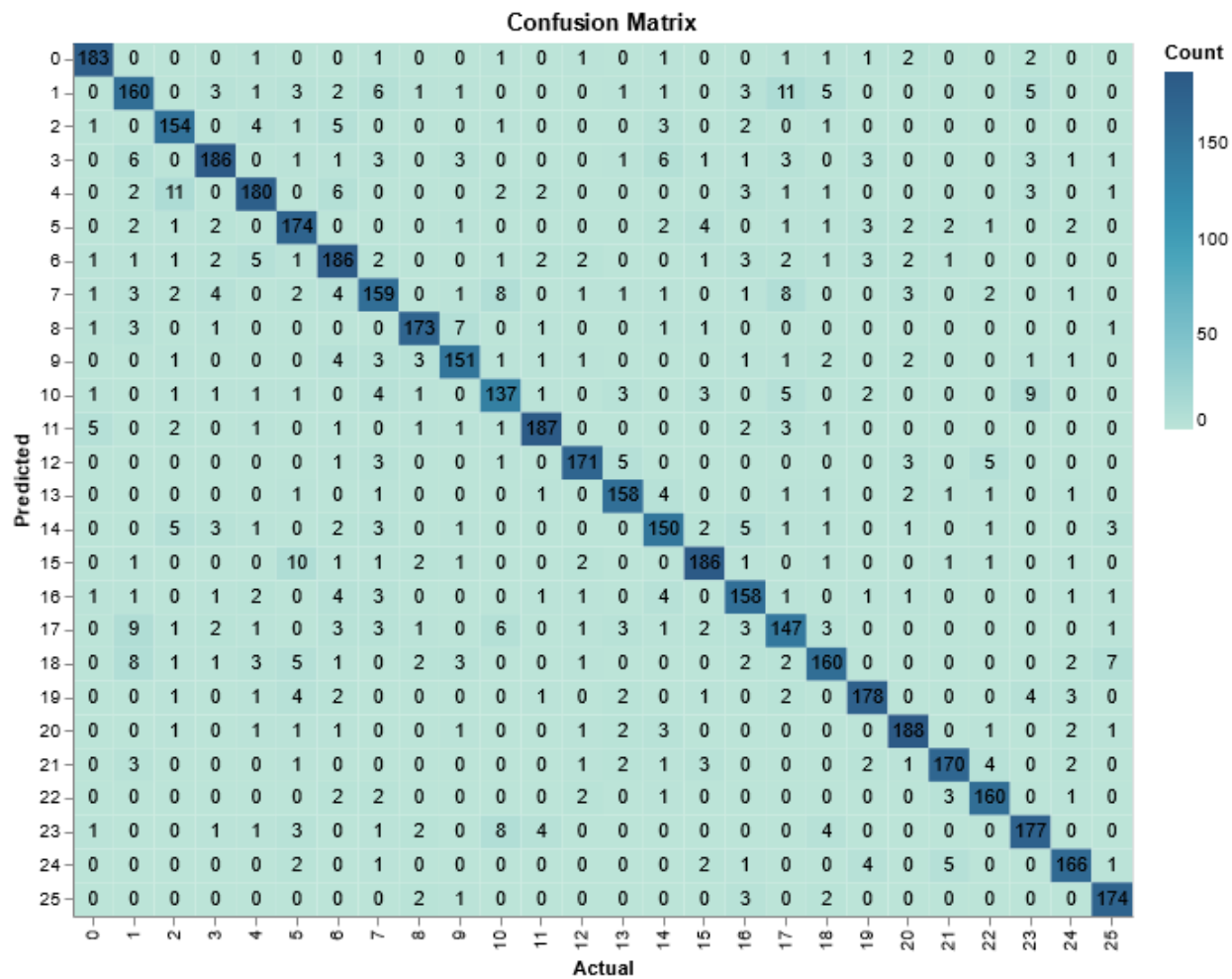


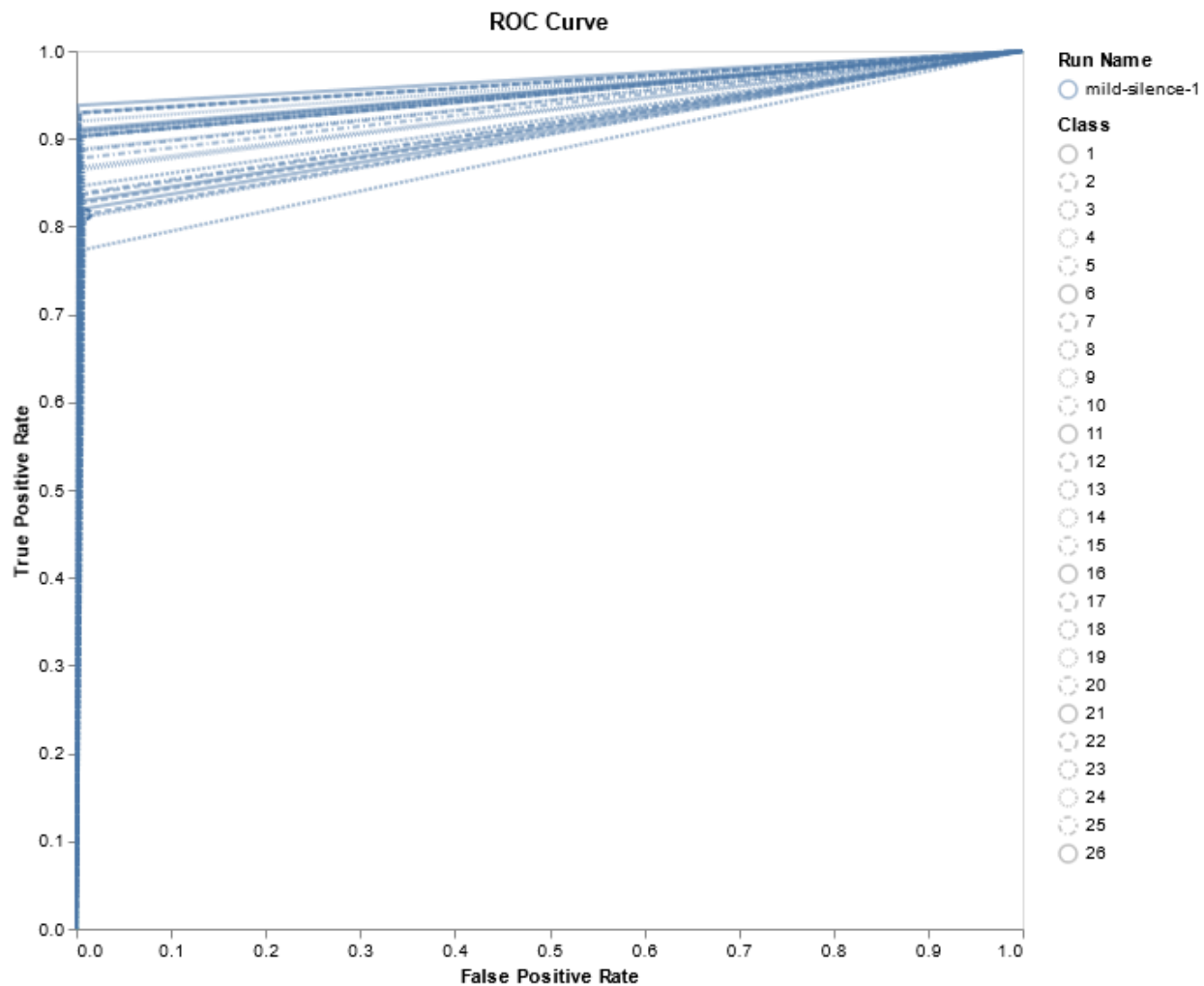
The training and testing accuracy are both stabilizing over the period.



KNN Classifier performance

Here are some plots for KNN: `n_neighbors = 5` (explained in a code):



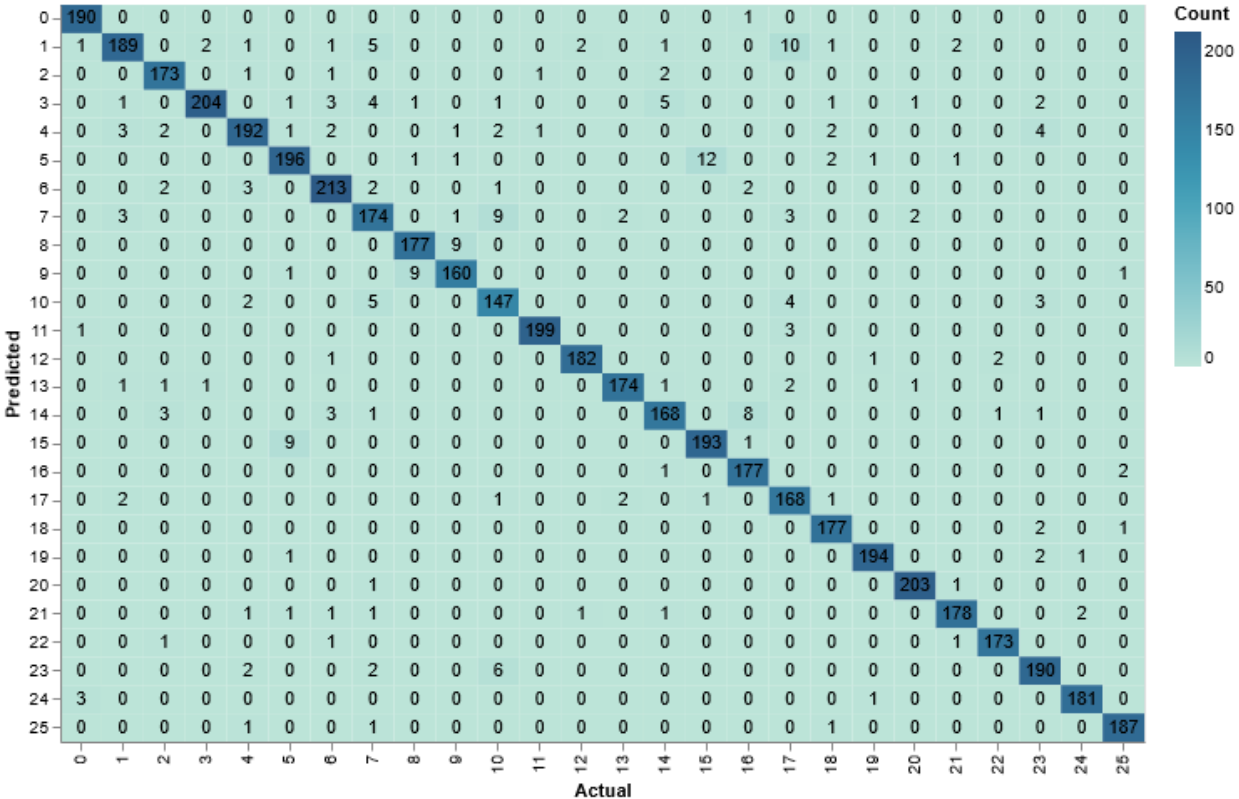


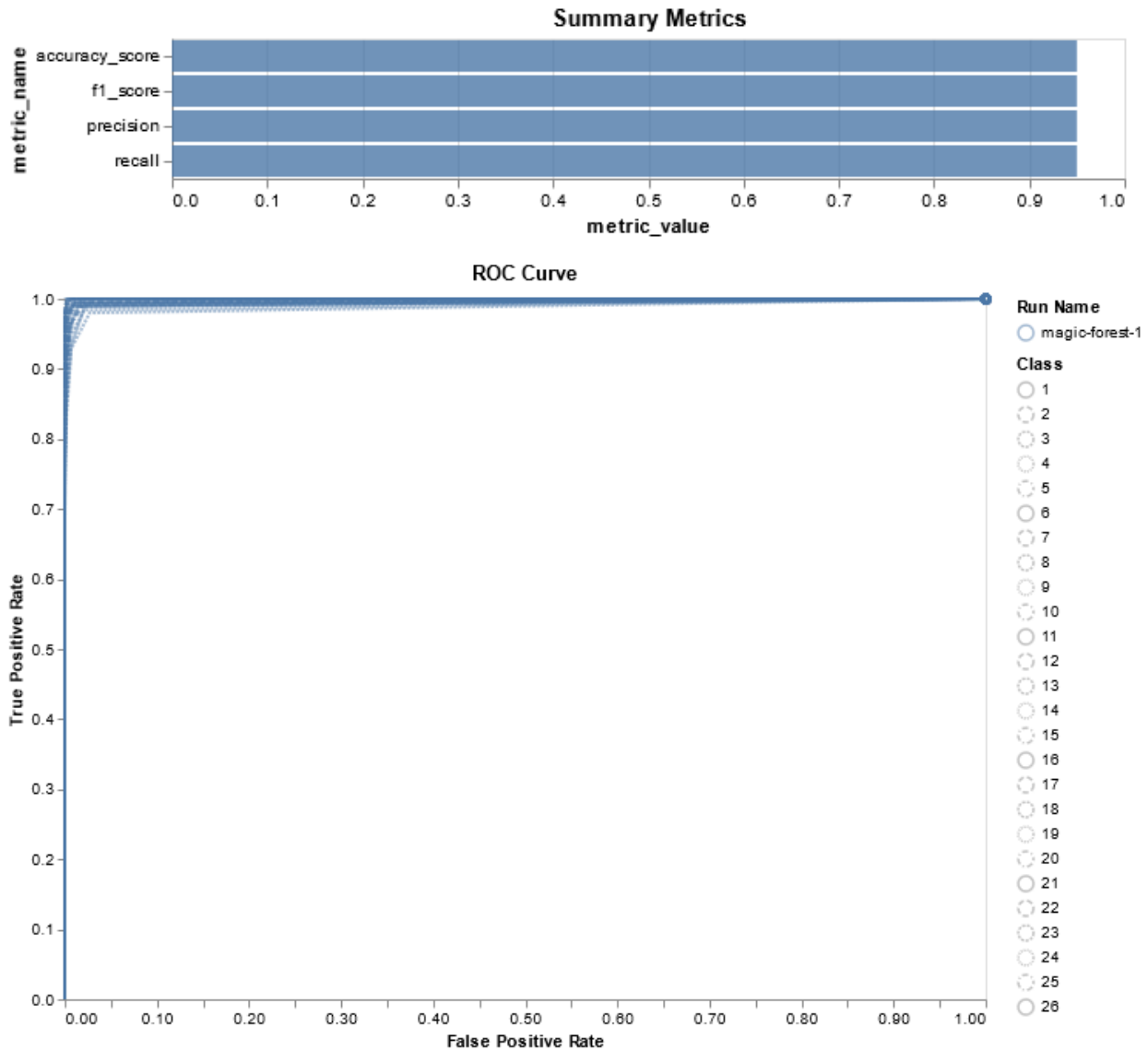
Hyperparameter tuning: best parameters are shown of knn, looks like $n_neighbors=5$ is the best choice

```
training_time(sec) for grid search 39.60739779472351
best accuracy:{ } 0.9415333333333334
best estimator:{ } KNeighborsClassifier(metric='manhattan', weights='distance')
best params:{ } {'metric': 'manhattan', 'n_neighbors': 5, 'weights': 'distance'}
```

Random Forest Classifier performance
Following are the plots for 100 trees in a random forest:

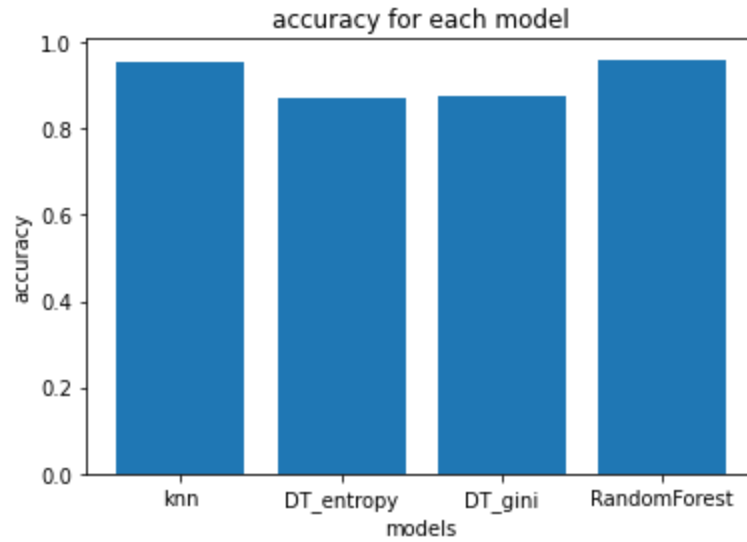
Confusion Matrix



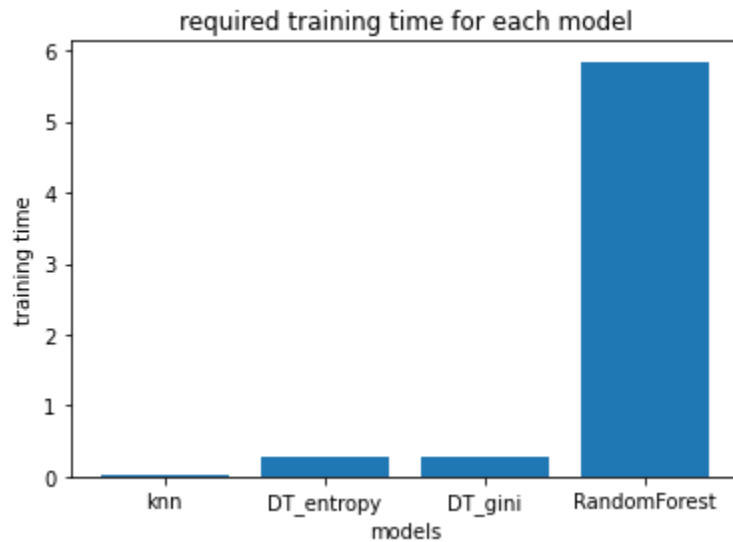


Global Observations:

1. Accuracy for each model:
 - a. Random forest has a greater accuracy but it takes time to train it.
 - b. Decision trees with both criteria has reasonable time and accuracy.
 - c. For this data, KNN will be a good choice because its performance is good training time and accuracy wise.



2. Training time for each model



Further notes:

I did not do hyperparameter tuning for RF because the given model gives reasonable accuracy for test data.

This model-centric approach can be further investigated:

1. Further hyperparameters tuning for each algorithm, and investigate how accuracy is varying with them.
2. By understanding data features in detailed (feature engineering)