Semir Salkić (63190409)

In this task, data is split in training data of 26% (first 130 samples) and test data is 74%. Misclassification rate is recorded for different minimal samples( Figure 1).
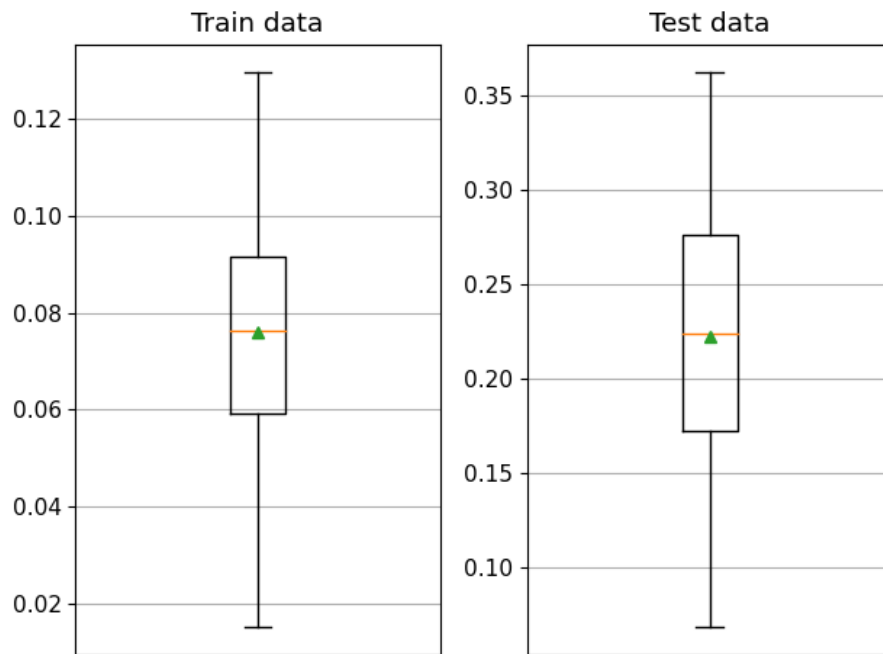


Figure 1: Misclassification rate on decision tree classifier

In Figure 1 misclassification rates samples are presented and divided into train and test data results where misclassification rates are presented. To quantify uncertainty bootstrap is used with 100 iterations for each case.
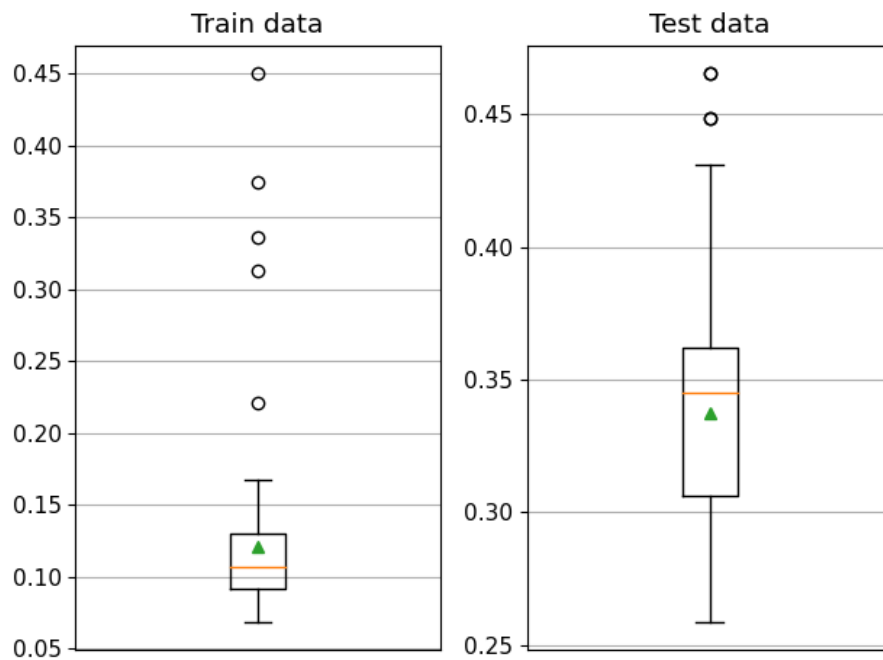


Figure 2: Misclassification rate on random forests classifier

In Figure 2 misclassification rates samples are presented and divided into train and test data results where misclassification rates are presented for random forests with 100 trees. To quantify uncertainty bootstrap is used

with 100 iterations for each case. In presented data we can see that decision tree is much more stable then random forests on given dataset due to multiple outliers produced by the trees in the forest. This can be due to overfitting, which impacted performance of RF ( small number of min. samples).

| Model | mcr – mean | SE – mean |
|-------|-----------|-----------|
| Tree  | 0.2396    | 0.0236    |
| RF    | 0.4965    | 0.0303    |

Table 1: Misclassification rates (**mcr**) and SE for random forests (**RF**) and decision tree (**Tree**)

In the Table 1 misclassification rates and standard errors are presented for case RF($n = 100$, $min.samples = 2$) and Tree ($min.samples = 2$). Using 10 bootstrap iterations samples we can see that Tree has better performance.
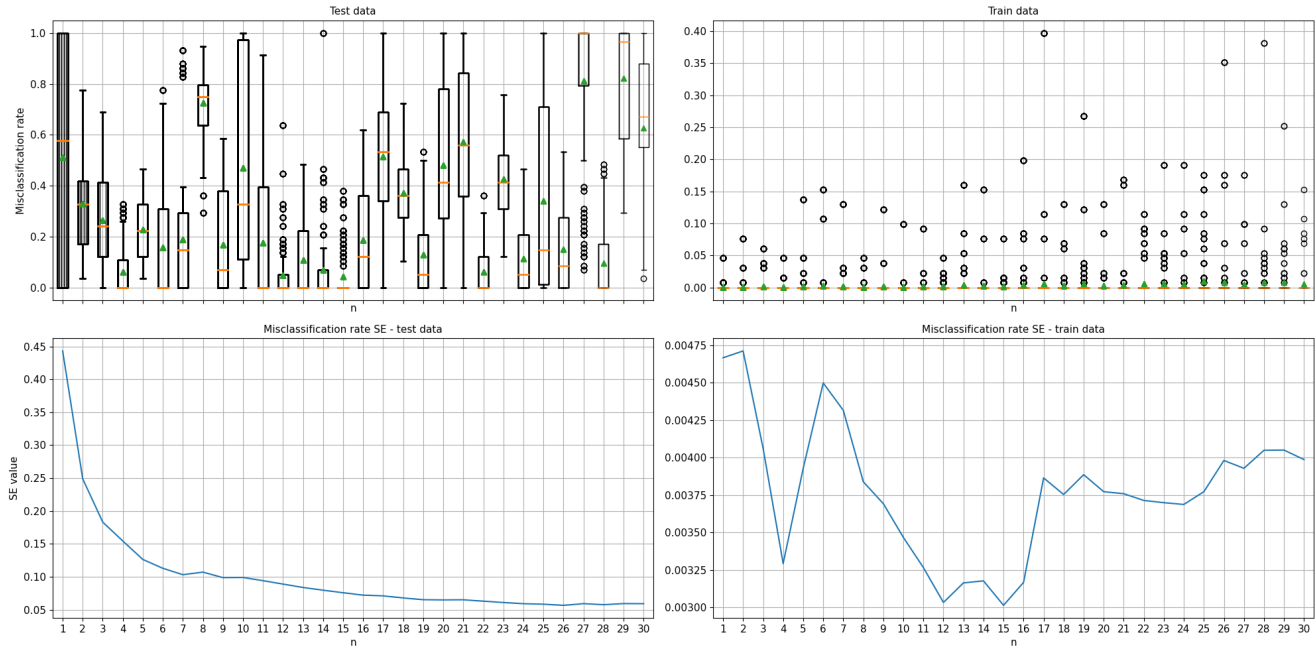


Figure 3: Misclassification rate and SE vs. number of trees n

In the Figure 3 standard error (**SE**) and misclassification rates are presented for random forests. We can see that data is presented for test and train data. We can see that SE of test data is stable,reducing with increasing the number of trees. For train data we can see that most of the classification rates median and means are close to 0.
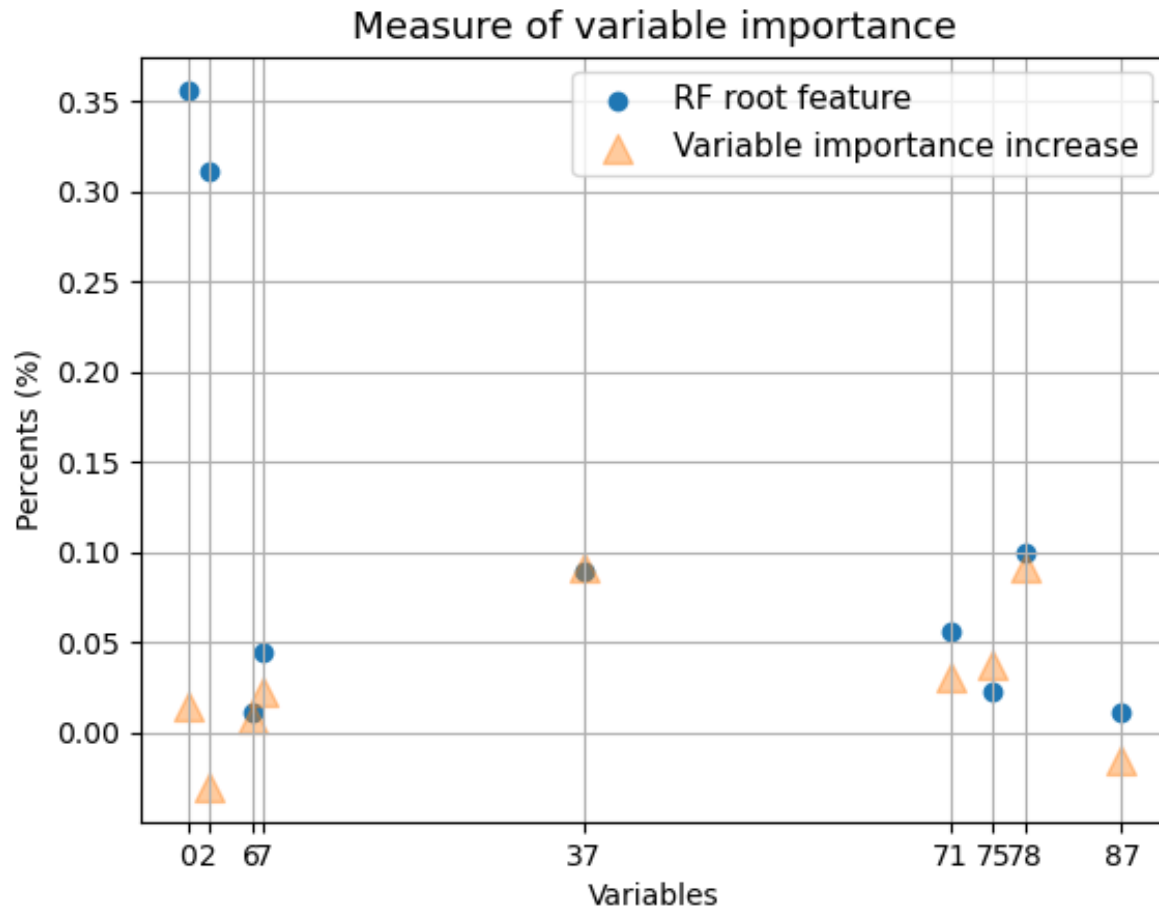
Figure 4: Variable importance of the RF

In Figure 4 feature importance with RF ($n = 100$, 5 minimal samples) is presented. Our dataset has 190 data points and it is consisted from 198 features. We can see that largest increase of mcr is present in features *37* and *78*. Also to compare and verify distribution of root features in trees is presented. In our case we can see that random forest's 36.24% of trees used the best split of *0* feature, while next most used feature is *2* feature is used. Having this in mind, we can see that RF model evaluated correctly variable importance.