

Variation of accuracy as function of maximum depth using 5-fold cross validation:

Maximum depth	Average accuracy
0	0.33333333333333
1	0.66666666666666
2	0.91333333333333
3	0.94666666666667
4	0.94666666666667
20	0.94666666666667

The data provided (150 samples) was partitioned into 5 different parts, with each part containing 30 samples (10 of each class to maintain balance). Then the tree was built (training) using 4 partitions and after that the prediction was validated with the remaining partition. This was done 5 different times and the accuracy is the average of those 5 times.

I assumed the maximum depth referred to the maximum level of the tree, so a depth of 0 would mean only the root, which is all the samples and as there's the same proportion for every sample, the accuracy is $\frac{1}{3}$ to guess "randomly". But after I allowed the tree to have more depth and splits, the accuracy improved until I got to level 3 which was where all the nodes became leaf nodes, either because they only contained one class or because there was not another way to reduce impurity. The reason for this is probably because the amount of samples was too low. The problem with using a decision tree alone is that it tries to find an optimal solution at each node with a greedy algorithm, which would probably find a local optimum in each step, but not a global optimum. It's also a supervised model which depends on the training data, so maybe it's representing well the whole population. On the other side, a random forest combines different independent classifiers and aggregates information, which would reduce variance. At the end, they limit overfitting without increasing that much error due to bias.