# Decision Trees

Maciej Małecki

August 24, 2019

## Introduction

The purpose of my project is to analyze the impact of different hyperparameters on obtained clustering solutions using Decision Trees from sklearn python package. The dataset was 'Social Network Ads', with attributes such as UserID, Gender, Age, EstimatedSalary, Purchased (if advertised item was purchased or not). Table presents first 5 of 400 observations of dataset and table shows basic descriptive statistics.

| UserID | Gender | Age | EstimatedSalary | Purchased |
|--------|--------|-----|-----------------|-----------|
| 15624510 | Male | 19 | 19000 | 0 |
| 15810944 | Male | 35 | 20000 | 0 |
| 15668575 | Female | 26 | 43000 | 0 |
| 15603246 | Female | 27 | 57000 | 0 |
| 15804002 | Male | 19 | 76000 | 0 |

Table 1: First 5 observations of dataset

|  | Age | EstimatedSalary | Purchased |
|--------|-----|-----------------|-----------|
| mean | 37.655 | 69742.5 | 0.3575 |
| std | 10.4829 | 34096.9603 | 0.4799 |
| min | 18 | 15000 | 0 |
| 25% | 29.75 | 43000 | 0 |
| median | 37 | 70000 | 0 |
| 75% | 46 | 88000 | 1 |
| max | 60 | 150000 | 1 |

Table 2: Summary of basic statistics dataset

# Comparing of influence of different hyperparameters

## max_depth=2

Firstly, we have chosen the max_depth = 2. Classification is based on two attributes – on Age and EstimatedSalary, and the criterion is 'enthropy'. The accuracy for this classification on the test set was 94% and actually it was the best result we have got. The decision tree we can see on 1 and the outputs on figures 2. As we state from these, The clusterisation is simple but yet very effective for data.
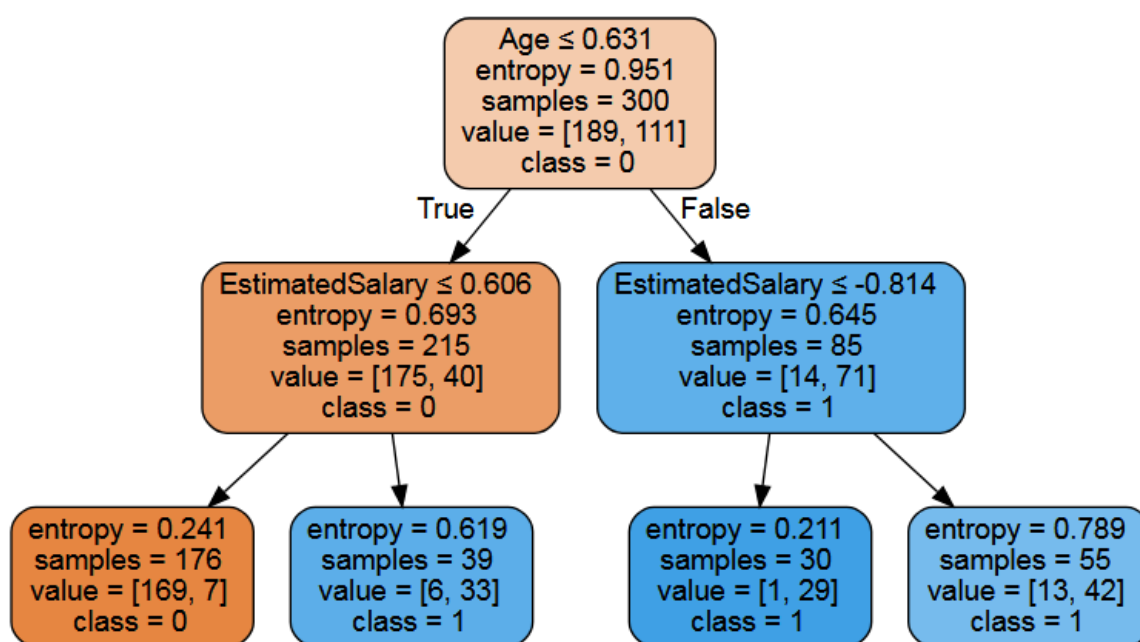


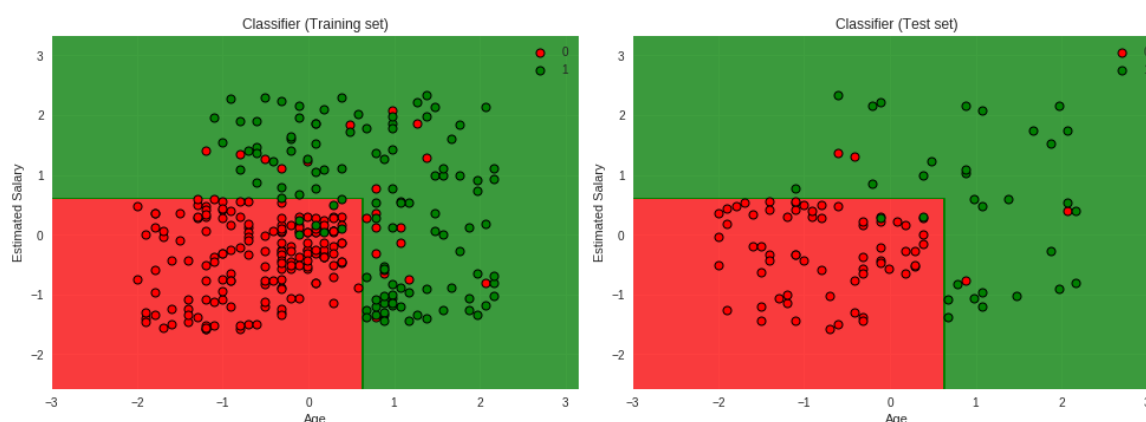Figure 1: Image of decision tree generated using Graphviz, max_depth=2



Figure 2: Training and validation visualization, max_depth=2

2

## max_depth=4

In case of max_depth=4, the accuracy was 93%. The decision tree from fig. 3 is starting to be complicated, but as we can see on fig. 4 it is maybe better suited to training data but not necessarily to test data, but accuracy is still satisfying.
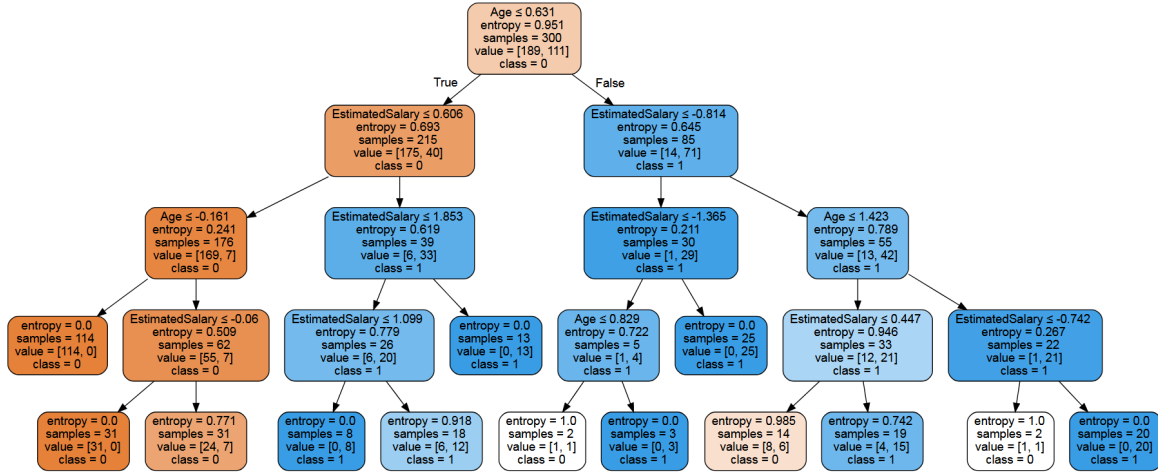


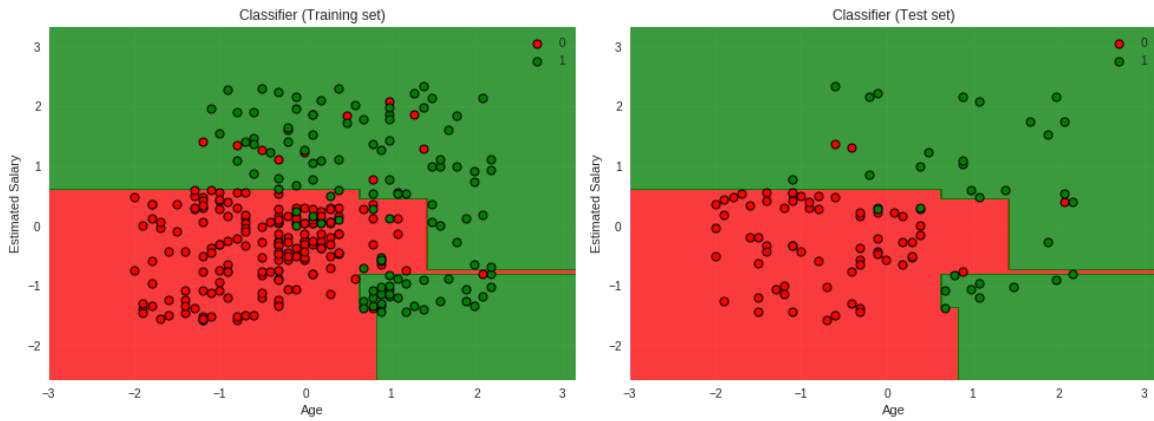Figure 3: Image of decision tree generated using Graphviz, max_depth=4



Figure 4: Training and validation visualization, max_depth=4

## max_depth=6

The decision tree from fig. 5 is very, maybe even unnecessarily, complicated. With only two attributes it starts to have the characteristics of overfitting. The accuracy is much lower than in previous cases and is equal 90%. In the fig. 6 we can see rectangular-shaped areas which badly classify both training set points and test points.
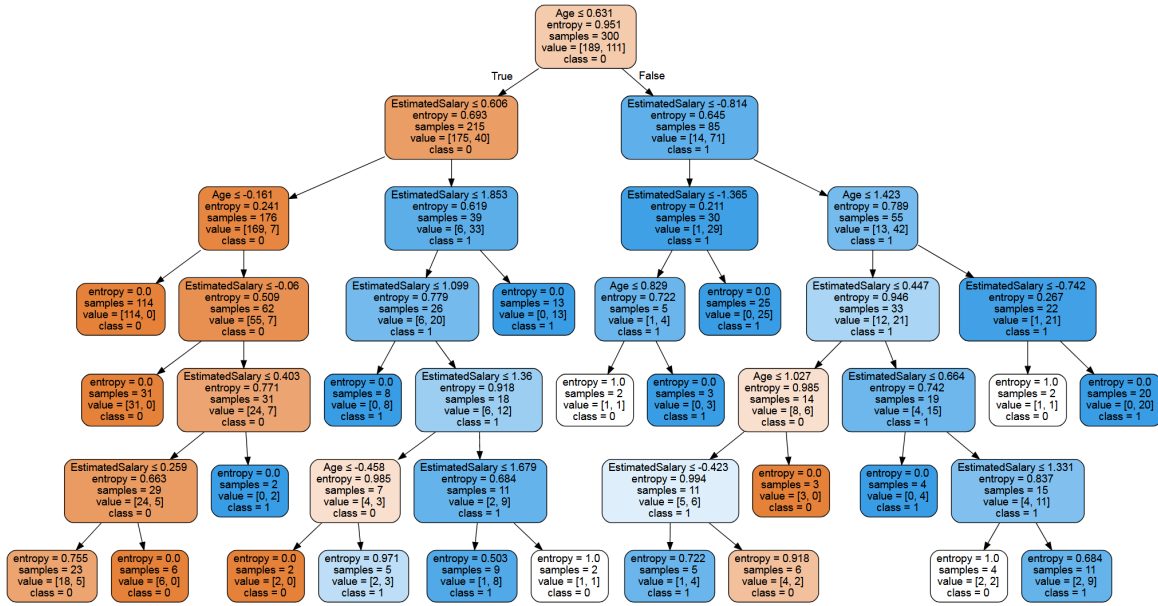


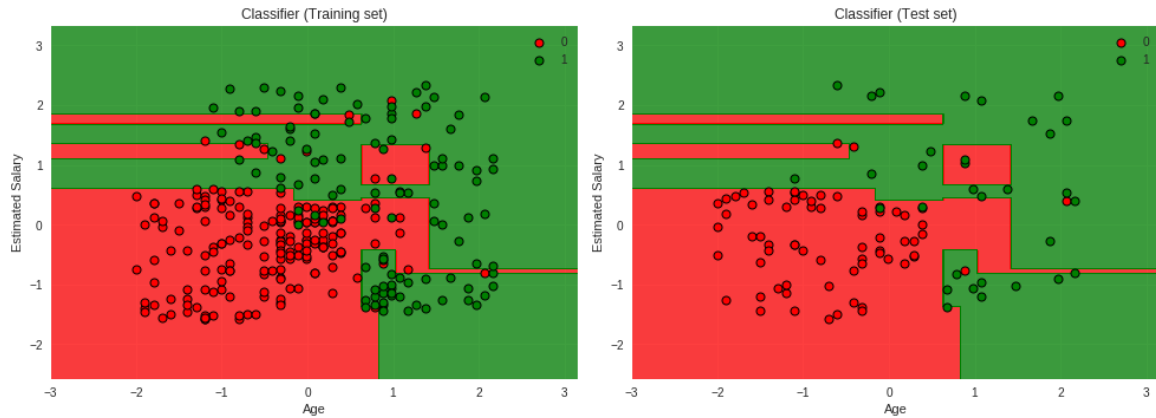Figure 5: Image of decision tree generated using Graphviz, max_depth=6



Figure 6: Training and validation visualization max_depth=6

## criterion = entropy

For this and the next section tree we have let the sklearn library function to pick the maximal depth of a tree and we wanted to check the influence of of criterion on the output. So the first one is 'entropy' and the tree can be seen on fig. 7. It is even more complicated than for previous trials. The accuracy was also equal 90% for this case and we see that actually the classifier tends to be overfitted – it is good for the training set, but have issues with test data.
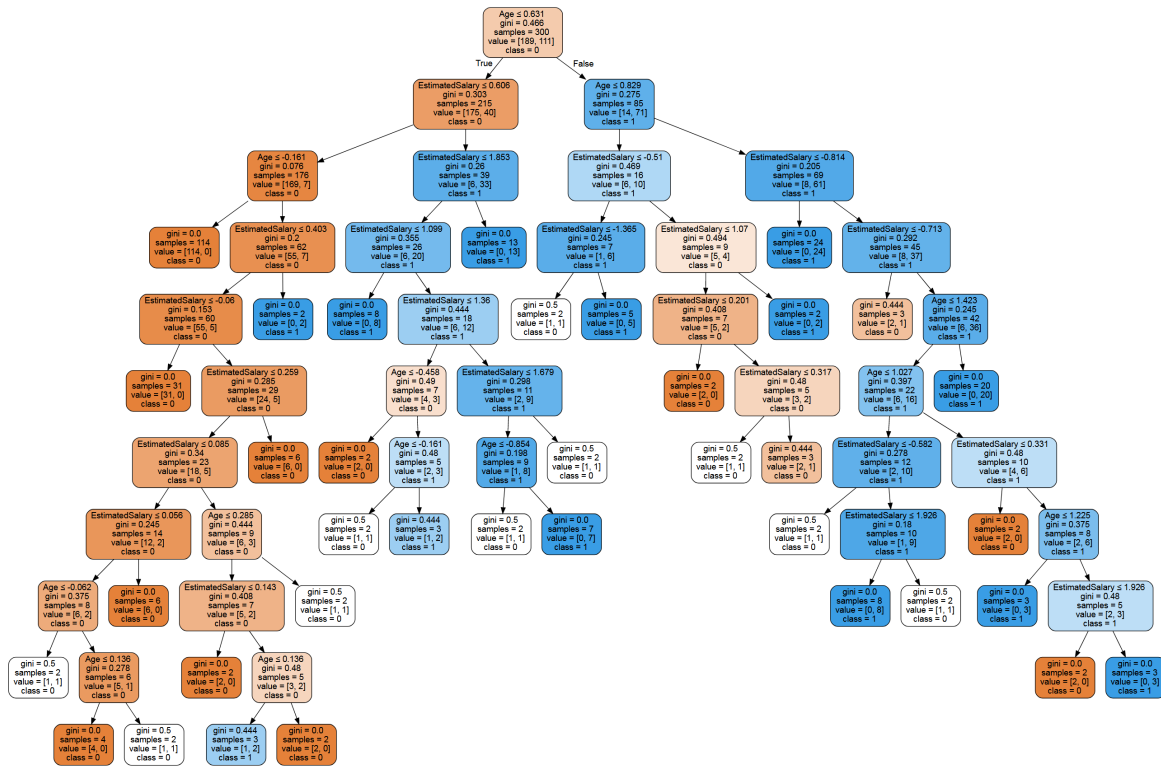


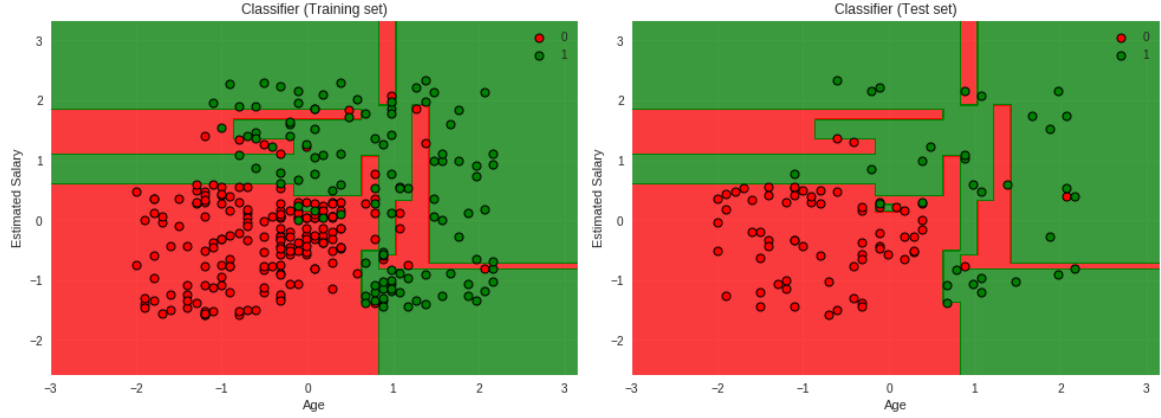Figure 7: Image of decision tree generated using Graphviz, criterion = entropy

Figure 8: Training and validation visualization, criterion = entropy

## criterion = gini

For this case we also let the library to choose depth of the tree with is also more complicated like we can see no fig. 9. Comparing 11 and plot 6 we can state that differences are actually not visible. The accuracy was, the same like in previous cases, 90%.
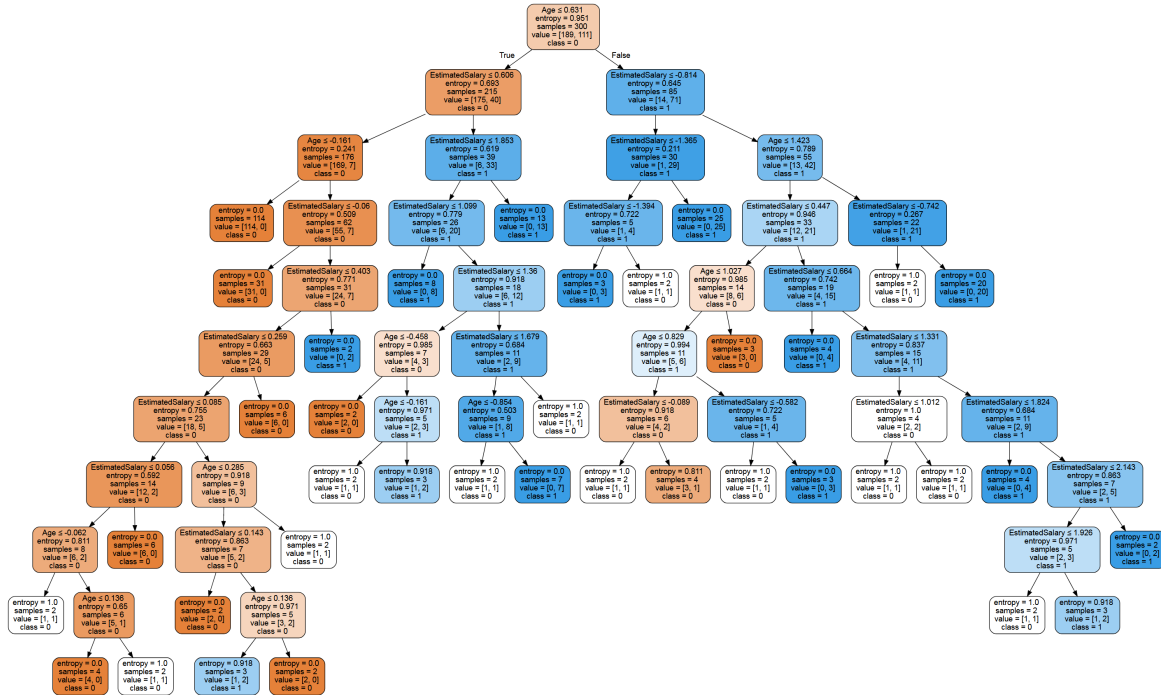


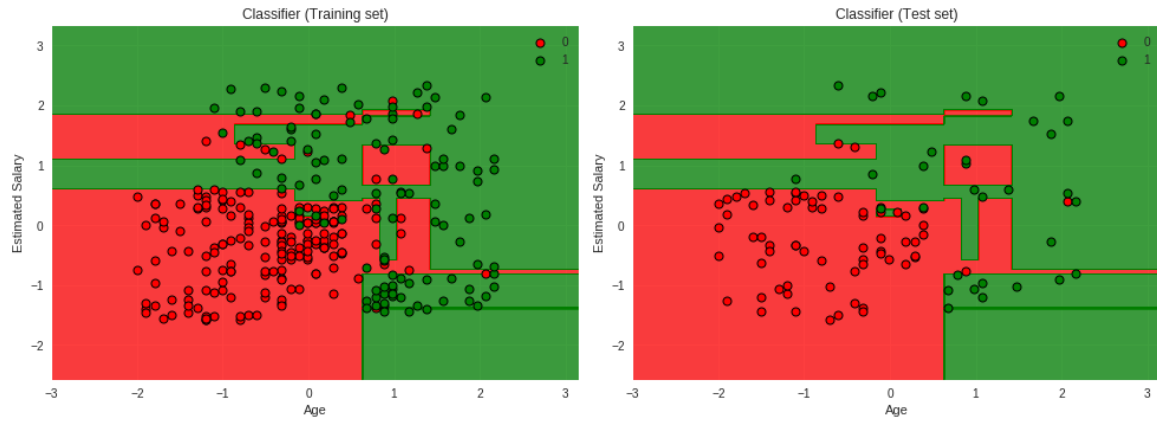Figure 9: Image of decision tree generated using Graphviz, criterion = gini

Figure 10: Training and validation visualization, criterion = gini

## splitter='random'

This time we wanted to check what is the influence of the splitter, because previously we always had splitter set to best. Now accuracy was accuracy = 87.0% for max_depth=4, comparing to 93% for the earlier tree for max_depth=4.
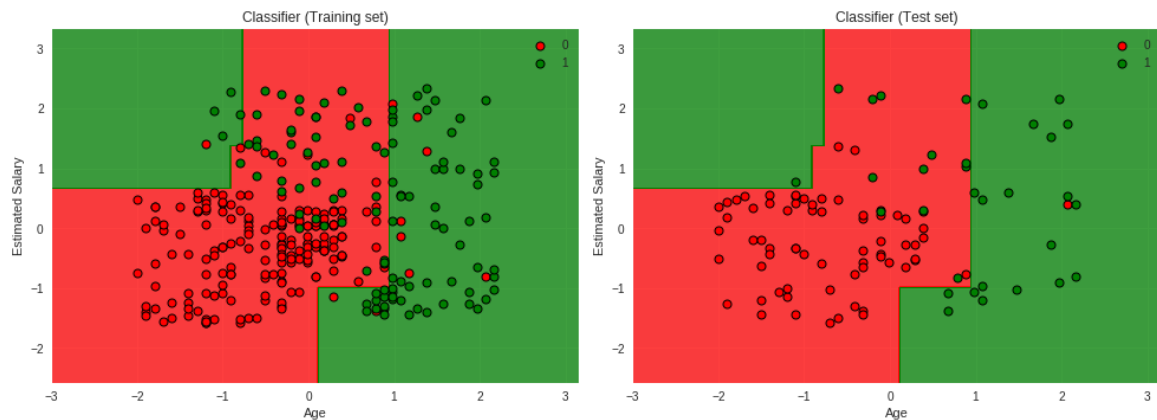


Figure 11: Training and validation visualization, splitter='random'

## Conclusions

As we could see on my examples not always bigger depth of a tree gives better results. The trees with too big depth can have better results on training sets, but they tend to be overfitted and gave poorer results on test set. It is also important to check available criterion functions which measure the quality of a split, because sometimes for the same depth, they can us give different results as everything depends on data that we have.