

Introduction

For this Home Assignment, I mostly stuck with the R functionalities discussed in the practice lessons, e.g. for generating the datasets or plotting. When I used functions from the practices, I marked it in this report in parentheses, e.g. saying "Practice1".

1 Distance Functions

This task was pretty straight forward, as discussed in the practice, all distance functions are combined in one function to avoid duplication and redundancy (Practice1). All distance functions compute the distance between two points with any number of dimensions, which is determined in the function.

2 Entropy and Fisher Score

The Entropy function accepts the label vector of a dataset as the only input parameter, since only the labels are needed for the calculation of the entropy. For the Fisher Score, the formula discussed in class was used, therefore the function accepts two parameters, one being the dataset along with the labels and the second being the feature for which the fisher score should be calculated for.

To test both functions, a 2D dataset with two classes was created (Practice2) which are separable by one feature and non-separable by the other feature. The Fig. 1 shows the used

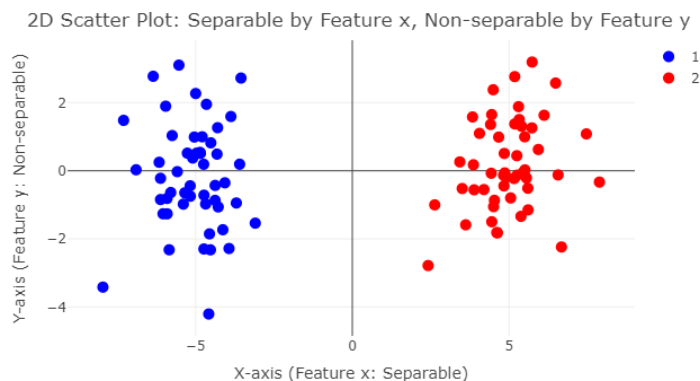


Figure 1: Plot of the Dataset used for testing

dataset. It's clearly visible, that classes 1 and 2 are separable by feature x and non-separable by feature y. Results of the test show, that the entropy of the entire dataset is clearly 1, since the dataset is split equally between class 1 and 2. For the Fisher Score, the results can be seen in this Table 1. As

Table 1: Fisher Score	
Feature	Fisher Score
Feature x	21.47333
Feature y	0.003366095

anticipated, the Fisher Score for feature "x" is significantly higher than the one for feature "y", indicating feature "x" should be used for classification.

3 Classification with Decision Tree

For the decision tree, essentially two functions are needed, one for building the tree and one for predicting the class labels of given data. To build a tree, the "train_tree" function is called with the dataset and the separate label vector along with the current depth of the tree (needed for recursion) and the max depth (set by the user, default is 4). The function then determines the feature with the highest information gain, which is calculated in the function "information_gain". This function uses the before implemented entropy function to determine both the feature and the best split for each feature, and returns the best feature along with the best possible split (threshold). It's inspired by the information gain function from Practice 3, yet different in the way that it performs the split based on the indices of the shuffled dataset. A possible improvement for this function could be an implementation of an early stopping functionality when determining the best split. Once the best split is determined, the "train_tree" function performs the split and calls the "train_tree" function recursively for the left branch with the left subset and the right branch with the right subset. All the datapoints in the right subset are greater and in the left subset less or equal to the determined split (threshold). When the "train_tree" function is called it first checks the stopping criteria, being the reach of max depth, the case of a Subdataset with the length of one (meaning it can't be split anymore) or the case that of the data points in the Subdataset are of the same class (there is no need for further splitting at this branch). When one of these cases occur, the function returns the label of the data points at this terminal node using a majority vote. This way the tree gets build in the depth first approach and the result is the linked list (using R "list" object) where every node consists of the used feature for splitting, the threshold for the split and a "link" to the left and right branch of the node. In case of a terminal node, all this data is missing and just a label is accessible in the list object. For first testing the algorithm, a 2D dataset was created and a tree of max depth equal to 4 was build, it's shown in Fig. 2. We haven't implemented goodness metrics yet, therefore no train-test-split was performed, so this is the next step to determine whether this tree is already overfitting on the training data. As seen in the Fig. 2 above, the tree performs multiple splits for the same feature when it results in the best gain.

3.1 Goodness Metrics

For the Goodness Metrics, the calculations for the True Positives, True Negatives, etc. are performed in a separate function, again to avoid redundancy. All Goodness Metrics are calculated in one function and can then be accessed separately.

3.2 Cross-validation

The different Goodness Metrics can now be used to evaluate the model with the further help of cross-validation. In the following, a cross-validation was performed on an example 2d Dataset for simplification unfortunately something didn't go quite as planned, and I couldn't investigate the reasons further due to time limitations. Nevertheless, the outcome can be seen in Fig. 3, it shows the accuracy of the testing data. Because the outcome doesn't seem very reliable, we can't make an assumption about the tree overfitting.

Decision Tree Visualization

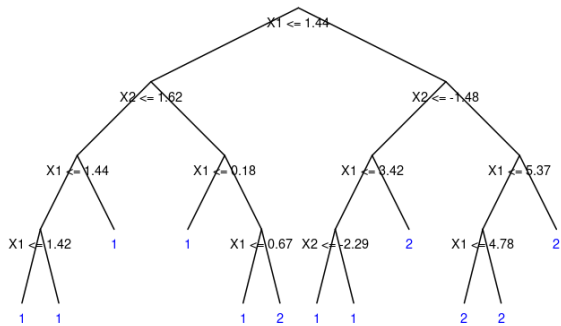


Figure 2: Example Decision Tree with max Depth of 4

4 Linear Regression

For visualizing the Regression, I've chosen the techniques we discussed in the labs. In Fig. 4 you can see the result for the linear Regression of a simple one feature dataset, but the concept stays the same for multiple features. For the sake of simplicity and because I'm forced to limit my time spend on this home assignment, the dataset has just one feature.

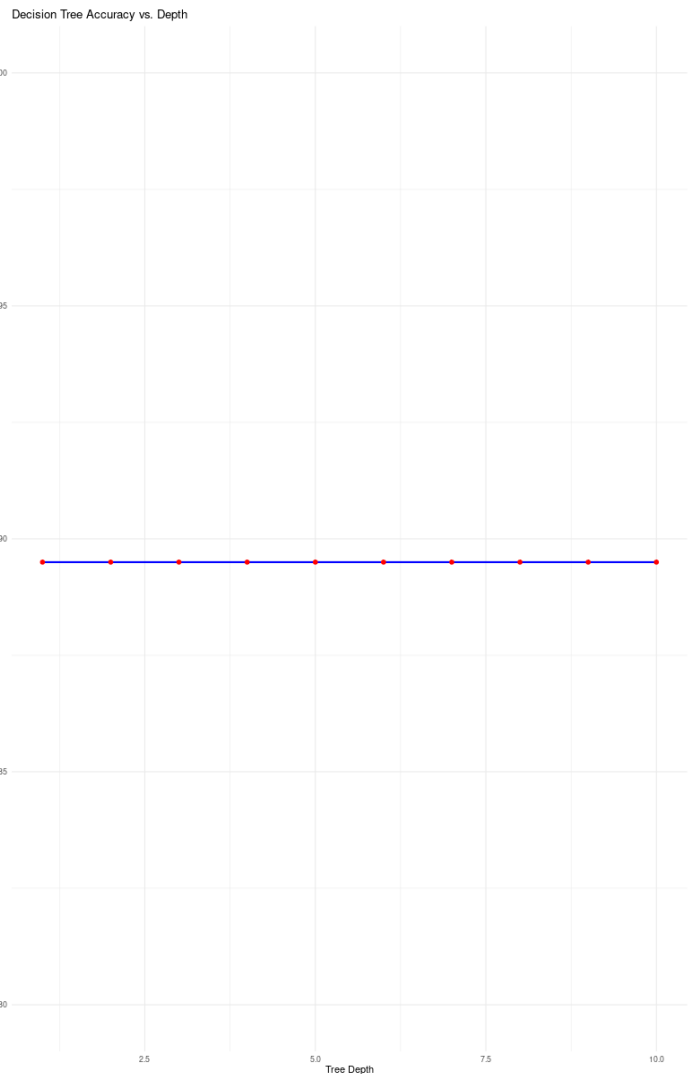


Figure 3: Accuracy of trees with different depths using cross validation

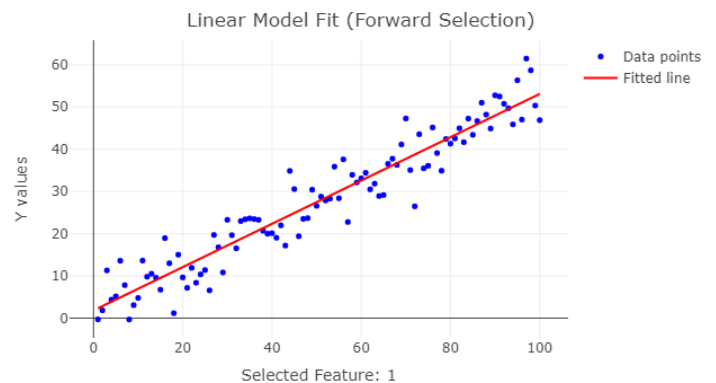


Figure 4: Linear Regression with forward selection