# Homework 02

Gabriele Lorenzo
s314913

December 2022

## Answer #1:

Using the whole dataset to train the Decision Tree, setting the default configuration and setting the minimal gain to 0.01, we can analyze that:

### (a)

As we can see in Figure 1:, the most discriminative attribute for class prediction is the *Node-caps* attribute. Said attribute has been divided in its three values: **"yes"**, **"no"** and **"?"**.

### (b)

The height of the generated Decision Tree is 6. We recall that: the height of a tree is defined as the maximum depth of any leaf node from the root node. That is, the length of the longest path (links traversed) from the root node *(Node-Caps)* to any leaf node.

### (c)

There are multiple examples of partial pure partitions, as we can see in Figure 2. In general there is a pure partition when all the data in the partition is correctly classified. In other words: if all the elements are accurately divided in different classes, that's a pure split.

# Answer #2:

In the default configuration, the Decision Tree uses the *gain ratio splitting criterion*. Using this criterion, the algorithm computes the gain value before each partition. If the calculated gain value is higher than the *minimal gain* value (configurated in the parameters of the Decision Tree), the node is partitioned.

Another important value is the *maximal depth* value. This value represents the maximal depth that the tree can grow to. In pratical terms, it represents the maximum number of "vertical" nodes of the tree.

Some examples of different parameters configurations: Figure 3, Figure 4, Figure 5, Figure 6, Figure 7.

# Answer #3:

Here are some examples of the confusion matrices with different configurations of the Decision Tree model: Figure 8, Figure 9, Figure 10, Figure 11, Figure 12, Figure 13.

As we can see, the model with the highest accuracy (Figure 13) is the one with balanced values of *minimal gain* and *maximal depth.*

# Answer #4:

Here are some examples of the confusion matrices with different configurations of the K-NN model: Figure 14, Figure 15, Figure 16, Figure 17, Figure 18, Figure 19.

As we can see, the model has an average accuracy above 70%, with the highest value for k=9.

Here is the confusion matrix for the Naive Bayes model: Figure 20.

If we compute the average accuracy value for the K-NN model with the different configurations, we get an average accuracy value of: 72.32, which is sligtly lower than the accuracy value of the Naive Bayes model (72.45). In conclusion, the Naive Bayes model perform better on average with respect to the K-NN model.

# Answer #5:

Here is the correlation matrix of the attributes: Figure 21.

## (a)

The Naïve independence assumption does not hold for the Breast dataset. In fact we can see that the highest absolute correlation value is 0.465, between the *inv-nodes* attribute and the *node-caps* attribute (quite high compared to the maximum theoretical value of 1).

## (b)

The pair of most correlated attributes are *inv-nodes* and *node-caps* with a correlation value of $-0.465$.

# Appendix:



Figure 1: Decision Tree



Figure 2: Example of left pure partition



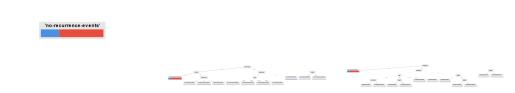Figure 3: Decision Tree with minimal gain=0.01 and maximal depth=1



Figure 4: Decision Tree with minimal gain=0.1 and maximal depth=10



Figure 5: Decision Tree with minimal gain=0.03 and maximal depth=4



Figure 6: Decision Tree with minimal gain=0.05 and maximal depth=10



Figure 7: Decision Tree with minimal gain=0.06 and maximal depth=3



Figure 8: Confusion Matrix with minimal gain=0.01 and maximal depth=1



Figure 9: Confusion Matrix with minimal gain=0.1 and maximal depth=10

Figure 10: Confusion Matrix with minimal gain=0.03 and maximal depth=4



Figure 11: Confusion Matrix with minimal gain=0.05 and maximal depth=10



Figure 12: Confusion Matrix with minimal gain=0.06 and maximal depth=3



Figure 13: Confusion Matrix with minimal gain=0.04 and maximal depth=3



Figure 14: Confusion Matrix with k=1



Figure 15: Confusion Matrix with k=3



Figure 16: Confusion Matrix with k=5



Figure 17: Confusion Matrix with k=7



Figure 18: Confusion Matrix with k=9



Figure 19: Confusion Matrix with k=11



Figure 20: Confusion Matrix with Naive Bayes model



Figure 21: Correlation Matrix