UNIVERSITY OF BERGEN

INTRODUCTION TO MACHINE LEARNING
INF264

# Project 1: Implementing decision trees

**Student :**
Arthur GAUTIER

September 18, 2023

# Contents

# 1   Implement a decision tree learning algorithm, from scratch

## 1.1   Node and Tree class

I chose to use object-oriented programming to implement the decision tree. In this implementation the tree is composed of two classes *TreeNode* and *DecisionTree*. The first class contained several attribute detailed below :

**TreeNode :**

- *feature* : to select a feature for the split.

- *threshold* : threshold value to split the feature.

- *left* : left subtree (lower or equal values).

- *right* right subtree (upper or equal values).

- *info_gain* : information gain of the split.

- *value* : value predicted by the leaf (for leaf nodes).

**DecisionTree :**

- *root* : related to the root node of the tree

- *max_depth* : recursion stopping condition

- *min_samples_split* : recursion stopping condition

## 1.2   Entropy and information gain

We know that

$$H(X) = -\sum_{x \in \mathcal{X}} P(x) \log_2 P(x)$$

Based on this formula and using pandas array properties, we get this entropy method.

```
def entropy(self, y):
    """Calculate the entropy of a variable Y"""
    p = y.value_counts() / y.shape[0]
    return -np.sum(p*np.log2(p))
```

Same thing for the information gain,

$$IG = H(Y) - H(Y|X)$$

```python
1   def information_gain(self, y, mask):
2       """Calculate the information gain of a variable Y using a mask"""
3       cardinal = mask.shape[0]
4       m_true = sum(mask)
5       m_false = cardinal - m_true
6
7       if m_true == 0 or m_false == 0:
8           ig = 0
9       else:
10          ig = self.entropy(y) - m_true/cardinal * self.entropy(y[mask])
11              - m_false/cardinal* self.entropy(y[-mask])
12          # y[mask] and y[-mask] are respectively two vectors containing
13          the values of y where the mask
14          is true and false.
15      return ig
```

## 1.3   Building the tree

We can now calculate the Information Gain for all variables. Next, we choose the split that generates the highest Information Gain as our splitting criterion. To do that, we can search the split that maximize the information gain using the method *max_information_gain*. We can apply this function to every feature of the dataset and find where is the best split using *get_best_split*. This means we select the variable that provides the most valuable information for making decisions at each node of the tree. Lastly, we repeat this process recursively until at least one of the conditions set by hyperparameters of the algorithm is not fulfilled. All the tree is store in the root attribute of the decision tree class.

```python
1   def build_tree(self, data, current_depth=0, depth=[]):
2       """recursive function to build the tree"""
3       # Print the current_depth
4       if current_depth not in depth:
5           depth.append(current_depth)
6           print(current_depth)
7       #  If all data points have the same label :
8       y = data.iloc[:, -1]
9       if y.nunique() == 1:
10          leaf_value = self.calculate_leaf_value(y)
11          #  return a leaf with that label
12          return TreeNode(value=leaf_value)
13
14      # Else, if all data points have identical feature values:
15      all_rows_same = data.duplicated().sum() == 0
16      if all_rows_same is True:
17          leaf_value = self.calculate_leaf_value(y)
18          return TreeNode(value=leaf_value)
```

```
19
20       # Else, choose a feature that maximizes the information gain :
21       num_samples, num_features = data.shape
22       if num_samples >= self.min_samples_split and current_depth <= self.max_depth:
23           best_split_label, best_split_value,
24                                     best_split_ig = self.get_best_split(y, data)
25           # We check if the split is pure
26           if best_split_ig is not None and best_split_ig > 0:
27               data_left, data_right = self.make_split(best_split_label,
28                                                 best_split_value, data)
29               left_subtree = self.build_tree(data_left, current_depth+1, depth)
30               right_subtree = self.build_tree(data_right, current_depth + 1, depth)
31               return TreeNode(best_split_label, best_split_value, left_subtree,
32                             right_subtree, best_split_ig)
33           else:
34               leaf_value = self.calculate_leaf_value(y)
35               return TreeNode(value=leaf_value)
36       # compute leaf nodel
37       leaf_value = self.calculate_leaf_value(y)
38       # return leaf node
39       return TreeNode(value=leaf_value)
```

## 2   Add Gini index

$$G(X) = \sum_{x \in \mathcal{X}} P(x)(1 - P(x))$$

```
1   def gini(self, y):
2       """Gini index"""
3       p = y.value_counts()/y.shape[0]
4       return np.sum(p*(1-p))
```

We can modify a bit the *information_ gain* method by replacing the *entropy* function by a variable *function* in order to choose which cost function to use in the decision tree.

## 3   Add reduced-error pruning

We recursively traverse the nodes of the tree, for each node we replace the node with a leaf. We calculate the precision on the plume dataset, if the precision is better than the accuracy compute on the training dataset we keep the leaf, otherwise we cancel the modifications while keeping the original node.

```
1   def prune(self, node, data, y, data_prune, y_prune, accuracy):
2       """Prune function to cut the useless branches of the tree"""
3       # If it's not a leaf :
```

```
4          if node.value is None:
5              # We save the two subtrees
6              node_save_left = node.left
7              node_save_right = node.right
8              # If the subtree are actually two leafs with the same value :
9              if node_save_left.value is not None and node_save_right.value is not None:
10                 if node_save_left.value == node_save_right.value:
11                     # We replace the node with a leaf
12                     return TreeNode(value=node.left.value)
13             data1 = data[data[node.feature] < node.threshold]
14             data2 = data[data[node.feature] >= node.threshold]
15
16             y1 = data1.iloc[:, -1]
17             y2 = data2.iloc[:, -1]
18
19             value1 = y1.mode().iloc[0]
20             value2 = y2.mode().iloc[0]
21             # We build to new leaf for the actual node :
22             node.left = TreeNode(value=value1)
23             node.right = TreeNode(value=value2)
24
25             # Compute the accuracy with this configuration
26             result = self.prediction(data_prune)
27
28             # if the new accuracy is not improved :
29             if accuracy_score(y_prune, result) < accuracy:
30                 # We continue to explore the tree with the original node
31                 node.left = self.prune(node_save_left, data1, y1, data_prune,
32                                        y_prune, accuracy)
33                 node.right = self.prune(node_save_right, data2, y2, data_prune,
34                                         y_prune, accuracy)
35             # otherwise we return the new node
36             return node
37         # if it's a leaf :
38         return node
```

# 4    Evaluate your algorithm

We start by importing the data from the file *wine_dataset.csv*. Then, we need to split the dataset in a training and test set. We proceed as below :

```
1   def load_data():
2       # Load the dataset :
3       data_types = {
4           "citric acid": float,
5           "residual sugar": float,
```

```
 6              "sugar": float,
 7              "sulphates": float,
 8              "alcohol": float,
 9              "type": int
10          }
11          columns = ["citric acid", "residual sugar", "pH", "sulphates", "alcohol", "type"
12          data = pd.read_csv("wine_dataset.csv", header=0, dtype=data_types, names=columns
13          return data
14
15      def trainTest(data):
16          """Devide the data in train and test datasets"""
17          y = data.iloc[:, -1]
18          data_train, data_test, y_train, y_test = model_selection.train_test_split(
19              data, y, test_size=0.30, random_state=42)
20          return data_train, data_test, y_train, y_test
```

Then we can train our tree using the previous explained functions and print the results.
Results with the decision tree implemented from scratch :



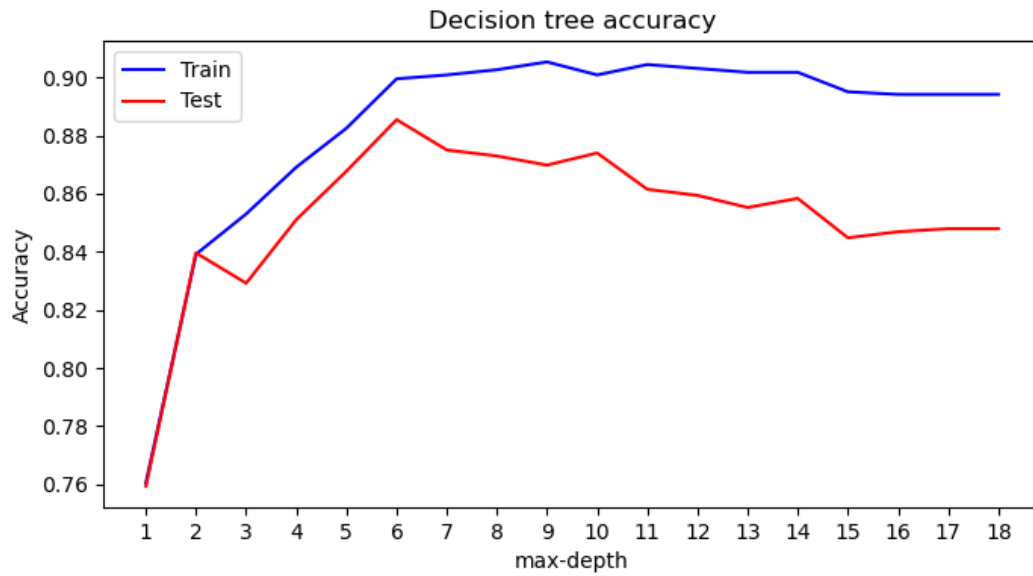Figure 1: Implemented decision tree accuracy, impurity_measure='entropy'

Figure 2: Implemented decision tree accuracy, impurity_measure='gini'
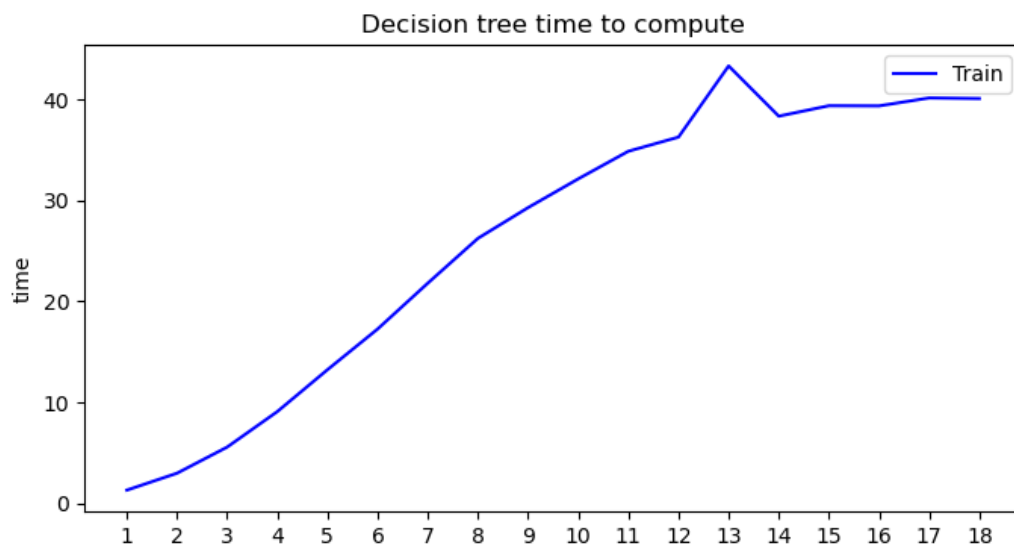


Figure 3: Time to compute in second for the implemented decision tree, with entropy (
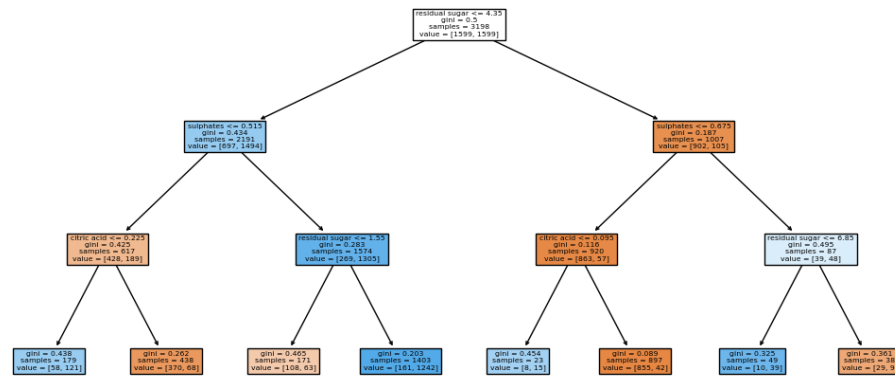
# 5   Compare to an existing implementation



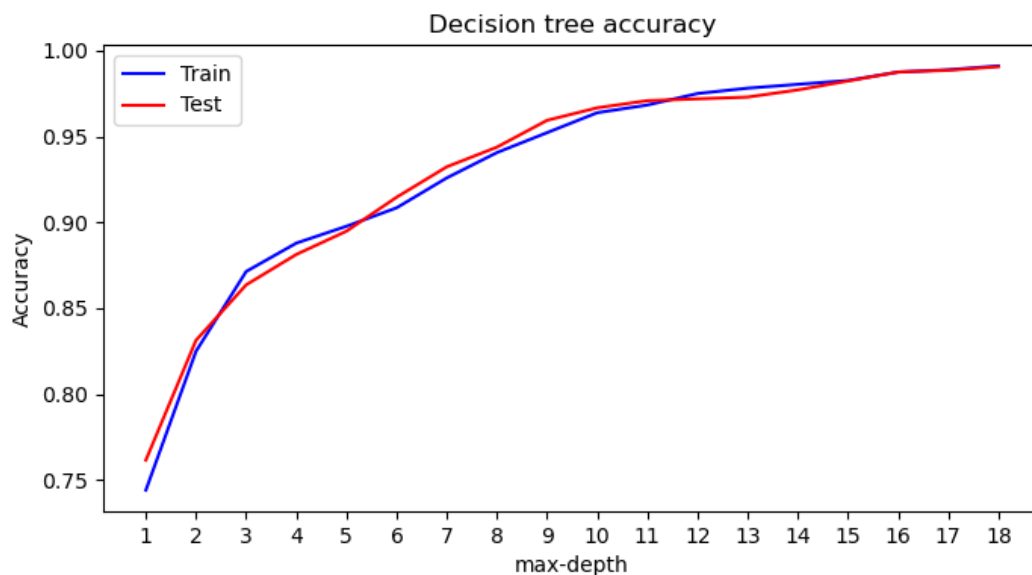Figure 4: Decision tree from *sklearn* with maxdepth=3



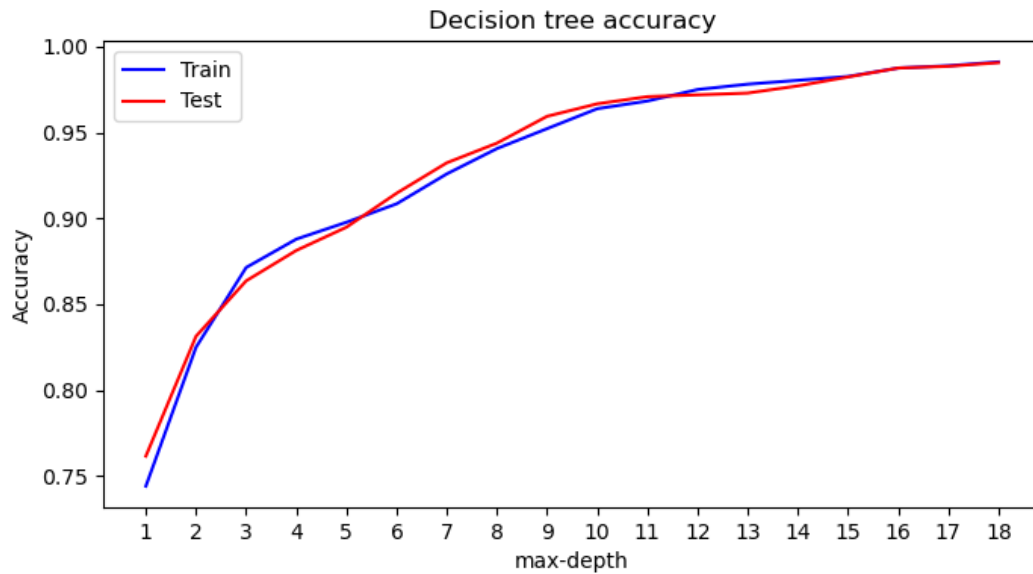Figure 5: Sklearn decision tree accuracy, impurity_measure='entropy'

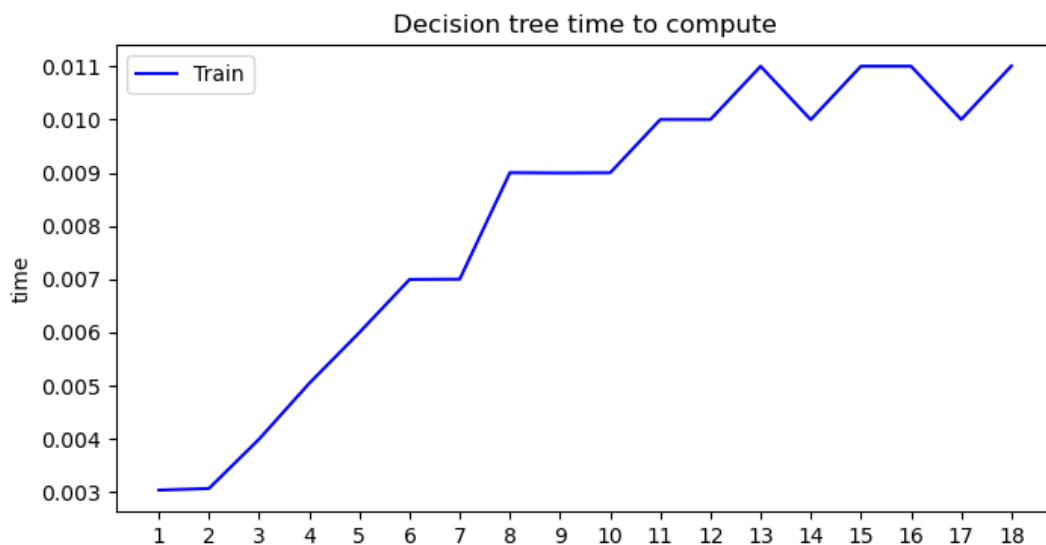Figure 6: Sklearn decision tree accuracy, impurity_measure='gini'



Figure 7: Time to compute in second for the sklearn decision tree, with entropy

As we can see my program is not efficient in comparison to the sklearn one. This difference can be attributed to scikit-learn's extensive codebase, which has undergone years of refinement and optimization by a community of developers.