

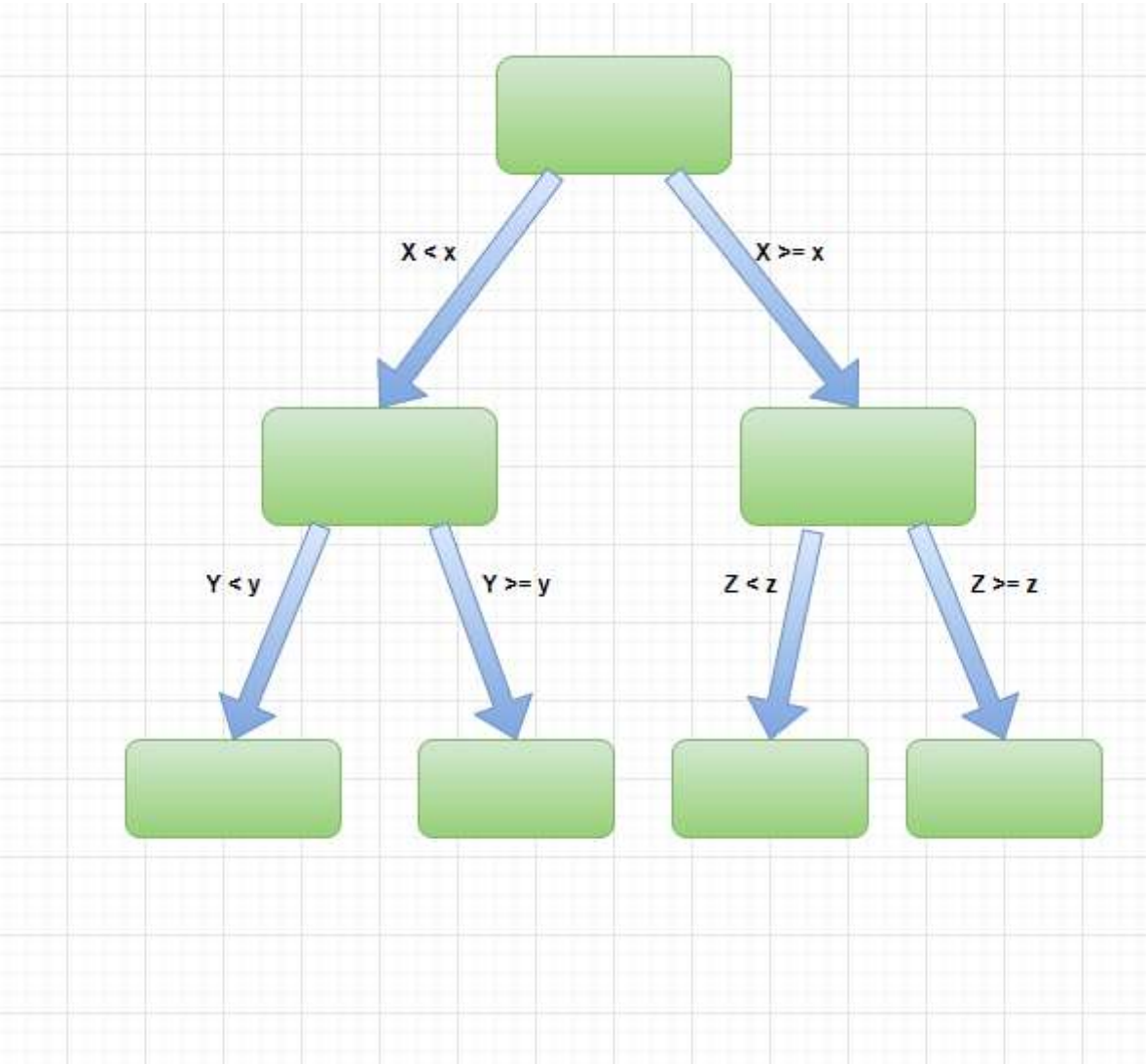


Decision Tree Algorithm in Python From Scratch

Coding the popular algorithm using just NumPy and Pandas in Python and explaining what's under the hood



Eligijus Bujokas Apr 14 · 8 min read



Decision tree schema; graph by author

. . .

The aim of this article is to make all the parts of a decision tree classifier clear by walking through the code that implements the algorithm. The code uses only NumPy, Pandas and the standard python libraries.

The full code can be accessed via <https://github.com/Eligijus112/decision-tree-python>

As of now, the code creates a decision tree when the target variable is binary and the features are numeric. This is completely sufficient to

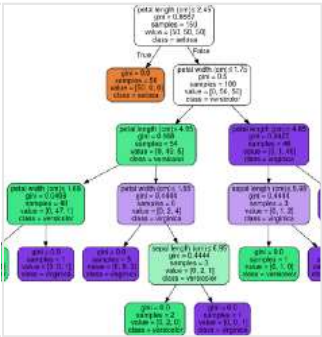
understand the algorithm.

The golden standard of building decision trees in python is the scikit-learn implementation:

1.10. Decision Trees - scikit-learn 0.24.1 documentation

Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression . The goal...

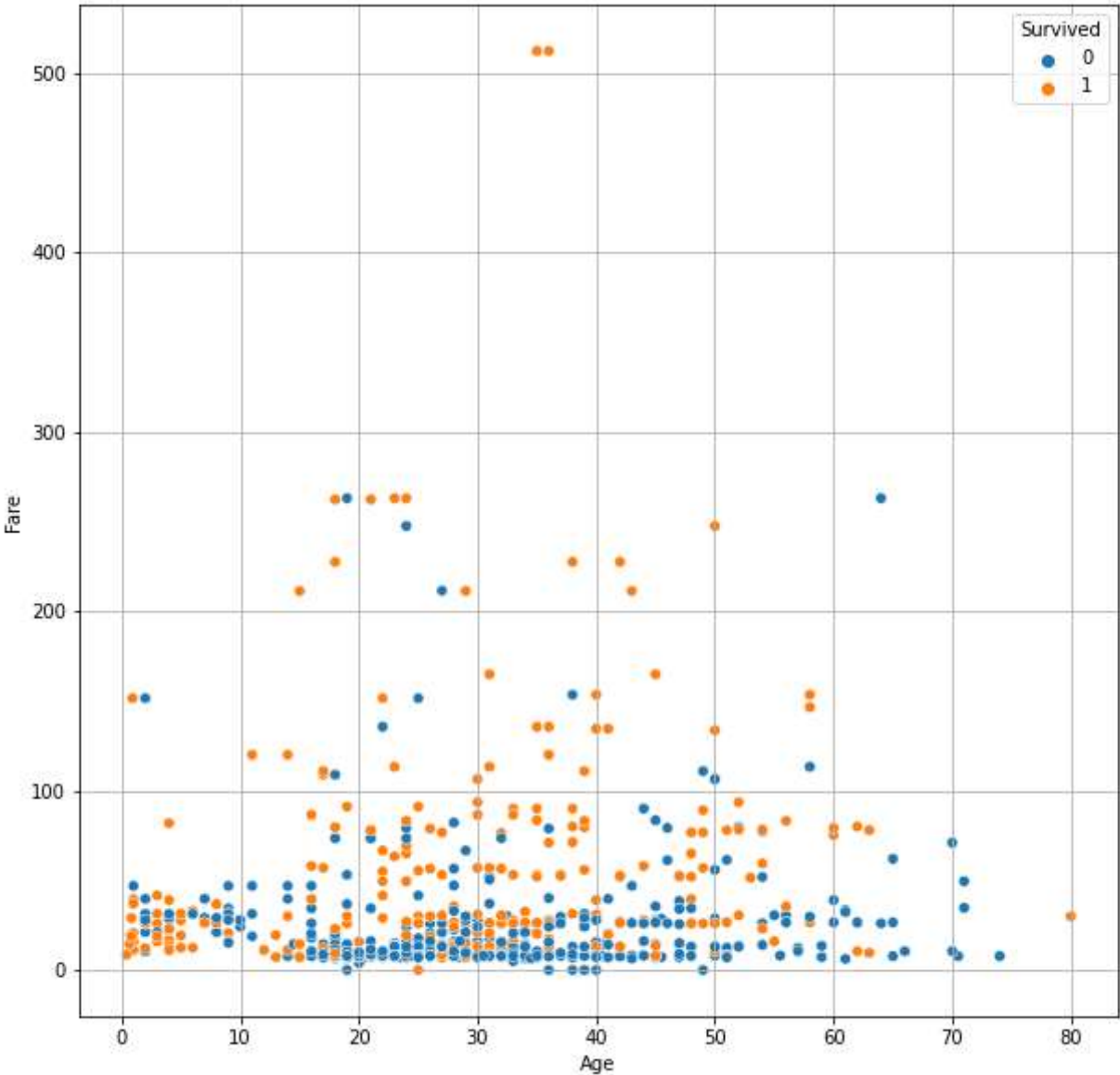
scikit-learn.org



When I tested out my code I wanted to make sure that the results are identical to the scikit-learn implementation.

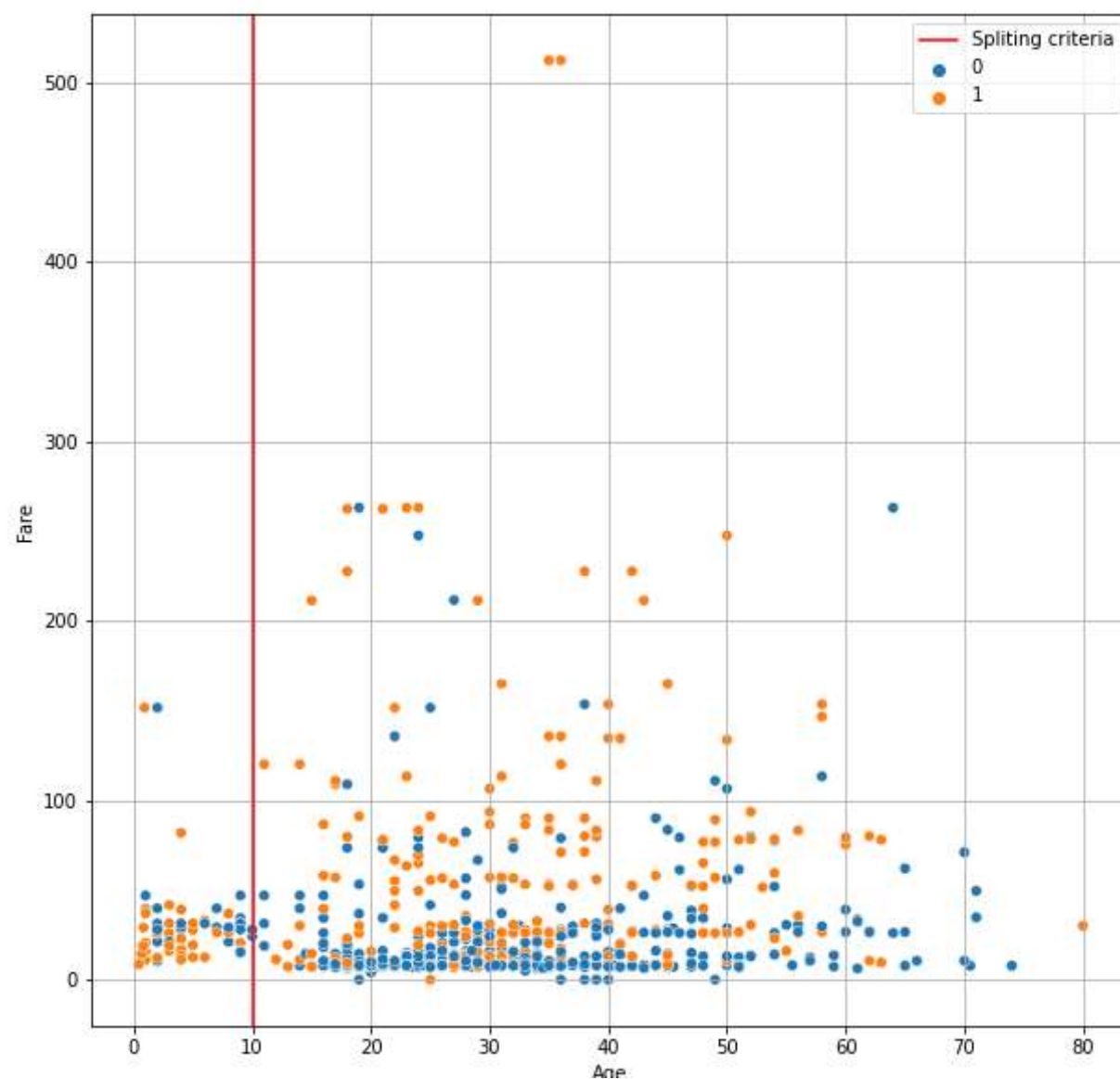
. . .

The data used in this article is the famous Titanic survivor dataset. We will use two numeric variables — Age of the passenger and the Fare of the ticket — to predicting whether a passenger survived or not.



Age + Fare ~ Survival; Graph by author

The goal is to create the “**best**” splits of the numeric variables. Just eyeballing the data, we could guess that one good split is to split the data into two parts: observations that have Age < 10 and observations that have Age ≥ 10:



Splitting the dataset; Graph by author

Now some immediate questions may rise:

Is this a good split?

Maybe a split at the Fare = 200 is a better one?

How do we quantify the “goodness” of a split?

How does a computer search for the best split?

All of these questions will be answered by the end of this article.

. . .

A decision tree algorithm (DT for short) is a machine learning algorithm that is used in classifying an observation given a set of input features. The algorithm creates a set of rules at various decision levels such that a certain metric is optimized.

The target variable will be denoted as $Y = \{0, 1\}$ and the feature matrix will be denoted as X .

Keywords to expand on:

Node

Gini impurity (a metric which we are optimizing)

Level

Splitting

. . .

A **node** is the building block in the decision tree. When viewing a typical schema of a decision tree (like the one in the title picture) the nodes are the rectangles or bubbles that have a downward connection to other nodes.

Each node has the following main attributes:

Gini impurity score

Number of observations

The number of observations belonging to each of the binary target classes.

The feature matrix X representing the observations that fall into the node.

The custom **Node** class in python (written by me):

```
1  # Data wrangling
2  import pandas as pd
3
4  # Array math
5  import numpy as np
6
7  # Quick value count calculator
8  from collections import Counter
9
10
11  class Node:
12      """
13      Class for creating the nodes for a decision tree
14      """
15      def __init__(
16          self,
17          Y: list,
18          X: pd.DataFrame,
19          min_samples_split=None,
20          max_depth=None,
21          depth=None,
22          node_type=None,
23          rule=None
24      ):
25          # Saving the data to the node
26          self.Y = Y
27          self.X = X
28
29          # Saving the hyper parameters
```

```

30     self.min_samples_split = min_samples_split if min_samples_split else 20
31     self.max_depth = max_depth if max_depth else 5
32
33     # Default current depth of node
34     self.depth = depth if depth else 0
35
36     # Extracting all the features
37     self.features = list(self.X.columns)
38
39     # Type of node
40     self.node_type = node_type if node_type else 'root'
41
42     # Rule for splitting
43     self.rule = rule if rule else ""
44
45     # Calculating the counts of Y in the node
46     self.counts = Counter(Y)
47
48     # Getting the GINI impurity based on the Y distribution
49     self.gini_impurity = self.get_GINI()
50
51     # Sorting the counts and saving the final prediction of the node
52     counts_sorted = list(sorted(self.counts.items(), key=lambda item: item[1]))
53
54     # Getting the last item
55     yhat = None
56     if len(counts_sorted) > 0:
57         yhat = counts_sorted[-1][0]
58
59     # Saving to object attribute. This node will predict the class with the most frequent c
60     self.yhat = yhat
61
62     # Saving the number of observations in the node
63     self.n = len(Y)
64
65     # Initiating the left and right nodes as empty nodes
66     self.left = None
67     self.right = None
68
69     # Default values for splits
70     self.best_feature = None
71     self.best_value = None
72
73     @staticmethod
74     def GINI_impurity(y1_count: int, y2_count: int) -> float:
75         """
76         Given the observations of a binary class calculate the GINI impurity
77         """
78         # Ensuring the correct types
79         if y1_count is None:
80             y1_count = 0
81
82         if y2_count is None:
83             y2_count = 0
84
85         # Getting the total observations
86         n = y1_count + y2_count
87
88         # If n is 0 then we return the lowest possible gini impurity
89         if n == 0:
90             return 0.0
91
92         # Getting the probability to see each of the classes
93         p1 = y1_count / n
94         p2 = y2_count / n
95
96         # Calculating GINI
97         gini = 1 - (p1 ** 2 + p2 ** 2)
98

```

```

99         # Returning the gini impurity
100         return gini
101
102     @staticmethod
103     def ma(x: np.array, window: int) -> np.array:
104         """
105         Calculates the moving average of the given list.
106         """
107         return np.convolve(x, np.ones(window), 'valid') / window
108
109     def get_GINI(self):
110         """
111         Function to calculate the GINI impurity of a node
112         """
113         # Getting the 0 and 1 counts
114         y1_count, y2_count = self.counts.get(0, 0), self.counts.get(1, 0)
115
116         # Getting the GINI impurity
117         return self.GINI_impurity(y1_count, y2_count)
118
119     def best_split(self) -> tuple:
120         """
121         Given the X features and Y targets calculates the best split
122         for a decision tree
123         """
124         # Creating a dataset for splitting
125         df = self.X.copy()
126         df['Y'] = self.Y
127
128         # Getting the GINI impurity for the base input
129         GINI_base = self.get_GINI()
130
131         # Finding which split yields the best GINI gain
132         max_gain = 0
133
134         # Default best feature and split
135         best_feature = None
136         best_value = None
137
138         for feature in self.features:
139             # Dropping missing values
140             Xdf = df.dropna().sort_values(feature)
141
142             # Sorting the values and getting the rolling average
143             xmeans = self.ma(Xdf[feature].unique(), 2)
144
145             for value in xmeans:
146                 # Splitting the dataset
147                 left_counts = Counter(Xdf[Xdf[feature]<value]['Y'])
148                 right_counts = Counter(Xdf[Xdf[feature]>=value]['Y'])
149
150                 # Getting the Y distribution from the dicts
151                 y0_left, y1_left, y0_right, y1_right = left_counts.get(0, 0), left_counts.get(1
152
153                 # Getting the left and right gini impurities
154                 gini_left = self.GINI_impurity(y0_left, y1_left)
155                 gini_right = self.GINI_impurity(y0_right, y1_right)
156
157                 # Getting the obs count from the left and the right data splits
158                 n_left = y0_left + y1_left
159                 n_right = y0_right + y1_right
160
161                 # Calculating the weights for each of the nodes
162                 w_left = n_left / (n_left + n_right)
163                 w_right = n_right / (n_left + n_right)
164
165                 # Calculating the weighted GINI impurity
166                 wGINI = w_left * gini_left + w_right * gini_right
167

```



```

167
168         # Calculating the GINI gain
169         GINIgain = GINI_base - wGINI
170
171         # Checking if this is the best split so far
172         if GINIgain > max_gain:
173             best_feature = feature
174             best_value = value
175
176         # Setting the best gain to the current one
177         max_gain = GINIgain
178
179     return (best_feature, best_value)
180
181 def grow_tree(self):
182     """
183     Recursive method to create the decision tree
184     """
185     # Making a df from the data
186     df = self.X.copy()
187     df['Y'] = self.Y
188
189     # If there is GINI to be gained, we split further
190     if (self.depth < self.max_depth) and (self.n >= self.min_samples_split):
191
192         # Getting the best split
193         best_feature, best_value = self.best_split()
194
195         if best_feature is not None:
196             # Saving the best split to the current node
197             self.best_feature = best_feature
198             self.best_value = best_value
199
200             # Getting the left and right nodes
201             left_df, right_df = df[df[best_feature]<=best_value].copy(), df[df[best_feature
202
203             # Creating the left and right nodes
204             left = Node(
205                 left_df['Y'].values.tolist(),
206                 left_df[self.features],
207                 depth=self.depth + 1,
208                 max_depth=self.max_depth,
209                 min_samples_split=self.min_samples_split,
210                 node_type='left_node',
211                 rule=f"{best_feature} <= {round(best_value, 3)}"
212             )
213
214             self.left = left
215             self.left.grow_tree()
216
217             right = Node(
218                 right_df['Y'].values.tolist(),
219                 right_df[self.features],
220                 depth=self.depth + 1,
221                 max_depth=self.max_depth,
222                 min_samples_split=self.min_samples_split,
223                 node_type='right_node',
224                 rule=f"{best_feature} > {round(best_value, 3)}"
225             )
226
227             self.right = right
228             self.right.grow_tree()
229
230 def print_info(self, width=4):
231     """
232     Method to print the information about the tree
233     """
234     # Defining the number of spaces
235     const = int(self.depth * width ** 1.5)

```

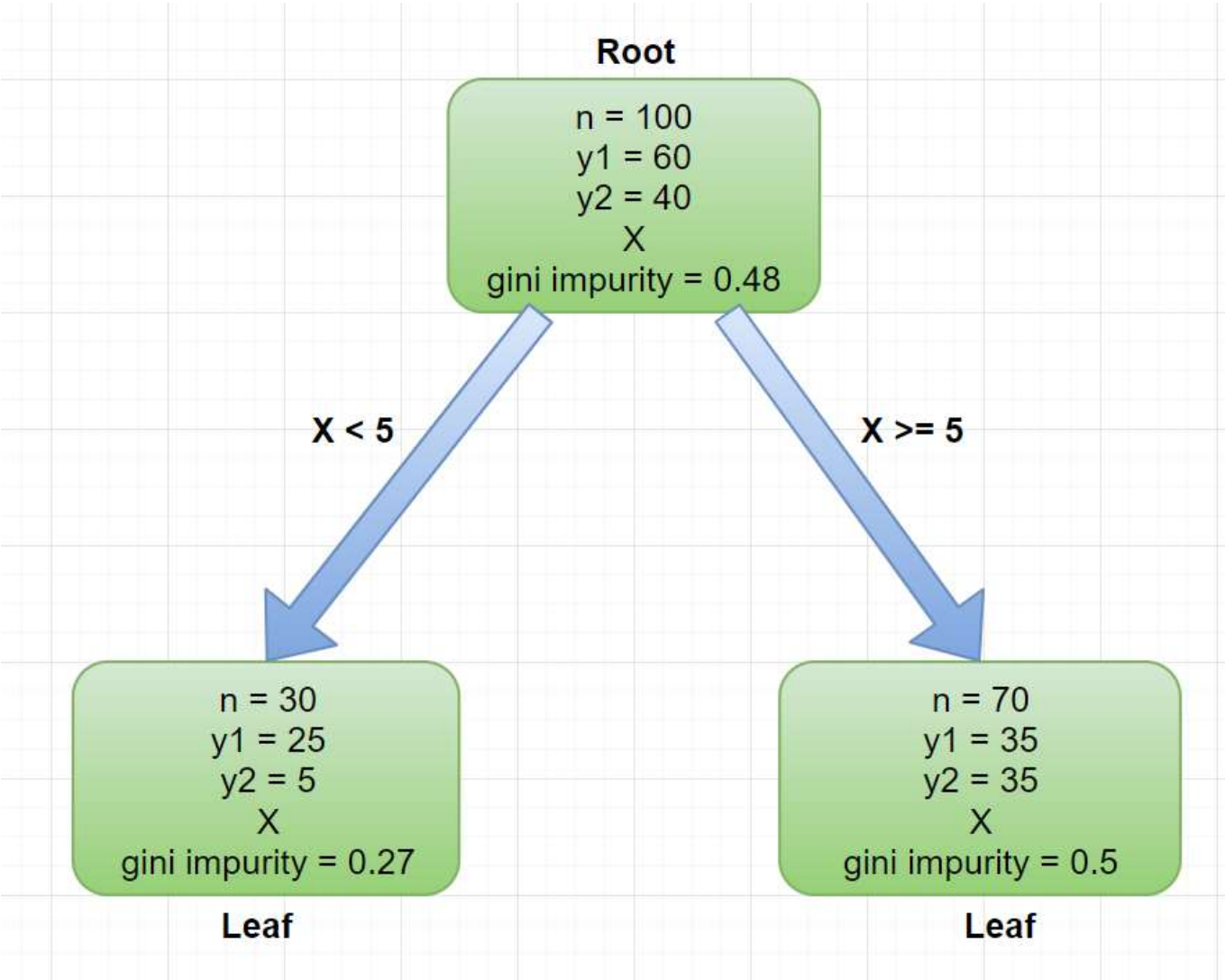
```

236     spaces = "-" * const
237
238     if self.node_type == 'root':
239         print("Root")
240     else:
241         print(f"|{spaces} Split rule: {self.rule}")
242         print(f"{' ' * const} | GINI impurity of the node: {round(self.gini_impurity, 2)}")
243         print(f"{' ' * const} | Class distribution in the node: {dict(self.counts)}")
244         print(f"{' ' * const} | Predicted class: {self.yhat}")
245
246     def print_tree(self):
247         """
248         Prints the whole tree from the current node to the bottom
249         """
250         self.print_info()
251
252         if self.left is not None:
253             self.left.print_tree()
254
255         if self.right is not None:
256             self.right.print_tree()
257
258     def predict(self, X:pd.DataFrame):
259         """
260         Batch prediction method
261         """
262         predictions = []
263
264         for _, x in X.iterrows():
265             values = {}
266             for feature in self.features:
267                 values.update({feature: x[feature]})
268
269             predictions.append(self.predict_obs(values))
270
271         return predictions
272
273     def predict_obs(self, values: dict) -> int:
274         """
275         Method to predict the class given a set of features
276         """
277         cur_node = self
278         while cur_node.depth < cur_node.max_depth:
279             # Traversing the nodes all the way to the bottom
280             best_feature = cur_node.best_feature
281             best_value = cur_node.best_value
282
283             if cur_node.n < cur_node.min_samples_split:
284                 break
285
286             if (values.get(best_feature) < best_value):
287                 if self.left is not None:
288                     cur_node = cur_node.left
289             else:
290                 if self.right is not None:
291                     cur_node = cur_node.right
292
293         return cur_node.yhat

```

The first node in a decision tree is called ***the root***. The nodes at the bottom of the tree are called ***leaves***.

If splitting criteria are satisfied, then each node has two linked nodes to it: the left node and the right node. For example, a very simple decision tree with one root and two leaves may look like this:



Example decision tree; Graph by author

n – number of observations

y1 – number of first class elements

y2 – number of second class elements

x – numeric feature for the observations

An input if the feature X less than 5 would go to the left node. A feature value of more or equal to 5 would go to the right node.

. . .

As mentioned above, each node has a GINI impurity score. To calculate the GINI impurity, all that is needed is the distribution of the target variable in the node or simply, how many Y=1 and Y=0 observations there are in a node.

The formal definition of GINI impurity is as follows:

Gini impurity is a measure of how often a randomly chosen element from the set would be incorrectly labelled if it was randomly labelled according to the distribution of labels in the subset.

The algebraic definition:

Suppose we have two classes in the dataset:

$$k_1, k_2$$

Each of the classes have n_1 and n_2 observations.

The probability of observing something from one of the k classes is:

$$p(i) = P(x_i \in k_i) = \frac{n_i}{n_1 + n_2}, i \in \{1, 2\}$$

The GINI impurity of such a system is calculated with the following formula:

$$G = 1 - \sum_{i=1}^2 p(i)^2$$

GINI impurity formula; Formulas are taken from the author's notebook

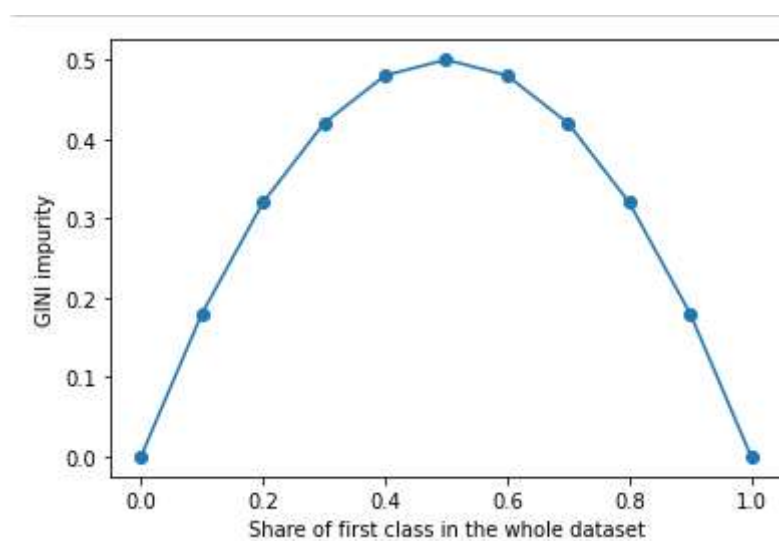
For example, if we have in a node 10 observations that are survivors and 5 observations that are not survivors, then:

```
p(survivors) = 10 / (10 + 5) = 0.66..
p(non-survivors) = 5 / (10 + 5) = 0.33..
```

Thus the GINI impurity can be calculated by squaring the two numbers, adding them up and subtracting from one:

```
gini impurity = 1 - (0.66..^2 + 0.33..^2) = 0.44..
```

In a binary case, the maximum Gini impurity is equal to 0.5 and the lowest is 0. The lower the value, the more “pure” the node is. No matter how many observations we have, we can calculate the share of one of the classes in terms of the whole dataset and draw a relationship:



Gini impurity vs share of class in a dataset; Image by author

We always search for a split where the GINI impurity will be the lowest.

If nothing else, the basic intuition is that the more one class observations there is in a node, the lower its impurity.

. . .

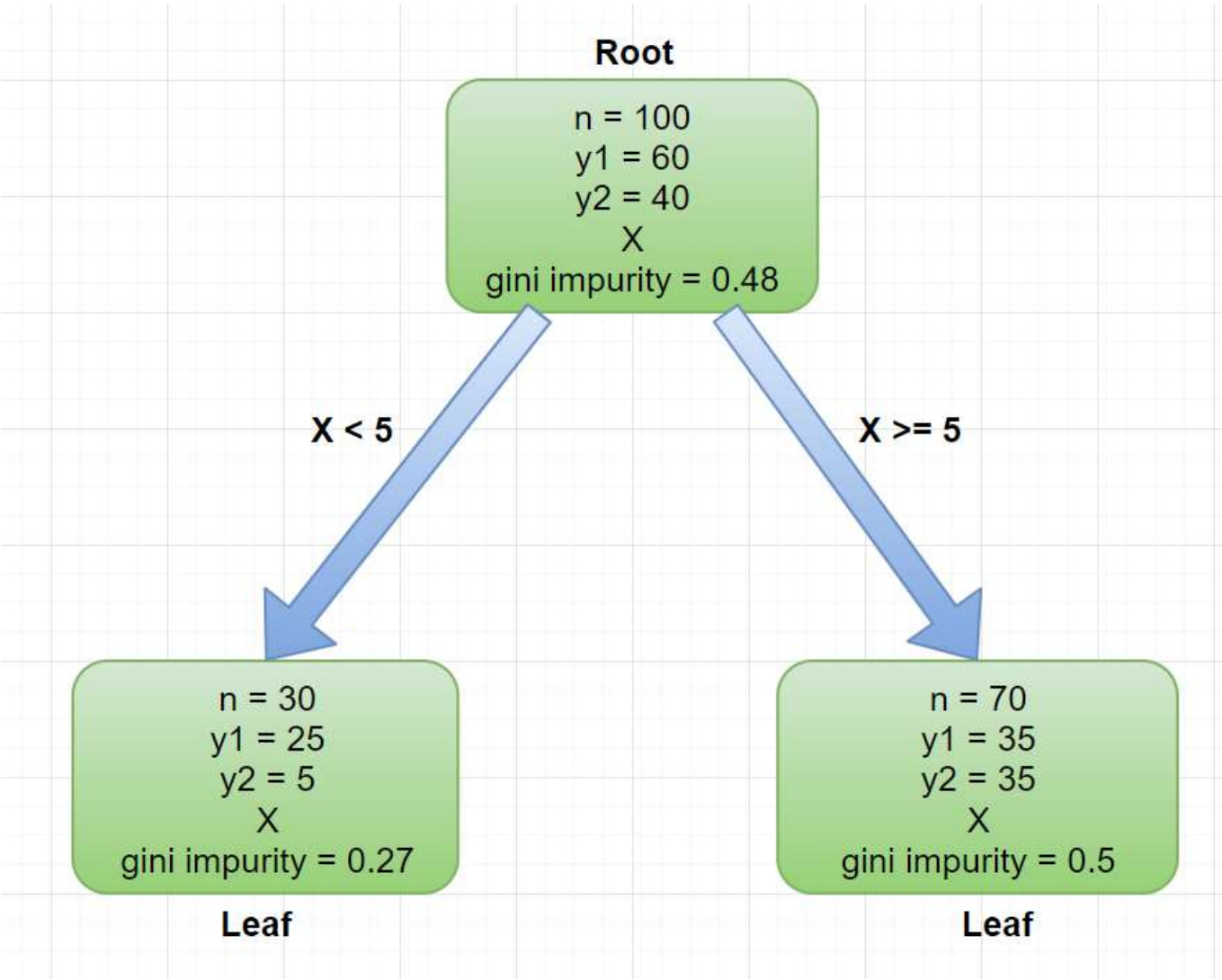
The level attribute defines how many splits were made before the node was created. So for example, the root node’s level is 0, then the left and the right nodes have a level of 1, and so on.

In the custom Node class, the maximum depth of the tree can be regulated with the hyperparameter **max_depth**.

. . .

The splitting procedure is a process where we search in each node which feature and which feature value is the best to divide the data into two smaller parts.

In the classification case, we want to maximize the **Gini gain**. Going back to the example before:



Example decision tree; Graph by author

The root node has a Gini score of 0.48. The left node has a score of 0.27 and the right node has 0.5. *How much Gini did we “gain”?*

The formula to calculate this is:

Gini gain = Parent node Gini impurity subtracted by the weighted average of the Gini impurities of the left and right nodes.

$$GINI_{gain} = \Delta Gini = Gini_{parent} - (Gini_{left} \frac{n_{left}}{n_{right} + n_{left}} + Gini_{right} \frac{n_{right}}{n_{right} + n_{left}})$$

Gini gain formula; From author’s notebook

In the example above, plugging in all the values we would get:

$$0.48 - (0.27 * 30 / 100 + 0.5 * 70 / 100) = \mathbf{0.049}$$

The GINI gain is equal to 0.049. **Any positive Gini gain is an improvement. This means that our decision would make the nodes more “pure”.**

How does the algorithm search for the best split among numeric columns?

For each of the feature, we sort the feature values and get the means of two neighbouring values.

For example, suppose our feature_1 is the following:

```
feature_1 = [7, 5, 9, 1, 2, 8]
```

Sorted:

```
feature_1 = [1, 2, 5, 7, 8, 9]
```

Means of neighbours:

```
feature_1 = [1.5, 3.5, 6.5, 7.5, 8.5]
```

We then check what is the GINI gain with each of the values from the above vector. And, we do this for all the features in our dataset.

The final split value and split feature is the one that has the highest GINI gain.

. . .

The hyperparameters of the custom Node object that grows the tree (more will be added in the future) are the **max_depth: int** and **min_samples_split: int** variables.

The `max_depth` integer defines how deep should the tree grow. At the depth of `max_depth`, the searching for the best split feature and split feature values stops.

The `min_samples_split` integer defines the minimal number of observations in a node for the best split search to start. So for example, if the node has 51 observations but the `min_samples_split = 55`, then the growth of the tree stops.

. . .

So, how does the code work?

First of all, read the data:

```
# Loading data
d = pd.read_csv('data/train.csv')

# Dropping missing values
dtree = d[['Survived', 'Age', 'Fare']].dropna().copy()

# Defining the X and Y matrices
Y = dtree['Survived'].values
X = dtree[['Age', 'Fare']]

# Saving the feature list
features = list(X.columns)
```

Then we define the dictionary of hyperparameters.

```
hp = {
    'max_depth': 3,
    'min_samples_split': 50
}
```

Then we initiate the root node:

```
root = Node(Y, X, **hp)
```

The main tree building function is the **grow_tree()** function.

```
root.grow_tree()
```

And that's it!

To view the results, we can invoke the `print_tree()` function.

```
root.print_tree()
```

The results:

```
Root
| GINI impurity of the node: 0.48
| Class distribution in the node: {0: 424, 1: 290}
| Predicted class: 0
|----- Split rule: Fare <= 52.277
|   | GINI impurity of the node: 0.44
|   | Class distribution in the node: {0: 389, 1: 195}
|   | Predicted class: 0
|----- Split rule: Fare <= 10.481
|   | GINI impurity of the node: 0.32
|   | Class distribution in the node: {0: 192, 1: 47}
|   | Predicted class: 0
|----- Split rule: Age <= 32.5
|   | GINI impurity of the node: 0.37
|   | Class distribution in the node: {0: 134, 1: 43}
|   | Predicted class: 0
|----- Split rule: Age > 32.5
|   | GINI impurity of the node: 0.12
|   | Class distribution in the node: {0: 58, 1: 4}
|   | Predicted class: 0
|----- Split rule: Fare > 10.481
|   | GINI impurity of the node: 0.49
|   | Class distribution in the node: {0: 197, 1: 148}
|   | Predicted class: 0
|----- Split rule: Age <= 6.5
|   | GINI impurity of the node: 0.41
|   | Class distribution in the node: {0: 12, 1: 30}
|   | Predicted class: 1
|----- Split rule: Age > 6.5
|   | GINI impurity of the node: 0.48
|   | Class distribution in the node: {0: 185, 1: 118}
|   | Predicted class: 0
|----- Split rule: Fare > 52.277
|   | GINI impurity of the node: 0.39
|   | Class distribution in the node: {1: 95, 0: 35}
|   | Predicted class: 1
|----- Split rule: Age <= 63.5
|   | GINI impurity of the node: 0.38
|   | Class distribution in the node: {1: 95, 0: 32}
|   | Predicted class: 1
|----- Split rule: Age <= 29.5
|   | GINI impurity of the node: 0.44
|   | Class distribution in the node: {0: 17, 1: 34}
|   | Predicted class: 1
|----- Split rule: Age > 29.5
|   | GINI impurity of the node: 0.32
|   | Class distribution in the node: {1: 61, 0: 15}
|   | Predicted class: 1
|----- Split rule: Age > 63.5
|   | GINI impurity of the node: 0.0
|   | Class distribution in the node: {0: 3}
|   | Predicted class: 0
```

Full decision tree; Snippet by author

The decision tree obtained from the scikit-learn implementation is identical:

```
--- Fare <= 52.28
|   --- Fare <= 10.48
|   |   --- Age <= 32.50
|   |   |   class: 0
|   |   |   Age > 32.50
```

• • •

As it turns out, the best first initial split is the Fare feature at a value of 52.28 and not at the proposed Age feature at value 10.

Feel free to create a pull request in this repo
<https://github.com/Eligijus112/decision-tree-python> if you see any bugs or
just want to add functionalities.

By Towards Data Science

 Get this newsletter

Gini Impurity