Eligijus Bujokas   ( Follow )

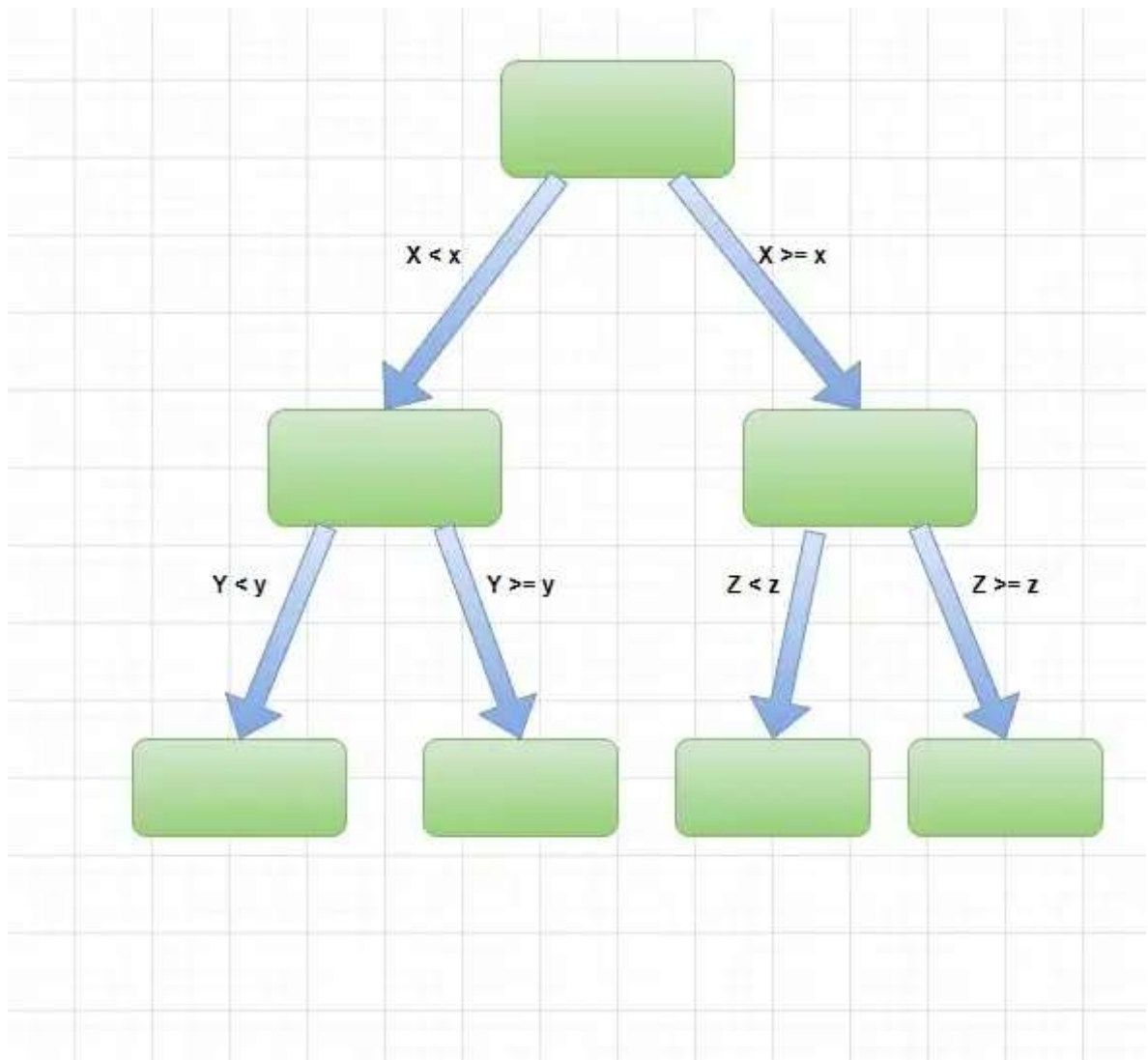Apr 14, 2021 · 8 min read · ✦ · ▶ Listen

⌐ Save    🐦    f    in    🔗

# 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

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.

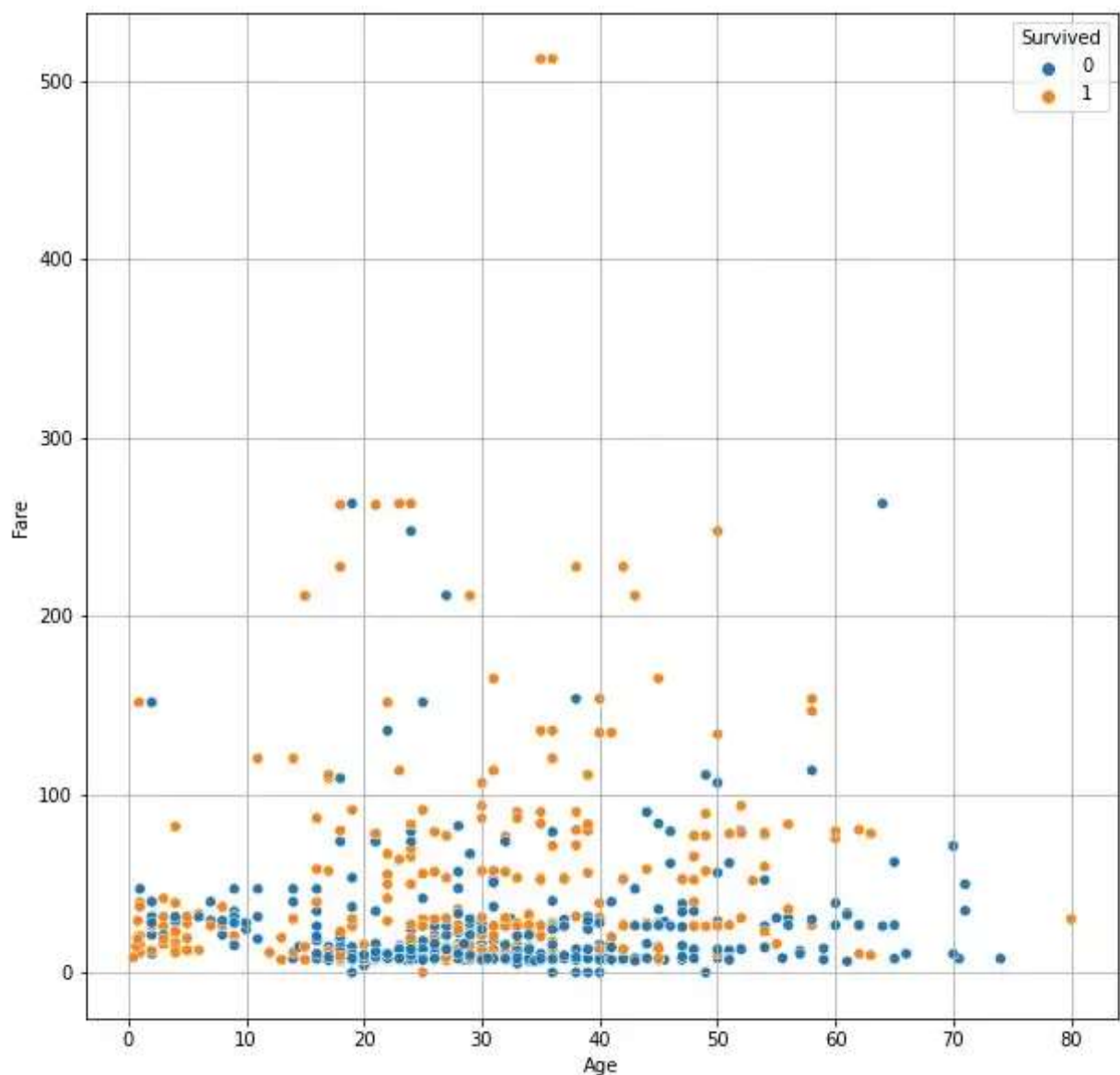The golden standard of building decision trees in python is the scikit-learn implementation:
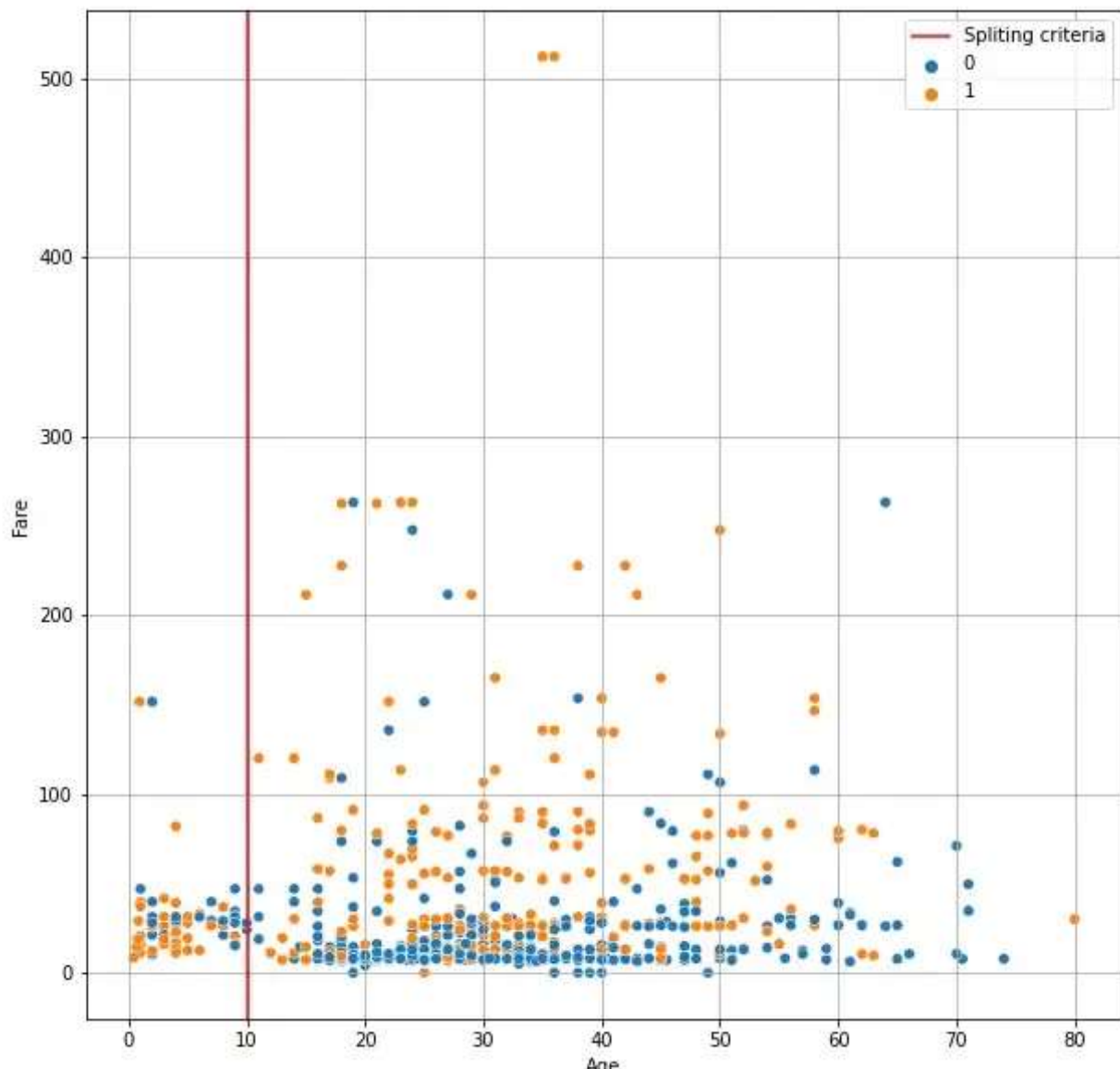
When I tested out my code I wanted to make sure that the results are identical to the scikit-learn implementation.

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.

*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.*

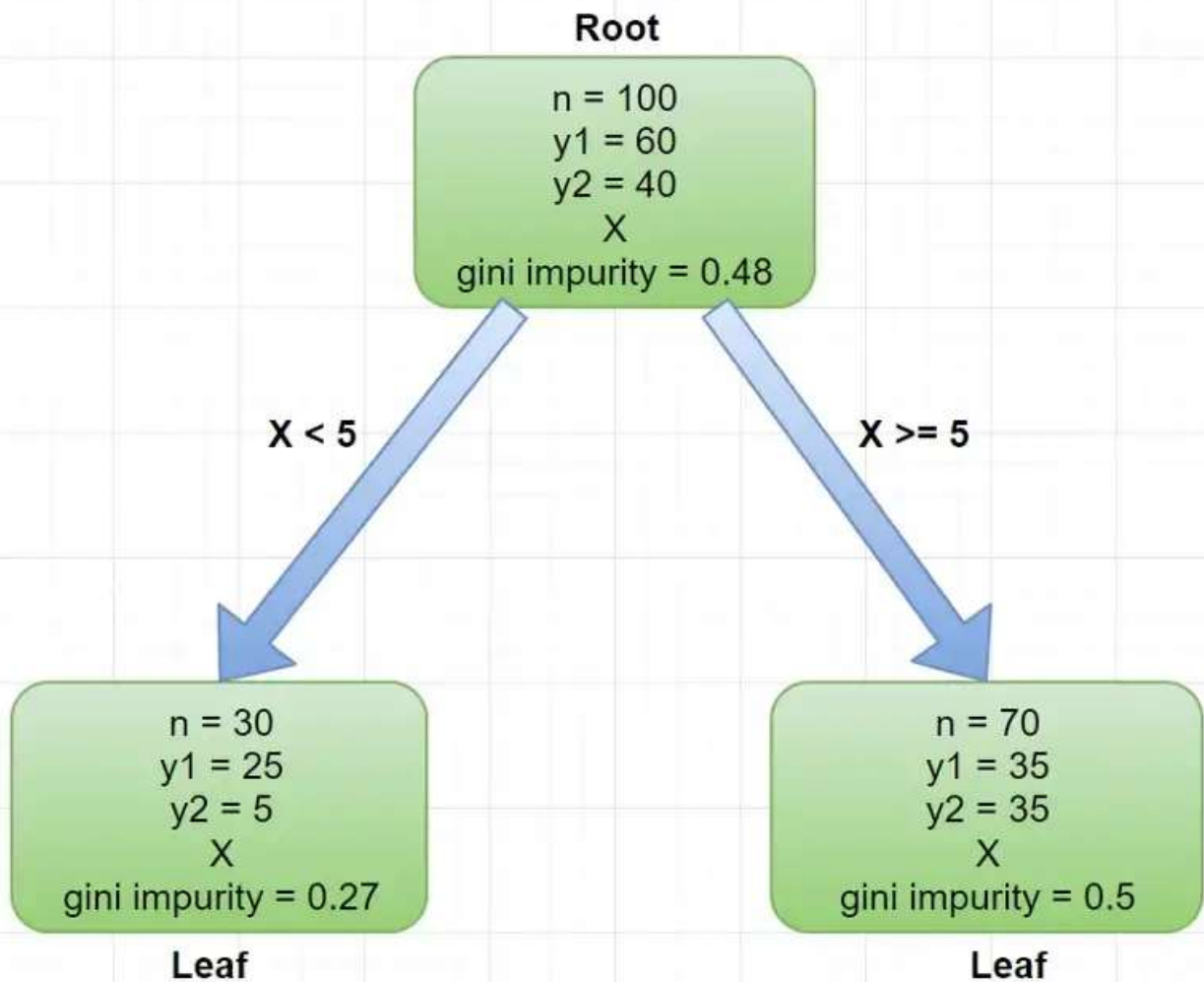The custom **Node** class in python (written by me):

```python
# Data wrangling
import pandas as pd

# Array math
import numpy as np

# Quick value count calculator
from collections import Counter


class Node:
    """
    Class for creating the nodes for a decision tree
    """
    def __init__(
        self,
        Y: list,
        X: pd.DataFrame,
        min_samples_split=None,
        max_depth=None,
        depth=None,
        node_type=None,
        rule=None
    ):
        # Saving the data to the node
        self.Y = Y
        self.X = X

        # Saving the hyper parame
        self.min_samples_split = min_samples_split if min_samples_split else 20
        self.max_depth = max_depth if max_depth else 5

        # Default current depth of node
        self.depth = depth if depth else 0

        # Extracting all the features
        self.features = list(self.X.columns)

        # Type of node
        self.node_type = node_type if node_type else 'root'

        # Rule for spliting
        self.rule = rule if rule else ""

        # Calculating the counts of Y in the node
        self.counts = Counter(Y)

        # Getting the GINI impurity based on the Y distribution
```

```
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:
```



```
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         """
```

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:

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

```
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             # Droping 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)
```
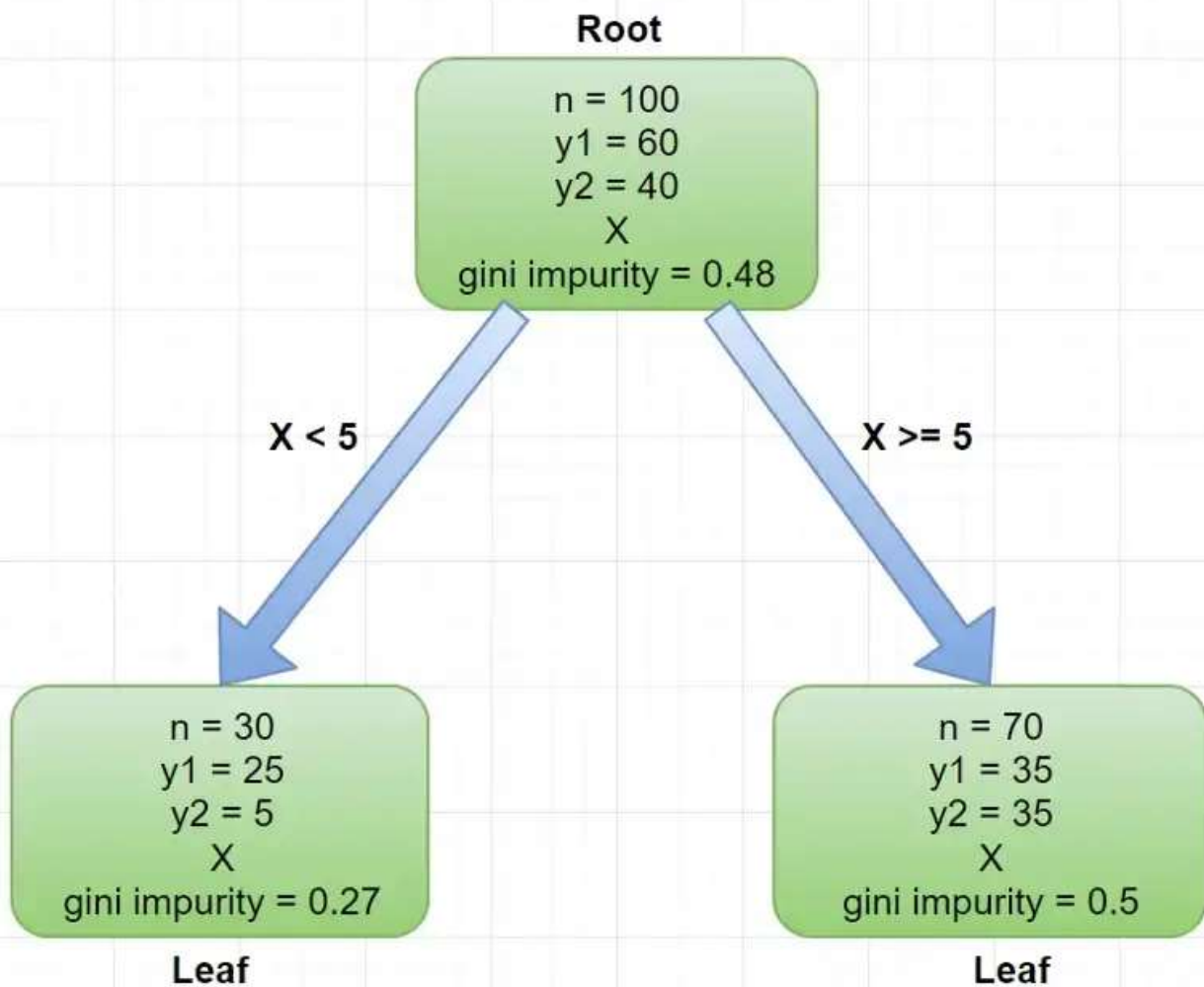
```python
            xmeans = self.ma(Xdf[feature].unique(), 2)

            for value in xmeans:
                # Spliting the dataset
                left_counts = Counter(Xdf[Xdf[feature]<value]['Y'])
                right_counts = Counter(Xdf[Xdf[feature]>=value]['Y'])

                # Getting the Y distribution from the dicts
                y0_left, y1_left, y0_right, y1_right = left_counts.get(0, 0), left_counts.get(

                # Getting the left and right gini impurities
                gini_left = self.GINI_impurity(y0_left, y1_left)
                gini_right = self.GINI_impurity(y0_right, y1_right)

                # Getting the obs count from the left and the right data splits
                n_left = y0_left + y1_left
                n_right = y0_right + y1_right

                # Calculating the weights for each of the nodes
                w_left = n_left / (n_left + n_right)
                w_right = n_right / (n_left + n_right)

                # Calculating the weighted GINI impurity
                wGINI = w_left * gini_left + w_right * gini_right

                # Calculating the GINI gain
                GINIgain = GINI_base - wGINI

                # Checking if this is the best split so far
                if GINIgain > max_gain:
                    best_feature = feature
                    best_value = value

                    # Setting the best gain to the current one
                    max_gain = GINIgain

        return (best_feature, best_value)

    def grow_tree(self):
        """
        Recursive method to create the decision tree
        """
        # Making a df from the data
        df = self.X.copy()
        df['Y'] = self.Y

        # If there is GINI to be gained, we split further
        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
```



```
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            """
```

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

```
234                # Defining the number of spaces
235                const = int(self.depth * width ** 1.5)
236                spaces = "-" * const
237
238                if self.node_type == 'root':
```

```python
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}")

246     def print_tree(self):
247         """
248         Prints the whole tree from the current node to the bottom
249         """
250         self.print_info()

252         if self.left is not None:
253             self.left.print_tree()

255         if self.right is not None:
256             self.right.print_tree()

258     def predict(self, X:pd.DataFrame):
259         """
260         Batch prediction method
261         """
262         predictions = []

264         for _, x in X.iterrows():
265             values = {}
266             for feature in self.features:
267                 values.update({feature: x[feature]})

269             predictions.append(self.predict_obs(values))

271         return predictions

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

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
```

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
|    |    |    |--- class: 0
|    |--- Fare >  10.48
|    |    |--- Age <= 6.50
|    |    |    |--- class: 1
|    |    |--- Age >  6.50
|    |    |    |--- class: 0
|--- Fare >  52.28
|    |--- Age <= 63.50
|    |    |--- Age <= 29.50
|    |    |    |--- class: 1
|    |    |--- Age >  29.50
|    |    |    |--- class: 1
|    |--- Age >  63.50
|    |    |--- class: 0
```

Although, the scikit-learn implementation prints out less information than my implementation.

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.

The code that I have written builds the same trees as scikit-learn implementation and the predictions are the same. But the training time for the scikit-learn algorithm is much faster. But my goal was not to grow the trees faster. My goal was to write an understandable code for any machine learning enthusiast to have a better grasp of what is happening under the hood.

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.

Recursion    Python    Decision Tree    Machine Learning    Gini Impurity

**Sign up for The Variable**

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. Take a look.

By signing up, you will create a Medium account if you don't already have one. Review our Privacy Policy for more information about our privacy practices.

✉⁺ Get this newsletter

Get the Medium app