

# Decision Tree classifier

Today your task is to get familiar with decision tree classifier - simple, but powerful case of discrete math usage.

---

## General idea

You are expected to write a quite simple, yet good core logic of decision tree classifier class. Additionally, you need to test your results and write down a report on what you've done, which principles used and explain the general process.

Hopefully, you have already learned what is decision tree classifier and how it work. For better understanding, and in case if something is still unclear for you, here are some useful links on basics of DTC:

- <https://towardsdatascience.com/decision-tree-from-scratch-in-python-46e99dfea775>
- <https://towardsdatascience.com/decision-tree-algorithm-in-python-from-scratch-8c43f0e40173>
- <https://machinelearningmastery.com/implement-decision-tree-algorithm-scratch-python/>
- <https://anderfernandez.com/en/blog/code-decision-tree-python-from-scratch/>

Also, for those interested to learn more about machine learning and particulary Desicion Trees - here is a great course on Coursera (you may be interested in the whole course or just this particular week):

- <https://www.coursera.org/learn/advanced-learning-algorithms/home/week/4>
- 

## Dataset

You can use Iris dataset for this task. It is a very popular dataset for machine learning and data science. It contains 150 samples of 3 different species of Iris flowers (Iris setosa, Iris virginica and Iris versicolor). Four features were measured from each sample: the length and the width of the sepals and petals, in centimeters.

Read more on this:

[https://scikit-learn.org/stable/auto\\_examples/datasets/plot\\_iris\\_dataset.html](https://scikit-learn.org/stable/auto_examples/datasets/plot_iris_dataset.html)  
[https://en.wikipedia.org/wiki/Iris\\_flower\\_data\\_set](https://en.wikipedia.org/wiki/Iris_flower_data_set)

P.S. In case you are not sure if your dataset is suitable, feel free to ask assistants :).

```
# install the required packages

# !pip install pandas
# !pip install numpy
# !pip install matplotlib
# !pip install graphviz
# !pip install scikit-learn

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

# scikit-learn package
from sklearn.datasets import load_iris
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier
from sklearn import tree
from sklearn.model_selection import train_test_split

iris = load_iris()
dir(iris)

['DESCR',
 'data',
 'data_module',
 'feature_names',
 'filename',
 'frame',
 'target',
 'target_names']

iris.target

array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0,
      0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0,
      0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1,
      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1,
      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2])
```

```
2,
    2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2,
    2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2])
```

This means that we have 150 entries (samples, infos about a flower). The columns being: Sepal Length, Sepal Width, Petal Length and Petal Width(features). Let's look at first two entries:

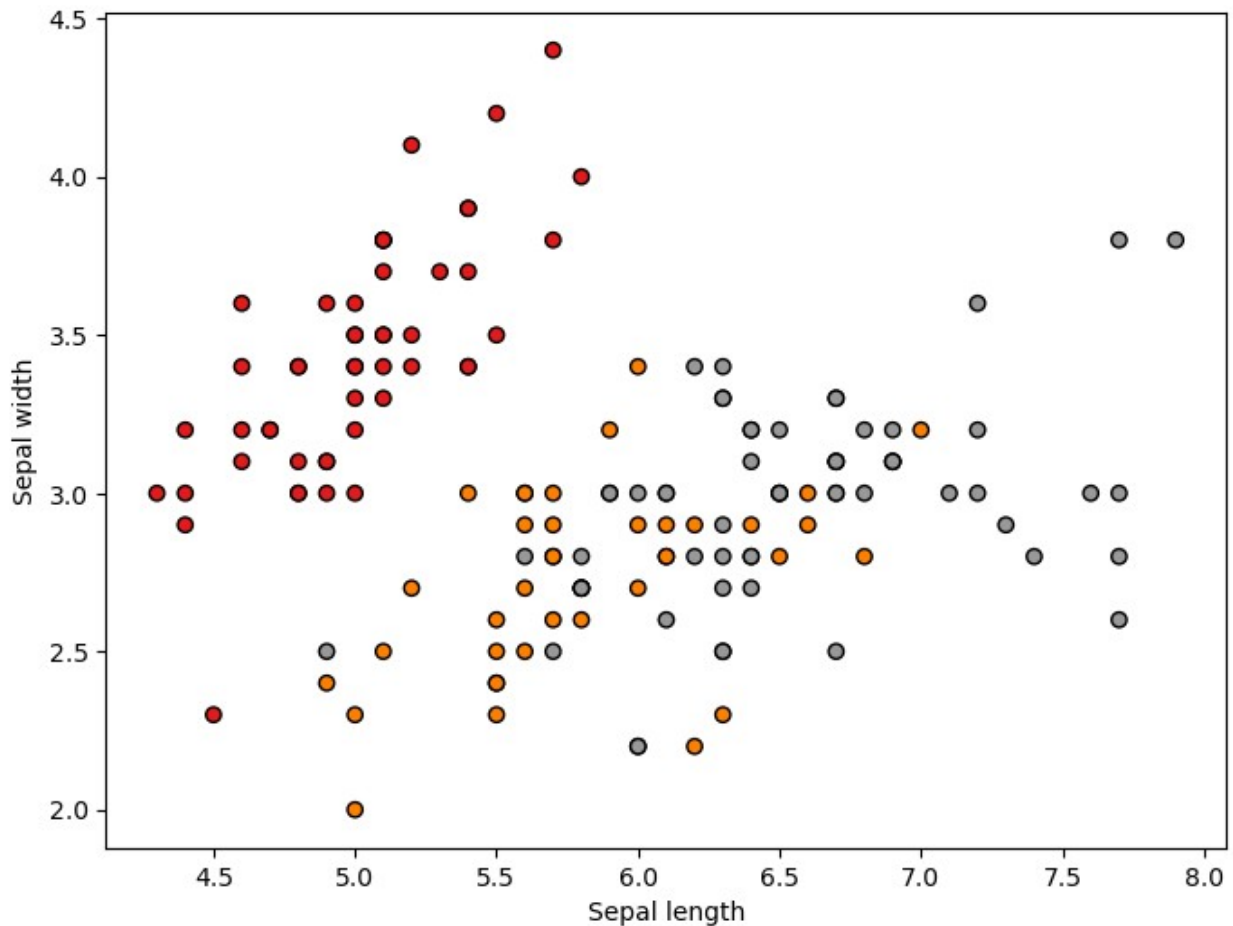
```
iris.data[0:2]
array([[5.1, 3.5, 1.4, 0.2],
       [4.9, 3. , 1.4, 0.2]])
```

## To understand data little bit better, let's plot some features

```
X = iris.data[:, :2] # we only take the first two features.
y = iris.target
plt.figure(2, figsize=(8, 6))
plt.clf()

# Plot the training points
plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1, edgecolor="k")
plt.xlabel("Sepal length")
plt.ylabel("Sepal width")

Text(0, 0.5, 'Sepal width')
```



From this we can clearly see, that even basing on those two parameters, we can clearly divide (classify) out data into several groups. For this, we will use decision tree classifier: <https://scikit-learn.org/stable/modules/tree.html#tree>

## Example of usage

**Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression.** The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. A tree can be seen as a piecewise constant approximation.

```
clf = DecisionTreeClassifier()
dir(clf)

['_abstractmethods__',
 '_annotations__',
 '_class__',
 '_delattr__',
 '_dict__',
```

```
'__dir__',
'__doc__',
'__eq__',
'__format__',
'__ge__',
'__getattr__',
'__getstate__',
'__gt__',
'__hash__',
'__init__',
'__init_subclass__',
'__le__',
'__lt__',
'__module__',
'__ne__',
'__new__',
'__reduce__',
'__reduce_ex__',
'__repr__',
'__setattr__',
'__setstate__',
'__sizeof__',
'__sklearn_clone__',
'__str__',
'__subclasshook__',
'__weakref__',
'_abc_impl',
'_build_request_for_signature',
'_check_feature_names',
'_check_n_features',
'_compute_missing_values_in_feature_mask',
'_doc_link_module',
'_doc_link_template',
'_doc_link_url_param_generator',
'_estimator_type',
'_fit',
'_get_default_requests',
'_get_doc_link',
'_get_metadata_request',
'_get_param_names',
'_get_tags',
'_more_tags',
'_parameter_constraints',
'_prune_tree',
'_repr_html_',
'_repr_html_inner',
'_repr_mimebundle_',
'_support_missing_values',
'_validate_X_predict',
```

```
'_validate_data',
'_validate_params',
'apply',
'ccp_alpha',
'class_weight',
'cost_complexity_pruning_path',
'criterion',
'decision_path',
'feature_importances_',
'fit',
'get_depth',
'get_metadata_routing',
'get_n_leaves',
'get_params',
'max_depth',
'max_features',
'max_leaf_nodes',
'min_impurity_decrease',
'min_samples_leaf',
'min_samples_split',
'min_weight_fraction_leaf',
'monotonic_cst',
'predict',
'predict_log_proba',
'predict_proba',
'random_state',
'score',
'set_fit_request',
'set_params',
'set_predict_proba_request',
'set_predict_request',
'set_score_request',
'splitter']
```

```
X, y = iris.data, iris.target
X.shape, y.shape
```

```
((150, 4), (150,))
```

## Train / test split

We train our model using training dataset and evaluate its performance basing on the test dataset. Reason to use two separate datasets is that our model learns its parameters from data, thus test set allows us to check its possibilities on completely new data.

```
X, X_test, y, y_test = train_test_split(X, y, test_size= 0.20)
```

## Model learning

It learns its parameters (where it should split data and for what threshold value) basing on the training dataset. It is done by minimizing some cost function (e.g. Gini impurity or entropy).

```
clf = clf.fit(X, y)
```

## Visualization of produced tree

You do not need to understand this piece of code :)

```
import graphviz
dot_data = tree.export_graphviz(clf, out_file=None)
graph = graphviz.Source(dot_data)
graph.render("iris")

'iris.pdf'

dot_data = tree.export_graphviz(clf, out_file=None,
                                feature_names=iris.feature_names,
                                class_names=iris.target_names,
                                filled=True, rounded=True,
                                special_characters=True)
graph = graphviz.Source(dot_data)
graph
```

```
X_test.shape  
(30, 4)
```

## Prediction step

Now we can use our model to predict which type has a flower, basing on its parameters.

This is conducted basically via traversing the tree that you can see above.

```
predictions = clf.predict(X_test)
```

We can also measure the accuracy of our model

```
sum(predictions == y_test) / len(y_test)  
0.9666666666666667
```



To get clearer intuition about prediction, let's look at those X, that should be labeled to some flower

```
y_test
array([1, 1, 1, 0, 0, 2, 2, 1, 0, 2, 2, 0, 0, 2, 0, 2, 1, 0, 0, 0, 1,
       2,
       1, 2, 0, 0, 0, 1, 0, 1])
```

Here you can traverse the tree above by yourself and make sure that prediction works

```
X_test[1]
array([5. , 2.3, 3.3, 1. ])
clf.predict([X_test[1]])
array([1])
```

Finally, it is your turn to write such classifier by yourself!

## Gini impurity

Decision trees use the concept of Gini impurity to describe how “pure” a node is. A node is pure ( $G = 0$ ) if all its samples belong to the same class, while a node with many samples from many different classes will have a Gini closer to 1.

$$G = 1 - \sum_{k=1}^n p_k^2$$

For example, if a node contains five samples, with two belonging to the first class (first flower), two of class 2, one of class 3 and none of class 4, then

$$G = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{2}{5}\right)^2 - \left(\frac{1}{5}\right)^2 = 0.64$$

```
class Node:
    def __init__(self, X, y, gini):
        self.X = X
        self.y = y
        self.gini = gini
        self.predicted_class = y
        self.feature_index = 0
        self.threshold = 0
        self.left = None
        self.right = None
```

```

class MyDecisionTreeClassifier:
    '''My decision tree'''
    def __init__(self, max_depth, classes):
        self.max_depth = max_depth
        self.classes = classes

    def gini(self, groups:list, classes)->float:
        '''
        A Gini score gives an idea of how good a split is by how mixed
the
        classes are in the two groups created by the split.

        A perfect separation results in a Gini score of 0,
        whereas the worst case split that results in 50/50
        classes in each group result in a Gini score of 0.5
        (for a 2 class problem).
        '''

        output = 1.0 - sum((n / groups) ** 2 for n in classes)
        return output
    def split_data(self, X, y):
        """Find the best split for a node.
        "Best" means that the average impurity of the two children,
weighted by their
        population, is the smallest possible. Additionally it must be
less than the
        impurity of the current node.
        To find the best split, we loop through all the features, and
consider all the
        midpoints between adjacent training samples as possible
thresholds. We compute
        the Gini impurity of the split generated by that particular
feature/threshold
        pair, and return the pair with smallest impurity.
        Returns:
            best_idx: Index of the feature for best split, or None if
no split is found.
            best_thr: Threshold to use for the split, or None if no
split is found.
        """
        m = y.size
        if m <= 1:
            return None, None

        # Count of each class in the current node
        num_parent = [np.sum(y == c) for c in range(self.n_classes_)]

        best_gini = self.gini(m, num_parent)
        best_idx, best_thr = None, None

```

```

# Loop through all features
for idx in range(self.n_features_):
    # Sort data along selected feature
    thresholds, classes = zip(*sorted(zip(X[:, idx], y)))

    num_left = [0] * self.n_classes_
    num_right = num_left.copy()
    for i in range(1, m):
        c = classes[i - 1]
        num_left[c] += 1
        num_right[c] -= 1
        gini_left = 1.0 - sum(
            (num_left[x] / i) ** 2 for x in
range(self.n_classes_))
        gini_right = 1.0 - sum(
            (num_right[x] / (m - i)) ** 2 for x in
range(self.n_classes_))
        gini = (i * gini_left + (m - i) * gini_right) / m

        if thresholds[i] == thresholds[i - 1]:
            continue

        if gini < best_gini:
            best_gini = gini
            best_idx = idx
            best_thr = (thresholds[i] + thresholds[i - 1]) / 2

    return best_idx, best_thr

def build_tree(self, X, y, depth=0):
    """
    create a root node

    recursively split until max depth is not exeeed
    """

    num_samples_per_class = [np.sum(y == i) for i in
range(self.n_classes_)]
    predicted_class = np.argmax(num_samples_per_class)
    node = Node(
        gini=self.gini(y.size, num_samples_per_class),
        X=num_samples_per_class,
        y=predicted_class,
    )

    # Split recursively until maximum depth is reached
    if depth < self.max_depth:
        idx, thr = self.split_data(X, y)

```

```

        if idx is not None:
            indices_left = X[:, idx] < thr
            X_left, y_left = X[indices_left], y[indices_left]
            X_right, y_right = X[~indices_left], y[~indices_left]
            node.feature_index = idx
            node.threshold = thr
            node.left = self.build_tree(X_left, y_left, depth + 1)
            node.right = self.build_tree(X_right, y_right, depth +
1)

    return node

def fit(self, X, y):
    """basically wrapper for build tree / train"""
    self.n_classes_ = len(set(y))
    self.n_features_ = X.shape[1]
    self.tree_ = self.build_tree(X, y)

def predict(self, inputs):
    """
    traverse the tree while there is a child
    and return the predicted class for it,
    note that X_test can be a single sample or a batch
    """
    node = self.tree_
    while node.left:
        if inputs[node.feature_index] < node.threshold:
            node = node.left
        else:
            node = node.right
    return node.predicted_class

def evaluate(self, X_test, y_test):
    '''return accuracy'''
    predictions = [self.predict(inputs) for inputs in X_test]
    return f'{sum(predictions == y_test) / len(y_test) * 100}%'

clf = MyDecisionTreeClassifier(max_depth=10, classes=y)
clf.fit(X, y)
print(clf.evaluate(X, y))

100.0%

```

For those who want to do it a little bit more complicated ;) (**optional**)

Consider also using some techniques to avoid overfitting, like pruning or setting a maximum depth for the tree. You can also try to implement some other metrics, to measure the quality of a split and overall performance. Also, you can try to implement some other algorithms, like proper CART, ID3 or C4.5. You can find more information about them here:

<https://scikit-learn.org/stable/modules/tree.html#tree>