# Decision Tree - Lab

## September 14, 2021

Modify the Decision Tree scratch code in our lecture such that: - Modify the scratch code so it can accept an hyperparameter max_depth, in which it will continue create the tree until max_depth is reached.

- Put everything into a class DecisionTree. It should have at least two methods, fit(), and predict()
- Load the iris data and try with your class

```
[20]: import numpy as np

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


class DecisionTree:
    def __init__(self, max_depth=None):
        self.max_depth = max_depth

    def fit(self, X, y):
        self.n_classes_ = len(set(y))
        self.n_features_ = X.shape[1]
        self.tree_ = self._grow_tree(X, y)

    def predict(self, X):
        return [self._predict(inputs) for inputs in X]

    def _best_split(self, X, y):
        m = y.size
        if m <= 1:
            return None, None
        num_parent = [np.sum(y == c) for c in range(self.n_classes_)]
        best_gini = 1.0 - sum((n / m) ** 2 for n in num_parent)
        best_idx, best_thr = None, None
```

```python
        for idx in range(self.n_features_):
            thresholds, classes = zip(*sorted(zip(X[:, idx], y)))
            num_left = [0] * self.n_classes_
            num_right = num_parent.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 _grow_tree(self, X, y, depth=0):
        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(predicted_class=predicted_class)
        if depth < self.max_depth:
            idx, thr = self._best_split(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._grow_tree(X_left, y_left, depth + 1)
                node.right = self._grow_tree(X_right, y_right, depth + 1)
        return node

    def _predict(self, inputs):
        node = self.tree_
        while node.left:
            if inputs[node.feature_index] < node.threshold:
                node = node.left
            else:
                node = node.right
```

```
        return node.predicted_class
```

```
[21]: import sys
      from sklearn.datasets import load_iris
      from sklearn.model_selection import train_test_split
      from sklearn.metrics import classification_report
      dataset = load_iris()
      X, y = dataset.data, dataset.target

      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
[22]: clf = DecisionTree(max_depth=5)
      clf.fit(X_train, y_train)
      y_pred = clf.predict(X_test)
      print(classification_report(y_test, y_pred))
```

```
                precision    recall  f1-score   support

             0       1.00      1.00      1.00        12
             1       1.00      0.71      0.83         7
             2       0.85      1.00      0.92        11

      accuracy                           0.93        30
     macro avg       0.95      0.90      0.92        30
  weighted avg       0.94      0.93      0.93        30
```

[ ]: