





```
class Node:
    def __init__(self, data: pd.DataFrame, max_depth=None, depth=None, node_type=None, question=None, min_split_data=
        self.outputCol = outputCol
        self.data = data
        self.min_split_data = min_split_data
        self.X = data.iloc[:, 0:]
        self.Y = data[self.outputCol]
        self.max_depth = max_depth
        self.depth = depth if depth else 0
        self.features = list(self.X.columns)
        self.question = question
        self.node_type = node_type if node_type else 'root'
        self.left = None
        self.right = None
        self.classes = self.Y.unique()
        self.classCount = self.Y.value_counts()
        self.classProbabilities = self.classCount/len(self.data)
        self.prediction = None
        self.infoGain = None
        self.previousProbabilities = previousProbabilities
        self.predictedProbabilities = None

    def split_impurity():
        pass

    return bestFeature, bestValue, typeOfFeature

def entropy(self, data: pd.DataFrame):
    Y = data[self.outputCol]
    classCount = Y.value_counts()
    classProbabilities = classCount/len(data)
    CE = 0
    for i in classProbabilities:
        CE = CE - i*math.log(i,2)
    return CE

def bestSplitInfoGain(self):
    infoGainBest = 0
    CE_data = self.entropy(self.data)
    for f in self.features:
        typeOfFeature = self.data[f].dtype
        columnData = self.data[f]
        if typeOfFeature == 'float64':
            columnData = columnData.sort_values()
            numData = columnData.unique()
            averageSplit = [(a + b) / 2 for a, b in zip(columnData[:2], columnData[1:2])] # taking average of
            for item in averageSplit:
                split1 = self.data[self.data[f] < item] # for each average value dividing dataframe into 2
                # left
                split2 = self.data[self.data[f] >= item]

                infoGain = CE_data - ((len(split1)/len(self.data))*(self.entropy(split1)) - ((len(split2)/len

                if (infoGain > infoGainBest): # best split with minimum mcr
                    infoGainBest = infoGain
                    bestFeature = f
                    bestValue = item

    print(infoGainBest)
    print('rule =', bestFeature, '<', bestValue)
    return bestFeature, bestValue, typeOfFeature, infoGainBest

def applySplit(self, bestFeature, bestValue, typeOfFeature):
    leftSplit = self.data.query(bestFeature + '<' + str(bestValue))
    rightSplit = self.data.query(bestFeature + '>=' + str(bestValue))
    return leftSplit, rightSplit

def learnDecisionTree(self, previousProbabilities=None):
    if self.node_type == 'root':
        self.data['previousProbabilities'] = previousProbabilities
        if self.node_type == 'firstTree':
            self.prediction = self.classCount.argmax()
            predictedProbabilities = []
            for i in self.Y:
                predictedProbabilities.append(self.classProbabilities[i])
            self.predictedProbabilities = predictedProbabilities
        elif (len(self.classes) > 1 and (self.depth < self.max_depth and len(self.data) > self.min_split_data
            bestFeature, bestValue, typeOfFeature, infoGainBest = self.bestSplitInfoGain()

            self.question = bestFeature + '<' + str(bestValue)
            self.infoGain = infoGainBest
            leftSplit, rightSplit = self.applySplit(bestFeature, bestValue, typeOfFeature)

            left = Node(leftSplit, self.max_depth, self.depth+1, 'left_node', None, 5, self.outputCol, self.previousP
            self.left.learnDecisionTree()

            right = Node(rightSplit, self.max_depth, self.depth+1, 'right_node', None, 5, self.outputCol, self.previous
            self.right.learnDecisionTree()

        else:
            numSum = 0
            dumSum = 0
            for index, row in self.data.iterrows():
                numSum = numSum + row['residual']
                dumSum = dumSum + row['previousProbabilities']

            self.prediction = numSum/dumSum

def breadth_first_search(self, root=None, width = 4 ):
    levelDict = {}
    root = self.root if root is None else root
    to_visit = [root]
    histList = [root]
    while to_visit:
        current = to_visit.pop(0)
        const = int(current.depth * width ** 1.5)
        if (levelDict.get(current.depth) != None ):
            nodes = levelDict.get(current.depth)
            nodes.append(current)
            levelDict.update({current.depth: nodes})
        else:
            nodes = []
            nodes.append(current)
            levelDict.update({current.depth: nodes})
        if (current.prediction == None):
            print(f"{'-' * const} {current.depth} | Split rule: {current.question}")
            if current.left:
                to_visit.append(current.left)
                histList.append(current.left)
            if current.right:
                to_visit.append(current.right)
                histList.append(current.right)
            else:
                print(f"{'-' * const} {current.depth} | Prediction: {current.prediction}")
    return histList, levelDict

def histogram(self, root=None, width = 4 ):
    root = self.root if root is None else root
    to_visit = [root]
    while to_visit:
        current = to_visit.pop(0)
        const = int(current.depth * width ** 1.5)
        if (current.prediction == None):
            axes = current.classProbabilities.plot.bar(title= current.question + ' level:' + str(current.dept
            if current.left:
                to_visit.append(current.left)
            if current.right:
                to_visit.append(current.right)
            else:
                axes = current.classProbabilities.plot.bar(title= current.prediction + ' level:' + str(current.de
        predict_class(self, X_test: pd.DataFrame):
            curNode = self
            if self.node_type == 'firstTree':
                return curNode.prediction
            best_feature, best_value = curNode.question.split('<')
            if (X_test[best_feature] < float(best_value)):
                curNode = curNode.left
            else:
                if self.right is not None:
                    curNode = curNode.right
            while curNode.depth < curNode.max_depth and (curNode.question != None):
                best_feature, best_value = curNode.question.split('<')
                if len(curNode.Y) < curNode.min_split_data:
                    break
                if (X_test[best_feature] < float(best_value)):
                    if self.left is not None:
                        curNode = curNode.left
                    else:
                        if self.right is not None:
                            curNode = curNode.right
            return curNode.prediction

def predict(self, d_test: pd.DataFrame, learningRate):
    if (self.node_type == 'firstTree'):
        X_test = d_test.iloc[:, 0:2]
        Y_test = d_test['Y']
        Y_hat = []
        for index, row in d_test.iterrows():
            Y_hat.append(self.predict_class(row))
        return (Y_hat)
    else:
        X_test = d_test.iloc[:, 0:2]
        Y_test = d_test[self.outputCol]
        predictedProbabilities = []
        for index, row in d_test.iterrows():
            logProb = row['previousProbabilities'] + (learningRate * self.predict_class(row))
            predictionProb = math.exp(logProb)/(1+math.exp(logProb))
            predictedProbabilities.append(predictionProb)
        self.predictedProbabilities = predictedProbabilities
        print(np.array(predictedProbabilities).sum())
        return (predictedProbabilities)
```

```
In [129]: listOffTrees = []
Out[129]: []
In [129]: firstTree = Node(train, 2, 1, 'firstTree', None, 5, 'Y', predictedProbabilities)
Out[129]: firstTree.learnDecisionTree(predictedProbabilities)
```

```
In [129]: firstTree.classCount
Out[129]: 0      36
          1      34
          Name: Y, dtype: int64
```

```
In [129]: len (firstTree.predictedProbabilities)
Out[129]: 70
```

```
In [129]: listOffTrees.append(firstTree)
Out[129]: Y_hat = firstTree.predict(train, 0.2)
In [129]: firstTree.classProbabilities
```

```
Out[129]: 0      0.514286
          1      0.485714
          Name: Y, dtype: float64
```

```
In [129]: residual = []
Out[129]: predictedProbabilities = []
for i,j in zip(train['Y'],Y_hat):
    residual.append(i-firstTree.classProbabilities[j])
train['predictedProbabilities'] = predictedProbabilities
train['residual'] = residual
```

```
In [130]: train['predictedProbabilities'] = predictedProbabilities
Out[130]: train
```

	x1	x2	Y	residual	predictedProbabilities	previousProbabilities
0	-0.096023	0.995379	0	-0.514286	0.514286	0.597095
1	1.672301	-0.240278	1	0.485714	0.514286	0.684788
2	0.991790	0.127877	0	-0.514286	0.514286	0.618081
3	0.050944	0.184892	1	0.485714	0.514286	0.669787
4	1.032052	-0.499486	1	0.485714	0.514286	0.684788
...	...	...	...	...	...	...
65	1.623490	-0.281831	1	0.485714	0.514286	0.684788
66	-0.462538	0.886599	0	-0.514286	0.514286	0.597095
67	1.572117	-0.320172	1	0.485714	0.514286	0.684788
68	0.623490	0.781831	0	-0.514286	0.514286	0.597095
69	1.801414	-0.098111	1	0.485714	0.514286	0.684788

70 rows x 6 columns

```
In [130]: numberoffree = 5
Out[130]: learningRate = 0.2
for i in range(1,5):
    tree = Node(train, 2, 1, 'root', None, 5, 'residual', listOffTrees[i-1].predictedProbabilities)
    tree.learnDecisionTree(predictedProbabilities)
    predictedProbabilities = tree.predict(train.data, learningRate)
    train['predictedProbabilities'] = predictedProbabilities
    residual = []
    for i,j in zip(train['Y'], predictedProbabilities):
        residual.append(i-j)
    train['residual'] = residual
    listOffTrees.append(tree)
```

```
0.40783647581206334
rule = x2 < 0.25
43.382939585236146
0.9994110647387533
rule = x2 < 0.25
44.20021570348916
0.9994110647387533
rule = x2 < 0.25
44.92416120060605
0.9994110647387533
rule = x2 < 0.25
44.8630740392363
```

```
In [130]: for i in listOffTrees:
Out[130]: print('-----Tree-----')
          histList, levelDict = l.breadth_first_search(i)
```

```
-----Tree-----
----- 1 | Prediction: 0
-----Tree-----
----- 1 | Split rule: x2<0.25
----- 2 | Prediction: 0.65849673020261441
----- 2 | Prediction: 0.7299382716049385
-----Tree-----
----- 1 | Split rule: x2<0.25
----- 2 | Prediction: 0.30000685858020293
----- 2 | Prediction: -0.7650123722212955
-----Tree-----
----- 1 | Split rule: x2<0.25
----- 2 | Prediction: 0.249731843047814
----- 2 | Prediction: -0.7714865805026665
-----Tree-----
----- 1 | Split rule: x2<0.25
----- 2 | Prediction: 0.2457543595677741
----- 2 | Prediction: -0.7728606934477926
```

```
In [124]:
```

	x1	x2	Y	residual	predictedProbabilities	previousProbabilities
0	-0.096023	0.995379	0	-0.598273	0.598273	0.598274
1	1.672301	-0.240278	1	0.314969	0.685031	0.685031
2	0.991790	0.127877	0	-0.618496	0.618496	0.618498
3	0.050944	0.184892	1	0.329840	0.670160	0.670160
4	1.032052	-0.499486	1	0.314969	0.685031	0.685031
...	...	...	...	...	...	...
65	1.623490	-0.281831	1	0.314969	0.685031	0.685031
66	-0.462538	0.886599	0	-0.598273	0.598273	0.598274
67	1.572117	-0.320172	1	0.314969	0.685031	0.685031
68	0.623490	0.781831	0	-0.598273	0.598273	0.598274
69	1.801414	-0.098111	1	0.314969	0.685031	0.685031

70 rows x 6 columns

```
In [ ]:
```