

# Flowchart for writing a decision tree algorithm using the ID3 algorithm :

1. Start.



2. Define a Node class to represent decision tree nodes.



3. Define Decision Tree class.



4. Initialize Decision Tree with hyperparameters such as max-depth and min-samples-split



5. Define an entropy() function to calculate the entropy of a given set of samples.



6. Define an information-gain() function to calculate the information gain of a given split.



7. Define a build-tree() function to recursively build the decision tree.



8. Define a ~~best~~ get-best-split() function to find the best split for a given set of samples.



9. Define a split() function to split a set of samples based on given features and threshold value.



10. In build-tree() function, check if stopping conditions are met (eg. min-samples-split or max-depth reached)



11. If stopping conditions are not met, use get-best-split() to find the best split.



12. If information-gain is positive, recursively call build-tree() on the left and right subtrees.

- ↓
13. Create a leaf node if stopping conditions are met or information gain is negative
  - ↓
  14. Return the root node of the decision tree
  - ↓
  15. End.

### Pseudo code

#### 1. Node class

START

CLASS Node.

FUNCTION \_\_init\_\_ (self, feature\_index = None, threshold = None, left = None, right = None, info\_gain = None, value = None)

SET self.feature\_index to feature\_index # Index of feature that node splits on

SET self.threshold to threshold # the threshold value used to split the data at this node.

SET self.right to right # right child of this node

SET self.left to left # left child " "

SET self.info\_gain to info\_gain # the information gain obtained by splitting at this node.

SET self.value to value # predicted value of target variable at this node (leaf node)

END FUNCTION

END CLASS.

END.

#### 2. Define Decision Tree class

START

CLASS Decision\_Tree\_Classifier.

FUNCTION \_\_init\_\_ (self, min\_samples\_split = 2, max\_depth = 2)

# initialize the root of the tree.

SET self.root to None.

# stopping conditions

SET self.min-samples-split to min-samples-split # min sample req to split the node.

SET self.max-depth to max-depth # max depth of tree.

END FUNCTION

END CLASS

END

3. Define Entropy()

$$\text{Entropy}(P) = - \sum_{i=1}^N P_i \log_2 P_i$$

START

FUNCTION entropy(y)

# compute the unique class in y

SET class-labels to the unique elements in y.

# initialize entropy to zero.

SET entropy to 0.

# iterate over each unique class-label.

FOR each cls in class-labels.

# compute the proportion of samples with class label cls (Probability)

SET p-cls to the number of samples in y with class label cls divided by the total no. of samples.

SET entropy to entropy +  $-p\text{-cls} \log_2(p\text{-cls})$

END FOR

class entropy

# return the computed entropy

RETURN entropy

END FUNCTION

END

4. Information-gain()

$$\text{Information gain} = \text{Entropy}(\text{Parent}) - \left( \frac{\text{Weighted}}{\text{average}} \right) \text{Entropy}(\text{children})$$

$$= I(D_p) - \frac{N_{\text{left}}}{N_p} I(D_{\text{left}}) - \frac{N_{\text{right}}}{N_p} I(D_{\text{right}})$$

START

FUNCTION information-gain(parent, l-child, r-child)

# calculate the entropy of the parent node

SET parent-entropy to the result of calling the entropy fn with parent.

# calculate the entropy of the left child node.

SET left-entropy to the result of calling the entropy function with l-child.

# calculate the entropy of right child node.

SET right-entropy to the result of entropy (r-child)

# calculate the weighted avg<sup>entropy</sup> of the child nodes.

SET child-entropy to  $(\text{len}(l\text{-child})/\text{len}(\text{parent})) * \text{left-entropy}$   
 $+ (\text{len}(r\text{-child})/\text{len}(\text{parent})) * \text{right-entropy}$

# calculate inform gain.

SET gain to parent-entropy - child-entropy

# return the computed inform gain

RETURN gain

END FUNCTION

END.

## 5. Define build-tree()

START

FUNCTION build-tree (dataset, curr-depth):

X, Y = get-features & target from dataset.

num-samples, num-features as shape of X. # get no. of samples & features

# check stopping condition.

IF num-samples  $\geq$  min-sample-split & curr-depth  $\leq$  max-depth

# find best split.

best-split = get-best-split (dataset, num-samples, num-features)

# check if inform gain is positive.

IF best-split ["inform-gain"]  $> 0$

left-subtree = build-tree (best-split ["dataset-left"], curr-depth+1)

right-subtree = build-tree (best-split ["dataset-right"], curr-depth+1)

# return decision node

RETURN Node

# compute ~~leaf~~ leaf node

leaf-value = calculate-leaf-value (Y)

return Node

END

6. get-best-split ()

START

FUNCTION get-best-split (dataset, num-samples, num-features)  
    <sup>Initialise</sup>  
    # ~~define~~ dictionary to store best split & variable to store maximum info gain.

    best-split = {}

    max-info-gain = ~~float~~ .float ()

    # loop over all the features & their possible threshold.

        FOR feature-index in <sup>range of</sup> num-feature

            feature-values = dataset (:, feature-index)

            possible threshold as unique values of feature-values

        # Split dataset into left & right subset using the split () method.

        FOR threshold in possible-threshold.

            dataset-left, dataset-right = split (dataset, feature-index, threshold)

        # compute information gain.

            current-information-gain as information-gain (y, left-y, right-y)

        # If the info gain is greater than the maximum information gain seen so far, update the best split with current split & new max info gain

        IF curr-info-gain > max-info-gain :

            best-split ["feature-index"] = feature-index.

            best-split ["threshold"] = threshold

            best-split ["dataset-left"] = dataset-left

            best-split ["dataset-right"] = dataset-right

            best-split ["Info-gain"] = curr-info-gain

            max-info-gain = curr-info-gain.

        RETURN best-split # returns the best split.

7. Split () :- split dataset into two subset.

START

FUNCTION split (dataset, feature-index, threshold)

# Initialise 2 ~~to~~ empty arrays

dataset-left = empty numpy array

dataset-right = empty numpy array

# Iterate/loop over each row in the dataset & append them to left subset if the feature value is less than or equal to the threshold, and to the right subset otherwise.

FOR ~~IF~~ each row in dataset

IF row[feature-index]  $\leq$  threshold  
append row to dataset-left

ELSE  
append row to dataset-right.

# return left & right subsets

RETURN dataset-left, dataset-right.

---

Commit