

## Assignment 3

Q.1 Compare decision tree classification with K-NN classification.

⇒

K-NN Classification	Decision tree
<ol style="list-style-type: none"> <li>1) Does not require specific training.</li> <li>2) Zero learning, that is why called lazy algorithm.</li> <li>3) Data must always be available for making predictions.</li> <li>4) Most interpretable algorithm.</li> <li>5) Takes much time in decision making, since it have to traverse the whole dataset to calculate the distance with each point.</li> <li>6) Since no training is required hence it zero(0) time.</li> <li>7) Suitable for small size datasets.</li> </ol>	<ol style="list-style-type: none"> <li>1) Require a proper training phase.</li> <li>2) Builds a model based on training.</li> <li>3) Once training is completed training data is not needed.</li> <li>4) It is also interpretable since we can see decision making in tree format.</li> <li>5) Once training is completed prediction time is less.</li> <li>6) Require initial training time to create decision node and branches.</li> <li>7) Used for small as well as large size dataset.</li> </ol>



Q.2 Compare decision tree classification with logistic regression classification.

Logistic Regression	Decision tree
<ol style="list-style-type: none"> <li>1) Logistic regression is a type of supervised learning algorithm used for classification tasks where the goal is to model probability that an instance belongs to given class or not.</li> <li>2) Coefficients can be interpreted to understand the influence of each feature on the outcomes.</li> <li>3) Less prone to overfitting compared to decision tree.</li> <li>4) It is not majorly affected by noise.</li> <li>5) Requires a large enough training dataset.</li> </ol>	<ol style="list-style-type: none"> <li>1) A decision tree is a type of supervised learning that is commonly used to predict outcomes based on input data.</li> <li>2) Highly interpreted, as the tree structure clearly shows the decision made at each node.</li> <li>3) Prone to overfitting.</li> <li>4) It is majorly affected by noise.</li> <li>5) Can be trained on small training set.</li> </ol>

Q.3 List down the attribute selection measures used by ID3 algorithm to construct a decision tree.

⇒

1) Attribute Selection Measures are

1) Gini Index

2) Information Gain

3) Entropy

4) Gain Ratio

1) Gini Index

⇒ All attributes are assumed to be continuous valued. It is assumed that there exist several possible (split) values for each attribute.

- Gini index method can be modified for categorical attributes.

- Gini index is used in Classification and Regression tree (CART).

If a dataset  $T$  contains example from  $n$  classes, gini index,  $gini(T)$  is defined as

$$gini(T) = 1 - \sum_{j=1}^n (p_j)^2$$

- In the above equation  $p_j$  represents the relative frequency of class  $j$  in  $T$ .



After splitting  $T$  into two subsets  $T_1$  and  $T_2$  with sizes  $N_1$  and  $N_2$ , gini index of split data is,

$$gini_{split}(T) = \frac{N_1}{N} gini(T_1) + \frac{N_2}{N} gini(T_2).$$

## 2) Information Gain (IG)

⇒ In this method all attributes are assumed to be categorical.

→ The amount of information required to decide if a random record  $s$  belongs

to  $A$  is defined as  $I(A)$

$$I(p/n) = - \left( \frac{p}{p+n} \right) \log_2 \left( \frac{p}{p+n} \right) - \left( \frac{n}{p+n} \right) \log_2 \left( \frac{n}{p+n} \right)$$

$$I(p/n) = - \left( \frac{p}{p+n} \right) \log_2 \left( \frac{p}{p+n} \right) - \left( \frac{n}{p+n} \right) \log_2 \left( \frac{n}{p+n} \right)$$

→ The higher the IG, the better the attribute is at classifying the examples.

## 3) Entropy

⇒ Entropy measures the impurity or disorder in the dataset. It is a measure of the uncertainty in the dataset's classification.

$$Entropy(S) = - \sum_{i=1}^n P_i \log_2 P_i$$

- To quantify how mixed or impure a set of data is.
- The higher the entropy, the more uncertain or impure the data is.

#### 4) Gain Ratio

⇒ Gain ratio is an extension of IG that takes into account the intrinsic information of a split, penalizing attribute that result in a large number of branches with small sets of data.

$$\text{Gain Ratio}(A) = \frac{\text{Information Gain}(S, A)}{\text{Split Information}(A)}$$

Q.4: Explain the properties of Gini Index

⇒

- The Gini index is a measure used in decision trees to evaluate quality of the split.
- It quantifies the degree of impurity or disorder in a dataset.
- The Gini is used in Classification and Regression Tree (CART).
- If a dataset  $T$  contains example from  $n$  classes, gini index,  $\text{gini}(T)$  is defined as-

$$\text{Gini}(T) = 1 - \sum_{j=1}^n (p_j)^2$$



- In the above equation,  $P_j$  represents the relative frequency of class  $j$  in  $T$ .
- After splitting  $T$  in two subsets  $T_1$  and  $T_2$  with sizes  $N_1$  and  $N_2$ , gini index of split data is:

$$\text{gini(split)} = \frac{N_1}{N} \text{gini}(T_1) + \frac{N_2}{N} \text{gini}(T_2)$$

- The attribute with smallest  $\text{gini(split)}(T)$  is selected to split the node.
- The gini index ranges from 0 to 0.5.
- 0 indicates perfect purity meaning all the instances belong to a single class.
- 0.5 represents maximum impurity, indicating that the instances are uniformly distributed across all classes.

$$s(i) = \frac{1}{n} \rightarrow 1 \quad \text{if } (T) \text{ is}$$

Q.3 Discuss in brief pruning in decision tree.

- ⇒ Pruning in a decision tree is a technique used to reduce the size of the decision tree by removing sections of tree that provide little to no additional predictive power.
- The main goal of pruning is to prevent overfitting and improve model's ability to generalize to unseen data.
  - Decision trees can grow very large and complex, capturing noise and outliers in the training data, which leads to overfitting.
  - Pruning helps to mitigate this by simplifying the tree.
  - By removing unnecessary branches, pruning reduces the complexity of the model, leading to better performance on new, unseen data.
  - There are two main types of pruning in decision tree :-

1) Pre-pruning.

- ⇒ This involves halting the growth of the tree, before it becomes too complex.
- It sets a limit on tree depth, maximum samples per leaf, or minimum information gain required for a split.



## 2) Post-pruning

⇒ This involves first allowing the tree to grow fully and then pruning back unnecessary branches.

The pruning techniques are:

### 1) Cost-complexity pruning

⇒ This method assigns a price to each subtree primarily based on its accuracy and complexity, then selects the subtree with the lowest free price.

### 2) Reduced Error Pruning

⇒ Removes the branches that do not significantly affect the overall accuracy.

*18/1/20*