

## Assignment No.2

Roll Number- 41449

Title: To find the decision based on a given scenario from a dataset using Decision Tree Classifier.

Problem Statement: A dataset collected in a cosmetics shop showing details of customers and whether or not they responded to a special offer to buy a new lip-stick is shown in table below. Use this dataset to build a decision tree, with Buys as the target variable, to help in buying lip-sticks in the future. Find the root node of decision tree. According to the decision tree you have made from previous training data set, what is the decision for the test data: [Age < 21, Income = Low, Gender = Female, Marital Status = Married]?

Objective: To understand how to decision tree classifier algorithm works on the give dataset

Outcome: To find the decision based on a given scenario of people with income, gender and marital status information from a dataset using Decision Tree Classifier.

Concept Related Theory:

1. **Nodes** : Test for the value of a certain attribute.
2. **Edges/ Branch** : Correspond to the outcome of a test and connect to the next node or leaf.
3. **Leaf nodes** : Terminal nodes that predict the outcome (represent class labels or class distribution).

There are two main types of Decision Trees:

1. Classification Trees.
2. Regression Trees.

### 1. Classification trees (Yes/No types) :

What we've seen above is an example of classification tree, where the outcome was a variable like 'fit' or 'unfit'. Here the decision variable is **Categorical/ discrete**. Classification tree methods (i.e., decision tree methods) are recommended when the data mining task contains classifications or predictions of outcomes, and the goal is to generate rules that can be easily explained and translated into SQL or a natural query language. A Classification tree labels, records, and assigns variables to discrete classes. A Classification tree can also provide a measure of confidence that the classification is correct. A Classification tree is built through a process known as binary recursive partitioning. This is an iterative process of splitting the data into partitions, and then splitting it up further on each of the branches.

### 2. Regression trees (Continuous data types) :

Decision trees where the target variable can take continuous values (typically real numbers) are called regression trees. (e.g. the price of a house, or a patient's length of stay in a hospital). All regression techniques contain a single output (response) variable and one or more input (predictor) variables. The output variable is numerical. The general regression tree building methodology allows input variables to be a mixture of continuous and

categorical variables. A decision tree is generated when each decision node in the tree contains a test on some input variable's value. The terminal nodes of the tree contain the predicted output variable values. A Regression tree may be considered as a variant of decision trees, designed to approximate real-valued functions, instead of being used for classification methods.

### Creation of Decision Tree :

In this method a set of training examples is broken down into smaller and smaller subsets while at the same time an associated decision tree get incrementally developed. At the end of the learning process, a decision tree covering the training set is returned. The key idea is to use a decision tree to partition the data space into cluster (or dense) regions and empty (or sparse) regions. In Decision Tree Classification a new example is classified by submitting it to a series of tests that determine the class label of the example. These tests are organized in a hierarchical structure called a decision tree. Decision Trees follow Divide-and-Conquer Algorithm.

### Divide and Conquer

Decision trees are built using a heuristic called **recursive partitioning**. This approach is also commonly known as **divide and conquer** because it splits the data into subsets, which are then split repeatedly into even **smaller subsets**, and so on and so forth until the process stops when the algorithm determines the data within the subsets are **sufficiently homogenous**, or another stopping criterion has been met.

### Basic Divide-and-Conquer Algorithm :

1. Select a test for root node. Create branch for each possible outcome of the test.
2. Split instances into subsets. One for each branch extending from the node.
3. Repeat recursively for each branch, using only instances that reach the branch.
4. Stop recursion for a branch if all its instances have the same class.

### Decision Tree Classifier

Using the decision algorithm, we start at the tree root and split the data on the feature that results in the **largest information gain (IG)** (reduction in uncertainty towards the final decision).

In an iterative process, we can then repeat this splitting procedure at each child node **until the leaves are pure**. This means that the samples at each leaf node all belong to the same class.

In practice, we may set a **limit on the depth of the tree to prevent overfitting**. We compromise on purity here somewhat as the final leaves may still have some impurity.

### Attribute Selection Measures

Attribute selection measure is a heuristic for selecting the splitting criterion that partition data into the best possible manner. It is also known as splitting rules because it helps us to determine breakpoints for tuples on a given node. ASM provides a rank to each feature(or attribute) by explaining the given dataset. Best score attribute will be selected as a splitting

attribute. In the case of a continuous-valued attribute, split points for branches also need to define. Most popular selection measures are Information Gain, Gain Ratio, and Gini Index.

### Information Gain

Shannon invented the concept of entropy, which measures the impurity of the input set. In physics and mathematics, entropy referred as the randomness or the impurity in the system. In information theory, it refers to the impurity in a group of examples. Information gain is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values. ID3 (Iterative Dichotomiser) decision tree algorithm uses information gain.

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

Where,  $P_i$  is the probability that an arbitrary tuple in  $D$  belongs to class  $C_i$ .

$$\text{Info}_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

Where,

$\text{Info}(D)$  is the average amount of information needed to identify the class label of a tuple in  $D$ .

$|D_j|/|D|$  acts as the weight of the  $j$ th partition.

$\text{Info}_A(D)$  is the expected information required to classify a tuple from  $D$  based on the partitioning by  $A$ .

The attribute  $A$  with the highest information gain,  $\text{Gain}(A)$ , is chosen as the splitting attribute at node  $N()$ .

### Gain Ratio

Information gain is biased for the attribute with many outcomes. It means it prefers the attribute with a large number of distinct values. For instance, consider an attribute with a unique identifier such as `customer_ID` has zero  $\text{info}(D)$  because of pure partition. This maximizes the information gain and creates useless partitioning. C4.5, an improvement of ID3, uses an extension to information gain known as the gain ratio. Gain ratio handles the issue of bias by normalizing the information gain using Split Info.

$$\text{SplitInfo}_A(D) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2 \left( \frac{|D_j|}{|D|} \right)$$

Where,

$|D_j|/|D|$  acts as the weight of the  $j$ th partition.

$v$  is the number of discrete values in attribute  $A$ .

The gain ratio can be defined as The attribute with the highest gain ratio is chosen as the splitting attribute.

$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}_A(D)}$$

## Gini index

Another decision tree algorithm CART (Classification and Regression Tree) uses the Gini method to create split points.

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

Where,  $p_i$  is the probability that a tuple in  $D$  belongs to class  $C_i$ . The Gini Index considers a binary split for each attribute. You can compute a weighted sum of the impurity of each partition. If a binary split on attribute  $A$  partitions data  $D$  into  $D_1$  and  $D_2$ , the Gini index of  $D$  is:

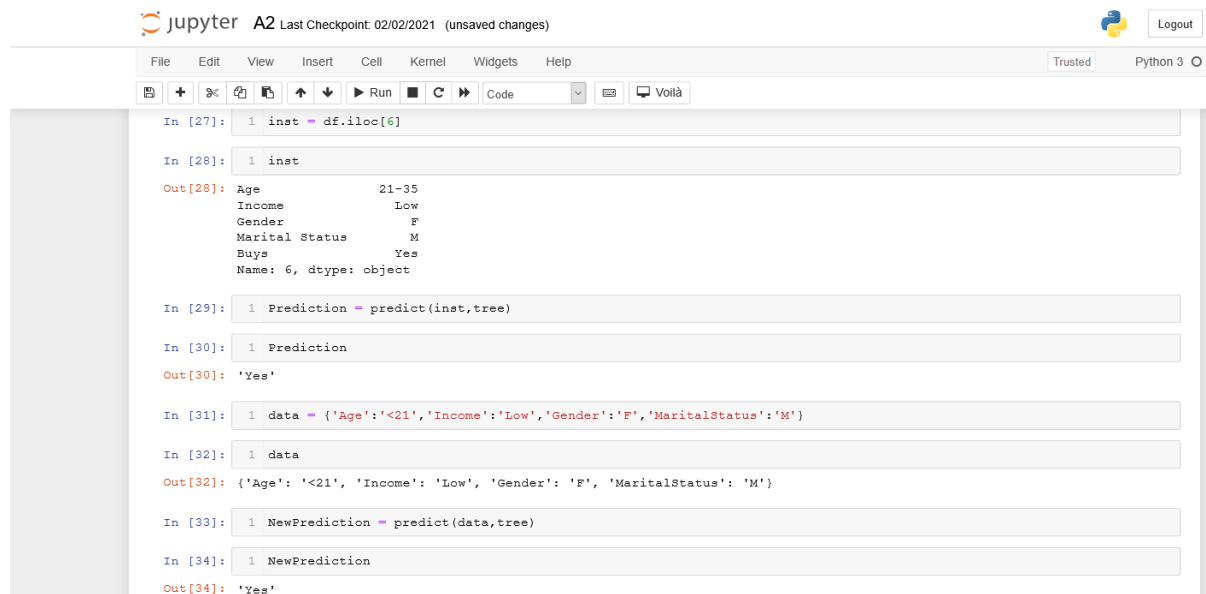
$$\text{Gini}_A(D) = \frac{|D_1|}{|D|} \text{Gini}(D_1) + \frac{|D_2|}{|D|} \text{Gini}(D_2)$$

In case of a discrete-valued attribute, the subset that gives the minimum gini index for that chosen is selected as a splitting attribute. In the case of continuous-valued attributes, the strategy is to select each pair of adjacent values as a possible split-point and point with smaller gini index chosen as the splitting point.

$$\Delta \text{Gini}(A) = \text{Gini}(D) - \text{Gini}_A(D).$$

The attribute with minimum Gini index is chosen as the splitting attribute.

## Results:



The image shows a Jupyter Notebook interface with the following code and output:

```

In [27]: 1 inst = df.iloc[6]

In [28]: 1 inst
Out[28]: Age      21-35
Income      Low
Gender      F
Marital Status  M
Buys      Yes
Name: 6, dtype: object

In [29]: 1 Prediction = predict(inst, tree)

In [30]: 1 Prediction
Out[30]: 'Yes'

In [31]: 1 data = {'Age': '<21', 'Income': 'Low', 'Gender': 'F', 'MaritalStatus': 'M'}

In [32]: 1 data
Out[32]: {'Age': '<21', 'Income': 'Low', 'Gender': 'F', 'MaritalStatus': 'M'}

In [33]: 1 NewPrediction = predict(data, tree)

In [34]: 1 NewPrediction
Out[34]: 'Yes'

```