## Q.

| department | status | age | salary | count |
|---|---|---|---|---|
| Sales | Senior | 31 – 35 | 46k – 50k | 30 |
| Sales | Junior | 26 – 30 | 26k – 30k | 40 |
| Systems | Junior | 31 – 35 | 31k – 35k | 40 |
| Systems | Junior | 21 – 25 | 46k – 50k | 20 |
| Systems | Senior | 31 – 35 | 66k – 70k | 5 |
| Systems | Junior | 26 – 30 | 46k – 50k | 3 |
| Systems | Senior | 41 – 45 | 66k – 70k | 3 |
| marketing | Senior | 36 – 40 | 46k – 50k | 10 |
| marketing | Junior | 31 – 35 | 41k – 45k | 4 |
| secretary | Senior | 46 – 50 | 36k – 40k | 4 |
| secretary | Junior | 26 – 30 | 26k – 30k | 6 |

**a]** How would you modify the basic decision tree algorithm to take into consideration the count of each generalized data tuple?

Ans → The basic decision tree algo should be modified as follows to take into consideration the count of each generalized data tuple —

- The count of each tuple must be integrated into the calculation of the attribute selection measure (such as information gain)
- Take the count into consideration to determine the most common class among the tuples.

**b]** use your algo to construct a decision tree from given data.

Ans  First we calculate gini (Gini) - Index for entire dataset

Total - count = 150
Sales - count = 110
System - count = 28
marketing - count = 14
Secretary - count = 10

$\rightarrow$ Gini - Total $= [1 - ((\text{sales\_count}/\text{total\_count}) \times 2$
$+ (\text{system\_count}/\text{total\_count}) \times 2$
$+ (\text{marketing\_count}/\text{total\_count}) \times 2$
$+ (\text{secretary\_count}/\text{total\_count}))]$

Gini - total $= [1 - (0.7332 + 0.1872 + 0.0932 + 0.0672)]$

$\rightarrow$ Now calculate Gini - Index for each $= \boxed{0.612}$

$\rightarrow$ Now calculate Gini - Index for each attribute →

**① Department**

$\rightarrow$ **Sales:** Senior $= 0$

Junior $= [1 - ((\frac{40}{70}) \times 2 + (\frac{30}{70}) \times 2] = \boxed{0.489}$

→ $\underline{System}$ : Senior $= \left[1 - \left(\left(\frac{5}{8}\right) \times 2 + \left(\frac{3}{8}\right) \times 2\right)\right] = \boxed{0.469}$

Junior $= \left[1 - \left(\left(\frac{23}{28}\right) \times 2 + \left(\frac{5}{28}\right) \times 2\right)\right] = \boxed{0.408}$

→ $\underline{Marketing}$ : Senior $= 0$    Junior $= 0$

→ $\underline{Secretary}$ : Senior $= \boxed{0.375}$   Junior $= \boxed{0.5}$

→ $\underline{Status}$ : Senior → Sales $= 0.469$     |    Junior → Sales $= 0.489$
                 Systems $= 0.375$                     Systems $= 0.408$
                 Marketing $= 0.5$                    Marketing $= 0$
                 Secretary $= 0$                      Secretary $= 0.5$

→ $\underline{Age}$ :   $(21-25) = 0$           $(31-35) →$         |   $(36-40) = 0$
        $(26-30) →$ Sales $= 0.489$    Sales $= 0.489$     $(41-45) = 0$
                   Systems $= 0.408$    Systems $= 0.469$   $(46-50) →$
                   Marketing $= 0$      Marketing $= 0$
                   Secretary $= 0.5$    Secretary $= 0$     Secretary $= 0.375$

→ $\underline{Salary}$ : $(26K-30K) →$   |   $(31K-35K) →$   |   $(36K-40K) →$
         Sales $= 0.489$      Sales $= 0.489$      Secretary $= 0.375$
         Systems $= 0$        System $= 0.469$   $(41K-45K) →$
         Marketing $= 0$     Marketing $= 0$      Marketing $= 0$
         Secretary $= 0.5$    Secretary $= 0$     $(46K-50K) →$

→ Attribute of Lowest Gini-Index is            Sales $= 0.489$
   Department with value $= \boxed{0.373}$        System $= 0.469$
   Split dataset with based on department    Marketing $= 0$
   attribute →   Sales : Status →        Secretary $= 0$

                              Senior $= 0$
                              Junior $= \left[1 - \left[\left(\frac{40}{70}\right) \times 2 + \left(\frac{30}{70}\right) \times 2\right]\right]$

                                     $= [1 - (1.14 + 0.85)]$

                                     $= \boxed{0.997}$



                                  ← — — — $[\underline{Decision-tree}]$

c] Given a data-tuple having the values "systems", "26-30" & "46-50K" for the attributes department / age / salary respectively what would be **a** naive Bayesian classification?

→ Given a data tuples with the values — "System"   "Junior"   "26---30" for the attribute department status / age respectively what would a naive Bayesian classification for salary tuple be →

$P(x | Senior) = 0$    |    $P(x | Junior) = 0.018$

                        — Thus a naive Bayesian classification
                         predicts Junior (Am)

**Q** Why is tree pruning useful in decision tree induction? What is a drawback of using a separate set of tuples to evaluate pruning?

**Ans,** The decision tree built may overfit the training data. There could be too many branches, some of which may reflected anomalies in the training data due to noise. Tree Pruning addresses this issue of overfitting the data by removing the least reliable branches. This generally result in a more compact and reliable decision tree that is faster and more accurate in it's classification.

The drawback of using a separate set of tuples to evaluate pruning is that it may not be representative of the training tuple used to create the original decision tree. If the separate set of tuples are skewed then using them to evaluate the pruned tree would not be a good indicator of the pruned tree's classification accuracy. Furthermore, using a separate set of tuple to evaluate pruning means there are less tuples to use for creation & testing of the tree. While this is also considered a drawback in machine learning. It may not be so in data mining due to the availability of larger data sets.

**Q** Briefly outline the major steps of decision tree classification?

**Ans,**
    Step-I : Determine the root of the tree
    Step-II: Calculate Entropy for the classes.
    Step-III: Calculate Entropy after split for each attribute.
    Step-IV: Calculate information gain for each split.
    Step-V: Perform further split.
    Step-VI: Complete the decision tree.

Calculation formula →

$$\text{Gini} = 1 - \sum_{i=1}^{n} p^2(c_i)$$

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

$\left[\begin{array}{l} P(c_i) \leftarrow \\ \text{probabity} / \\ \text{Percentage of class} \\ (c_i) \leftarrow \text{In a node} \end{array}\right]$

$$E(s) = -(P_+) * \log_2(P_+) - (P_-) * \log_2(P_-)$$

Entropy (S)

$\left[\begin{array}{l} \text{where } (P_+) \text{ positive sample} \\ (P_-) \text{ negative sample} \\ (s) \text{ sample of attribution} \end{array}\right]$

**Q.** Why is naive Bayesian classification called "naive"? Briefly outline the major ideas of naive Bayesian classification?

**Ans,**

Naive Bayesian classification is called naive cause, it assumes class conditional independence.

- That is, the effect of an attribute value on a given class is independent of the value of other attributes.

- The assumption is made to reduce computational costs, and hence is considered "_naïve_"

- The major idea behind "naïve" Bayesian classification is to try and classify data by maximizing

$P(X/c_i)P(c_i)$ [where, $i$ = index of the class] using the Baye's theorem of posterior probability.

- We are given a set of unknown data tuples, where each tuple is represented by an n-dimensional vector. $X = (X_1, X_2 \dots X_n)$ depicting n-measurement made on the tuple from n-attribute, respectively $(A_1, A_2 \dots A_n)$. also given a set of m-classes $(C_1, C_2, \dots C_m)$

- using Bayes theorem, the naive Bayesian classifier calculates the posterior probability of each class conditioned on X. ( $X \leftarrow$ assigned the class label of the class with max posterior )

try to maximize $P(c_i/X) = P(X/c_i)P(c_i)/P(X)$

However since, $P(X)$ is constant for all classes. only the $P(X/c_i)P(c_i)$ need be maximized. If the —

- class prior probability are not unknown. then it's common assumed that the classes are equally likely —

$$P(c_1) = P(c_2) = \dots P(c_m)$$

Therefore, should be maximize $P(X/c_i)$

- otherwise,

Maximize $P(X/c_i)P(c_i)$ the class prior probabilities may be estimated by $P(c_i) = s_i$

( where '$s_i$' ∈ num of training tuples of class $c_i$ )

$s \in$ total num of training tuples

- In order of reduce computation in evaluating $P(X/c_i)$

If $A_k$ is categorical attribute then $P(x_k/c_i)$ equal to the num of training tuples in '$c_i$'. that have '$x_k$' as the value for that attribute, divided by total num of training tuples in $(c_i)$.

- If $A_k$ is continuous attribute then $P(x_k/c_i)$ can be calculated using Gaussian density function.

**Q** what is information gain, Gain ratio, Gini-Index?

**Ans,** Information gain : →

    (i) Information gain is used for determining the best features / attribute that render information about a class.

    (ii) If it follows the concept of entropy while aiming at decreasing the level of entropy, begining from root node to the leaf node.

$$\left[\text{Information Gain} = \left(\begin{array}{c}\text{Entropy before}\\ \text{Splitting}\end{array}\right) - \left(\begin{array}{c}\text{Entropy after}\\ \text{Splitting}\end{array}\right)\right]$$

Gain ratio : →

    (i) First, determine the information gain of all the attributes, and then compute the avg of info gain.

    (ii) Second, calculate the gain ratio of all of the attribute whose calculated information Gain is larger / equal to the computed avg information Gain, then pick to the attribute of higher gain ratio.

$$\left[\text{Gain ratio} = \left(\frac{\text{Information gain}}{\text{Entropy}}\right)\right]$$

Gini-Index : →

    (1) Gini-Index computes the degree of probability of a specific variables that is wrongly being classified when chosen randomly and a variation of the gini-coefficient. It works on categorical variables provides outcome.

    (11) It varies from 0 → 1 where –

        – 0 depicts that all the elements be allied to a certain class or any / only one class exists there.

        – Gini-Index of values as 1 signifies that all the elements are randomly distribute across various classes.

        – value of (0.5) denotes the elements are uniformly distributed into some classes

**Q** Generate decision-tree algorithm?

**Ans,** Input:
- Data partition, D which is a set of training tuples and their associated class labels.
- attribute-list, the set of candidate attribute.
- Attribute-selection-method, a procedure to determine the spliting criterion that "best" partitions the data tuples into individual classes. The criterion consists of a spliting-attribute and posibly, either a spliting subnet.

Output: A decision tree

Method:

1. create a node N
2. If tuples in 'D' are all of the same class C, then return N as a leaf node labelled with class 'C'.
3. If attribute-list is empty then
4. return N as a leaf node labelled with the majority class in D // majority voting.
5. Apply attribute-selection-method (D, attribute-list) to find the best spliting-criterion;
6. Label node N with spliting criterion;
7. If spliting-attribute is discrete-valued and muuiway splits allowed then
8. attribute-list ← attribute-list - spliting attribute
9. for each outcome 'j' of spliting-creation.
   // partition the tuples & grow subtree
10. let 'Dj' be the set of data tuples in D satisfying outcome 'j'; // a partition.
11. If 'Dj' is empty then
12. attach a leaf labeled with the majority class in D to node N;
13. else attach a leaf labeled with the majority class return by Generate-decision-tree (D, attribute-list) to node N;

14. return N;

**Q =** EXPLAIN Bayesian classification?

**Ans,**

"Bayesian classification"

→ Baye's theorem →

    i) Posterior probability $[P(H/X)]$ — where, $X \leftarrow$ data tuple, $H \leftarrow$ hypothesis

    ii) Prior probability $[P(H)]$

    According to Baye's theorem —

$$P(H/X) = P(X/H)\,P(H) \big/ P(X)$$

→ Bayesian belief network →

    i) A belief network allows class conditional independencies to be defined between subset of variables.

    ii) It Provides a graphical model of casual relationship on which learning can be performed.

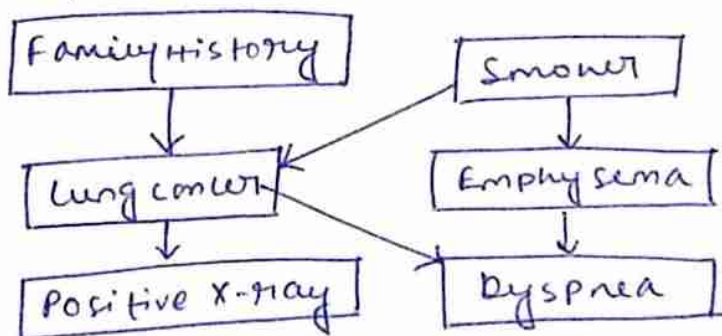    iii) we can use a trained Bayesian Network for classification.

→ Components of Baye'sian network —

    ① Directed acyclic graph

    ② A set of conditional probability table

→ Directed acyclic graph →

    i) Each node in a directed acyclic graph represents a random variable.

    ii) These variable may be discrete valued.

    iii) These variable may correspond to the actual attribute given in the data.

→ conditional probability table →

Family History → Lung cancer → Positive X-ray

Smoker → Lung cancer, Emphysema

Lung cancer → Positive X-ray, Dysprea

Emphysema → Dysprea

The conditional probability table for the values of the variable LC lung cancer (LC) showing each possible comb of values of it's parent nodes Family History (FH) | Smoker (S)

|      | (FH, S) | (FH, -S) | (-FH, S) | (-FH, -S) |
|------|---------|----------|----------|-----------|
| LC   | 0.8     | 0.5      | 0.7      | 0.1       |
| -LC  | 0.2     | 0.5      | 0.3      | 0.9       |

**Q** what is the use of regression? what may be the reason for not using the linear regression model to estimate the output data?

**Ans** → Regression analysis is a statistical method that is used to estimate the relationship between a dependent variable and one/more independent variables.

The primary use of regression analysis is to predict/estimate the value of the dependent variable based on the values of independent variable.

For example, researcher may use regression analysis to estimate sales of product.

→ Linear regression popularly & widely used for →

i) Non linear relationships
ii) outliers
iii) Multicollinearity
iv) categorical variables
v) Appropriateness of linear regression.

**Q 8.3** → If pruning a subtree, we would remove the subtree completely with method (b). However with method (a) If pruning a rule, we may remove any prediction of it. The latter is less restrictive.

**Q 8.4** → The worst case scenario occurs when we have to use as many attributes as possible before being able to classify each group of tuples. The Max depth of the tree $= \log(|D|)$. At each level we will have to compute the attribute selection measure $= O(n)$ times. The total num of tuples at each level of tree $O(n \times |D|)$ summing overall of the levels obtain $O(n \times |D| \times \log(|D|))$.

**Q 8.5** → We will use the rainforest algo for this problem, assumes there are 'c' class labels. Most memory required will be for AVC-set for the root of the tree. To compute the AVC-set for the root node we scan the database once & construct AVC-List $(100 \times c)$. Total site of AVC-set $= (100 \times c \times 50)$ & which will easily fit into 512 MB of memory for reasonable c. The computation of other AVC set is done in a smaller way but they will be smaller cause will be less attribute present. To reduce the num of scans we can compute the AVC-set for nodes at the same level of the tree in Parallel with such small AVC set.