

# MODULE III

## CHAPTER 3

### Classification

**University Prescribed Syllabus w.e.f Academic Year 2021-2022**

Basic Concepts; Classification methods : 1. Decision Tree Induction: Attribute Selection Measures, Tree pruning.  
2. Bayesian Classification: Naïve Bayes Classifier. Prediction : Structure of regression models; Simple linear regression, Accuracy and Error measures, Precision, Recall, Holdout, Random Sampling, Cross Validation, Bootstrap, Introduction of Ensemble methods, Bagging, Boosting, AdaBoost and Random forest.

**Self-learning Topics :** Multiple linear regression, logistic regression, Random forest, nearest neighbour classifier, SVM.

3.1	Introduction.....	3-1
3.2	Basic Concepts .....	3-2
3.2.1	What is Classification ? .....	3-3
3.2.2	How Does Classification Work ? .....	3-3
3.2.3	Classification Issues.....	3-3
3.2.4	Comparison of Classification Methods.....	3-4
3.3	Decision Tree Induction .....	3-4
UQ.	Define and explain: (i) Information Gain (ii) Entropy (iii) Gini Index (MU - Dec. 2019) .....	3-4
3.3.1	Decision Tree Induction Algorithm.....	3-4
3.3.2	Tree Pruning .....	3-5
GQ.	Given a decision tree, you have the option of (a) converting the decision tree to rules and then pruning the resulting rules, or (b) pruning the decision tree and then converting the pruned tree to rules. What advantage does (a) have over (b)?.....	3-5
3.3.3	Cost Complexity .....	3-6
3.3.4	Classification using Information Gain (ID3).....	3-6
3.3.5	Classification using CART Algorithm .....	3-17
3.3.6	ID3 Vs CART.....	3-22

### 3.1 INTRODUCTION

- There are two forms of data analysis that can be used for extracting models describing important classes or to predict future data trends.
- These two forms are as follows :
  - (i) Classification      (ii) Prediction
- Classification is a type of data analysis in which models defining relevant data classes are extracted.
- Classification models, called classifier, predict categorical class labels; and prediction models predict continuous valued functions.
- For example, we can build a classification model to categorize bank loan applications as either safe or risky, or a prediction model to predict the expenditures in dollars of potential customers on computer equipment given their income and occupation.

### 3.2 BASIC CONCEPTS

In this section we will discuss what classification is, working of classification, issues in classification and criteria for comparing the methods of classification.

#### 3.2.1 What is Classification ?

- Following are the examples of cases where the data analysis task is Classification –
  - (i) A bank loan officer wants to analyze the data in order to know which customer (loan applicant) are risky or which are safe.

name	age	income	loan_decision
Albin	youth	low	risky
Bryan	youth	low	risky
Carol	middle_aged	high	safe
Rikin	senior	medium	safe
Sahil	middle_aged	high	safe
Dhiren	middle_aged	low	risky
...	...	...	...

- (ii) A marketing manager at a company needs to analyze a customer with a given profile, who will buy a new laptop.
- In both of the above examples, a model or classifier is constructed to predict the categorical labels. These labels are "risky" or "safe" for loan application data and "yes" or "no" for marketing data.

#### 3.2.2 How Does Classification Work ?

With the help of the bank loan application that we have discussed above, let us understand the working of classification. The Data Classification process includes two steps –

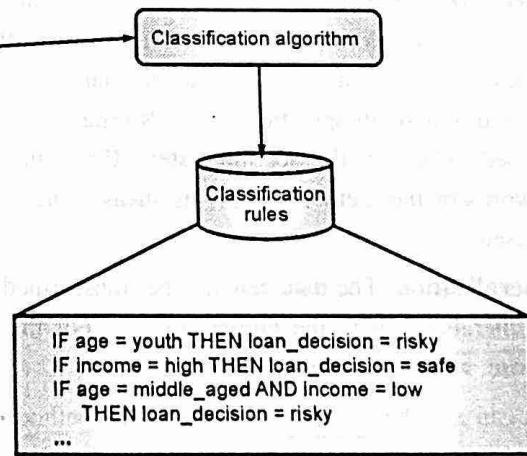
1. Building the Classifier or Model
2. Using Classifier for Classification

#### 1. Building the Classifier or Model

- This step is the learning step or the learning phase.
- In this step the classification algorithms build the classifier.
- The classifier is built from the training set made up of database tuples and their associated class labels.
- Each tuple that constitutes the training set is referred to as a category or class. These tuples can also be referred to as sample, object or data points.

Module

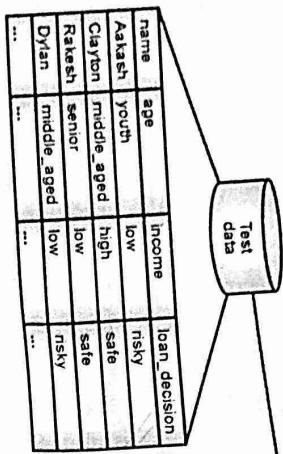
3



(1c) Fig. 3.2.1 : Building a Classifier

## 2. Using Classifier for Classification

- In this step, the classifier is used for classification. Here the test data is used to estimate the accuracy of classification rules. The classification rules can be applied to the new data tuples if the accuracy is considered acceptable.



(Fig. 3.2.2 : Testing a Classifier)

### 3.2.3 Classification Issues

The major issue is preparing the data for Classification. Preparing the data involves the following activities:

- Data Cleaning :** Data cleaning involves removing the noise and treatment of missing values. The noise is removed by applying smoothing techniques and the problem of missing values is solved by replacing a missing value with most commonly occurring value for that attribute.
- Relevance Analysis :** Database may also have the irrelevant attributes. Correlation analysis is used to know whether any two given attributes are related.
- Data Transformation and Reduction :** The data can be transformed by any of the following methods.
  - (i) Normalization : The data is transformed using normalization. Normalization involves scaling all values for given attribute in order to make them fall within a small specified range. Normalization is used when in the learning step, the neural networks or the methods involving measurements are used.
  - (ii) Generalization: The data can also be transformed by generalizing it to the higher concept. For this purpose, we can use the concept hierarchies.
  - (iii) Data can also be reduced by some other methods such as - wavelet transformation, - binning, histogram analysis and clustering.

- Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label.

The topmost node in the tree is the root node.

The following decision tree is for the concept buys\_laptop that indicates whether a customer at a company is likely to buy a computer or not. Each internal node represents a test on an attribute. Each leaf node represents a class (either buys\_laptop = yes or no).



(Fig. 3.3.1: Representation of a Decision Tree

The benefits of having a decision tree are as follows:

1. It does not require any domain knowledge.
2. It is easy to comprehend.
3. The learning and classification steps of a decision tree are simple and fast.

### 3.3.1 Decision Tree Induction Algorithm

A machine learning researcher named J. Ross Quinlan in 1980 developed a decision tree algorithm known as ID3 (Iterative Dichotomiser). Later, he presented C4.5, which was the successor of ID3. ID3 and C4.5 adopt a greedy approach.

In this algorithm, there is no backtracking; the trees are constructed in a top-down recursive divide-and-conquer manner.

**Algorithm :** Generate\_decision\_tree. Generating a decision tree from training tuples of data partition D

**Input**

UQ Define and explain: (i) Information Gain (ii) Entropy (MU - Dec. 2019)

Gini Index

MU - Dec. 2019

- Decision tree induction is the learning of decision trees from class-labeled training tuples.
- A decision tree is a structure that includes a root node, branches, and leaf nodes.

- Attribute selection method, a procedure to determine the splitting criterion that best partitions the data tuples into individual classes. This criterion includes a splitting attribute and either a splitting point or splitting subset.

**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 labeled with class, C;
3. If attribute\_list is empty then return N as a leaf node labeled with majority class in D; //majority voting
4. If attribute\_list is empty then apply the best splitting criterion.
5. If attribute\_list is not empty then
  - 5.1 return N as a leaf node with labeled with majority class in D;
  - 5.2 attach a leaf labeled with the majority class in D to node N;
  - 5.3 attach the node returned by Generalize\_decision\_tree for each outcome j of splitting criterion // partition the tuples and grow subtrees for each partition
  - 5.4 remove attribute\_list - splitting attribute;
  - 5.5 for each outcome j of splitting criterion
    - 5.5.1 attach a node Nj to node N;
    - 5.5.2 let Dj be the set of data tuples in D satisfying outcome j; // a partition
    - 5.5.3 If Dj is empty then attach a leaf labeled with the majority class in Dj to node Nj;
    - 5.5.4 else attach the node returned by Generalize\_decision\_tree for Dj to node Nj;
  - 5.6 return N;
6. apply the best splitting criterion.
7. label node N with splitting\_criterion;
8. If splitting\_attribute is discrete-valued and majority splits allowed then// no restricted to binary trees
  - 8.1 remove attribute\_list - splitting attribute;
  - 8.2 for each outcome j of splitting criterion
    - 8.2.1 attach a node Nj to node N;
    - 8.2.2 let Dj be the set of data tuples in D satisfying outcome j; // a partition
    - 8.2.3 If Dj is empty then attach a leaf labeled with the majority class in Dj to node Nj;
    - 8.2.4 else attach the node returned by Generalize\_decision\_tree for Dj to node Nj;
  - 8.3 return N;
9. If attribute\_list is not empty then
  - 9.1 attach a node Nj to node N;
  - 9.2 let Dj be the set of data tuples in D satisfying outcome j; // a partition
  - 9.3 If Dj is empty then attach a leaf labeled with the majority class in Dj to node Nj;
  - 9.4 else attach the node returned by Generalize\_decision\_tree for Dj to node Nj;
10. If splitting\_attribute is discrete-valued and majority splits allowed then// no restricted to binary trees
  - 10.1 remove attribute\_list - splitting attribute;
  - 10.2 for each outcome j of splitting criterion
    - 10.2.1 attach a node Nj to node N;
    - 10.2.2 let Dj be the set of data tuples in D satisfying outcome j; // a partition
    - 10.2.3 If Dj is empty then attach a leaf labeled with the majority class in Dj to node Nj;
    - 10.2.4 else attach the node returned by Generalize\_decision\_tree for Dj to node Nj;
  - 10.3 return N;
11. If attribute\_list is not empty then
  - 11.1 attach a node Nj to node N;
  - 11.2 let Dj be the set of data tuples in D satisfying outcome j; // a partition
  - 11.3 If Dj is empty then attach a leaf labeled with the majority class in Dj to node Nj;
  - 11.4 else attach the node returned by Generalize\_decision\_tree for Dj to node Nj;
12. If Dj is empty then attach a leaf labeled with the majority class in Dj to node N;
13. attach a leaf labeled with the majority class in D to node N;
14. else attach the node returned by Generalize\_decision\_tree for each outcome j of splitting criterion // partition the tuples and grow subtrees for each partition
15. return N;

### 3.3.2 Tree Pruning

The decision tree built may overfit the training data. There could be too many branches, some of which may reflect anomalies in the training data due to noise or outliers.

- Tree pruning addresses this issue of overfitting the data by removing the least reliable branches (using statistical measures).
- This generally results in a more compact and reliable decision tree that is faster and more accurate in its classification of data.

(Classification)...Page no. (3-6)

Data Mining & Business Intelligence (MU-Sem 6-IT)

**Data Mining & Business Intelligence (MU-Sem 6-IT)**

- There are two approaches to prune a tree :

(i) **Pre-pruning** : The tree is pruned by halting its construction early.

(ii) **Post-pruning** : This approach removes a sub-tree from a fully grown tree.

**Drawback of using a separate set of tuples to evaluate pruning**

- If a separate set of tuples are used to evaluate pruning is that it may not be representative of the training tuples 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 tuples to evaluate pruning means there are less tuples to use for creation and testing of the tree. While this is considered a draw back in machine learning, it may not be so in data mining due to the availability of larger data sets.

Q. Given a decision tree, you have the option of

- converting the decision tree to rules and then pruning the resulting rules, or (b) pruning the decision tree and then converting the pruned tree to rules. What advantage does (a) have over (b)?

**Ans.:**

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 precondition of it. The latter is less restrictive.

### 3.3.3 Cost Complexity

The cost complexity is measured by the following two parameters –

- Number of leaves in the tree, and
- Error rate of the tree

### 3.3.4 Classification using Information Gain

#### (ID3)

- ID3 stands for Iterative Dichotomiser 3 and is named such because the algorithm iteratively (repeatedly) dichotomizes(divides) features into two or more groups at each step.

(Classification)...Page no. (3-7)

**ID3 Steps**

- Calculate the Information Gain of each feature.
- Considering that all rows don't belong to the same class, split the dataset D into subsets using the feature for which the Information Gain is maximum.
- Make a decision tree node using the feature with the maximum Information gain.
- If all rows belong to the same class, make the current node as a leaf node with the class as its label.

Repeat the above steps for the remaining features until we run out of all features, or the decision tree has all leaf nodes.

Ex. 3.3.1 : Apply ID3 algorithm on the following training dataset and extract the classification rule from the tree.

Day	Outlook	Temp.	Humidity	Wind	Play_Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

$$\text{Gain}(\text{Outlook}) = H(D) - H(\text{Outlook}) \\ = 0.940 - 0.694 = 0.246$$

$$\text{Gain}(\text{Temperature}) = H(D) - H(\text{Temperature}) \\ = 0.940 - 0.911 = 0.029$$

$$\therefore H(\text{Temperature}) = \frac{4}{14} \times H(\text{Hot}) + \frac{6}{14} \times H(\text{Mild}) \\ + \frac{4}{14} \times H(\text{Cool})$$

$$= \frac{4}{14} \times 1 + \frac{6}{14} \times 0.92 + \frac{4}{14} \times 0.81 = 0.911$$

$$\text{Gain}(\text{Temperature}) = H(D) - H(\text{Temperature}) \\ = 0.940 - 0.911 = 0.029$$

#### (iv) Choosing Humidity as Splitting Attribute

Humidity	C1	C2	Play_Tennis = Yes	Play_Tennis = No	H
High	3	4	4	0.935	
Normal	6	1	1	0.592	

$$\therefore H(\text{Humidity}) = \frac{7}{14} \times H(\text{High}) + \frac{7}{14} \times H(\text{Normal}) \\ = \frac{7}{14} \times 0.985 + \frac{7}{14} \times 0.592 = 0.789$$

$$\text{Gain}(\text{Humidity}) = H(D) - H(\text{Humidity}) \\ = 0.940 - 0.789 = 0.151$$

**Data Mining & Business Intelligence (MU-Sem 6-IT)**

**(v) Choosing Wind as Splitting Attribute**

Wind	C1	C2	Entropy
Strong	3	3	1
Weak	6	2	0.811

$$\therefore H(Wind) = \frac{6}{14} \times H(Strong) + \frac{8}{14} \times H(Weak)$$

$$= \frac{6}{14} \times 1 + \frac{8}{14} \times 0.811 = 0.892$$

$$Gain(Wind) = H(D) - H(Wind)$$

$$= 0.940 - 0.892 = 0.048$$

**Summary**

$$Gain(OutlookID) = 0.246$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

Temperature	C1	C2	Entropy
- Hot	0	2	0
Mild	1	1	0
Cool	1	0	0

$$\therefore H(Temperature) = \frac{2}{3} \times H(Hot) + \frac{2}{3} \times H(Mild)$$

$$+ \frac{1}{3} \times H(Cool)$$

$$= \frac{2}{3} \times 0 + \frac{2}{3} \times 1 + \frac{1}{3} \times 0 = 0.4$$

$$= 0.971 - 0.4 = 0.571$$

$$Gain(Temperature) = H(D) - H(Temperature)$$

$$= 0.971 - 0.4 = 0.571$$

Humidity	C1	C2	Entropy
Play_Tennis = Yes Sunny	3	0	0
Play_Tennis = No Sunny	0	0	0

$$Gain(Humidity) = \frac{3}{5} \times H(High) + \frac{2}{5} \times H(Normal)$$

$$= \frac{3}{5} \times 0 + \frac{2}{5} \times 0 = 0$$

$$Gain(Humidity) = H(D) - H(Humidity)$$

$$= 0.971 - 0 = 0.971$$

$$Gain(Humidity) = H(D) - H(Humidity)$$

$$= 0.971 - 0 = 0.971$$

Outlook	C1	C2	Entropy
Sunny	0	1	0
Overcast	1	1	0

$$= \frac{1}{2} \times 0 + \frac{1}{2} \times 1 = \frac{1}{2} = 0.5$$

$$Gain(Outlook) = H(D) - H(Outlook)$$

$$= 0.971 - 0.5 = 0.471$$

$$Gain(Outlook) = H(D) - H(Outlook)$$

$$= 0.971 - 0.5 = 0.471$$

$$Gain(Outlook) = H(D) - H(Outlook)$$

$$= 0.971 - 0.5 = 0.471$$

**Data Mining & Business Intelligence (MU-Sem 6-IT)**

**(i) Entropy before split for the given database D1 :**

$$H(D_1) = \sum_{i=1}^5 P_i \log_2 \left( \frac{1}{P_i} \right) = \frac{2}{5} \log_2 \frac{5}{2} + \frac{3}{5} \log_2 \frac{5}{3} = 0.971$$

$$Gain(TemperatureID) = 0.571$$

$$Gain(HumidityID) = 0.971$$

$$Gain(WindID) = 0.02$$

$$Humidity' attribute has the highest gain; therefore, it is placed below Outlook = "Sunny".$$

$$Since, Humidity has two possible values, the Humidity' node has two branches (High, Normal).$$

$$From dataset D1, we find that when Humidity = High, Play_Tennis = No and when Humidity = Normal, Play_Tennis = Yes.$$

$$Let the class label attributes be as follows:$$

$$C1 = Play_Tennis = Yes|Rain = 3 Samples$$

$$C2 = Play_Tennis = No|Rain = 2 Samples$$

$$Therefore, P(C1) = 3/5 and P(C2) = 2/5$$

$$(i) Entropy before split for the given database D3:$$

$$H(D_3) = \sum_{i=1}^4 P_i \log_2 \left( \frac{1}{P_i} \right)$$

$$\therefore H(D_3) = \frac{3}{5} \log_2 \frac{5}{3} + \frac{2}{5} \log_2 \frac{5}{2} = 0.971$$

$$Gain(Temperature) = H(D_3) - H(Temperature)$$

$$= 0.971 - 0.4 = 0.571$$

$$Gain(Humidity) = H(D_3) - H(Humidity)$$

$$= 0.971 - 0.4 = 0.571$$

$$Gain(Wind) = H(D_3) - H(Wind)$$

$$= 0.971 - 0.951 = 0.02$$

$$(i) Summary$$

$$Gain(OutlookID) = 0.246$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

$$Gain(OutlookID) = 0.048$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

Day	Outlook	Temp.	Humidity	Wind	Play_Tennis
4	Rain	Mild	High	Weak	Yes
5	Rain	Cloudy	Normal	Weak	Yes
6	Rain	Cloudy	Normal	Strong	No
10	Rain	Mild	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

**Data Mining & Business Intelligence (MU-Sem 6-IT)**

**(ii) Choosing Temperature as Splitting Attribute**

Temperature	C1	C2	Entropy
Play_Tennis = Yes Rain	0	0	0
Play_Tennis = No Rain	1	1	0.913

$$Gain(Temperature) = \frac{1}{3} \times H(Hot) + \frac{2}{3} \times H(Normal)$$

$$+ \frac{2}{3} \times H(Cool)$$

$$= \frac{0}{3} \times 0 + \frac{2}{3} \times 0.913 + \frac{2}{3} \times 1 = 0.951$$

$$Gain(Temperature) = H(D_3) - H(Temperature)$$

$$= 0.971 - 0.951 = 0.02$$

$$Gain(Humidity) = H(D_3) - H(Humidity)$$

$$= 0.971 - 0.951 = 0.02$$

$$Gain(Wind) = H(D_3) - H(Wind)$$

$$= 0.971 - 0 = 0.971$$

$$(iii) Choosing Wind as Splitting Attribute$$

Wind	C1	C2	Entropy
Play_Tennis = Yes Rain	0	2	0
Play_Tennis = No Rain	3	0	0

$$\therefore H(Wind) = \frac{2}{3} \times H(Strong) + \frac{1}{3} \times H(Weak)$$

$$= \frac{2}{3} \times 0 + \frac{1}{3} \times 0 = 0$$

$$Gain(Wind) = H(D_3) - H(Wind)$$

$$= 0.971 - 0 = 0.971$$

$$(iv) Choosing Wind as Splitting Attribute$$

Wind	C1	C2	Entropy
Play_Tennis = Yes Sunny	1	1	0
Play_Tennis = No Sunny	2	0	0

$$\therefore H(Wind) = \frac{2}{3} \times H(Strong) + \frac{1}{3} \times H(Weak)$$

$$= \frac{2}{3} \times 1 + \frac{1}{3} \times 0.918 = 0.951$$

$$Gain(Wind) = H(D_3) - H(Wind)$$

$$= 0.971 - 0.951 = 0.02$$

$$(v) Summary$$

$$Gain(OutlookID) = 0.246$$

$$Gain(TemperatureID) = 0.029$$

$$Gain(HumidityID) = 0.151$$

$$Gain(WindID) = 0.048$$

Day	Outlook	Temp.	Humidity	Wind	Play_Tennis
1	Sunny	Hot	High	Weak	Yes
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	Yes
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

Wind	C1	C2	Entropy
Play_Tennis = Yes Sunny	1	1	0
Play_Tennis = No Sunny	2	0	0

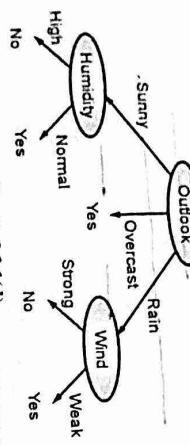
Wind	C1	C2	Entropy
Play_Tennis = Yes Sunny	1	1	0
Play_Tennis = No Sunny	2	0	0

Wind	C1	C2	Entropy

**Summary**  
 $\text{Gain}(\text{Temperature}(D)) = 0.02$   
 $\text{Gain}(\text{Wind}(D)) = 0.971$

Wind attribute has the highest gain; therefore, it is placed below Outlook = "Rain". Since, Wind has two possible values, the Wind node has two branches (Strong, Weak).

From dataset D3, we find that when Wind = Strong, Play\_Tennis = Yes.



Ex. 3.3.2 : A simple example from the stock market involving only discrete ranges has profit as categorical attribute with values \{Up, Down\} and the training data is :

Age	Competition	Type	Profit
Old	Yes	Software	Down
Old	No	Software	Down
Old	No	Hardware	Down
Mid	Yes	Software	Down
Mid	Yes	Hardware	Down
Mid	No	Software	Up
Mid	No	Hardware	Up
Mid	Yes	Software	Up
Mid	No	Software	Up
New	Yes	Software	Up
New	No	Hardware	Up
New	No	Software	Up

(i) Fig. P. 3.3.1(d)

The decision tree can also be expressed in rule format as :

IF Outlook = Sunny AND Humidity = High THEN Play\_Tennis = No  
 IF Outlook = Sunny AND Humidity = Normal THEN Play\_Tennis = Yes  
 IF Outlook = Overcast THEN Play\_Tennis = No

IF Outlook = Overcast THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Strong THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Weak THEN Play\_Tennis = YES

Apply decision tree algorithm and show the generated rules.

Soln. :

Let the class label attributes be as follows:

C1 = Profit = Down = 5 Samples

C2 = Profit = Up = 5 Samples

Therefore,  $P(C1) = 5/10$  and  $P(C2) = 5/10$

(i) Entropy before split for the given database D :

$$H(D) = \sum_{i=1}^5 p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$\therefore H(D) = \frac{5}{10} \log_2 \frac{10}{5} + \frac{5}{10} \log_2 \frac{10}{5} = 0.5 + 0.5 = 1$$

∴ HD = 1

(iii) Choosing Competition as the Splitting Attribute.

Competition	C1	C2	Entropy
Yes	3	1	0.8113
No	2	4	0.9183
Old	No	Software	Down
Old	No	Hardware	Down
Old	Yes	Software	Down
Mid	Yes	Hardware	Down
Mid	Yes	Software	Up
Mid	No	Software	Up
Mid	No	Hardware	Up
New	Yes	Software	Up
New	No	Hardware	Up
New	No	Software	Up

$$\therefore H(\text{Competition}) = \frac{4}{10} \times H(\text{Yes}) + \frac{6}{10} \times H(\text{No})$$

$$= \frac{4}{10} \times 0.8113 + \frac{6}{10} \times 0.9183 = 0.8755$$

Gain(Competition) =  $H(D) - H(\text{Competition})$

$$= 1 - 0.8755 = 0.1245$$

(ii) Choosing Type as the Splitting Attribute.

Type	C1	C2	Entropy
Software	3	1	0.8113
Software	2	2	0.9183
Hardware	0	3	0

$$H(\text{Type}) = \frac{6}{10} \times H(\text{Software}) + \frac{4}{10} \times H(\text{Hardware})$$

$$= \frac{6}{10} \times 1 + \frac{4}{10} \times 1 = 1$$

∴ H(Type) = 1

Now, consider Age = Mid and count the number of tuples from the original dataset D. Let us denote it as D1.

Age	Competition	Type	Profit
Mid	Yes	Software	Down
Mid	Yes	Hardware	Up
Mid	No	Software	Up
Mid	No	Hardware	Up

Let the class label attributes be as follows:

C1 = Profit = Down = 2 Samples

Therefore,  $P(C1) = 2/4$  and  $P(C2) = 2/4$

$$\therefore H(D1) = \sum_{i=1}^3 p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$\therefore H(D1) = \frac{2}{4} \log_2 \frac{1}{2} + \frac{2}{4} \log_2 \frac{1}{2} = 0.5 + 0.5 = 1$$

(i) Choosing Competition as the Splitting Attribute.

Competition	C1	C2	Entropy
Yes	2	0	0
No	0	2	0

$$\therefore H(\text{Competition}) = \frac{2}{4} \times H(\text{Yes}) + \frac{2}{4} \times H(\text{No})$$

$$= \frac{2}{4} \times 0 + \frac{2}{4} \times 0 = 0$$

∴ H(Competition) =  $H(D1) - H(\text{Competition})$

$$= 1 - 0 = 1$$

(ii) Choosing Type as the Splitting Attribute.

Type	C1	C2	Entropy
Software	1	1	1
Hardware	1	1	1

$$\therefore H(\text{Type}) = \frac{2}{4} \times H(\text{Software}) + \frac{2}{4} \times H(\text{Hardware})$$

$$= \frac{2}{4} \times 1 + \frac{2}{4} \times 1 = 1$$

∴ H(Type) =  $H(D1) - H(\text{Type})$

$$= 1 - 1 = 0$$

(iii) Choosing Age as the Splitting Attribute.

Age	C1	C2	Entropy
Old	3	0	0
Mid	2	2	1
New	0	3	0

$$\therefore H(\text{Age}) = \frac{3}{10} \times H(\text{Old}) + \frac{4}{10} \times H(\text{Mid}) + \frac{3}{10} \times H(\text{New})$$

$$= \frac{3}{10} \times 0 + \frac{4}{10} \times 1 + \frac{3}{10} \times 0 = 0.4$$

∴ H(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

∴ Gain(Age) =  $H(D1) - H(\text{Age}) = 1 - 0.4 = 0.6$

Sohn :  
Let the class label attributes be as follows:

C1 = Open House = Yes = 7 Samples  
C2 = Own House = Rented = 5 Samples

Therefore,  $P(C1) = 7/12$  and  $P(C2) = 5/12$

~~Data Mining & Business Intelligence (MU-Sem 3-I)~~  
~~Gain(Age) =  $H(D) - H(Age)$~~   
~~= 0.930 - 0.944 = 0.076~~

~~(ciii) Fig. P.3.3.3(d)~~

~~Ans. 12~~

~~Fig. 3.3.3~~

~~Ans. 12~~

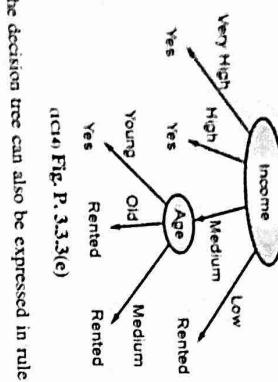


Fig. P. 3.3.3(e)

The decision tree can also be expressed in rule format

- IF Income = Very High THEN Own House = Yes
- IF Income = High THEN Own House = Yes
- IF Income = Low THEN Own House = Rented
- IF Income = Medium AND Age

- = Young THEN Own House = Yes
- IF Income = Medium AND Age
- = Medium THEN Own House
- = Rented
- IF Income = Medium AND Age

- = Old THEN Own House = Rented

**Ex. 3.3.4 :** The following table consists of training data from an employee database. The data have been generalized. For example, "31 ... 35" for age represents the age range of 31 to 35. For a given row entry, count represents the number of data tuples having the values for department, status, age, and salary given in that row.

department	status	Age	Salary	count
Sales	senior	31...35	46K...50K	30
Sales	junior	26...30	26K...30K	40
Sales	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

$$= \frac{110}{165} \times 0.8454 + \frac{31}{165} \times 0.8238 + \frac{14}{165} \times 0.8631 + \frac{10}{165} \times 0.9709 = 0.8504$$

$$\text{Gain(department)} = H(D) - H(\text{department})$$

$$= 0.899 - 0.8504 = 0.0486$$

Soln. :

- (a) The basic decision tree algorithm 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.

- Use ID3 algorithm to construct a decision tree from the given data.

Let the class label attributes be as follows:

- C1 = status = junior = 113 Samples

- C2 = status = senior = 52 Samples

Therefore,  $P(C1) = 113/165$  and  $P(C2) = 52/165$

$$H(D) = \sum_{i=1}^s p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$\therefore H(D) = \frac{113}{165} \log_2 \frac{165}{113} + \frac{52}{165} \log_2 \frac{165}{52}$$

$$= 0.3740 + 0.5250 = 0.899$$

(i) Choosing department as the Splitting Attribute.

Let the class label attributes be as follows:

- C1 = Status=junior

- C2 = Status=senior

Entropy before split for the given database D:

$$H(D) = \sum_{i=1}^s p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$= \frac{20}{165} \times 0 + \frac{49}{165} \times 0 + \frac{79}{165} \times H(31...35) + \frac{10}{165} \times H(36...40)$$

$$+ \frac{3}{165} \times H(41...45) + \frac{4}{165} \times H(46...50)$$

$$+ \frac{20}{165} \times 0 + \frac{49}{165} \times 0 + \frac{79}{165} \times 0.9906 + \frac{10}{165} \times 0.4743$$

Gain(age) =  $H(D) - H(\text{age})$

$$= 0.899 - 0.4743 = 0.4247$$

$$\therefore H(\text{age}) = \frac{110}{165} \times H(\text{Sales}) + \frac{31}{165} \times H(\text{systems}) + \frac{14}{165} \times H(\text{marketing}) + \frac{10}{165} \times H(\text{secretary})$$

(ii) Choosing department at the Splitting Attribute.

Choosing salary as the Splitting Attribute.

$$H(D) = \sum_{i=1}^s p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$= \frac{80}{165} \log_2 \frac{165}{80} + \frac{8}{165} \log_2 \frac{165}{8} + \frac{4}{165} \log_2 \frac{165}{4} + \frac{6}{165} \log_2 \frac{165}{6}$$

$$= 0.8454 + 0.8238 + 0.9709 + 0.4743$$

$$\therefore H(\text{department}) = \frac{110}{165} \times H(\text{Sales}) + \frac{31}{165} \times H(\text{systems}) + \frac{14}{165} \times H(\text{marketing}) + \frac{10}{165} \times H(\text{secretary})$$

(iv) Choosing salary as the Splitting Attribute.

Choosing salary as the Splitting Attribute.

$$H(D) = \sum_{i=1}^s p_i \log_2 \left( \frac{1}{p_i} \right)$$

$$= \frac{80}{165} \log_2 \frac{165}{80} + \frac{8}{165} \log_2 \frac{165}{8} + \frac{4}{165} \log_2 \frac{165}{4} + \frac{6}{165} \log_2 \frac{165}{6}$$

$$= 0.8454 + 0.8238 + 0.9709 + 0.4743$$

$$\therefore H(\text{salary}) = \frac{110}{165} \times H(\text{Sales}) + \frac{31}{165} \times H(\text{systems}) + \frac{14}{165} \times H(\text{marketing}) + \frac{10}{165} \times H(\text{secretary})$$

Summary

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375

Gain(department) = 0.0486

Gain(age) = 0.4247

Gain(salary) = 0.5375



### Data Mining & Business Intelligence (MU-Sem 6-IT)

(Classification)...Page no. (3-8)

- Gini index :** Gini index is a metric for classification tasks in CART. It stores sum of squared probabilities of each class. We can formulate it as illustrated below.

$$\text{Gini (Attribute = value)} = 1 - \sum (P_i)^2 \quad \text{for } i=1 \text{ to number of classes}$$

$\text{Gini (Attribute)} = \sum_{\text{values}} P_i \times \text{Gini (Attribute = value)}$

**Note : Select the attribute as the node of decision tree whose Gini Index is Low**

- Advantages of CART**
  - Decision trees can inherently perform multiclass classification.
  - They provide most model interpretability because they are simply series of if-else conditions.
  - They can handle both numerical and categorical data.
  - Nonlinear relationships among features do not affect the performance of the decision trees.

#### Disadvantages of CART

- A small change in the dataset can make the tree structure unstable which can cause variance.
- Decision tree learners create underfit trees if some classes are imbalanced. It is therefore recommended to balance the data set prior to fitting with the decision tree.

**Ex. 3.3.5 : Apply CART algorithm on the following training dataset and construct the decision tree.**

Day	Outlook	Temperature	Humidity	Wind	Play_Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	No
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	Weak	No	
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

**Soln. :** There are 14 instances of tennis playing factors.

**Outlook**

Outlook is a nominal feature. It can be sunny, overcast or rain. We will summarize the final decisions for outlook feature.

Outlook	Yes	No	Number of Instances
Sunny	2	3	5
Overcast	4	0	4
Rain	3	2	5

Humidity	Yes	No	Number of Instances
High	3	4	7
Normal	6	1	7

**Wind**

Wind	Yes	No	Number of Instances
Weak	6	2	8
Strong	3	3	6

Humidity	Yes	No	Number of Instances
High	0	2	2
Normal	1	1	2

**Temperature**

Then, we will calculate weighted sum of gini indexes for outlook feature.

$\text{Gini}(\text{Outlook}=\text{Rain}) = 1 - (3/5)^2 - (2/5)^2$

$= 1 - 0.36 - 0.16 = 0.48$

$\text{Gini}(\text{Outlook}=\text{Sunny}) = 1 - (2/5)^2 - (3/5)^2$

$= 1 - 0.16 - 0.36 = 0.48$

$\text{Gini}(\text{Outlook}=\text{Overcast}) = 1 - (4/4)^2 - (0/4)^2 = 0$

$= 1 - 1 - 0 = 0$

$\text{Gini}(\text{Humidity}=\text{Normal}) = 1 - (3/7)^2 - (4/7)^2$

$= 1 - 0.183 - 0.326 = 0.489$

$\text{Gini}(\text{Humidity}=\text{High}) = 1 - (3/7)^2 - (4/7)^2$

$= 1 - 0.183 - 0.326 = 0.489$

$\text{Gini}(\text{Humidity}=\text{Low}) = 1 - (6/7)^2 - (1/7)^2$

$= 1 - 0.734 - 0.02 = 0.244$

Weighted sum for humidity feature will be calculated next

$\text{Gini}(\text{Humidity}) = (7/14) \times 0.489 + (7/14) \times 0.244 = 0.367$

**Temperature**

Similarly, temperature is a nominal feature and it could have 3 different values : Cool, Hot and Mild. Let's summarize decisions for temperature feature.

We've calculated weighted sum of gini index for each feature.

The winner will be temperature feature because its cost is the lowest.

**Wind**

We've calculated weighted sum of gini index for each feature.

The winner will be wind feature because its cost is the lowest.

**Outlook**

We've calculated weighted sum of gini index for each feature.

The winner will be outlook feature because its cost is the lowest.

**Humidity**

We've calculated weighted sum of gini index for each feature.

The winner will be humidity feature because its cost is the lowest.

**Wind**

We've calculated weighted sum of gini index for each feature.

The winner will be wind feature because its cost is the lowest.

**Temperature**

We'll put outlook decision at the top of the tree.

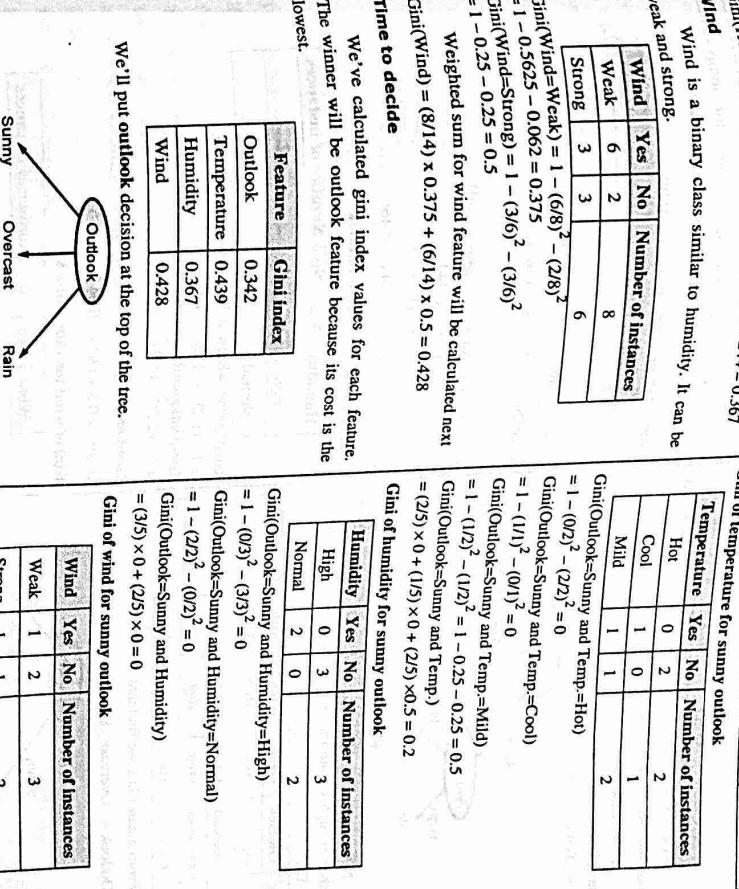


Fig. P.3.3.5(a)

Now, consider Outlook = Sunny and count the number of tuples from the original dataset D. Let us denote it as D1.

$$= 1 - (1/2)^2 - (1/2)^2 = 0.2$$

Humidity is a binary class feature. It can be high or normal.

$Gini(Outlook=Sunny \text{ and } Wind) = (3/5) \times 0.265 + (2/5) \times 0.2 = 0.445$

**Decision for sunny outlook**  
We've calculated gini index scores for feature when outlook is sunny. The winner is humidity because it has the lowest value.

Feature	Gini index
Temperature	0.2
Humidity	0
Wind	0.466

Since, Humidity has two possible values, the Humidity node has two branches (High, Normal).

From dataset D1, we find that when Humidity = High, Play\_Tennis = No and when Humidity = Normal,

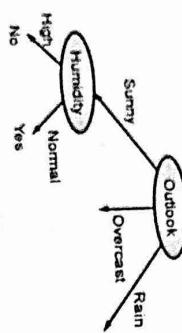


Fig. P. 3.3.5(b)

Now, consider Outlook = Overcast and count the number of tuples from the original dataset D. Let us denote it as D2.

Day	Outlook	Temp.	Humidity	Wind	Play_Tennis
3	Overcast	Hot	High	Weak	Yes
7	Overcast	Cool	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes

From dataset D2, we find that for all values of Outlook = "Overcast", Play\_Tennis = Yes.

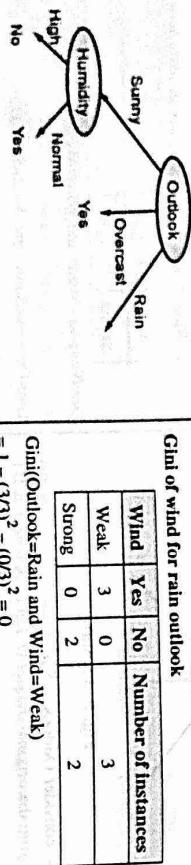


Fig. P. 3.3.5(c)

Now, consider Outlook = Rain and count the number of tuples from the original dataset D. Let us denote it as D3.

Temperature	Yes	No	Number of Instances
Cool	1	1	2
Mild	2	1	3

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Temp}=\text{Cool}) &= 1 - (1/2)^2 - (1/2)^2 = 0.5 \\ Gini(\text{Outlook}=\text{Rain and Temp}=\text{Mild}) &= 1 - (2/3)^2 - (1/3)^2 = 0.444 \\ Gini(\text{Outlook}=\text{Rain and Temp}=\text{Mild}) &= (2/5) \times 0.5 + (3/5) \times 0.444 = 0.466 \end{aligned}$$

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) &= 1 - (0/2)^2 - (2/2)^2 = 0 \\ Gini(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) &= 1 - (0/2)^2 - (2/2)^2 = 0 \end{aligned}$$

We'll calculate gini index scores for temperature, humidity and wind features when outlook is rain.

Now, consider Outlook = Rain and count the number of tuples from the original dataset D. Let us denote it as D3.

Humidity	Yes	No	Number of instances
High	1	1	2
Normal	2	1	3

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Humidity}=\text{High}) &= 1 - (1/2)^2 - (1/2)^2 = 0.5 \\ Gini(\text{Outlook}=\text{Rain and Humidity}=\text{Normal}) &= 1 - (1/3)^2 - (2/3)^2 = 0.444 \end{aligned}$$

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Humidity}=\text{High}) &= 1 - (2/3)^2 - (1/3)^2 = 0.444 \\ Gini(\text{Outlook}=\text{Rain and Humidity}=\text{Normal}) &= 1 - (2/3)^2 - (0/3)^2 = 0 \end{aligned}$$

Now, consider Outlook = Rain and count the number of tuples from the original dataset D. Let us denote it as D3.

Wind	Yes	No	Number of instances
Weak	3	0	3
Strong	0	2	2

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Wind}=\text{Weak}) &= 1 - (3/3)^2 - (0/3)^2 = 0 \\ Gini(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) &= 1 - (0/2)^2 - (2/2)^2 = 0 \end{aligned}$$

Now, consider Outlook = Rain and count the number of tuples from the original dataset D. Let us denote it as D3.

Humidity	Yes	No	Number of instances
High	1	1	2
Normal	2	1	3

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Humidity}=\text{High}) &= 1 - (1/2)^2 - (1/2)^2 = 0.5 \\ Gini(\text{Outlook}=\text{Rain and Humidity}=\text{Normal}) &= 1 - (1/3)^2 - (2/3)^2 = 0.444 \end{aligned}$$

$$\begin{aligned} Gini(\text{Outlook}=\text{Rain and Humidity}=\text{High}) &= 1 - (2/3)^2 - (1/3)^2 = 0.444 \\ Gini(\text{Outlook}=\text{Rain and Humidity}=\text{Normal}) &= 1 - (2/3)^2 - (0/3)^2 = 0 \end{aligned}$$

$Gini(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) = 1 - (0/2)^2 - (2/2)^2 = 0$

**Decision for rain outlook**  
The winner is wind feature for rain outlook because it has the minimum gini index score in features.

Feature	Gini index
Temperature	0.466
Humidity	0.466
Wind	0

Put the wind feature for rain outlook branch and monitor the new sub data sets.

Since, Wind has two possible values, the Wind node has two branches (Strong, Weak).

From dataset D3, we find that when Wind = Strong, Play\_Tennis = No and when Wind = Weak, Play\_Tennis = Yes.

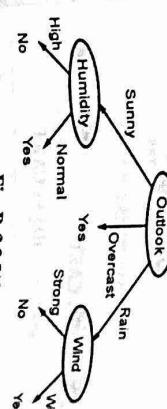


Fig. P. 3.3.5(d)

The decision tree can also be expressed in rule format as:

IF Outlook = Sunny AND Humidity = High  
THEN Play\_Tennis = No

IF Outlook = Sunny AND Humidity = Normal  
THEN Play\_Tennis = YES

IF Outlook = Overcast THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Strong  
THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Weak  
THEN Play\_Tennis = NO

THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Weak  
THEN Play\_Tennis = NO

THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Strong  
THEN Play\_Tennis = NO

THEN Play\_Tennis = YES

IF Outlook = Rain AND Wind = Weak  
THEN Play\_Tennis = NO

THEN Play\_Tennis = YES

$Gini(\text{Outlook}=\text{Rain and Wind}=\text{Strong}) = 1 - (0/2)^2 - (2/2)^2 = 0$

**Decision for rain outlook**  
The winner is wind feature for rain outlook because it has the minimum gini index score in features.

Tuples	Age	Salary	Jobs	Performance	Score
1	Young	High	Programmer	Average	No
2	Young	High	Programmer	Excellent	Yes
3	Middle-Age	High	Programmer	Average	No
4	Old	Medium	Programmer	Excellent	Yes
5	Old	Low	Government	Average	No
6	Old	Low	Government	Excellent	Yes
7	Middle-Age	Low	Government	Average	No
8	Young	Medium	Private	Average	No
9	Young	Low	Government	Average	Yes
10	Old	Medium	Government	Average	Yes
11	Young	Medium	Government	Excellent	Yes
12	Middle-Age	Medium	Private	Excellent	Yes
13	Middle-Age	High	Government	Average	Yes
14	Old	Medium	Private	Excellent	No

**Soln.:**  
There are 14 instances of selection decisions based on age, salary, job and performance factors.

Age is a nominal feature. It can be young, old or middle-aged. We will summarize the final decisions for age feature.

Age	Yes	No	Number of instances
Young	2	3	5
Middle-Age	4	0	4
Old	3	2	5

$$\begin{aligned} Gini(\text{Age}=\text{Young}) &= 1 - (2/5)^2 - (3/5)^2 \\ &= 1 - 0.16 - 0.16 = 0.48 \\ Gini(\text{Age}=\text{Middle-Age}) &= 1 - (4/4)^2 - (0/4)^2 = 0 \\ Gini(\text{Age}=\text{Old}) &= 1 - (3/5)^2 - (2/5)^2 \\ &= 1 - 0.36 - 0.36 = 0.48 \end{aligned}$$

Then, we will calculate weighted sum of gini indexes for age feature:

$$\begin{aligned} Gini(\text{Age}) &= (5/14) \times 0.48 + (4/14) \times 0 + (5/14) \times 0.48 \\ &= 0.171 + 0 + 0.171 = 0.342 \end{aligned}$$

### 3.4 NAIVE BAYESIAN CLASSIFICATION

Data Mining & Business Intelligence (MU-Sem 6-IT)

#### Salary

Similarly, salary is a nominal feature and it could have 3 different values: High, Medium and Low. Let's summarize decisions for salary feature.

Salary	Yes	No	Number of Instances
High	2	2	4
Low	3	1	4
Medium	4	2	6

$$\text{Gini}(\text{Salary}=\text{High}) = 1 - (2/4)^2 - (2/4)^2$$

$$\text{Gini}(\text{Salary}=\text{Low}) = 1 - (3/4)^2 - (1/4)^2$$

$$\text{Gini}(\text{Salary}=\text{Medium}) = 1 - (4/6)^2 - (2/6)^2$$

$$= 1 - 0.5625 - 0.0625 = 0.375$$

$$\text{Gini}(\text{Salary}=\text{Medium}) = 1 - (4/6)^2 - (2/6)^2$$

$$= 1 - 0.444 - 0.111 = 0.445$$

$$\text{Gini}(\text{Salary}) = (4/14) \times 0.5 + (4/14) \times 0.375 + (6/14) \times 0.445$$

$$= 0.142 + 0.107 + 0.190 = 0.439$$

We'll calculate weighted sum of gini index for salary feature.

$$\text{Gini}(\text{Salary}=\text{High}) = 1 - (2/4)^2 - (2/4)^2$$

$$\text{Gini}(\text{Salary}=\text{Low}) = 1 - (3/4)^2 - (1/4)^2$$

$$\text{Gini}(\text{Salary}=\text{Medium}) = 1 - (4/6)^2 - (2/6)^2$$

$$= 1 - 0.5625 - 0.0625 = 0.375$$

$$\text{Gini}(\text{Salary}=\text{Medium}) = 1 - (4/6)^2 - (2/6)^2$$

$$= 1 - 0.444 - 0.111 = 0.445$$

$$\text{Gini}(\text{Salary}) = (4/14) \times 0.5 + (4/14) \times 0.375 + (6/14) \times 0.445$$

$$= 0.142 + 0.107 + 0.190 = 0.439$$

Job

Humidity is a binary class feature. It can be private or government.

Job	Yes	No	Number of Instances
Private	3	4	7
Govt.	6	1	7

$$= 1 - 0.734 - 0.02 = 0.244$$

Weighted sum for job feature will be calculated next

$$\text{Gini}(\text{Job}=\text{Private}) = (7/14) \times 0.489 + (7/14) \times 0.244 = 0.367$$

Performance

Performance is a binary class similar to job. It can be average or excellent.

Performance	Yes	No	Number of Instances
Average	6	2	8
Excellent	3	3	6

$$\text{Gini}(\text{Performance}=\text{Average}) = 1 - (6/8)^2 - (2/8)^2$$

$$= 1 - 0.5625 - 0.0625 = 0.375$$

$$\text{Gini}(\text{Performance}=\text{Excellent}) = 1 - (3/6)^2 - (3/6)^2$$

$$= 1 - 0.25 - 0.25 = 0.5$$

Weighted sum for performance feature will be calculated next

$$\text{Gini}(\text{Performance}) = (8/14) \times 0.375 + (6/14) \times 0.5 = 0.428$$

$$\text{Gini}(\text{Performance}) = 0.428$$

**Time to decide**  
We've calculated gini index values for each feature. The winner will be age feature because its cost is the lowest.

Feature	Gini Index
Age	0.342
Salary	0.439
Job	0.367
Performance	0.428

We'll put age decision at the top of the tree. From the dataset it is clear that for age = middle-aged, select = yes.

```

graph TD
    Root[Age] --> Young[Young]
    Root --> Old[Old]
    Root --> Middle[Middle-aged]
    Young --> Yes[Yes]
    Old --> No[No]
    Middle --> Middle
    Middle --> Fair[Fair]
    Middle --> Excellent[Excellent]
  
```

Fig. P. 3.3.6

### 3.3.6 ID3 Vs CART

Table 3.3.1: ID3 Vs CART

ID	Age	Income	Student	Credit rating	Buys Computer
1	<=30	High	No	Fair	No
2	<=30	High	No	Excellent	No
3	<=30	High	No	Fair	Yes
4	>30	Medium	No	Fair	Yes
5	>30	Low	Yes	Fair	Yes
6	>30	Low	Yes	Excellent	No
7	>30	Low	Yes	Excellent	Yes
8	>30	Medium	No	Fair	Yes
9	>30	Low	Yes	Fair	Yes
10	>30	Medium	Yes	Fair	Yes
11	>30	Medium	Yes	Excellent	Yes
12	>30	Medium	No	Excellent	Yes
13	>30	High	Yes	Fair	Yes
14	>30	Medium	No	Excellent	No

### 3.4.2 Bayesian Interpretation

In the Bayesian interpretation, probability determines a "degree of belief."

Bayes theorem connects the degree of belief in a hypothesis before and after accounting for evidence.

For example, let's us consider an example of the coin. If we toss a coin, then we get either head or tail, and the percent of occurrence of either head and tail is 50%. If the coin is flipped numbers of times, and the outcomes are observed, the degree of belief may rise, fall, or remain the same depending on the outcomes.

For proposition X and evidence Y,

- (a)  $P(X)$ , the prior probability, is the primary degree of belief in X
- (b)  $P(X|Y)$ , the posterior probability, is the degree of belief having accounted for Y.
- (c) The quotient  $\frac{P(X|Y)}{P(Y)}$  represents the support Y provides for X.

Bayes theorem can be derived from the conditional probability :

**Soln.:**  
Class label attribute is Buys\_Computer.  
C1: Buys\_Computer = Yes = 9 samples  
C2: Buys\_Computer = No = 5 samples

$$\therefore P(C_1) = \frac{9}{14} \text{ and } P(C_2) = \frac{5}{14}$$

Let event X1 be Age = "<=30", event X2 be Income = "Medium", event X3 be Student = "Yes", and event X4 be Credit rating = "Fair"

$$P(X_1 \cap X_2) = \frac{P(X_1) \cdot P(X_2)}{P(X_1 \cup X_2)}, \text{ if } P(X_1 \cup X_2) \neq 0$$

$$P(X_1 \cap X_2) = \frac{P(X_1) \cdot P(X_2)}{P(X_1) + P(X_2)}, \text{ if } P(X_1 \cup X_2) \neq 0$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{Or, } P(X_1 \cap X_2) = \frac{P(X_1) \cdot P(X_2)}{P(X_1) + P(X_2) + P(X_3) + P(X_4)}$$

$$\text{being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P(X_1 \cap X_2) = P(X_1) \cap P(X_2)$$

$$\text{where } P(X_1 \cap X_2) \text{ is the joint probability of both X and Y being true, because } P$$





Homeowner	Marital Status	Job experience (in years)	Defaulted
No	Married	4	No
Yes	Divorced	2	No
No	Married	3	Yes
No	Married	3	No
Yes	Single	3	Yes

Soln. :

From the above table, it is clear that there are 3 tuples classified as

Defaulted = Yes and 7 tuples classified as Defaulted = No.

From the given training data, we estimate

$$P(\text{Defaulted} = \text{Yes}) = 3/10$$

$$P(\text{Defaulted} = \text{No}) = 2/10$$

Probabilities associated with attributes is given in table below.

Attribute	Value	Count		Probabilities	
		Defaulted = Yes	Defaulted = No	Defaulted = Yes	Defaulted = No
Homeowner	Yes	1	3	1/3	3/7
	No	2	4	2/3	4/7
	Total	3	7		
Marital Status	Single	1	2	1/3	2/7
	Married	1	4	1/3	4/7
	Total	3	7		
Job experience (in years)	2	2	1	2/3	1/7
	3	1	2	1/3	2/7
	4	0	3	0	3/7
Total	5	0	1	0	1/7
	Total	3	7		

The unknown tuple is  $t = \langle \text{Homeowner: No}, \text{Marital Status: Married}, \text{Job experience: 3} \rangle$

$$P(\text{Defaulted} = \text{Yes}) \times P(\text{Defaulted} = \text{Yes})$$

$$= P(\text{Homeowner: No}) \times P(\text{Marital Status: Married}) \times P(\text{Defaulted} = \text{Yes})$$

$$= (\text{Yes}) \times P(\text{Job experience: 3}) \times P(\text{Defaulted} = \text{Yes})$$

$$= \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} \times \frac{1}{3} = 0.011$$

$$P(\text{Defaulted} = \text{No}) \times P(\text{Defaulted} = \text{No})$$

$$= P(\text{Homeowner: No}) \times P(\text{Marital Status: Married}) \times P(\text{Defaulted} = \text{No})$$

$$= (\text{No}) \times P(\text{Defaulted} = \text{No})$$

$$= \frac{4}{7} \times \frac{4}{7} \times \frac{2}{7} \times \frac{7}{10} = 0.032$$

Based on the above probabilities, we classify the new tuple as Defaulted = No because it has the highest probability.

- Rule-based classifier makes use of a set of IF-THEN rules for classification.
- We can express a rule in the following form
- If condition THEN conclusion
- If we consider a rule R1,

R1: If age = youth AND student = yes THEN buy\_computer = yes

The IF part of the rule is called rule antecedent or precondition.

The THEN part of the rule is called rule consequent.

The antecedent part the condition consists of one or more attribute tests and these tests are logically ANDed.

The consequent part consists of class prediction.

To extract a rule from a decision tree –

1. One rule is created for each path from the root to the leaf node.

2. To form a rule antecedent, each splitting criterion is logically ANDed.

3. The leaf node holds the class prediction, forming the rule consequent.

### M 3.6 REGRESSION

UQ: Explain Regression. Explain linear regression with example. [MU - May 2019]

- Regression is a data mining technique used to predict a range of numeric values (also called continuous values), given a particular dataset. For example, regression might be used to predict the cost of a product or service, given other variables.
- Data smoothing can also be done using regression.

Regression can be classified as:

- Linear Regression :

It finds the best line to fit two variables or attributes so that one variable can be used to predict the other variable.

$$\text{Linear Regression : } y = \alpha + \beta x$$

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Here,  $y$  is the variable that is to be predicted,  $x$  is the variable used to predict  $y$ ,  $\alpha$  is the intercept and  $\beta$  is

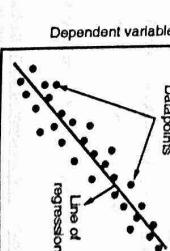


Fig. 3.6.1: Linear Regression Model

$$\beta = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

where,  $\bar{x}$  and  $\bar{y}$  are the mean of variables  $x$  and  $y$  respectively.

Linear regression is one of the easiest and most popular Supervised Machine Learning algorithms.

It is a statistical method that is used for predictive analysis.

Linear regression makes predictions for continuous/real or numeric variables such as sales, salary, age, product price, etc.

Linear regression algorithm shows a linear relationship between a dependent ( $y$ ) and one or more independent ( $x$ ) variables, hence called as linear regression.

Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

The linear regression model provides a sloped straight line representing the relationship between the variables.



$$\bar{x} = \frac{1}{12} \sum_{i=1}^{12} x_i = 72.17, \bar{y} = \frac{1}{12} \sum_{i=1}^{12} y_i = 74$$

Midterm Exam (x)	Final Exam (y)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
72	84	-0.17	10	0.03	-1.67
50	63	-22.17	-11	491.36	243.83
81	77	8.83	3	78.03	26.50
74	78	1.83	4	3.36	7.33
94	90	21.83	16	476.69	349.33
86	75	13.83	1	191.36	13.83
59	49	-13.17	-25	173.36	329.17
83	79	10.83	5	117.36	54.17
65	77	-7.17	3	51.36	-21.50
33	52	-39.17	-22	1534.03	861.67
88	74	15.83	0	250.69	0.00
81	90	8.83	16	78.03	141.33
				$\Sigma = 3445.67$	$\Sigma = 2004$

$$\beta = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = 0.58$$

$$\alpha = \bar{y} - \beta \bar{x} = 74 - 0.58 \times 72.17 = 32.14$$

∴ The linear regression line becomes :  $y = 32.14 + 0.58 \times x$

The final exam grade of a student who received 86 on the midterm exam :

$$y = 32.14 + 0.58 \times 86 = 82.02 = 82 \text{ marks}$$

**Ex. 3.6.3 :** Say we're given data on student exam results and our goal is to predict whether a student will pass or fail based on number of hours slept and hours spent studying. We have two features (hours slept, hours studied) and two classes: passed (1) and failed (0).

Studied (x <sub>1</sub> )	Slept (x <sub>2</sub> )	Passed (y)
4.85	9.63	0
8.62	3.23	1
5.43	8.23	0
9.21	6.34	1

Let  $b_0 = -0.406$ ,  $b_1 = 0.8525$ ,  $b_2 = -1.105$ . Determine the predicted value of y using Logistic Regression.

Soln. :

**Logistic regression equation**

Linear regression  $y = b_0 + b_1 \times x_1 + b_2 \times x_2 + \dots + b_k \times x_k$

$$\text{Sigmoid Function : } P = \frac{1}{1 + e^{-y}}$$

$$\text{For 1<sup>st</sup> input tuple, } x_1 = 4.85 \text{ and } x_2 = 9.63,$$

$$y = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + b_2 x_2)}} = \frac{1}{1 + e^{-( -0.406 + 0.8525 \times 4.85 - 1.105 \times 9.63)}} = 0.0522$$

$$\text{For 2<sup>nd</sup> input tuple, } x_1 = 8.62 \text{ and } x_2 = 3.23,$$

$$y = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + b_2 x_2)}} = \frac{1}{1 + e^{-( -0.406 + 0.8525 \times 8.62 - 1.105 \times 3.23)}} = 0.5076$$

$$\text{For 3<sup>rd</sup> input tuple, } x_1 = 5.43 \text{ and } x_2 = 8.23,$$

$$y = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + b_2 x_2)}} = \frac{1}{1 + e^{-( -0.406 + 0.8525 \times 5.43 - 1.105 \times 8.23)}} = 0.5076$$

$$\text{For 4<sup>th</sup> input tuple, } x_1 = 9.21 \text{ and } x_2 = 6.34,$$

$$y = \frac{1}{1 + e^{-(b_0 + b_1 x_1 + b_2 x_2)}} = \frac{1}{1 + e^{-( -0.406 + 0.8525 \times 9.21 - 1.105 \times 6.34)}} = 0.6082$$

In Logistic Regression,

$P \geq 0.5$ , class = 1

$P < 0.5$ , class = 0

The predicted and actual value are shown in table below:

Studied (x <sub>1</sub> )	Slept (x <sub>2</sub> )	Passed (y)	Predicted Value	Prediction
4.85	9.63	0	0.0022	0
8.62	3.23	1	0.5076	1
5.43	8.23	0	0.5076	1
9.21	6.34	1	0.6082	1

### 3.6.1 Linear Regression Vs Logistic Regression

Sr. No.	Linear Regression	Logistic Regression
1	Linear regression is used to predict the continuous dependent variable using a given set of independent variables.	Logistic Regression is used to predict the categorical dependent variable using a given set of independent variables.
2	Linear Regression is used for solving Regression problem.	Logistic regression is used for solving Classification problems.
3	In Linear regression, we predict the value of continuous variables.	In logistic Regression, we predict the values of categorical variables.
4	In linear regression, we find the best fit line, by which we can easily predict the output.	In Logistic Regression, we find the S-curve by which we can classify the samples.
5	Least square estimation method is used for estimation of accuracy.	Maximum likelihood estimation method is used for estimation of accuracy.
6	The output for Linear Regression must be a continuous value, such as price, age, etc.	The output of Logistic Regression must be a value such as 0 or 1, Yes or No, etc.
7	In Linear regression, it is required that relationship between dependent variable and independent variable must be linear.	In Logistic regression, it is not required to have the linear relationship between the dependent and independent variable.
8	In linear regression, there may be collinearity between the independent variables.	In logistic regression, there should not be collinearity between the independent variable.



- The data is trained on the training set and the mean square error (MSE) is obtained from the predictions on the test set.
- This method is not recommended because the MSE would depend on the split. So a new split can give you a new MSE and then you don't know which to trust.
- The overall accuracy is calculated by taking the average of the accuracies obtained from each iteration.

$$E = \frac{1}{k} \sum_{i=1}^k E_i$$

Total number of examples

example

- FIG. 3.8.2 : Random Subsampling**
- Experiment 1
  - Experiment 2
  - Experiment 3

- 3.8.3 Cross Validation**
- When only a limited amount of data is available, to achieve an unbiased estimate of the model performance we use k-fold cross-validation.

- In k-fold cross-validation, we divide the data into k subsets of equal size.
- We build models k times, each time leaving out one of the subsets from training and use it as the test set.

- If k equals the sample size, this is called "leave-one-out".

#### 3.8.4 Bootstrapping

- Bootstrapping is a technique used to make estimations from data by taking an average of the estimates from smaller data samples.
- The bootstrap method involves iteratively resampling a dataset (with replacement). Instead of only estimating our statistic once on the complete data, we can do it many times on a re-sampling (with replacement) of the original sample. Repeating this re-sampling multiple times allows us to obtain a vector of estimates.
- We can then compute variance, expected value, empirical distribution, and other relevant statistics of these estimates.

- Ensemble methods are techniques that aim at improving the accuracy of results in models by combining multiple models instead of using a single model.
- The combined models increase the accuracy of the results significantly. This has boosted the popularity of ensemble methods in machine learning.

- Ensemble methods are ideal for regression and classification, where they reduce bias and variance to boost the accuracy of models.
- using majority voting to find the final prediction.

- AdaBoost** is a more general version of the Boosting algorithm.

- Random Forest**
- Random Forest is proposed particularly for trees.
- It is the combination of bagging and random subspace method for inducing the tree.
- It is similar to bagging except that each model is a random tree rather than a single model and each tree is grown according to the bootstrap sample of the training set to N.
- Another random step is used to split each node.
- A small subset of feature m is selected randomly ( $m < M$ ) rather than considers all possible splits M, and the best split is chosen from this subset.
- Final classification is done by majority vote across trees.

- In AdaBoost, the algorithm only makes a node with two leaves, and this is known as Stump. These stumps are weak learners, and boosting is important in AdaBoost. The order of stumps is very important in AdaBoost. The error of the first stump influences how the other stump is made. So, in AdaBoost algorithm also, the majority of votes take place between the stumps, the same as in random trees. Final classification is done by majority vote across trees.

#### 3.10 NEAREST NEIGHBOR CLASSIFIER

- K-Nearest Neighbor** is one of the simplest algorithms based on Supervised Learning technique.
- K-NN algorithm assumes the similarity between the new case and available cases and put the new case into the category that is most similar to the available categories.
- K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suited category by using K-NN algorithm.
- Except for the first, each subsequent learner is grown from previously grown learners. In simple words, weak learners are converted into strong ones.
- AdaBoost algorithm also works on the same principle as boosting, but there is a slight difference in working.