

Classification

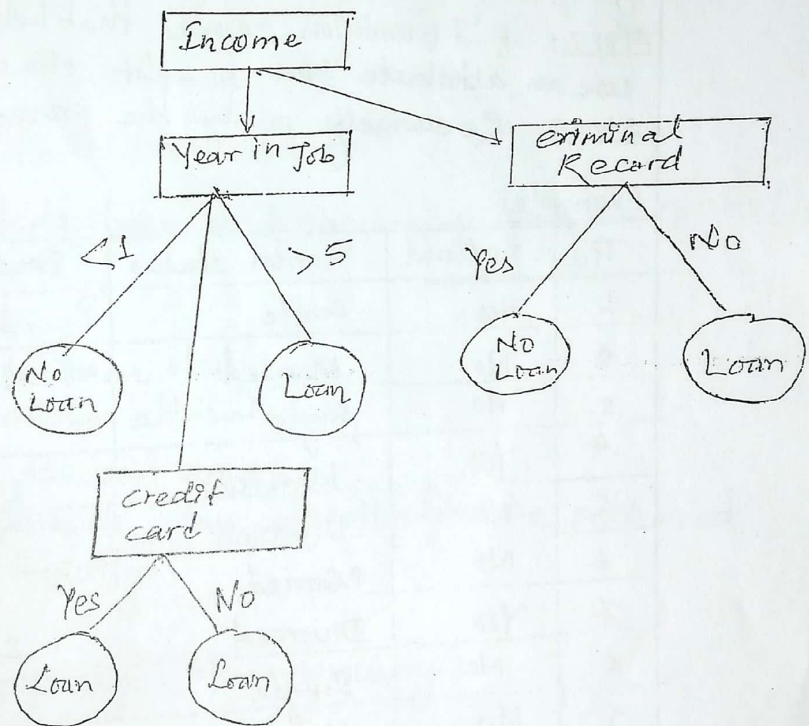
Grouping(1) Decision Tree:

→ A decision tree is a structure in which each branch represents a choice between a number of alternatives and each leaf node represents a classification or decision.

→ Decision tree is a classifier in the form of a tree structure where each node is either a leaf node, indicating a classification for instances or a decision node that represents some test to be carried out on a single attribute value with one branch or subtree for each possible outcome of the test.

→ It can be used to classify an instance by starting at the root of the tree and moving through it until a leaf node, which provides the classification of the instance.

Example:-

Advantages:-

- Extremely fast for classifying unknown records
- Easy to interpret
- High accuracy for simple data.

Some Common Decision Tree Algorithms Are:

- Hunt Algorithm
- CART
- ID3, J48, C45
- SLIQ, SPRLT

Hunt Algorithm:

Let D_t be the set of training records that reach at node T .

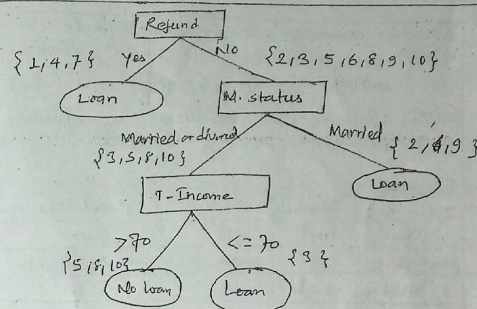
Algorithm:

- Step 1: If D_t contains records that belongs to the same class then the node T is a leaf node
- Step 2: If D_t is an empty set then T is a leaf node labelled by default class
- Step 3: If D_t contains records that belongs to more than one class use an attribute test to split the data into smaller subsets
- Step 4: Recursively apply the procedure to each subset.

Example:

Ti	Refund	Marital status	Taxable income	Class
1	Yes	Single	125K	Loan
2	No	Married	100K	Loan
3	No	Single	70K	Loan
4	Yes	Married	120K	Loan
5	No	Divorced	95K	no loan
6	No	Married	60K	Loan
7	Yes	Divorced	220K	Loan
8	No	Single	85K	No Loan
9	No	Married	75K	Loan
10	No	Single	90K	No Loan

Java Programs for Implementing Some of the Decision Tree Algorithms



Tree Induction:-

Greedy Strategy:-

→ split the record based on an attribute test that optimizes certain.

Issues

- (i) How to split record?
- (ii) How to specify attribute test condition?

→ depends on attribute type and number of ways to split (i.e. multiple way split)

(iii) when to stop splitting?

- All records belong to same class.
- All records have similar attributes.

(iv) How to determine the best split?

- Nodes with homogeneous class distribution are preferred
- Measure of node impurity

Measure of Node Impurity (Homogeneity)

- Gini Index
- Entropy calculation.

Gini Index:

the Gini Index measures the impurity of dataset (D) as,

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

where $p_i \rightarrow$ is the probability that a tuple in D belongs to class C_i .

→ Consider a binary split for each attribute
 → when D is partitioned into D_1 and D_2 then,

$Gini(D) = \frac{D_1}{D} Gini(D_1) + \frac{D_2}{D} Gini(D_2)$
 → The attribute that minimizes the reduction in impurity is selected as the splitting attribute.

Example:-

ID	Age	Income	student	credit-rating	class (buy-no)
1	Youth	high	No	fair	No
2	Youth	high	No	fair	No
3	middle age	high	No	Excellent	No
4	senior	medium	No	fair	No
5	senior	low	No	fair	Yes
6	senior	low	Yes	fair	Yes
7	M-age	low	Yes	Excellent	No
8	Youth	medium	Yes	Excellent	Yes
9	Youth	low	Yes	fair	No
10	senior	medium	Yes	fair	Yes
11	Youth	medium	Yes	fair	Yes
12	M-age	medium	Yes	Excellent	Yes
13	M-age	high	No	Excellent	Yes
14	senior	medium	No	fair	Yes

Solution:-

$$(buys-computer = Yes) = 9$$

$$(buys-computer = No) = 5$$

Impurity D , $Gini(D) = 1 - \left(\frac{9}{14}\right)^2 - \left(\frac{5}{14}\right)^2$
 Also calculating $Gini$ index for each attribute

Attribute: Income

Subset 1, $D_1 = \{low, medium\}$, $D_2 = \{high\}$
 $Gini(Income) = \frac{10}{14} Gini(D_1) + \frac{4}{14} Gini(D_2)$

$$= \frac{10}{14} \left(1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$$

[Here, 7 buys computer from set D_1 , & 3 don't buy computer]

$$= 0.556$$

subset 2: $D_1 = \{low, high\}$, $D_2 = \{medium\}$

$$Gini(Income) = \frac{8}{14} \left(1 - \left(\frac{5}{8}\right)^2 - \left(\frac{3}{8}\right)^2\right) + \frac{6}{14} \left(1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2\right)$$

$$= 0.4582$$

subset 3: $D_1 = \{medium, high\}$, $D_2 = \{low\}$

$$Gini(Income) = \frac{4}{14} \left(1 - \left(\frac{6}{4}\right)^2 - \left(\frac{4}{4}\right)^2\right) + \frac{10}{14} \left(1 - \left(\frac{1}{10}\right)^2 - \left(\frac{9}{10}\right)^2\right)$$

$$= 0.450$$

Since subset 3 has minimum value so selected as binary split for income.

Attribute: Age

subset 1: $D_1 = \{youth, senior\}$, $D_2 = \{M-age\}$

$$Gini(Age) = \frac{10}{14} \left(1 - \left(\frac{5}{10}\right)^2 - \left(\frac{5}{10}\right)^2\right) + \frac{4}{14} \left(1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2\right)$$

$$= 0.357$$

subset 2: $D_1 = \{youth, M-age\}$, $D_2 = \{senior\}$

$$Gini(Age) = \frac{9}{14} \left(1 - \left(\frac{3}{9}\right)^2 - \left(\frac{6}{9}\right)^2\right) + \frac{5}{14} \left(1 - \left(\frac{1}{5}\right)^2 - \left(\frac{4}{5}\right)^2\right)$$

$$= 0.3678$$

subset 3: $D_1 = \{youth\}$, $D_2 = \{M-age, senior\}$

$$Gini(Age) = \frac{5}{14} \left(1 - \left(\frac{5}{5}\right)^2 - \left(\frac{0}{5}\right)^2\right) + \frac{9}{14} \left(1 - \left(\frac{2}{9}\right)^2 - \left(\frac{7}{9}\right)^2\right)$$

$$= 0.3936$$

Attribute: student

$D_1 = \{Yes\}$, $D_2 = \{No\}$

$$Gini(student) = \frac{7}{14} \left(1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2\right) + \frac{7}{14} \left(1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2\right)$$

$$= 0.3775$$

Attribute: credit-rating

$D_1 = \{fair\}$, $D_2 = \{Excellent\}$

$$Gini(credit-rating) = \frac{8}{14} \left(1 - \left(\frac{3}{8}\right)^2 - \left(\frac{5}{8}\right)^2\right) + \frac{6}{14} \left(1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2\right)$$

$$= 0.4285$$

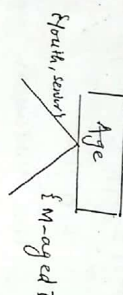
Impurity reduction for age = $0.459 - 0.357 = 0.102$

Impurity reduction for income = $0.459 - 0.450 = 0.009$

Impurity reduction for student = $0.459 - 0.3775 = 0.0815$

Impurity reduction for credit-rating = $0.459 - 0.4285 = 0.0305$

Since, Age has maximum impurity reduction it is selected as root node.



Nearest Neighbour Classifier:

→ It uses K-closest points for performing classification. K nearest of record 'X' are data points that have the K-smallest to X.

→ Classification based on learning by analogy i.e. by comparing given test tuple with training tuple that are similar.

→ Training tuples are described by an attribute.

→ When given an unknown tuple, a nearest neighbour classifier searches the pattern space for K-training tuples - are closest to the unknown tuple.

→ Nearest neighbour classifier requires:

(i) set of stored record.

(ii) distance matrix to compute the distance between records.

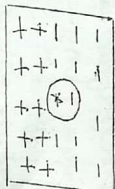
(iii) The value of K i.e. the number of nearest neighbours to

→ To classify an unknown records:

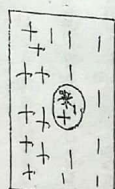
(i) Compute the distance to other training record

(ii) Identify the K-nearest neighbour

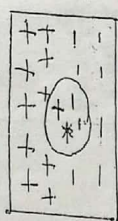
(iii) Use class labels of nearest neighbour to determine class of unknown record by using majority vote.



K = 1
Class = -



K = 2
Class = -/+



K = 3
Class = -

Classification Issues:

(a) Choosing the value of K:

If K is too small it is sensitive to noise point. If K is too large neighbours may include points from other class.

(b) Scaling Issue:

Attributes may have to be scaled to prevent distance measure from being dominated by one of attributes (eg. Height, weight).

(c) Distance computing for non-numeric data:
Use distance as 0 (minimum) for same data 1 (maximum) for different data set.

③ missing values:
use minimum possible distance.

Disadvantages:

- poor accuracy when data have noise and irrelevant attributes.
- classifying unknown records are relatively expensive.
- slow when classifying test tuple.

Rule based classifier:

It classifies record by using a collection of 'if...then...' rule. i.e. Rule based classifier uses a set of 'if...then' rules of classification, the 'if' part or left hand side of rule is known as the rule antecedent where as the 'then' part or right hand side of rule is known as rule consequent. In rule antecedent the condition consists of one or more attribute test.

Eg. $R_1: (Age = Youth) \wedge (Student = Yes) \Rightarrow (buy_computer = Yes)$

If the condition in the rule antecedent holds true for a given tuple the rule antecedent is satisfied and the rule covers the tuple i.e. coverage of rule. If the rule fraction of record that satisfy the antecedent of rule,

$$Coverage = \frac{n \text{ covers}}{n \text{ total dataset}} \times 100\%$$

The accuracy of rule is the fraction of record that satisfy both the antecedent and consequent of rule. i.e.

$$Accuracy = \frac{n \text{ correct}}{n \text{ covers}} \times 100\%$$

How does rule based classifier work?

1. If a rule is satisfied by a tuple the rule is said to be trigger.
2. If only one rule is satisfied, then the rule fire by returning the class prediction for the tuple.
3. Triggering doesn't always means firing because there may be more than one rule that can be satisfied.
4. If more than one rule is satisfied then the rule fires by returning the class by applying conflict resolution strategy to find which rule is fired.

Conflict Resolution Strategy:

- (1) when more than one rule is triggered rule ordering or ranking is used. the rule ordering may be class based or rule based. when rule is used the rule set is given as a decision class.
- (2) when no rule is satisfied by tuple (unknown) then, a default rule can be setup to specify a default class.

Example:

Name	blood type	give birth	can fly	lives in water	class
lemon	warm	yes	no	no	mammal
Turtle	cold	no	no	sometimes	Reptile
shark	cold	yes	no	yes	Aquatic

Rules:

$R_1: (give \text{ birth} = yes) \wedge (blood \text{ type} = warm) \Rightarrow \text{mammals}$

$R_2: (give \text{ birth} = no) \wedge (can \text{ fly} = no) \Rightarrow \text{Reptile}$

$R_3: (lives \text{ in water} = sometimes) \Rightarrow \text{Amphibian}$

- A lemon triggers R_1 , so, it is a mammal
- A turtle triggers R_2 and R_3 i.e. conflict
- A shark triggers none of the rule, use default class
- If conflict, apply rule ranking.

Characteristics of Rule based Classifier.

1. Mutually Exclusive rule: (ordering)

Classifier contains mutually-exclusive rule. If rules are independent to each other.

(i) Every record is covered by almost one rule.

(ii) Rules are no longer mutually exclusive if record may trigger more than one rule.

(iii) To make mutually exclusive we apply ordering.

2. Exhaustive rule: (default)

- Classifier has exhaustive coverage if it accounts for every possible combinations of attribute values.
- Each record is covered by almost one rule
- Rules are no longer exhaustive if a record may not trigger any rule
- To make exhaustive we apply default class

Building Classification Rules:

1. Direct Method: Sequential & Inductive Approach

→ Extract rule direct from data.

eg. RIPPER, CN2

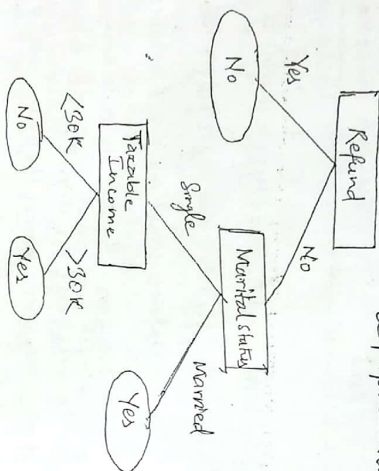
2. Indirect Method:

→ Extract rule from other classification.

eg. Decision Tree, ANN.

Rule extraction from decision tree.

leaf node always class



Rules:

R1: (Refund = Yes) = No

R2: (Refund = No) ∧ (Marital Status = Married) = Yes

Rule simplification:

R2 can be simplified as

(Marital Status = Married) = Yes.

Advantages of Rule-based classifier:

- Highly expressive
- Easy to generate and interpret
- can classify new instances rapidly.
- High performance.

Rules simplification

Posterior

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

← evidence

Bayesian Classifier:

- It is statistical classifier which predicts class membership probabilities for a membership class.
- It has high accuracy and speed for large databases.
- It has minimum error rate in comparison to other class.
- Based on Bayes Law.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Types:

- (1) Bayesian Belief Networks (Graphical Models)
 - It allows dependencies among subset of attribute i.e. has conditional dependencies among the attribute.
 - Has complex computational complexity.

(2) Naive Bayesian Classifier

- It assumes that the effect on an attribute value on a class is independent of the value of other attributes i.e. class is independent.
- It has simple computational complexity.
- Let 'x' be the training setup tuples and c_1, c_2, \dots, c_n their associated classes.

Given a tuple 'x', the classifier will predict that x to the class having highest posterior probability condition i.e. the naive Bayesian classifier predicts that the tuple belongs to the class c_i , only if

$$P(c_i/x) > P(c_j/x) \text{ for } 1 \leq i \leq m, j \neq i$$

$$\text{i.e. } P(c_i/x) = \frac{P(x/c_i) \cdot P(c_i)}{P(x)} \quad (\text{maximum})$$

Here, $P(x)$ = constant

$$P(c_i) = P(c_j) = \dots = P(c_m)$$

$$\text{So, we need to maximize } P(x/c_i)$$

Since, for naive assumption is class conditional independence

$$P(x/c_i) = \prod_{k=1}^n P(x_k/c_i)$$

$$= P(x_1/c_i) * P(x_2/c_i) * \dots * P(x_n/c_i)$$

These probability can be calculated from training tuple.

Example: 1

ID	Age	Income	Student	Credit Rating	Class buy computer
1	Youth	High	No	Fair	No
2	Youth	High	No	Excellent	No
3	Male	High	No	Fair	Yes
4	Senior	Medium	No	Fair	Yes
5	Senior	Low	Yes	Excellent	No
6	Senior	Low	Yes	Excellent	Yes
7	Male	Low	Yes	Fair	No
8	Youth	Medium	No	Fair	Yes
9	Youth	Low	Yes	Fair	Yes
10	Senior	Medium	Yes	Fair	Yes
11	Youth	Medium	Yes	Excellent	Yes
12	Male	Medium	No	Excellent	Yes
13	Male	High	Yes	Fair	Yes
14	Senior	Medium	No	Excellent	No

Find class of X:

$X = (\text{age} = \text{Youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit rating} = \text{fair})$

Solution:

$$\text{Let, } C_1 (\text{buys-computer} = \text{Yes}) = 9$$

$$C_2 (\text{buys-computer} = \text{No}) = 5$$

$$\text{Prior probability of } C_1 = 9/14 = 0.642$$

$$\text{Prior Probability of } C_2 = 5/14 = 0.357$$

$$P(\text{Age} = \text{Youth} / \text{buys-computer} = \text{Yes}) = 2/9 = 0.222$$

$$P(\text{Age} = \text{Youth} / \text{buys-computer} = \text{No}) = 3/5 = 0.6$$

$$P(\text{Income} = \text{medium} / C_1) = 4/9 = 0.444$$

$$P(\text{Income} = \text{medium} / C_2) = 2/5 = 0.4$$

$$P(\text{Student} = \text{Yes} / C_1) = 6/9 = 0.666$$

$$P(\text{Student} = \text{Yes} / C_2) = 3/5 = 0.6$$

$$P(\text{Credit-rating} = \text{Fair} / C_1) = 6/9 = 0.666$$

$$P(\text{Credit-rating} = \text{Fair} / C_2) = 2/5 = 0.4$$

$$P(X/C_1) = 0.222 \times 0.444 \times 0.666 \times 0.666 = 0.644$$

$$P(X/C_2) = 0.6 \times 0.4 \times 0.6 \times 0.4 = 0.096$$

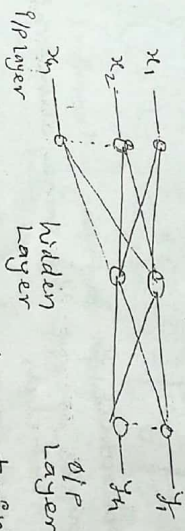
Since, $P(X/C_1) > P(X/C_2)$ buys-computer = yes.

X is classified as buys-computer = yes.

(Note: when count = 0 (Apply minimum probability) replace count with 1)

Artificial Neural Network Classifier (ANN)

- ANN is a set of connected 5/10 units in which each connection has a weight associated with it.
- During the learning phase, the network learns by adjusting the weights, so as to be able to predict the correct class label.
- ANN is also referred as connectionist learning due to connections between different units.
- It has long training time and poor interpretability but has tolerance to noisy data.
- It can classify pattern on which they haven't been trained.
- It is well suited for continuous value input. It has parallel processing.



for any classification task, before training, the topology (structure) of the network, number of input units, number of hidden units, number of output units, learning rate and bias value.

Back Propagation Algorithm:-

1. Initialization:
 - set all the weights and threshold levels inside a small range.
2. Activation:
 - Activate the network by applying the inputs and desired outputs.
 - calculate the actual output of the neurons (nodes) in the hidden layer and hence the output layers.
3. Weight Training:
 - update weight in the network by propagating the errors associated with the output neurons by calculating error gradients and

hence to the hidden layers
4. Repeat step 2 and 3 until selected error gradient is

Measure of Node Impurity:

- Gini Index

Entropy calculation

IDS Algorithm (Iterative Dichotomiser 3)

IDS is an algorithm used in decision tree to generate a decision tree. IDS builds the tree from the top down with no back tracking
→ Information gained is used to select the most useful attribute for classification
→ Information gained is calculated from entropy calculation

$$\text{Entropy}(H) = - \sum_{i=1}^m P(x_i) \cdot \log_2 P(x_i)$$

where, m = no of classes

$$P(x_i) = \text{Probability of 'x' in class 'i'}$$

Algorithm:

Step 1: Create a root node for a tree

Step 2: If all examples are of the same class return a single node with that class

Step 3: If all examples are not of same class, calculate entropy and information gain to select root node and branch into nodes.

Step 4: Partition the examples into subsets.

Step 5: Repeat the process until all examples are classified.

Example:

Person	Hair length	Weight	Age	Class
P ₁	0"	250	56	M
P ₂	10"	150	34	F
P ₃	2"	90	10	M
P ₄	6"	18	8	F
P ₅	4"	20	1	F
P ₆	1"	170	70	M
P ₇	8"	160	41	F
P ₈	10"	180	58	M
P ₉	6"	100	45	M

$$\text{Hans entropy of data (H)} = - \sum_{i=1}^m P(x_i) \log_2 P(x_i)$$

$$= - \left[\frac{4}{9} \log_2 \left(\frac{4}{9} \right) + \frac{5}{9} \log_2 \left(\frac{5}{9} \right) \right] = 0.9911$$

[5-4-

Entropy for hair length:

$$\left(\leq 5'' \right) = 4 \times \frac{1}{5} = - \left[\frac{4}{5} \log_2 \left(\frac{4}{5} \right) + \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right] = 0.8113$$

$$\left(\geq 5'' \right) = 5 \times \frac{5}{9} = - \left[\frac{5}{9} \log_2 \left(\frac{5}{9} \right) + \frac{4}{9} \log_2 \left(\frac{4}{9} \right) \right] = 0.9710$$

$$\text{Information gain} = 0.9911 - \left(\frac{4}{9} \times 0.8113 + \frac{5}{9} \times 0.9710 \right)$$

$$= 0.03107$$

Entropy for weight:

$$\left(\geq 160 \right) = - \left[\frac{0}{4} \log_2 \left(\frac{0}{4} \right) + \frac{4}{4} \log_2 \left(\frac{4}{4} \right) \right] = 0$$

$$\left(\leq 160 \right) = - \left[\frac{4}{5} \log_2 \left(\frac{4}{5} \right) + \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right] = 0.7219$$

$$\text{Information gain} = 0.9911 - \left(\frac{4}{5} \times 0.7219 + \frac{1}{5} \times 0 \right) = 0.59$$

Entropy for age:

$$\left(\leq 40 \right) = - \left[\frac{3}{6} \log_2 \left(\frac{3}{6} \right) + \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right] = 1$$

$$\left(> 40 \right) = - \left[\frac{1}{3} \log_2 \left(\frac{1}{3} \right) + \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] = 0.9183$$

$$\text{Information gain} = 0.9911 - \left(\frac{3}{6} \times 1 + \frac{3}{6} \times 0.9183 \right) = 0.0163$$

to choose maximum information gain

Since, weight has maximum information gain so selected as root node.

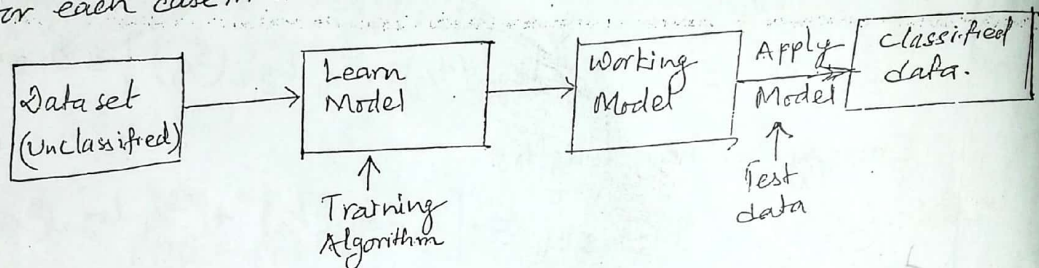
Advantages of ID3:

- Easy to construct and interpret for small sized data
- Higher accuracy for simple data.
- Fast for classifying unknown records

What is classification?

It is a data mining technique used to predict group-membership of data instances.

- classification assigns data in a collection of target category or class
- The goal of classification is to accurately predict the target class for each case in the data.



Model comparison:

- ✓ (1) Confusion Matrix (contingency Table)
- X (2) ROC