# Module 3
## CLASSIFICATION & PREDICTION

* There are two forms of data analysis that can be used for extracting models describing important classess or to predict future data trends
* These are the two forms: classification prediction
* classification models predict categorical class labels, and prediction models predict continuous valued function

• eg: we can build a classification model to categorize bank loan applications as either safe/risky. or a prediction model to predict the expenditures in dollers

## Classification vs prediction

### classification
• predicts categorical class labels[discrete or norminal]
• classifies data [constructs a model] based on the Training set and the values [class labels] in a classifying attribute and uses it in classifying new Data.

### prediction
models continuous-valued functions, ie predicts unknown or missing value.

## Typical Application
• Credit Approval

- Target marketing
- medical diagonosis
- Fraud detection

## what is classification

Following are the examples of cases where the Data analysis task is classification :-

- A bank loan officer wants to analyze the data in order to know which customor (loan application) are risky or which are safe
- A marketing manager at a company needs to analyze a customor with a given profile, who will buy a new Computer.

In both of the above examples, a model or classifier is Constructed to predict the categorical label. These labels are risky or safe for loan application Data and yes or no form marketing data.

→ How does classification works ?

The data classification process includes 2 type

i) Building The classifier or model
ii) Using Classifier For classification.

i, Building the classifier or model

The 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 model 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

## Supervised learning Vs unsupervised learning

### Supervised learning [classification]

- Supervised supervision:- The training data [observations, measurements etc] are accompanied by labels indicating the class of the observations.

- New data is classified based on the training set

### Unsupervised learning (clustering)

- These class labels of training data [observations] is unknown.

   a set of measurement, observations. ects with the aim of establishing the existance of classes or clusters in the data.

## Issues :- Data preparation

→ Preparing the data for classification & prediction
Data cleaning
   • preprocess the data in order to reduce noise & handle missing values
      • Relevance analysis [Feature selection]

- ~~Remove analysis~~
- Remove the irrelevant or redundant attrib
→ Data transformation & reduction
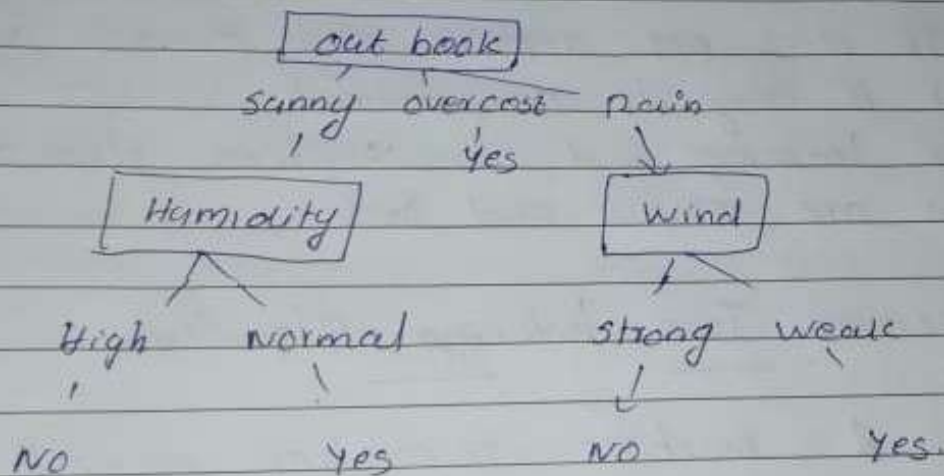    - Generalize and / or normalize data.

## Issues: Evaluating classification methods.

- Accuracy:
    - classifier accuracy: predicting class label
    - Predictor accuracy: guessing value of predict
      attributes.

- speed:
    - time to construct the model [training time]
    - time to use the model [classification /prediction
      time]

- ~~&~~ Robustenss: handling noise & missing values

- scalability: efficiency in disk. ~~resd~~ resident
                      data bases

- ~~Interpor~~ Interpretability
    - Understanding & insight provided by
      The model

- Other measure
    - eg: goodness of rules. such as decision
          tree size or compactness of
          classification rules.

# Decision Trees Induction

**Decision tree:-** Decision tree is the most powerful and popular tool for classification and prediction. A Decision tree is a flowchart like tree structure, where each Internal node denotes a test on an attribute, each branch represents the outcome of the test and each leaf node [terminal node] holds a class labels.
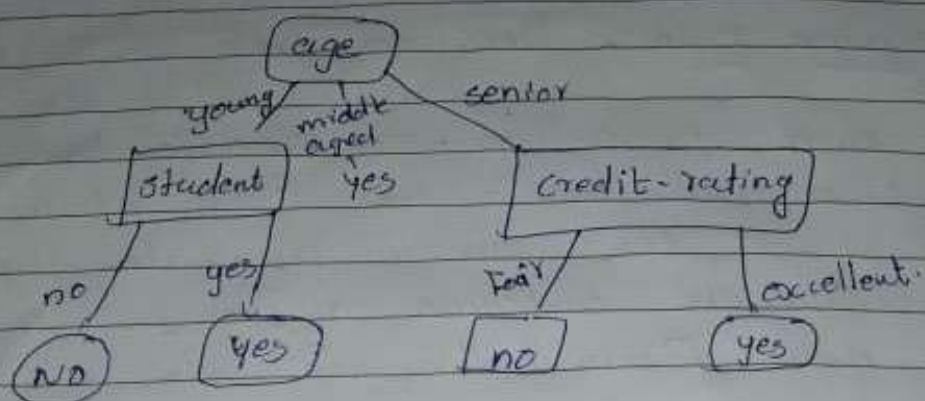
## Decision Tree for [play tennis]

```
                [out book]
           sunny  overcast  rain
             /       |        \
                    Yes
      [Humidity]              [Wind]
        / \                    / \
     High  normal          strong  weak
       |      \              |       \
      NO      Yes           NO       Yes.
```

A decision tree is a structure that Includes root node branches, and leaf nodes. Each Internal node denotes a test on attribute, each branch denotes the outcome of a test, and each leaf node holds a class labels. The topmost node In the tree is the root node

The following decision tree is for the concept buy computer that indicates whether a customer or not Each Internal node represents a test on an

attribute. Each leaf node represents a class



The benefits of having a decision tree are as follows:-

• It does not require any domain knowledge
• It is easy to comprehened
• The learning and classification steps of a decision tree are simple and fast

Decision Tree Inducing Algorithm.

A machine researcher named J. Ross Quinlan in 1980 developed a decision tree algorithm known as ID3 [Iterative Dichotomiser]. Later, he presented C45 adopt a greedy approach approach. In this algorithm, there is no back tracking, the trees are constructed in a top-down recusive divide-and-conquer manner

Generating a decision tree from training tuples of data partition D

Algorithms : General - decision tree

Input :

Data partition. D, which is a set of training
tuples and their associated class labels.
Attribute - list, the set of Candidate attributes.
Attribute selection method, a procedure to determine
The splitting criterion that best partitions that
The data tuples into individual classes. This crited
Criterion includes a splitting - attribute and.
either a splitting point or splitting subset

output :

A decision tree

method

Create a node N;

IF tuples in D are all of the same class, c
then return N as leaf node labeled with class.
c,

IF attribute list is empty then return N as leaf
with labeled with majority class in D, II major
majority voting.

apply attribute ,selection -method [D, attribute - list)
to find the best splitting criterions
label node N with splitting - criterion.

If splitting - attribute is discrete valued and multiway splitting allowed then all no restricted to binary trees.

Attribute list - spliting attributes. // remove spliting.

attribute for each outcome j of spliting criterion

// partition the types and grow subtrees for each partition let Dj be the set of data tuples in D satisfying outcome j. // a partition

if Dj is empty then
        attach the node a leaf labeled with the majority class in D to node n;
else
        attach the node returned by Generate dessi Decision tree [Dj, attribute list] to Node Nj
End for
return N;


Tree Pruning

Tree Pruning is performed in order order to remove anomalies in the training data due to noise or outliers. The pruned trees are smaller and less complex.


Tree Pruning Approaches:-


• Pre - Pruning - The tree is pruned by halting its

Construction early

- Post - pruning - This approach removes a sub-tree from a fully grown tree.

Cost Complexity

Cost Complexity is measured by the following two parameters :-

i) Number of leaves in the tree
ii) Error rate of the tree

strengths and Weakness of Decision Tree approach

strength :-

- Decision trees are able to generate understandable set rules
- Decision trees perform classification without requiring much computation
- Decision trees are able to handle both continuous and categorical variabels.
- Decision trees provide a clear indication of which fields are important for prediction or classification

weakness :-

- Decision trees are less appropriate for estimation tasks where the goal is to product the value of a

continuous attribute

• Decision trees are prone to errors in classification problems with many class relatively small number of training examples.

• Decision tree can be computationally expensive to train. The process of growing a decision tree is computationaly expensive. Pruning algorithm can be also expensive since many candidate sub-trees must be formed and compared.

→ classifier accuracy :-

The accuracy of a classifier refers to the ability of a given classifier to correctly predict the class label of a new or previously unseen data.

→

The accuracy of a predictor refers to how well a given predictor can guess the value of the predicted attribute for new previous unseen data.

Attributes selection measure.

• Attribute selection measure is a heuristic for selecting the splitting criterion that "best" separates a given data partition, D, of a class-labeled training tuples into individual classess. The Information gain is used to select the splitting attribute in each node in the tree.

- select the attribute with the highest information gain
  - Expected Information
  - Information needed
  - Information gained

## Computing Information Gain for Continuous value attributes.

- Let attribute A be a Continuous-valued attribute
- must determine the best split point for A
- Sort the value A in increasing order
- Typically, the midpoint between each pair of adjacent value is considered as a possible split point.
  $(a_i + a_{i+1})/2$ is the midpoint b/w values of $a_i$ & $a_{i+1}$
- The point with the minimum expected information requirement for A is selected as the split-point for A.

## Overfitting and Tree pruning

**overfitting:** An induced tree may overfit the Training data.
  - Too many branches, some may reflect anomalies due to noise or outliers.
  - poor accuracy for unseen samples.

→ Two Approaches for avoid overfitting :-

- **Porpruning :** Halt tree construction early - do not split a node if this would result.

In the goodness measure falling below a threshold
- Difficult to choose appropriate threshold

- Postpruning :- Remove branches from a fully
  grown" tree - get a sequence of
  progressively pruned trees.

- use a set of data different from training
  data to decide which is the " best pruned-tree"

Scalable Decisi8ion Tree Induction methods.

SLIQ

Builds an index for each attribute and only
class list and the current attribute list reside
in memory.

SPRINT

Constructs an attribute list data structure

PUBLIC

Integrates tree splitting and tree pruning. step
growing the tree earlier

- Rain Forest
  Builds an AVC - list (Attribute. value, class
  labes]

- BOAT
  uses bootstrapping to create several small
  samples.

# Bayesian classification

Bayesical classification is Based on Bayes Theorom Bayesian classifiers are the Stactistical classifiers Bayesian classfiers Can predict class membership probabilities such as the probability That a given tuple belongs to a particular class.

## Baye's Theorom

Baye's, Theorom is named after Thomas Bayes. There are two types of probabilities-
- Posterior probability $[p(H/x)]$
- Prior Probability $[p(H)]$

Where $x$ is data tuple and $H$ is some hypothesis According to Baye's Theorom,

$$P(H/x) = P(x/H) \, p(H) / p(x)$$

Bayesian classification uses Bayes Theorom to predict the occurrence of any event. Bayesian Classifiers are the Statisticcal classifiers with The Bayesian probability understanding. The theory expressess how a level of belief expressed as a probibility.

Bayes theorem come into existence offer Thomas Bayes who first utilized conditional probability to provide an algorithm that uses evidence to calculate limits on an unknown

parameter.

Baye's theorem is expressed mathematically by the following equation i'c given below

$$P(x|y) = \frac{p(y/x) \, p(x)}{p(y)}$$

where x and y are the events and p(y)≠

P(x/y) is a conditional probability that descri the occurence of event x is given that y is i's tr

P(y/x) is a conditional probability that descri The occurence of event y is given that x is tr

P(x) and p(y) are the probabilities of observi x and y independently of each other. This is kno as the marginal probability.

## Bayesian interpretation

In the bayesian interpretation, probability determines a "degree of belief." Bayes Theorem connects the degree of belief in a hypothesis be fore and after accounting for evidence. equ let us consider an example of the coin, if we toss a coin, then we got either heads or tails and The percent of occurence of either heads tail is 50%. If the coin is flipped numbers times and the outcomes are observed, the

degree of belief may rise, fall or remain remain the same depending on the outcomes.

For proposition X and evidence Y,

- $P(x)$, The prior, is the primary degree of belief in x.

- $P(x/y)$, The posterior is the degree of belief having accounted for Y

- The quotient $\dfrac{P(y/x)}{P(y)}$ represents the supports Y provides for x.

Bayes theorem can be derived from the conditional probability:-

$$P(x/y) = \frac{P(x \wedge y)}{P(y)}, \text{ if } P(y) \neq 0$$

$$P(y/x) = \frac{P(y \cap x)}{P(x)}, \text{ if } P(y) \neq 0$$

$$P(y/x) = \frac{P(y \cap x)}{P(x)}, \text{ if } P(x) \neq 0$$

where $P(x \cap y)$ is the joint probability of both x and y being true because

$$P(y \cap x) = P(x \cap y)$$
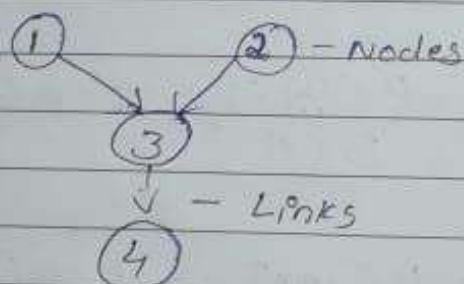$$\text{or } P(x \cap y) = P(x/y) \, P(y) = P(y/x) \, P(x)$$

or $P(x|y) = \dfrac{P(y|x) \, P(x)}{P(y)} \quad P(y) \neq 0$

## Bayesian n/w :-

A Bayesian n/w falls under the classification of probablistic Graphical modelling [PGM] procedure is utilized to compute uncertainties by utilizing the probability concept. Generally known as Belif n/ws, Bayesian n/ws are used to show uncertainities using directed Acyclic graph (DAG)

A Directed Acyclic graph is used to show a Bayesian n/w, and like some other statistical graph. DAG consists of a set of nodes and links, where the links Signify the connection between the nodes.



```
  ①           ②  - Nodes
     ↘      ↙
       ③
       ↓  - Links
       ④
```

The nodes here represent random variables and the edges define the relationship b/w these variables.

A DAG models the uncertainty of an event Taking place based on the Conditional

Probability Distribution [CDP] of each random variable. A conditional probability table [CPT] is used to represent the CPD of each variable in a nlw.

Bayesian classification: way 2

A statistical Classifier: performs probablistic prediction ie predicts class membership probabilities.

Foundation:- Based on Bayes Theorem

Performance: A simple Bayesian classifier, native Bayesian classifier, has comparable Performance with decision tree and selected neural network classifiers.

Incremental: Each training example can incrementally increase / decrease the Probability that a hypothesis is correct - Prior knowledge can be combined with observed data

Standard: Even when Bayesian methods are computationally intractable, they can provide a standard of optimal decision making against which other methods can be measured.