

# UNIT-IV

## Classification

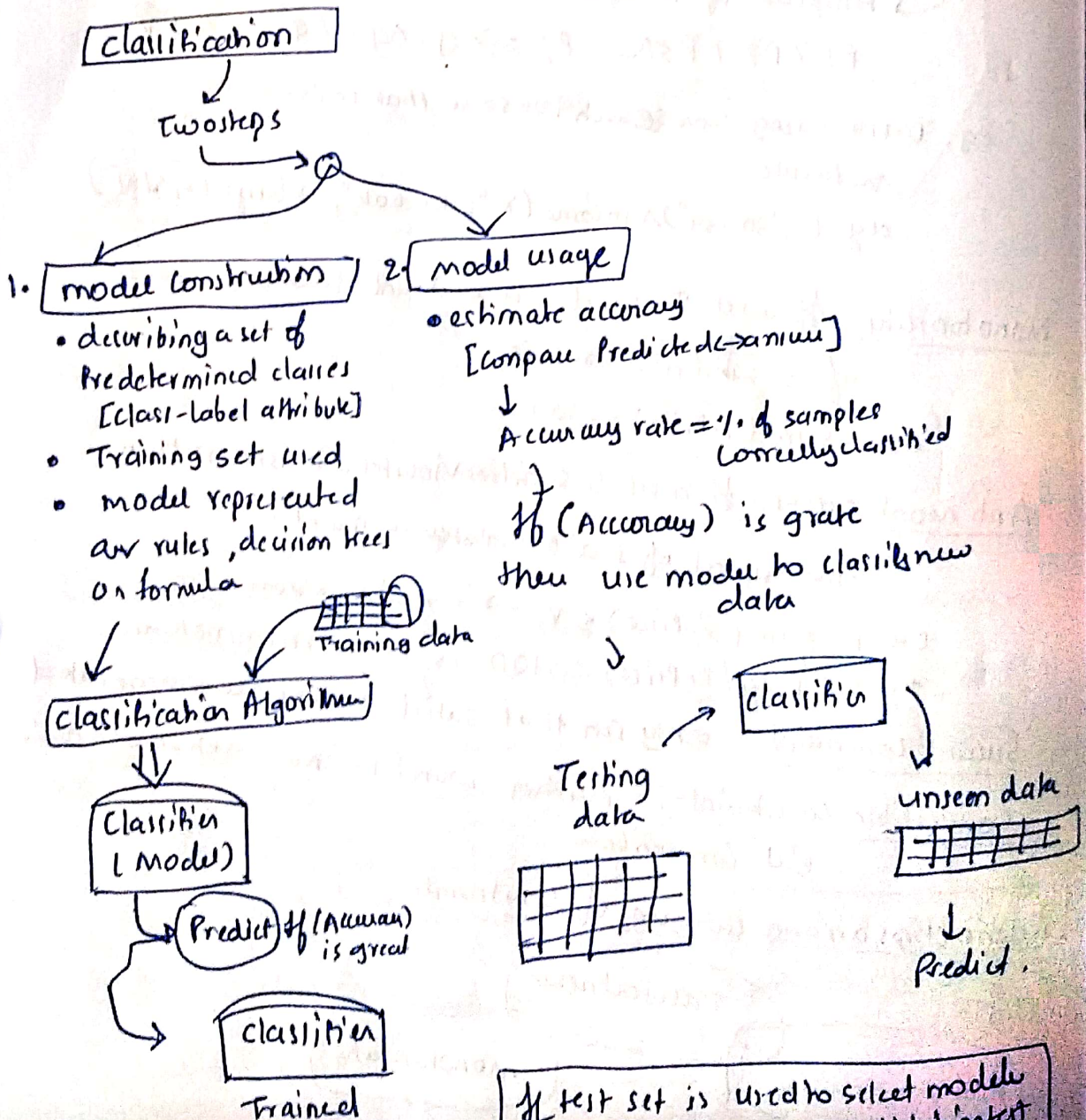
- Supervised learning (Classification)
  - New data classified based on training set
- Unsupervised learning (Clustering)
  - Establishing classes or clusters to identify unknown data

### Classification vs Numeric prediction

- Classification (discrete, nominal data) to predict class for new data
  - Numeric prediction (models continuously valued functions)
    - Predict unknown or missing values.
- Regression analysis is used

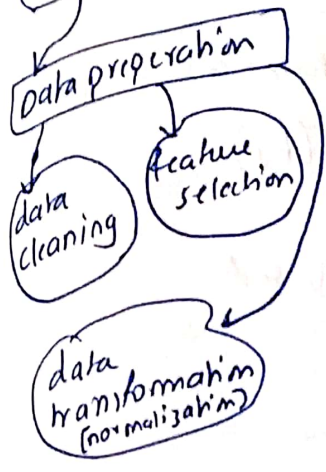
#### Applications:

- fraud detection
- medical diagnosis



If test set is used to select models then it is called validation test

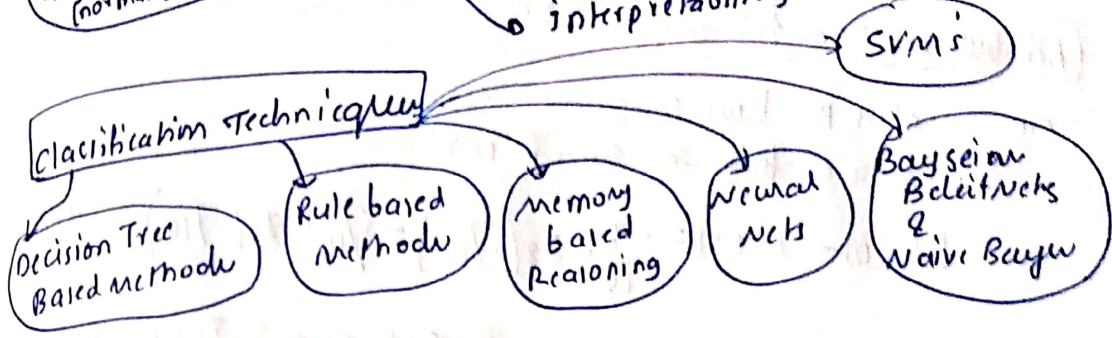
# Issues with classification & prediction



# Evaluating classification methods

- Prediction accuracy
- Speed & scalability
  - time to construct
  - LR use mode
  - efficiency
- Robustness
- Interpretability

# Classification techniques



# Decision Tree induction

- Basic algorithm:
  - Tree is constructed in Topdown, recursive divide & conquer manner
  - Starting, Training examples at root (all)
  - Attributes are categorical (if continuous valued, discretized in advance)
  - examples are partitioned recursively based on selected attributes
  - Test attributes are selected on basis of Statistical Measure (Information gain)
- Conditions for stopping partitioning:
  - All samples of a given node belong to same class
  - no remaining attributes. (Majority voting applied for classifying rest)
  - no samples left

**Entropy** [A measure of uncertainty]

- high entropy → high uncertainty
- low entropy → low uncertainty

Conditional entropy

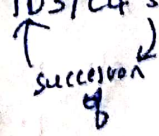
$$H(Y/X) = \sum_x P(x) H(Y/X=x)$$

**Attribute selection measure** : [Information gain 103/44.5]

- select attribute with highest information gain

→ expected information (entropy)

let  $D \rightarrow$  Tuple,  $P \rightarrow$  Probability  
 $C \rightarrow$  class



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

→ Information need after splitting  $D$  into  $V$  partitions using  $A$

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

→ Information gained by branching on attribute  $A$ :

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

**Attribute-selection: Information Gain**

eg:- class 'P': buys-computer = "yes"  
class 'N': buys-computer = "no"

$$\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,1) = 0.694$$

age	$p_i$	$n_i$	$I(p_i, n_i)$
$\leq 30$	2	3	0.971
31..40	4	0	0
$> 40$	3	2	0.971

→ 5 samples out of 14  
for  $\text{age} \leq 30$   
• 2 yes & 2 no.

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.246$$

∴ similarly

$$\text{Gain}(\text{income}) = 0.029$$

$$\text{Gain}(\text{student}) = 0.151$$

$$\text{Gain}(\text{credit-rating}) = 0.048$$

schema

age	income	stud	credit-rating

**Computing Information Gain for continuous valued attributes**

• let  $A$  → continuous valued attribute

steps:-

1. sort  $A$  in increasing order

2. find mid  $(a_1 + a_{n-1})/2$  and select it as  
Best split point

3. split tuple  $D$  into  $D_1$  &  $D_2$

$A(\leq \text{midpoint}) \rightarrow D_1$

$A(> \text{midpoint}) \rightarrow D_2$

## Gain ratio for attribute selection

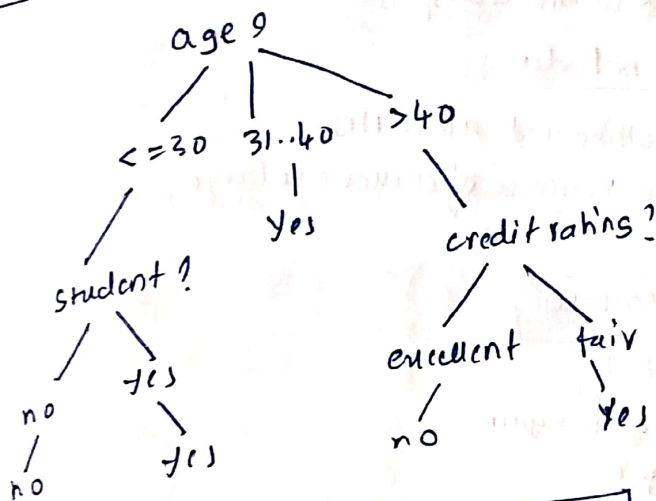
Information gain measure is biased towards attributes.  
 e.g. 4.5 (a succession of 103) uses gain ratio to overcome

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

$$\text{Gain Ratio} = \text{Gain}(A) / \text{SplitInfo}(A)$$

→ attribute with max gain ratio is selected as splitting attribute.

## Decision tree - after information gain



## Gini Index (CART, IBM Intelligent Miner)

• If a dataset  $D$  contains examples from  $n$  classes then gini index

$$\text{gini}(D) = 1 - \sum_{j=1}^n p_j^2 \quad p_j = \text{relative frequency of class } j \text{ in } D$$

If  $D$  on  $A$  is split into  $D_1$  &  $D_2$

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

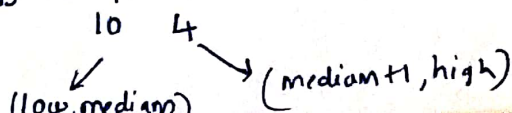
• reduction in impurity

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

eg:  $D \rightarrow$  9 tuples in buys-computer = "yes" & 5 in "no"

$$\text{gini}(D) = 1 - \left( \frac{9}{14} \right)^2 - \left( \frac{5}{14} \right)^2 = 0.459$$

Suppose  $D \rightarrow D_1$  &  $D_2$



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

$$= 0.443$$

$$= \text{Gini}(low, high)^D$$

•  $G_{in}(low, high)$  is 0.458,  $G_{in}(medium, high)$  is 0.450

∴ split the  $\{low, medium\}$  &  $\{high\}$

disadvantages of ~~GINI~~ Gini Ratio

- tends to prefer unbalanced splits

disadvantages of minitree:

- biased to multivalued attributes

- difficult when number of classes are large

### Other popular Selection Measures

- CHAID: based on  $\chi^2$  test
- C-SEP: better than gini & gain
- G-statistic: closest to  $\chi^2$
- MDL (Minimum description length) → Prefers simple solution
- CART: find multivariate split based on linear combination of attributes.

### Overfitting & Tree pruning

- Overfitting: - A induced tree may overfit the training data
  - Too many branches, noise
  - Poor accuracy for unseen samples
- Two approaches to avoid overfitting
- Pre-pruning: Halt tree construction early & remove branches falling below a goodness threshold
- Post-pruning: - Remove branches from fully grown trees

## Enhancements to Basic Decision Tree Induction

- Allow for continuous valued attributes.
- Handle missing attribute values
- attribute construction.
  - create new attributes based on existing
  - reduces fragmentation, repetition & replication

## Classification in large databases → Why decision trees so popular?

- can classify millions of examples & hundreds of attributes using reasonable speed in decision trees.
- Scalable
- relatively faster running speed.
- convertible to simple rules.

## Scalability framework for Rainforest

- Builds on Avc list Avc (Attribute, value, class-label)
- Avc-set (of attribute X)
  - Projection of training data onto attribute X
- Avc-group (of a node N)
  - Set of Avc sets of projection attributes at node N

eg: Avc set on income

income	Buys computer	
	yes	no
high	2	2
medium	4	2
low	3	1

## BOAT (Bootstrapped Optimistic Algorithm for Tree construction)

- uses a statistical technique called bootstrapping to create several smaller samples (subsets), each fit in memory.
  - each subset creates a new tree, resulting in several trees
- Adv: requires only two scans of DB, an incremental alg.

## Bayesian classification why?

- A statistical classifier → does probabilistic prediction
- works on Bayes Theorem
- good performance relatively along with Neural Nets & Decision trees
- Incremental: improves with data
- standard: Bayesian methods are intractable & Provides decision making

## Bayes Theorem: Basics

- Total probability theorem:

$$P(B) = \sum_{i=1}^M P(B/A_i) P(A_i)$$

- Bayes theorem:

$$P(H/x) = \frac{P(x/H) P(H)}{P(x)}$$

$x$  → data sample

$H$  → hypothesis that  $x$  belong to class

$C$

$$= P(x/H) \times P(H) / P(x)$$

→ classification is to determine  $P(H/x)$  [posteriori probability]

$P(H)$  → initial probability

$P(x)$  → Probability that sample data is observed

$P(x/H)$  → likelihood.

Theorem:-

[[[

Given training data  $x$ , posteriori probability of a hypothesis  $H$ ,  $P(H/x)$  follows Bayes theorem

$$P(H/x) = \frac{P(x/H) P(H)}{P(x)}$$

]]]

Simply ⇒ Posteriori = likelihood × Prior / evidence

→ let attributes  $x = \{x_1, x_2, \dots, x_n\}$

→ classes  $C_1, C_2, \dots, C_M$

→ classification → max Posteriori

$$P(c_i/x) = \frac{P(x/c_i) P(c_i)}{P(x)}$$

as  $P(x)$  is constant for all classes

$$P(c_i/x) = P(x/c_i) P(c_i)$$

### Naive Bayes classifier

If attributes are conditionally independent, this greatly reduces the computational cost.

$$P(x/c_i) = \prod_{k=1}^n P(x_k/c_i) = P(x_1/c_i) \times \dots \times P(x_n/c_i)$$

### Dataset

age	income	student	credit-rating	buys-computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31..40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31..40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31..40	medium	no	excellent	yes
31..40	high	yes	fair	yes
>40	medium	no	excellent	no



eg:-

•  $P_i(c) :=$

$$P(\text{buys-computer} = \text{"yes"}) = 9/14 = 0.643$$

$$P(\text{buys-computer} = \text{"no"}) = 5/14 = 0.357$$

• Compute  $P(x_i|c_i)$  for each class

$$P(\text{age} \leq 30 | \text{buys-computer} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{age} \leq 30 | \text{buys-computer} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{income} = \text{"medium"} | \text{buys-computer} = \text{"yes"}) = 4/9 = 0.444$$

$$P(\text{income} = \text{"medium"} | \text{buys-computer} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{student} = \text{"yes"} | \text{buys-computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{student} = \text{"yes"} | \text{buys-computer} = \text{"no"}) = 1/5 = 0.2$$

$$P(\text{credit} = \text{"fair"} | \text{buys-computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{credit} = \text{"fair"} | \text{buys-computer} = \text{"no"}) = 2/5 = 0.4$$

•  $X = \{\text{age} \leq 30, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit} = \text{fair}\}$

$$P(x|c_i): P(x | \text{buys-computer} = \text{"yes"})$$

$$= 0.222 \times 0.444 \times 0.667 \times 0.667$$

$$= 0.044$$

$$\therefore P(x | \text{buys-computer} = \text{"no"})$$

$$= 0.6 \times 0.4 \times 0.2 \times 0.4 = 0.019$$

$$P(x|c_i) * P(c_i): P(x | \text{buys-computer} = \text{"yes"}) *$$

$$P(\text{buys-computer} = \text{"yes"})$$

$$= 0.044 \times 0.643 = 0.028$$

$$P(x | \text{buys-computer} = \text{"no"}) * P(\text{buys-computer} = \text{"no"}) = 0.007$$

$\therefore X$  belongs to class ("buys-computer" = "yes")

## Problem in Bayesian

- Naive Bayesian prediction requires each condition prob be non-zero, otherwise the predicted probability will be zero.

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

- use Laplacian correction

## Advantages of Naive Bayes classifier

- easy to implement
- good results in most of cases

## Disadvantage

- If there is class conditional independence accuracy is lost

## Rule-Based classification

- representing the knowledge in the form of if-then rules

R: If age=youth and student=yes then buys-computer=yes

Rule antecedent /

Precondition

Condition

- Assessment of Rule is done by Accuracy & Coverage

$n_{\text{covers}}$  = number of tuples covered by R

$n_{\text{correct}}$  = number of tuples correctly classified by R

$$\text{Coverage}(R) = \frac{n_{\text{covers}}}{|D|} \quad |D| \rightarrow \text{Training data set}$$

$$\text{Accuracy}(R) = \frac{n_{\text{correct}}}{n_{\text{covers}}}$$

- If more than one rule are triggered, we need Conflict Resolution.

- size ordering
- class based ordering: decreasing order of prevalence
- Rule-based ordering (decision/priority list).

## Rules extraction from decision tree

- Rules are easier to understand
- one rule for each path from root to leaf  
- leaf holds class prediction
- Rules are mutually exclusive & exhaustive

## Rule Induction: Sequential Covering method

- Sequential covering algorithm:-  
extracts rules directly from training data.

eg: Algorithms are: FOIL, AQ, CN2, RIPPER

Rules are learned sequentially, each for given class  $c$ , will cover many tuples of  $C \cup D$

### steps:-

1. Rules are learned one at a time
2. each time a rule is learned, the tuples covered by rules are removed
3. Repeat the process until no more training examples left

while (enough target tuples left)

generate a rule

remove positive target tuples satisfying this rule

### - Rule generation

while (true)

find best predicate  $P$

if  $\text{foil-gain}(P) > \text{Threshold}$  then add  $P$  to current rule  
else break.

## How to learn one rule?

- start with most general rule possible (condition = empty)

- add new attributes by adopting a greedy depth-first strategy
- Rule quality measure (in FOIL & RIPPER)
  - new info-gain

$$\text{FOIL-Prune}(R) = \frac{\text{Pos} - \text{neg}}{\text{Pos} + \text{neg}}$$

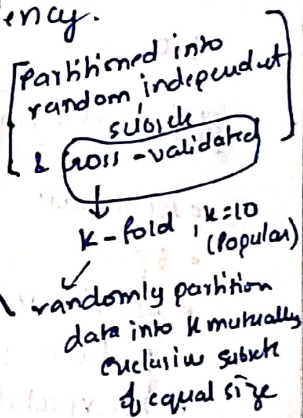
+ve examples  
-ve examples

$$\text{FOIL-Gain}(R) = \text{Pos}' \times \left( \log_2 \frac{\text{Pos}'}{\text{Pos}' + \text{neg}'} - \log_2 \frac{\text{Pos}}{\text{Pos} + \text{neg}} \right)$$

### Model evaluation & selection

- use validation test set of class-labelled tuples instead of Training set when assessing accuracy
- method for estimating a classifier's efficiency.

- Holdout method, random subsampling
- cross-validation
- Bootstrap (.632 bootstrap)  $K = \text{number of tuples}$
- Comparing classifiers:
  - Confidence intervals
  - Cost-benefit analysis & ROC curves



### Classifier evaluation metrics

- precision: % of tuples labelled by classifier as +ve are really +ve

$$\text{Precision} = \frac{TP}{TP + FP} \quad [\text{exactness}]$$

- Recall: % of what tuples did classifier labelled as +ve

$$\text{Recall} = \frac{TP}{TP + FN}$$

Perfect score is 1.0

- F measure (F-score): harmonic mean of precision & recall

$$F = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- $F_\beta$  (weighted measure of Precision & Recall)

$$F_\beta = \frac{(1 + \beta)^2 \times \text{Precision} \times \text{Recall}}{\beta^2 \times \text{Precision} + \text{Recall}}$$

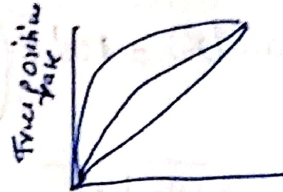
## Estimating Confidence Interval : null hypothesis

- If we can reject null hypothesis ( $M_1$  &  $M_2$  are classifiers & are the same). then the difference between  $M_1$  &  $M_2$  is statistically significant [choose lower error rate] one

→ t-test → Pairwise Comparison

## Model Selection: Roc curve

- Receiver operating characteristic curve for visual comparison of classifiers
- A model with perfect accuracy will have an area 1.0



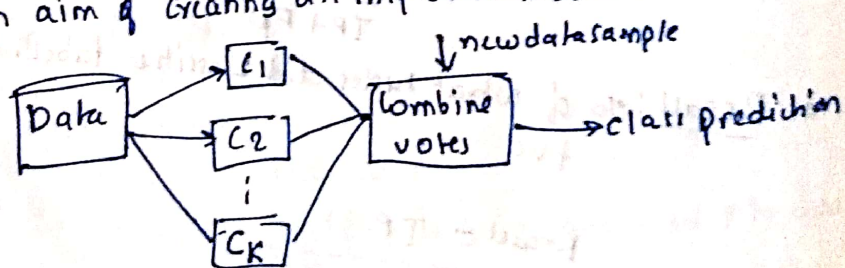
## Issues affecting model selection

- Accuracy
- Speed
- Robustness
- Scalability
- Interpretability

## Techniques to improve classification Accuracy

### ensemble methods

- Use a combination of models to increase accuracy
- Combine a series of  $K$  learned models  $M_1, M_2, \dots, M_K$  with aim of creating an improved model  $M^*$



### Ensemble methods

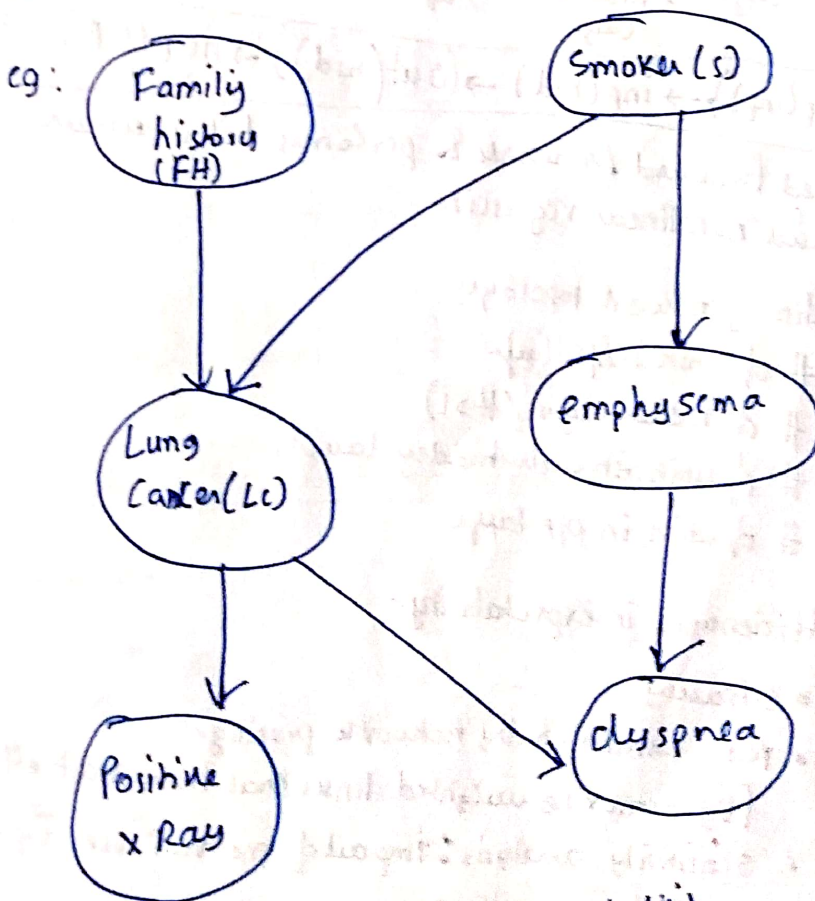
- ↳ Bagging [diagnosis based on multiple doctors majority vote]
- ↳ Ada Boost
- ↳ Boosting [consult several doctors, based on combination of weighted diagnoses - weight assigned on previous diagnosis accuracy]
- ↳ ensemble [Random forest → Forest- $R_1$  (random input selection) → Forest- $R_c$  (random linear combination)]

## Handling different kinds of cases in classification

- Traditional methods ~~assume~~ assumes a balanced distribution of data.
- Typical methods for imbalance data in 2-class classification
  - oversampling: - re-sampling of data from positive class
  - under sampling: - randomly eliminate tuples from negative class
  - Threshold moving: moves the decision threshold so that rare-class tuples are easy to classify
  - ensemble methods: ensemble multiple classifiers as discussed before.

## Bayesian Belief Network / Bayesian network

- Bayesian belief network, probabilistic networks allow class conditional independencies between subsets of variables.
- uses a directed acyclic graph. ~~graph~~



→ uses conditional probability.

**Classification by backpropagation** [Connectionist learning]

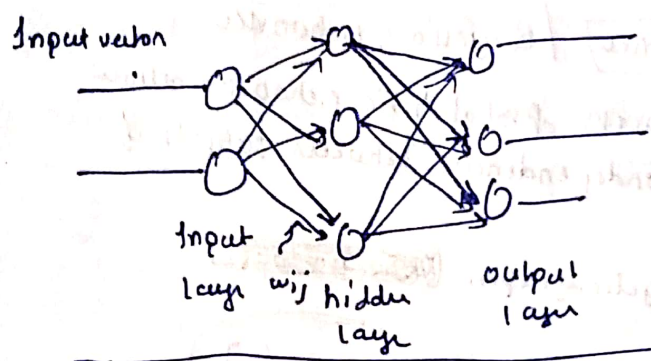
- A neural network learning algorithm
- A set of connected ip/op units where each connection has a weight associated with it.
- During learning phase, the network learns by adjusting weights to minimize mean squared error.
- Modifications are made in backward direction

**Neural Network as a classifier**

↳ Advantages:-

- high tolerance to noisy data
- Ability to classify unknown patterns.
- Algorithms are parallel.

dis:-  
long training time



$$[O_p(inp)] \rightarrow inp(hid) \rightarrow (out(hid)) \rightarrow inp(OP) \rightarrow OP$$

↳ Feed forward Network & performs better when used non linear regression

- defining network topology

- # of unit in ip layer
- # of hidden layers (If > 1)
- # of units in each hidden layer.
- # of units in op layer

- efficiency & interpretability

- efficiency
- Rule extraction by network pruning [By removing weighted links that have least effect]
- Sensitivity analysis: impact of one variable input on network output.

- discriminating classifiers

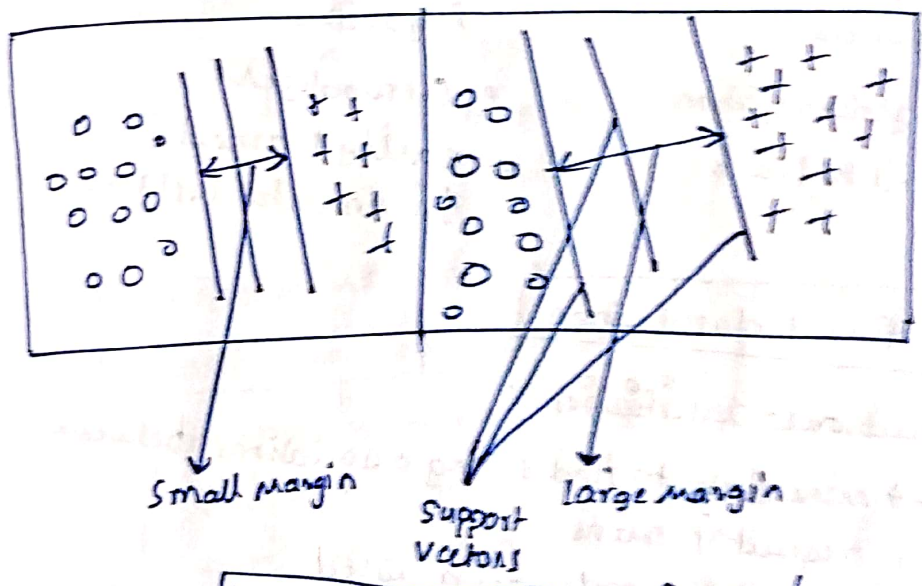
- high accuracy v/s long training time

# Support Vector Machines

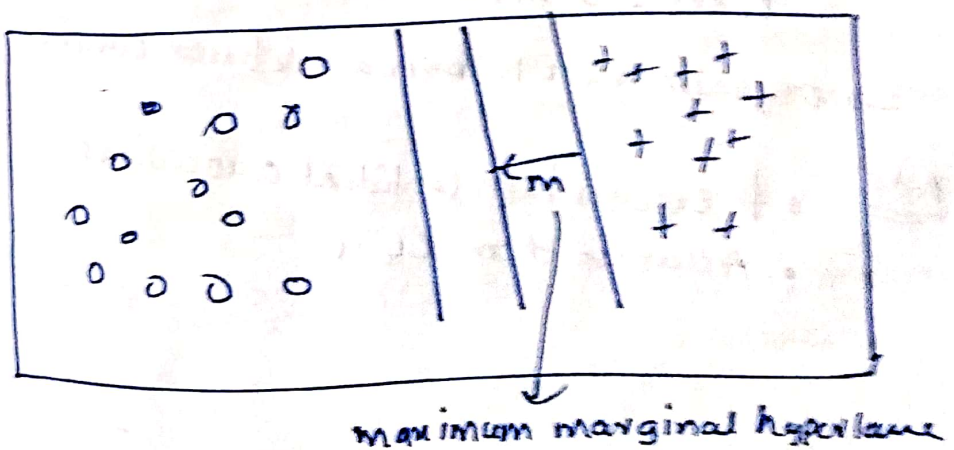
- A relatively new classification method for both linear & non-linear data
- It uses a non-linear mapping to transform the original training data into a higher dimension.
- With the new dimension, it searches for linear optimal separating hyperplane (decision boundary)
- With appropriate non-linear mapping, data from two classes can always be separated by a hyperplane.
- SVM finds this hyperplane using support vectors ("essentially training tuples") & margins (defined by support vectors)

Applications : handwritten digit recognition, object recognition etc.

(Classification & numeric predictions)



⇒ when data is linearly separable ≡≡≡ Concept





## Why SVM effective on high dimensional data

- The complexity of trained classifier is characterized by # of support vectors rather than dimensionality of data.
- With small # of support vectors, we can have good generalization

## SVM Linearly Inseparable

- Transform the original input data into a higher dimensional space.
- Search for linear separating hyperplane in the new space.

## SVM v/s Neural net

### SVM

- deterministic Algorithm
- Nice generalization
- Hard to learn

### NN

- Non-deterministic Algorithm
- Generalizer
- easily learned in incremental fashion.

## Pattern Based classification

• Association classification:

→ mine data to find strong associations between frequent patterns

→ Association rules are generated

$P_1, P_2, \dots, P_i \rightarrow "A \text{ class} = C" (\text{conf}, \text{sup})$

→ Organise rules to form a rule based classifier

Adv: • It explores high confident associations  
• Accurate than CLF

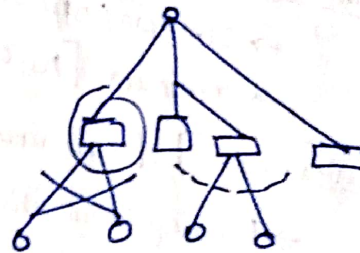
## Typical Associative classification methods

- CBA (classification based on Associations)
- CMAR (classification based on Multiple Association Rules)
- CPAR (classification based on predictive Association Rules)

→ Discriminative frequent Pattern - Based classification

- **DDP Mine** (Branch & Bound search)
- FPtree pruning

$$\begin{aligned} \text{Sup}(\text{child}) &\leq \text{Sup}(\text{Parent}) \\ \text{Sup}(b) &\leq \text{Sup}(a) \\ \text{maximize } \ln(C/b) \end{aligned}$$



## Lazy v/s eager learning

- lazy learning :- (instance-based learning)

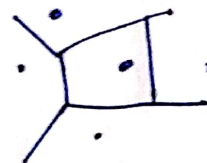
Simply stores data for minin processing & waits until it gives a test tuple

- early learning : constructs a classifier before receiving data to classify.

→ Instance Based methods:

- ① K-nearest neighbour approach :-

• Instances represented as points in Euclidean space



Voronoi diagram is used

- nearest neighbour defined in terms of euclidean distance  $(x_1, x_2)$

- Target function can be real or discrete valued

- for discrete valued, K-NN returns most common value among K training examples.

- ② locally weighted regression (uses local approximation)

- ③ case-based-reasoning { uses symbolic representations & inference }