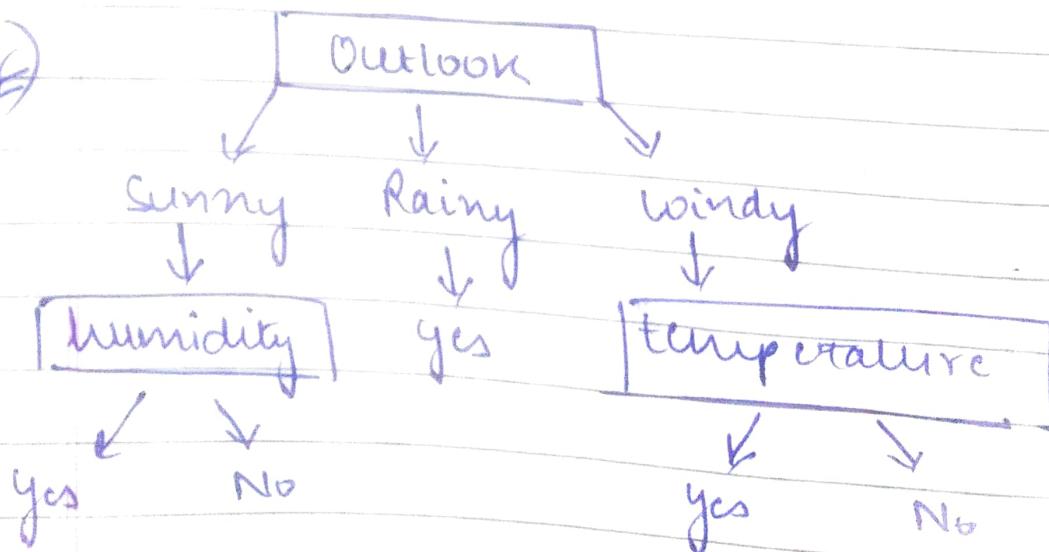


- Day
- Decision Tree :-
- ① A decision tree has a structure that consists of a root node, internal node, branches & terminal or leaf node.
 - ② Decision tree is a non linear representation. It is a classifier in the form of a tree.
 - ③ It is also known as tree structure classifier.
 - ④ It classifies instances by sorting them from root node to some leaf node.
 - ⑤ A decision tree is a graphical representation of all the possible solutions to a decision based on certain conditions.

Eg)



→ Appropriate problem for ITEM
Decision tree

- 1) Instances may be represented by ~~value~~ attribute value pairs
- 2) The target function has discrete output values
 - 3) The training data may contain errors
 - 4) The training data may contain missing attribute values

Which attribute is best classifier?

A statistical value called information gain measures how well a given attribute separates

b. Information Gain -

I_G is the decrease in entropy after a dataset is split on the basis of an attribute.

Constructing decision tree is all about finding the attribute that returns highest information gain.

$$\text{Gain} = \text{Entropy}(S) - I(\text{average})$$

Entropy → removing unwanted data

Entropy → degree of randomness tell how random is our data

Entropy = $-p \log_2 p - q \log_2 q$

→ Entropy - metric which measures the impurity - 1st step to some decision tree
 Entropy is the measure for uncertainty, impurity and information content.

It is an average optimal no of bits to encode information about certainty or uncertainty in a given example dataset.

$$\text{Entropy} = -\frac{p}{p+n} \log_2 \left(\frac{p}{p+n} \right) - \frac{n}{p+n} \left(\log_2 \left(\frac{n}{p+n} \right) \right)$$

Algorithms

CART Algo
 uses gini index

ID3 (Iterative Dichotomiser 3)
 uses entropy and information gain to calculate homogeneity of the sample

→ Strengths of Decision Tree:-

1) DT are able to generate understandable rules

SAND

2) Simple Calculation → $p \times q \times h$

They perform calculation without requiring much computation

3) Attribute Types -

DT are able to handle both continuous & categorical values

4) Not include all attributes → 2022 (6/16)

TENN

NCEEN
NEEN
N E E N C

→ Weakness

① No continuous predictions

less appropriate for tasks where we have to make continuous predictions

② ~~cost~~ can be computationally expensive to train

③ DT algo handles one attribute at a time for a node

④ Error prone

⑤ NO sequence → problematic for time series data

→ Applications

Speech Recognition
Character Recognition

Traffic Risk Analysis

Remote Sensing

Medical image Applications

→ Issues → MAHIA

→ ID3 algo takes
+ to discrete
values

① avoid overfitting the data

② incorporating continuous valued attributes

③ alternative measures for selecting attributes

④ handling training examples with missing attribute values

⑤ handling attributes with different costs

Instance Based Learning

Instance Based learning methods are conceptually straightforward approach

- ① Sometimes referred to as lazy learning
- ② In this we will not process the data to learn and develop a model
- ③ Instead we just store the example.
- ④ We will not immediately learn a model

Comprises of 2 phases

- ① Testing Phase → when a new query is encountered, a set of similar related learning instances
- ② Training Phase → we do not perform any processing on data. We simply store the data

is retrieved from memory and used to classify or predict the target function value.

Lazy learning (k - nearest neighbour)

Radial Based functions (RBF)

Case Based Reasoning (CBR)

⑥ can use more

Adv complex symbolic Disadv Most of

- ① target function can be estimated weakly and differently for each new instance to be classified
- ② new function can be learned

Model form
stores training data
requires less storage
scoring for new instance
Model Based is fast learning

Store the model in suitable form

Generalize the rules in form of model

Predict for unseen scoring instance using model

Can throw away training data after training

Train model from training data to estimate model parameters

may not have explicit model form
⑦ more storage
scoring for new instance may be slow

Instance Based Learning

① no model to store

no generalization beyond scoring

② - - - using training data

Training data must be kept

Do not train data

→ KNN → Adv

disadv

This algorithm classifies a new instance by determining k vote similar instances & summarizing the output
use k instances

If the target variable is discrete then it is a classification problem but selects the most common value among k instances by a majority vote.

If target function variable is continuous then it is a regression problem and hence the mean of k instances

Alg :-

(a) Continuous attribute -

Euclidean distance b/w 2 planes is (with co-ordinates)

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

(b) Categorical attribute

by the value of d instances

① if the distance $d = 0$ otherwise $d \neq 0$

② sort the distance in ascending order

& select 1st K nearest training data instance

③ predict the class of test instance by majority vote or mean of K selected nearest neighbour

→ Input :-

Training data set = T

distance matrix = d

Test instance = t

nearest neighbour = K

Output \Rightarrow predict class or category

for test instance T

compute the distance b/w

test instance and every other instance in training

dataset using a distance matrix

weighted KNN nearest neighbour :-

- ① KNN algorithm has some limitations as its performance depends on
 - (i) choosing nearest neighbour
 - (ii) Distance matrix
 - (iii) decision rule.
- ② The principle of weighted KNN is that n closest neighbors that are far away from the ^{test} instances. Here the weights are inversely proportional to distance
- ③ The selected K nearest neighbours can be assigned any way that means

Algo :-

training dataset = T

distance matrix = d

weighted function = $w(i)$

test instance t

no of n nearest neighbours K

~~Prediction → for test instance t~~

- ④ Compute the inverse of each distance of n selected nearest instance
- ⑤ find sum of inverse
- ⑥ Compute the weight by dividing each inverse distance by sum
- ⑦ add the weights of the same class
- ⑧ Predict the class by choosing the class with maximum votes

→ Locally Weighted Regression :-

Generalized form of KNN

① Local → ② Local → because the function is approximated based only on the data near the query point.

③ Weighted → because the contribution of each training example is weighted by its distance from query point.

④ Regression - because this term is widely used in statistical learning in the problem of approximating real valued functions

LWR uses nearly all distance weighted training example to form the local approximation function

LWR is also referred to as memory based method as it requires training data while prediction but uses only the training data instances locally around the point of interest.

T = training data set

Hypothesis function $H(n)$

the predicted target value output is

a linear function where,

b_0 is the intercept

2022-23-6-11-16

β_1 is the coefficient of x

$$H(n) = \beta_0 + \beta_1 x \quad \text{--- } ①$$

Cost function is such that it minimizes the error difference b/w the predicted value $H(n)$ & true value or the actual value y .

cost function \rightarrow

$$J(\beta) = \frac{1}{2} \sum_{i=1}^m (H_\beta(n) - y_i)^2 \quad \text{--- } ②$$

m = no of instances in the training dataset

Now the cost function is modified for locally weighted linear regression

Cost function is

$$J(\beta) = \frac{1}{2} \sum_{i=1}^m w_i (H_\beta(n) - y_i)^2 \quad \text{--- } ③$$

w_i = weight associated with each instance n_i

→ Case Based Reasoning

Case Based Reasoning is a problem solving technique that matches a new instance with the previous solved instance and its solution is stored in the database.

In design and implementation of any CBR application four RE / steps are involved.

- (i) Retrieve
- (ii) Reuse
- (iii) Revise
- (iv) Retain / Store

A CBR has 2 properties :-

- a) Lazy learning
- b) Classification is different for each instance

CBR only works on or is represented as symbols not values

there are 3 components of CBR

① Similarity functions or
distance measure

SAS

② Approximation or adjustment of
instances

③ Symbolic Representation

To measure or implement CBR

we use CADET system

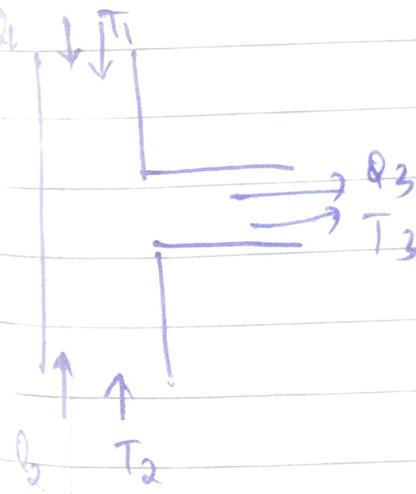
ISAPS

CADET system stands for

case Based Design Tools. It has
75 pre defined libraries.

These libraries will be utilized
to find the solution of new instance

④ A stored case Tjunction pipe



Q = water flow

T = temperature

Function

$$Q_1 \xrightarrow{+} Q_3$$

$$Q_2 \xrightarrow{+}$$

$$T_1 \nearrow +$$

$$T_3$$

$$T_2 \swarrow +$$

limitation

- ① It has limited capabilities for combining & adapting multiple retrieved cases to form the final design (and rely heavily on the user for adaptation stage of process.)
- ② Does not have the range of analysis algo that is needed to refine the abstract conceptual design into the final design.

→ Comparison with RNN :-

→ Application of CBR :-

- ① Reasoning about new legal cases based on previous cases and their solution
- ② come up with design of mechanical arms based on stored library of previous designs
- ③ solving, planning & scheduling problems by receiving & combining portion of previous solutions to similar problems.

Model Based Learning -

Model based learning describes all assumptions about the problem domains in the form of model.

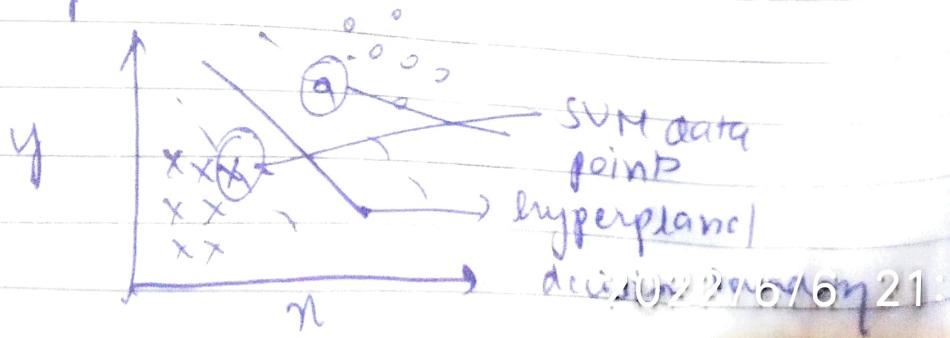
supervised mostly in learning used in classification

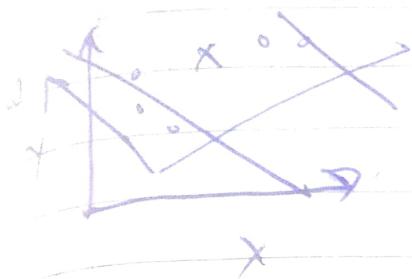
→ Support Vector Machine :-

- ① Support vectors are data points that lie closest to the decision surface or hyperplane
- ② They involve creation of hyperplanes (decision boundaries) that segregate data into classes
- ③ The aim of SVM is to produce a decision plane that defines the boundary b/w the classes to classify the data points
- ④ The decision boundary is drawn with following characteristics:-
 - ① The distance should be maximized between the line & nearest datapoints
 - ② Should avoid misclassification of data

Margin → amount of space b/w 2 classes as defined by the hyperplane.

Marginal distance → sum of shortest distance to the closest positive & negative points

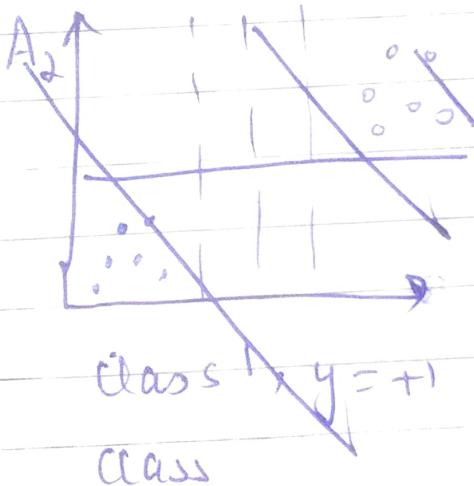




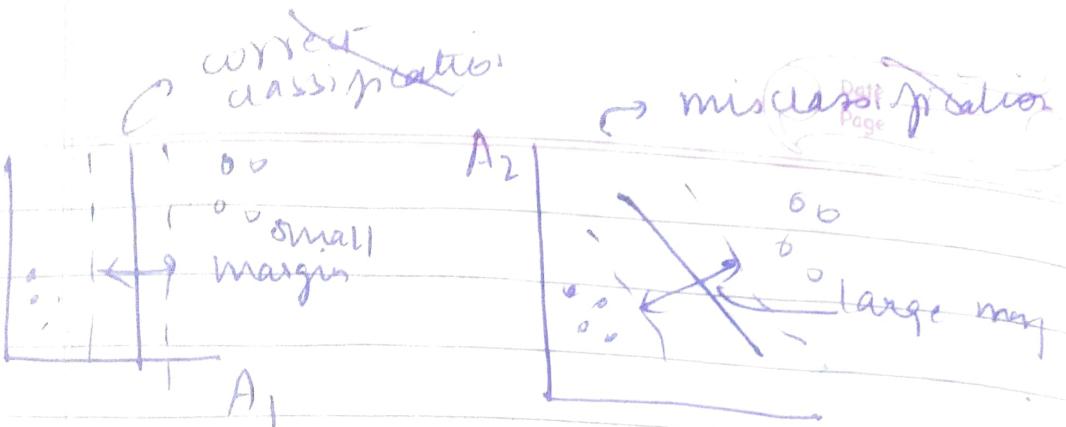
① SVM linearly separable

2D training data set are linearly separable. There are infinite no of separating hyperplanes \Rightarrow decision boundaries

In this we have to find the "best" one that will have minimum classification error



$y = b$



~~the one with larger margin
should have greater generalization
accuracy~~

A separating hyperplane can be written as

$$w \cdot x + b = 0$$

w = $[w_1, w_2, \dots, w_n]$ weight vector

b → scalar (bias)

w.n=0 for SD

$$w_0 + w_1 n_1 + w_2 n_2 = 0$$

hyperplanes -

$$H_1: w_0 + w_1 n_1 + w_2 n_2 \geq 1 \text{ for } y_1 = +1$$

$$H_2: w_0 + w_1 n_1 + w_2 n_2 \leq -1$$

$$H_3: w_0 + b \geq 0 \text{ for } y = +1$$

$$H_4: w_0 + b < 0 \text{ for } y = -1$$

any training tuples that

fall under hyperplane

H1 & H2 are support vectors.

⑥ Non linear SVM :

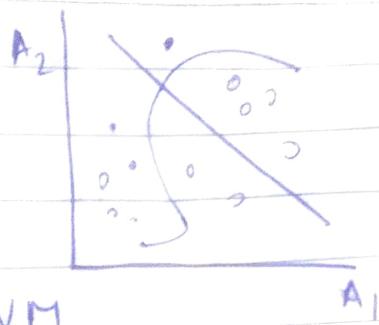
Kernel trick

This 2D

figure show

linearly

unseparable SVM



Kernel Trick means replacing
the dot product in the mapping
function with a kernel function

Kernel function returns the
inner product b/w 2 pts in a
suitable feature space.

Types :-

① Linear Kernel

if x, y are
data points
we want to
do classification

$$K(x, y) = (x^T y)$$
$$K(x, y) = \phi(x) \phi(y)$$

② Polynomial Kernel

$d=2 \Rightarrow$ quadratic
kernel $K(x, y) = (x^T y)^d$

$d=1 \Rightarrow$ linear
kernel
 $d=\text{degree of polynomial}$

more
generalised
form of non
linear

③

③ RBF / Gaussian Kernel

~~It gives SVM classifier in
higher dimensions which is
not possible to visualize naturally~~

It can map input space in
higher dimensional space.

$$k(\mathbf{x}, \mathbf{x}') = \text{Exp}(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2})$$

$$Y = \frac{1}{2\sigma^2}$$

→ Bayesian Belief Network

Bayesian Belief Network describes
the probability distribution over
the set of variables.

Let us consider random variables
($\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$)

Let us consider

$(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n) \Rightarrow$ set of random
variables

Each variable has a set of
values say $V(y_i)$

Here we can define the
joint space by cross product
 $V(y_1) \times V(y_2) \times \dots \times V(y_n)$

Joint probability distribution over the joint space is known as joint probability distribution.

BBN describes joint probability distribution for set of variables

Application

used in diagnosis

Classification

Decision making

Prediction

Joint Probability Distributions

Joint probability distribution for

set of values $\{y_1, y_2, \dots, y_n\}$

to the tuples of network variables
 (y_1, y_2, \dots, y_n) can be computed by

$$P(y_1, y_2, \dots, y_n) = \prod_{i=1}^n P(y_i | \text{parent}(y_i))$$

→ EM Algorithm -

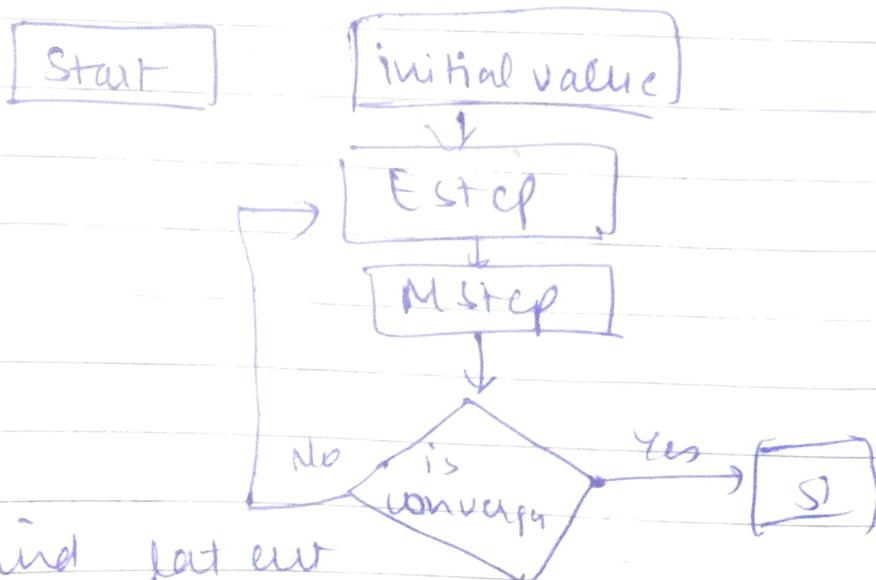
Expectation Maximization
Algorithm

It is used to find the latent variable.

It is basic for unsupervised clustering algorithm.

Use :-

- ① find missing date
- ② basic for unsupervised learning algo



- ③ find latent var
- ④

Adv

- ① guarantees that likelihood will increase each iteration
- ② E step and M-step are pretty easy & can often solve complex problem
- ③ Sol^h to M-step often exists in closed form

Disadv

- ① slow convergence
- ② requires both probabilities
- ③ makes convergence to local optima only