

## Motivation for ML

25/07/2024

To find patterns among large volumes of data -  
Not possible through simple query languages.

- (Insight Extraction)

- Automation

- Predictions & Forecasting

- Personalisation

- Optimisation

- Scalability & Adaptability

↓  
work with changing environments.

larger data

ML: discipline of AI which provides machines with the ability to automatically learn from data & past experiences (dataset) with minimal human intervention, for recognizing patterns/decision making.

### steps

- Pre-processing

- Training

- Post processing

Data is extracted from database & the knowledge extracted is stored in knowledge ~~base~~ base.

- Data warehouse :- Data integration technique to store all data in one place.

- Normalisation :-

Domain → Datatype of feature

Eg:- Age:- NUMBER(3).

$$\text{z-score normalisation} = \frac{\text{Value} - \mu}{\sigma}$$

To prevent priority of attributes having larger range of values.

Normalisation prevents bias of attributes having larger values.

### Feature extraction:-

Features can be used to distinguish objects from each other.

corresponding to every object, we need to extract features, based on domain knowledge. Most important features required to define the model.

feature selection:- Important subset of features required for the dataset.

Dimension reduction:- generates new features with the help of the existing features. (diff. to selection)

e.g:- Principal Component Analysis (PCA)

The reduced dataset is called Experience.

### Learning

- Supervised
- Unsupervised
- Semi-supervised
- Reinforcement

### Performance metrics

Accuracy, Precision, Recall, F-score value, Sensitivity.

Ensemble:- Random forest, Bagging, Boosting. Techniques

### Compare classification Algos

- Efficiency (Time & Space) & Effectiveness

If performance metrics are not satisfactory, we need to return to learning stage.

### Feature set types

- conditional features (other features)
- Decision features / class variable (single feature)

class variable having 2 classes is  $\frac{eg}{(Yes/No)}$  called binary classifier.

If class variables are not present, it is unsupervised learning.

every ~~MLP~~ is a ~~data~~ but not vice versa.

### Types of ML problems

#### 1) Supervised

Classification  
(DISCRETE)

Regression  
(CONTINUOUS)

Response/  
Random/  
Dependent  
variable

Regressor/  
Independent  
variable.

learning based on  
objective function, say  
Sum of Squares Error (SSE)  
Objective: To minimise SSE.

$$y_{pred} = \hat{m}x + \hat{c}$$

- Training based on  $(x, y)$  value pairs.
- Prediction on  $x$ -values to obtain  $y$ -values

If two features A & B are dependent.

A  $\rightarrow$  dependent on B

If A has missing values, it can be imputed/estimated from values that are present.

missing values may be replaced through

- Mean of remaining objects
- mean of remaining objects of a similar class (groupby or a feature having class values)

### Discretization

Continuous to discrete values.

- Results in loss of Information
- Sometimes it provides better decisions
- Actual values are being replaced.

~~Rough set~~  $\downarrow$  Discretization is read for soft computing  
example of soft computing :- ANN, rough-set theory

### Simple linear regression

$$Y = A + BX$$

### Multiple linear regression

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{n-1} x_{n-1}$$

### Unsupervised learning

Association  
(Market basket analysis)

Hard partitioning  $\rightarrow$  one object in one cluster

soft clustering  $\rightarrow$  one object may be present in multiple clusters with different membership values

Density based clustering  $\rightarrow$  DBSCAN

Hierarchy based clustering  $\rightarrow$  DBSHIFT, agglomerative

partition based clustering  $\rightarrow$  K-means

partition :- a set (group) of objects having similar characteristics.

Dividing Universal set into many subsets

$A_1, A_2, \dots, A_j$  such that

$$A_i \cap A_j = \emptyset$$

$$A_1 \cup A_2 \cup \dots \cup A_j = U$$

### Mapping :-

$$A \rightarrow B$$

mapping from A to B

$$f(x) = y$$

Salary (employee-object) = salary-value

salary :  $U \rightarrow \text{Domain(salary)}$

Prev. DBMS models  
Hierarchical (Tree structure)  
Network (Graph structure)

Relation R: subset of cartesian product.

Table is the subset of the cartesian products of the domains of the features

Table)  $A \times B \times C$ .  
 $D \rightarrow \text{domain}$   
 $\times \rightarrow \text{C.P.}$

$$T_1 \subseteq D(A) \times D(B) \times D(C)$$

30.07.2024

### Semi-Supervised learning

Training on initial labelled data is specific to the region of samples. In order to generalise it, we need to further train it upon unlabelled data.

Conditional Feature

Target feature - /  
Class Variable

A model is said to have learned, if it gives minimum error on the objective function for the training data itself.

Training performance: Measure of performance on training data should be high (near 100%)

### process

- Exploring the data - to find relationships b/w data
- Pre-processing - Normalising / Impudring missing values / Smoothening / Discretisation

clustering  
Convex - k-means

Concave - Density based

More features independent to each other  $\Rightarrow$   
provide better accuracy.

Feature selection algo. is used to select relevant features. Small subset of features is chosen.

[Study extra]

→ Data pre-processing concepts

Standardisation

Normalisation

Discretisation

Missing value estimation

[Long Ans Qs]

- what is ML?

- Motivation behind ML

- Types of ML

- How semi-supervised works?

Q1) Diff components of reinforcement learning.  
[Define each component]

# In Simple linear Regression

01/08/2024

$$y = ax + c + e$$

random error component which may arise due to error in the instrument.

- Minimising sum of square error (SSE) to estimate the values of  $a$  &  $c$ .
- Best-fitted line:- the curve which gives the minimum SSE.

SSR: Sum of square Residual

$$\sum (y_i - \hat{y}_i)^2$$

$$\left\{ \begin{array}{l} \text{Residual} = \text{Actual} - \text{Predicted} \\ \text{Error} = \text{Actual} - \text{Mean} \end{array} \right.$$



$$SS_{\text{Res}} = \sum_{i=1}^n (y_i - \hat{a} - \hat{c}x_i)^2$$

$$\frac{\partial SS_{\text{Res}}}{\partial a / \hat{a}, \hat{e}} = 0 \Rightarrow \sum (y_i - \hat{a} - \hat{c}x_i) = 0$$

$$\Rightarrow \sum y_i - \sum \hat{a} - \sum \hat{c}x_i = 0$$

$$\Rightarrow \sum y_i - n\hat{a} - \hat{c}\sum x_i = 0$$

$$\Rightarrow \hat{a} = \frac{\sum y_i - \hat{c}\sum x_i}{n}$$

$$\Rightarrow \hat{a} = \bar{y} - \hat{c}\bar{x} \quad (1)$$

~~( $y_i = \hat{a} + \hat{c}x_i$ )~~

$$\Rightarrow -2 \sum x_i (y_i - \hat{a} - \hat{c}x_i) = 0$$

$$\Rightarrow \sum x_i (y_i - \bar{y} + \hat{c}\bar{x} - \hat{c}x_i) = 0$$

$$\Rightarrow \sum x_i (y_i - \bar{y}) + \hat{c} \sum x_i (\bar{x} - x_i) = 0$$

Substituting value of  $\hat{a}$

$\bar{y}, \bar{x}$  = mean

$$\Rightarrow \sum x_i (y_i - \bar{y}) - \hat{c} \sum x_i (x_i - \bar{x}) = 0$$

$$\Rightarrow \frac{1}{\hat{c}} = \frac{\sum x_i (x_i - \bar{x})}{\sum x_i (y_i - \bar{y})}$$

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

$$= \frac{\sum x_i^2 - \frac{\sum x_i}{n} \sum x_i}{\sum x_i y_i - \frac{\sum y_i}{n} \sum x_i}$$

$$\Rightarrow \hat{c} = \frac{\sum x_i y_i - \frac{\sum y_i}{n} \sum x_i}{\sum x_i^2 - \frac{\sum x_i}{n} \sum x_i}$$

$$\boxed{\hat{c} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}} \quad (2)$$

Putting (2) in (1)

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

$$= \frac{\sum y_i}{n} - \frac{\sum x_i}{n} \left[ \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \right]$$

$$= \frac{\sum y_i (n \sum x_i^2 - (\sum x_i)^2) - \sum x_i [n \sum x_i y_i - \sum x_i \sum y_i]}{n \{ n \sum x_i^2 - (\sum x_i)^2 \}}$$

$$= \frac{n \sum x_i^2 \sum y_i - (\sum x_i)^2 \sum y_i - n \sum x_i \sum x_i y_i + (\sum x_i)^2 \sum y_i}{n \{ n \sum x_i^2 - (\sum x_i)^2 \}}$$

$$\hat{a} = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n[\sum x_i^2 - (\sum x_i)^2]}$$

$$\hat{a} = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$\epsilon_i \rightarrow$  Random variable from normal distribution  
i.e. mean=0 & variance =  $\sigma^2$

$$y_i = a + cx_i + \epsilon_i$$

$$\begin{aligned} E(y_i) &= E(a) + E(cx_i) + E(\epsilon_i) \\ &= a + cx_i \quad \text{Mean} \\ &= a + cx_i \end{aligned}$$

$$\begin{aligned} \text{var}(y_i) &= \text{var}(a + cx_i + \epsilon_i) \\ &= \text{var}(a) + \text{var}(cx_i) + \text{var}(\epsilon_i) \\ &= 0 + 0 + \sigma^2 \\ &= \sigma^2 \quad \text{Variance} \end{aligned}$$

### Multiple linear Regression

If one variable is not sufficient to explain the variation of the dependent variable, then we may require other <sup>independent</sup> variables -

Dimensions of matrices:-

$$Y \text{ } 2 \times n = n \times 1$$

$$\beta = K \times 1$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1K} \\ 1 & x_{21} & x_{22} & \dots & x_{2K} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nK} \end{bmatrix} \quad \text{dimensions } n \times K$$

$$Y = X\beta + \epsilon$$

$$SS_{\text{Res}} = \sum e_i^2 = \sum (y_i - \hat{y}_i)$$

$$\hat{y}_i = \hat{b}_0 + \hat{b}_1 x_{i1} + \dots + \hat{b}_{K-1} x_{i,K-1}$$

Upon partial derivative, we get  
 $K$  no. of equations, having  $K$  unknowns.

$$SS_{\text{Res}} = \sum e_i^2 = e^T e$$

$$[e_1 \ e_2 \ \dots \ e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

$$= (Y - \hat{Y})^T (Y - \hat{Y})$$

$$= (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

$$= (Y^T - \hat{\beta}^T X^T) (Y - X\hat{\beta})$$

$$= Y^T Y - Y^T X \hat{\beta} - \hat{\beta}^T X^T Y + \hat{\beta}^T X^T X \hat{\beta}$$

There is no epsilon term

$$Y = X\hat{\beta} + \epsilon$$

$$Y^T = 1 \times n$$

$$Y = n \times 1$$

$$X^T = n \times K$$

$$X = K \times n$$

all are  $1 \times 1$   
so all are scalar q'tys.

$$\begin{aligned} Y^T \times \beta &= (X^T Y)^T \hat{\beta} \\ &= (\hat{\beta}^T (X^T Y))^T \\ &= (\hat{\beta}^T X^T Y)^T = \hat{\beta}^T X^T Y \end{aligned}$$

2nd & 3rd terms are same.

$$SS_{\text{Res}} = Y^T Y - 2 \hat{\beta}^T X^T Y + \beta^T X^T X \hat{\beta}$$

$$\frac{\partial SS_{\text{Res}}}{\partial \hat{\beta}} = 0.$$

$$\Rightarrow Y^T Y - 2 X^T Y + 2 X^T X \hat{\beta}$$

$$\begin{aligned} \Rightarrow X^T Y &= X^T \hat{\beta} \\ \hat{\beta} &= (X^T X)^{-1} X^T Y \\ &\quad \text{estimated values} \\ &= \underbrace{(K \times 1)}_{(K \times K)(1 \times K)} \underbrace{(K \times K)}_{= K \times K} - K \times n \\ &= (K \times n) - n \times 1 \\ &= (K \times 1) \end{aligned}$$

~~SS Res~~

For every model, find the degree of freedom.

No. of eqns =  $K$  (constants)

No. of observations =  $n$ .

- Values that can be chosen with freedom =  $n - k$
- Values that need to be selected in

accordance to the constraints =  $k$   
(i.e. in order to satisfy the constraints)

Degree of freedom of SS<sub>Res</sub> model =  $\underline{n-k}$   
(residual)

$$SS_{\text{Res}} = \sum (y_i - \bar{y}_i)^2$$

$$SS_T = SS_{\text{Total}} = \sum_{i=1}^n (y_i - \bar{Y})^2$$

Taking partial derivative

$$\Rightarrow 2 \sum (Y_i - \bar{Y}) = 0$$

$$\Rightarrow \boxed{\sum (Y_i - \bar{Y}) = 0} \leftarrow \text{constraint}$$

Degree of freedom of  $SS_T = n-1$

$n-1$  values can be chosen with freedom. 1 value needs to be chosen in such a way so that it satisfies the constraint.

$$\begin{aligned} SS_T &= \sum (Y_i - \cancel{2 \bar{Y} \cdot \bar{Y} + \bar{Y}^2}) \\ &= \sum Y_i^2 - 2 \bar{Y} \cdot \sum Y_i + n \bar{Y}^2 \\ &= \sum Y_i^2 - 2n \bar{Y}^2 + n \bar{Y}^2 \\ &= \sum Y_i^2 - n \bar{Y}^2 \end{aligned}$$

$$SS_T = SS_{\text{Reg}} + SS_{\text{Res}}$$

Deg. of freedom (DF)

(n-1)

(cancel)

(n-k)

$$DF(SS_{\text{Reg}}) = k-1$$

$\left. \begin{array}{l} SS_{\text{Reg}} = \text{Regressor model} \\ = SS_{\text{Regressor}} \end{array} \right\}$

where  $K$  is the no. of coefficients.

SS<sub>reg</sub> explains the variability in the response variable

SS<sub>reg</sub> → error term → which does not explain variation in response variable, i.e. not expressed through regressor variable

SS<sub>r</sub> → ~~total~~ explains total variability in the regressor & residual variables.

∴ NO. of regressor variables = Degree of freedom. (SS<sub>reg</sub>)

6/8/24

SLR } Linear Regression

MLR

Non-linear Regression

1. Polynomial reg.
2. Linearisation of non-linear problems
3. Nonlinear estimation by Lsm

LSM → least square method.

POLYNOMIAL REG.

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_{k-1} X^{k-1} + \epsilon$$

X	Y
$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_i$	$y_i$
$\vdots$	$\vdots$
$x_n$	$y_n$

$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k + \epsilon_i \quad (1)$$

$\forall i = 1, 2, \dots, n$

$$x_i^j = z_{ij}$$

then (1) can be written as

$$Y_{i,j} = \beta_0 + \beta_1 z_{i,1} + \beta_2 z_{i,2} + \dots + \beta_j z_{i,j} + \beta_{k+1} z_{i,k+1} + \epsilon_i$$

which is similar to mLR.

$$\Rightarrow Y = \beta Z + \epsilon$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{k+1} \end{bmatrix}_{(k+1) \times 1}$$

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}_{n \times 1}$$

$$Z = \begin{bmatrix} 1 & z_{1,1} & z_{1,2} & \dots & z_{1,k+1} \\ 1 & z_{2,1} & z_{2,2} & \dots & z_{2,k+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & z_{n,1} & z_{n,2} & \dots & z_{n,k+1} \end{bmatrix}_{n \times (k+1)}$$

~~No. of terms~~ parameter

$$Y = Z \beta + \epsilon$$

$n \times 1$        $n \times k$        $k+1$        $n \times 1$

$\hat{Y} = Z \hat{\beta}$  (estimated model)

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_{k+1} \end{bmatrix}_{(k+1) \times 1}$$

$$\hat{\beta} = (Z^T Z)^{-1} Z^T Y$$

$\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$   $\downarrow$

$k \times n \quad n \times k \quad k \times k \quad k \times n \quad n \times 1$

$= (k+1) \times 1$

~~1~~

updated Z matrix

$$Z = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^{K-1} \\ 1 & x_2 & x_2^2 & \cdots & x_2^{K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_i & x_i^2 & \cdots & x_i^{K-1} \\ 1 & x_n & x_n^2 & \cdots & x_n^{K-1} \end{bmatrix}_{n \times K}$$

## 2 Linearisation of Non linear eqn

Dataset & Non linear eqn given

$$y = \frac{px}{q+x}$$

Dataset

x	y
x <sub>1</sub>	y <sub>1</sub>
x <sub>2</sub>	y <sub>2</sub>
⋮	⋮
x <sub>n</sub>	y <sub>n</sub>

$$\frac{1}{y} = \frac{q+px}{px}$$

$$= \frac{x}{px} + \frac{q}{px}$$

$$\Rightarrow \frac{1}{y} = \frac{1}{p} + \frac{q}{p} \cdot \frac{1}{x}$$

$$\text{let } Y = \frac{1}{y}$$

$$Y = \beta_0 + \beta_1 X \text{ where}$$

$$Y = \frac{1}{y}, X = \frac{1}{x}, \beta_0 = \frac{1}{p}, \beta_1 = \frac{q}{p}$$

$$\hat{Y} = X \hat{\beta} ; K=2$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y,$$

where

$$X = \begin{bmatrix} 1 & x_{10} \\ 1 & x_{20} \\ \vdots & \vdots \\ 1 & x_{n0} \end{bmatrix}_{n \times K} = n \times 2$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$$

$$2 \times n \quad n \times 2 \quad 2 \times n \quad n \times 1$$

$$= 2 \times 1$$

$$X = \begin{bmatrix} 1 & y_{x1} \\ 1 & y_{x2} \\ \vdots & \vdots \\ 1 & y_{xn} \end{bmatrix}_{n \times 2}$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} y_{y1} \\ y_{y2} \\ \vdots \\ y_{yn} \end{bmatrix}_{n \times 1}$$

## 3 Least square method

$$y = a + ce^{-bx}$$

$$S_{\text{Res}} = \sum_{i=1}^n (Y_i - \hat{Y})^2$$

$$= \sum \left\{ Y_i - (\hat{a} + \hat{c} e^{-\hat{b} x_i}) \right\}^2$$

$$= \sum (Y_i - \hat{a} - \hat{c} e^{-\hat{b} x_i})^2$$

Three parameters.

$$\frac{\partial S_{\text{Res}}}{\partial \hat{a}} = \frac{\partial S_{\text{Res}}}{\partial \hat{b}} = \frac{\partial S_{\text{Res}}}{\partial \hat{c}} = 0 \quad \left. \begin{array}{l} \text{to obtain} \\ \hat{a}, \hat{b} \text{ & } \hat{c} \end{array} \right\}$$

## Logistic Regression

- Supervised learning
- Binary classification (2 classes)
- Probability based classification.

Sigmoid function

$$h_{\theta}(x) = \sigma(z) = \frac{1}{1 + e^{-z}} \quad -\infty \leq z \leq \infty \quad 0 \leq \sigma(z) \leq 1$$

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} & x_2^{(i)} & \dots & x_n^{(i)} \end{bmatrix}^T \quad F_1, F_2, \dots, F_n \quad \Delta \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \quad D \rightarrow \text{decision variable}$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}_{n \times 1} \quad \text{parameters}$$

$$x^{(i)} = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}_{n \times 1}$$

$$\theta^T x^{(i)} = \sum_{j=1}^n \theta_j x_j^{(i)}$$

If  $x_0^{(i)} = 1$  then  
MLR

08/08/2024

$$h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}} \quad \text{Range} = 0 \text{ to } 1$$

$h_{\theta}(x^{(i)}) = \text{Prob}(y^{(i)} = 1/x^{(i)}; \theta)$   
 $\text{prob}(y^{(i)} = 0/x^{(i)}; \theta) = 1 - h_{\theta}(x^{(i)})$   
 $\text{prob}(y^{(i)} = k/x^{(i)}; \theta) = [h_{\theta}(x^{(i)})]^{y^{(i)}} \cdot [(1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}]$   
 where  $k = 0 \text{ or } 1$   
 probability that the sample  $x^{(i)}$  of class  $k$

Developing objective function using likelihood method

$$P_i = [h_{\theta}(x^{(i)})]^{y^{(i)}} \cdot [(1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}]$$

probabilities corresponding to samples  $x_1, x_2, \dots, x_m$  are  $P_1, P_2, \dots, P_m$ .

$$L = \text{likelihood of the model} = \prod_{i=1}^m P_i$$

$$\log L = \sum_{i=1}^m \log P_i$$

we have to maximize obj. fn.

$$\max_{\theta} \log L = \sum_{i=1}^m \log P_i$$

$$\min_{\theta} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \log P_i$$

cost fn

### Hypothesis

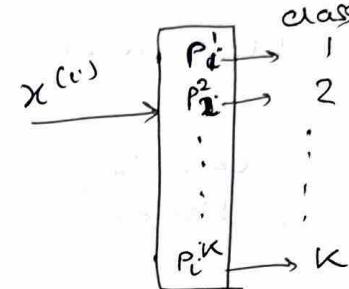
$$h_{\theta}(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}}$$

$x^{(t)}$   
unknown sample.

Training  
Estimates  
the value  
of  $\theta$ .

### Softmax Regression

Extended version of logistic regression to handle  
more than 2 (multi) classes (K).



Every class has  
parameters  $\theta_K$ .

if  $h_{\theta}(x^{(t)}) \geq 0.5 \Rightarrow y^{(t)} = 1$   
else  $y^{(t)} = 0$

minimise

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))]$$

$h_{\theta}(x^{(i)}) \rightarrow \text{sigmoid fn.}$

- There are n parameters:  $\theta_0, \theta_1, \dots, \theta_n$

Minimise  $J(\theta)$  using Gradient Descent approach

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Repeat

- until  $\theta$  values converge
- No. of iterations completed.

To prevent overfitting of train dataset,  
we use regularisation (2 types)

- L1 (Lasso)  $\lambda \sum_{j=1}^p |w_j|$

- L2 (Ridge)  $\lambda \sum_{j=1}^p w_j^2$

### Hypothesis

$$h_{\theta}(x) = \begin{bmatrix} p(y=1|x; \theta) \\ p(y=2|x; \theta) \\ \vdots \\ p(y=K|x; \theta) \end{bmatrix}$$

consider maximum  
probability for  
determination of  
class.

### Cross-entropy loss

$$-\sum_{c=1}^C y_c \log(p_c)$$

$\uparrow$  true class label       $\uparrow$  probability value

### Cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(p(y=k|x^{(i)}; \theta))$$

For each  $i$ ,  $y_k^{(i)} \geq 1$  only once.

Total of  $m$  values are summed.

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

Decision boundary  $\rightarrow$  Non-linear

[L1  $\rightarrow$  single summation  $\rightarrow$  as it takes mod value of each vector]

[L2  $\rightarrow$  double summation  $\rightarrow$  for each class, each term. (for 1 to m; for 1 to k)]

### One vs Rest

category 1  $\leftarrow$  1  
 $(M_1)$

Category 0  $\leftarrow$  Rest.

↓

class 2  $\rightarrow$  class 2 as Category 1  
 $(M_2)$

Rest = 0

No of models = K

$(M_1, M_2, \dots, M_K)$

$$x^{(t)} \rightarrow \begin{cases} M_1 = 0.4 \\ M_2 = 0.3 \\ \vdots \\ M_K = 0.7 \end{cases}$$

} check which has max. probability.  
(say  $M_5 = 0.8$ )

### Inference

it belongs to class 5,  
else  $x_i \notin$  class 5.

If  $x_i \notin$  class 5, remove all objects corresponding to class 5 and redo with  $K-1$  classes.

### liblinear solver

- Optimisation method.
- for sparse datasets

Hyperparameters - tuned using validation set (10)  
parameters - tuned using train set (80)

[model = LogisticRegression (solver = "name")]

### Exam

Q) what is logistic Regression (LR)?

Q) what are the diff types of LR:

binary, softmax (give all steps for constructing the models)

Q) How overfitting problem can be solved in LR?  
(L1 & L2)

Q) what are the diff solvers & compare them.

MCQs also

## Bayesian classification

13/08/2024

$$P(H|E) = \frac{P(H) * P(E|H)}{P(E)}$$

↑ likelihood      ↓ Prior  
↑ hypothesis      E → evidence  
↓ posterior      ↓ probability of sample

$P(E)$  is always constant wrt all classes.

$P(G_i|x)$  for all  $i=1$  to  $m$ . ( $m$  classes).

since  $P(E)$  is constant, we can ignore it.

max  $P(H|E)$  involves maximising

$$P(H) * P(E|H)$$

$X$  ←  $n$  dimensional matrix corresponding to  $n$  different features.

$$P(X, G_i) = \prod_{n=1}^N P(x_n, G_i)$$

ATTRIBUTES  $A_1, A_2, \dots, A_k, \dots, A_m$

$$X \rightarrow [x_1 \ x_2 \ \dots \ (x_k) \ \dots \ x_m]$$

$$P(x_k | G_i) \quad \forall k = 1, 2, \dots, m$$

No. of classes =  $G_m$ .

$A_k$	$x_k > 10$	$G_i$
	$x_k = 10$	$G_i$
	$x_k = 0$	$G_i$

Discrete data

$|C_{i,n}| \leftarrow$  No. of tuples in  $C_i$ , which belong in  $D$

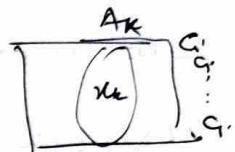
$P(x_k | C_i) =$  Given class  $C_i$ , how many values have  $A_k$  (Attribute value) equal to  $x_k$  (say 10)

Continuous data

→ Gaussian (Normal distribution)

$\mu_{G_i}$  ← mean

$\sigma_{G_i}$  ← standard deviation



Need to compute probability of each attribute, given  $C_i$   
then product of all gives  $P(x_k | C_i)$

Demerit of Bayesian classifier

- Each attribute must be independent to each other wrt class variable.
- if one class = 0, then probability calculation would yield zero.  
Hence, each class has 1 added to it through Laplacean correction.

so not too large  
but not zero

test on some sample & make it work  
for validation purpose & then go to next chapter

## KNN classifier

Supervised learning algorithm

→ lazy learning algo (No training set)

There is no early training.

- Training is during computation of class value

- No need to assume anything within dataset.

KNN can be used both for regression & classification



majority vote  
of nearest  
neighbours.



avg. of all  
values ~~to be~~ is  
predicted

## Disadv of KNN

1) Computationally expensive

2) Curse of dimensionality

→ Increasing search space (exponentially)

- For high dimensions, distance becomes equal for all pairs, i.e.

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2}$$

all equal distance metric becomes invalid.

- Increased search space means that it may not be representative of the actual data.

This requires Dimensionality reduction.

- High dimensionality also includes many irrelevant features. Overfitting of model takes place. Increased time complexity.

feature selection selects a subset of the existing features.

Dimensionality reduction reduces dimensions by generating new dimensions from existing features. The data of features are changed.

$$k = \sqrt{\text{NO. of objects}}$$

← determining optimal value of k.

## Decision Tree Model

20/08/2024

Supervised learning

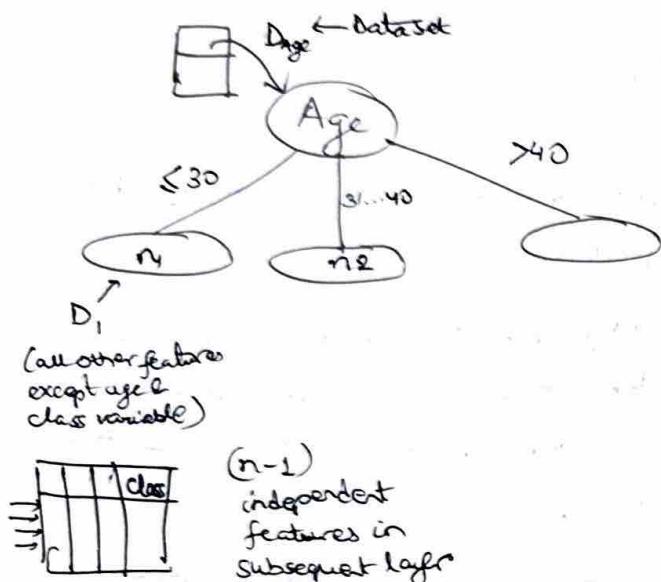
Quinlan Iterative Dichotomizer 3 (ID3)

Constructing a tree is exponential due to varying configurations of the nodes.  
Greedy heuristic is used to solve the problem.

### Hunt's Algorithm

First, all records  $D_t$  are associated with root node. The dataset is broken into subsets with each child node.

Based on values of the attribute in the node in question, each value would be treated as a different branch.



### Limitation of Decision Tree

- All the attributes need to be of discrete values.  
All ~~cate~~ continuous variables need to be discretised.
- If any node has only 1 feature in its associated dataset, then majority voting of class variable is considered.
- In the absence of multiple samples (only one), then we use that class value.
- minimise entropy, maximise info gain

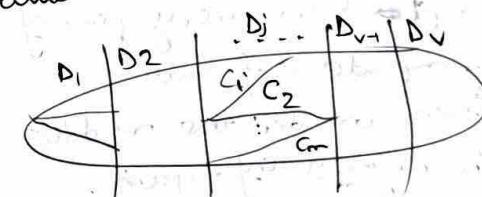
If probability values of each class are similar, it leads to higher uncertainty.

Let corresponding to database

	class	
Yes	9	5
No		
Total no of samples in dataset		

$$P_i = \frac{\text{No. of samples corresponding to class } i}{\text{Total no of } \cancel{\text{class}} \text{ samples in dataset}}$$

If any attribute A has v no. of values, then the dataset is divided into v no. of subsets.



No. of classes = m

Each such partition may have further internal partitions, corresponding to class variable.

$$\text{Info}(D) = - \sum_{i=1}^m P_i \log_2(P_i)$$

$$\forall j = 1 \dots v$$

$\text{Info}(D)$  based on attribute A is a weighted average.

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

Info gain for a particular attribute A:

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

If uncertainty is more then gain is less & vice versa.

More gain indicates a better split.

After construction of decision tree, the rules are derived.

Discretizing continuous valued attributes



Consider mid pt b/w every pair of values & compute info gain for each. Then consider the middle value with highest info gain.

for every no of requirement of discrete values, the loop is repeated that many times.

22/08/2024

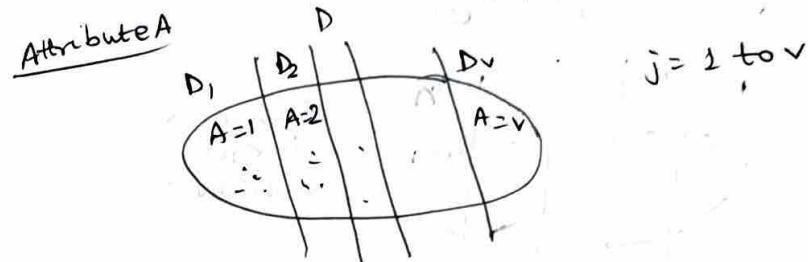
### Disadv of ID3

In ID3, if an attribute has many different discrete values then it will give more info gain (bias) towards that attribute.

Each class becomes correctly classified, in such a scenario & purity is increased. This bias is undesirable.

### C4.5 (ID3<sup>↑</sup>)

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



$$\text{GainRatio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)}$$

If Attr. A has large no. of diff. values, then splitInfo value is large & its importance reduces.

### Gini Index.

$$\text{gini}(D) = 1 - \sum_{j=1}^n p_j^2$$

calculated corresponding to target variable.

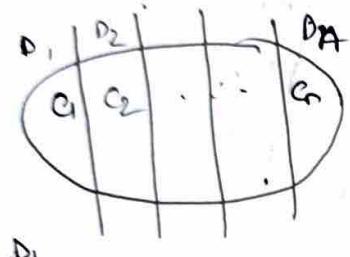
$p_i$  = probability that a randomly selected variable belongs to class  $C_i$ .

Gini index is the measure of inequality.

Joint probability of selecting 2 samples such that both belong to the same class (Independent) =  $P_i \cdot P_i = P_i^2$

$\sum p_j^2$  is thus a measure of purity or equality.

Hence  $1 - \sum p_j^2$  gives the measure of impurity or inequality.



Based on target variables



Based on attributes

Entropy  $\rightarrow$  dataset split based on feature

Attribute with highest gini index is chosen for further splits.

### Different types of overfitting

Overfitting  $\rightarrow$  model is more specific to train dataset.

Inclusion of noise in train dataset is learnt by the machine.

Two techniques to avoid overfitting

### Pruning

- Setting threshold value of maximum information gain for any attribute at any level.  
(Value needs to be set experimentally)

### Postpruning

- After construction of entire tree, some branches from any level are removed.

Replace the node by a leaf node & check performance against new dataset.

On old dataset, post pruning would give low accuracy for it.

If the measured performance is similar or better than the change is made.

Hyperparameter (ht. of decision tree) is being set here.

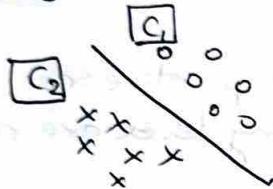
## Support Vector Machine

→ Linear Classifier → Supervised Model

Separates the objects of diff. classes.

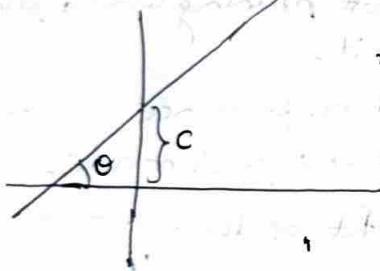
Let us consider 2 classes  $C_1$  and  $C_2$

Line/Curve by which we can discriminate the objects



$$y = mx + c$$

$m = \tan \theta$  orientation  
 $c \rightarrow$  position of classifier  
 $\theta \rightarrow$  orientation of classifier



$$Y = w^T x + b$$

$$w_0 x_0 + w_1 x_1 + \dots + w_n x_n + b$$

$$\begin{matrix} w_1 \\ \downarrow \\ x_1 \\ \vdots \\ b \end{matrix}$$

$$\begin{matrix} w_0 \\ \downarrow \\ x_0 \end{matrix}$$

bias term  
2 diff. features

$$X = [x_0 \ x_1 \ x_2]$$

$$w = [w_0 \ w_1 \ w_2]$$

$$w^T = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$$

; weightage of the features.

$$\begin{aligned} w^T x &= w_0 x_0 + w_1 x_1 + w_2 x_2 \\ &= \sum_{i=0}^2 w_i x_i \end{aligned}$$

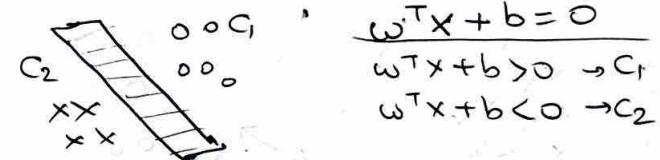
$x \leftarrow$  feature vector

$$g(x) = w^T x + b$$

if there are 2 features (st. line)

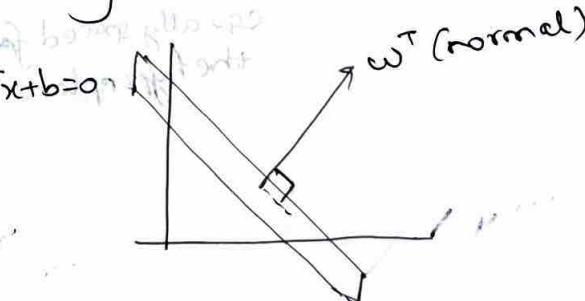
$$w^T x + b = 0$$

if there are 3 features (plane)



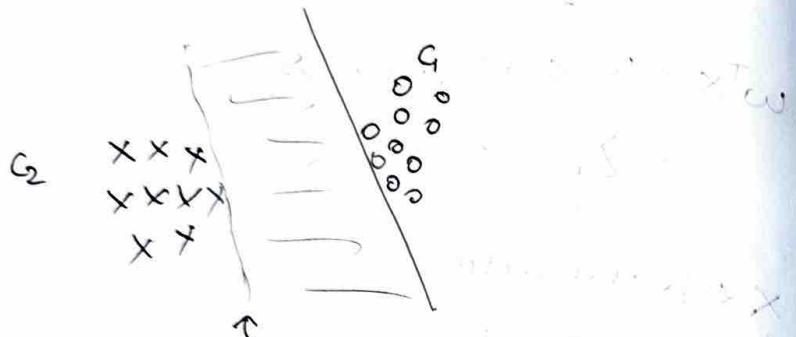
If dimensions are more than 3 (hyperplane)

- for any sample  $x \in C_1, w^T x + b > 0$
- for any sample  $x \in C_2, w^T x + b < 0$



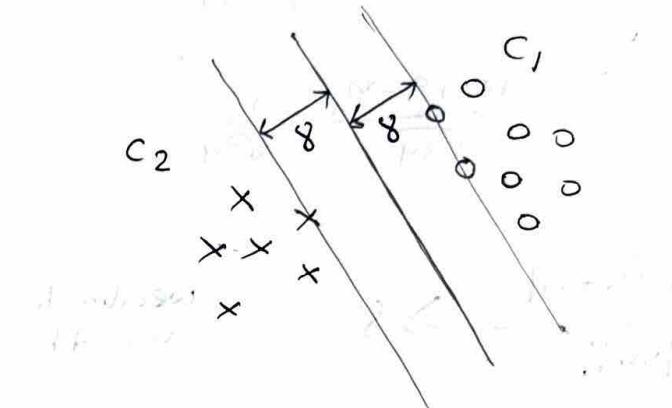
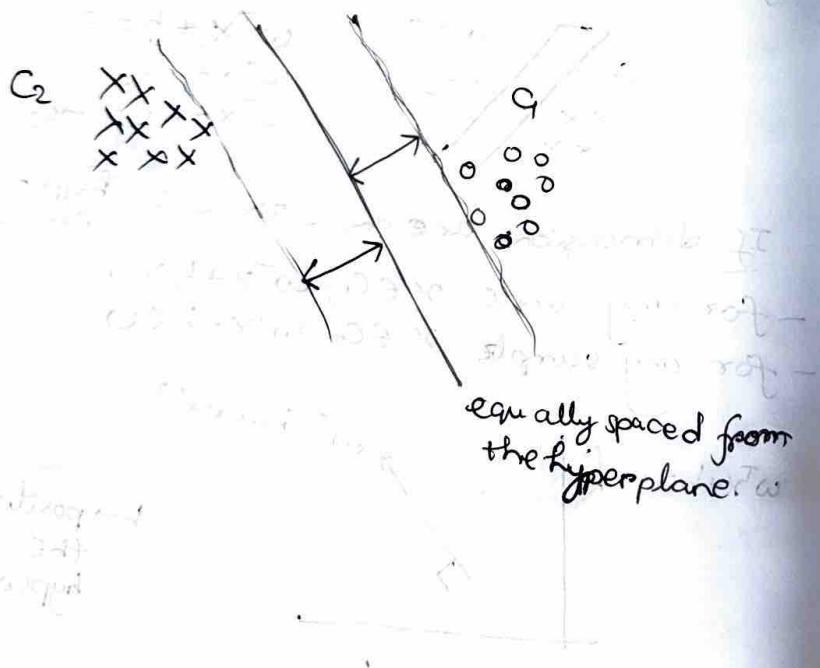
$b \rightarrow$  position of the hyperplane.

Many such hyperplanes can be obtained for correct classification.



max. margin to objects of  $C_2$  &  
min margin to obj of  $C_1$  &  
& vice versa.

So, we need to give max. margin to both.



The margin depends on the extreme objects, because removal of them would affect the marginal space (increases)  
Hence, they are called the Support Vectors.

$$\text{let, } C_1 \Rightarrow y_i = +1$$

$$C_2 \Rightarrow y_i = -1$$

$$\left\{ \begin{array}{l} w^T x_i + b > 0 \text{ if } x_i \in C_1 \\ \Rightarrow y_i = +1 \end{array} \right. \quad (1)$$

$$\left\{ \begin{array}{l} w^T x_i + b < 0 \text{ if } y_i = -1 \end{array} \right. \quad (2)$$

$$y_i(w^T x_i + b) > 0$$

Combining (1) and (2)

In order to fix the position of the hyperplane.

$$w^T x_i + b \geq 8 \text{ if } y_i = +1$$

$$w^T x_i + b \leq -8 \text{ if } y_i = -1$$

eqn of  
hyperplane

$$(2, 3)$$

$$5x + 3y - 2 = 0$$

Dist of pt from  
plane.

$$\frac{5 \times 2 + 3 \times 3 - 2}{\sqrt{5^2 + 3^2}} = \frac{10 + 9 - 2}{\sqrt{34}} = \frac{17}{\sqrt{34}}$$

$$d = \frac{|w^T x + b|}{\|w\|} \geq 8$$

we estimate  
 $w$  and  $b$

$$\frac{|w^T x + b|}{\|w\|} \geq 8$$

Left hand side is  
length of perpendicular  
distance from point to line

$$-8 \geq \frac{w^T x + b}{\|w\|} \geq 8$$

Left hand side is  
length of perpendicular distance  
from point to line

$$w^T x + b \geq \|w\| 8 \geq 1$$

$$\text{if } 8 > \frac{1}{\|w\|}$$

$$\frac{w^T x + b}{\|w\|} \leq -8$$

$$\Rightarrow w^T x + b \leq \|w\| 8 \leq -1$$

$$\text{Left hand side is length of perpendicular distance from point to line}$$

$$w^T x + b \leq \|w\| 8 \leq d + x^T w$$

$$w^T x + b \leq \|w\| 8 \leq d + x^T w$$

To solve  
optimization

$$\begin{cases} w^T x + b \geq 1 & \text{if } x \in C_1 \\ w^T x + b \leq -1 & \text{if } x \in C_2 \end{cases} \quad \begin{array}{l} y_i = +1 \\ y_i = -1 \end{array}$$

$$8 \geq \frac{1}{\|w\|}; \quad \|w\| 8 \geq 1$$

$$y_i (w^T x_i + b) \geq 1 \quad (\text{equal to 1 at support vectors})$$

$x_i \rightarrow$  feature vector  
 $y_i \rightarrow$  class label of feature vectors

~~Maximise~~  $\frac{2}{\|w\|}$  ← total margin

$$\text{Maximise } \frac{2}{\|w\|}$$

$\equiv$  Minimise  $\frac{1}{2} \|w\|^2$

for simplicity:

$$\equiv \text{Minimise } \frac{1}{2} \|w\|^2 = \frac{1}{2} w^T w$$

Optimization problem

$$\text{Minimise } \frac{1}{2} w^T w$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1$$

$$\forall i = 1, 2, \dots, m$$

No. of constraints = m.

We cannot apply Gradient Descent here because it is a constrained problem. It can be done in unrestricted unconstrained prob. only.

we cannot apply LPP (linear program)  
because obj fn. is quadratic

Instead, we use QPP (quadratic program)  
we use Lagrangian multiplier.



### Techniques

- 1) Obj fn  $\rightarrow$  quadratic  
constraints  $\rightarrow$  linear (one) equality

eg)  $\min f(x, y) = 2x^2 + 2y^2$

s.t.  $g(x, y) : x + y - 1 = 0$

Applying Lagrangian multiplier ( $\alpha$ )

$$\min L(x, y, \alpha) = f(x, y) - \alpha g(x, y)$$

(unconstrained)

- 2) Obj fn  $\rightarrow$  quadratic  
constraints many

$$\begin{aligned} & g_1(x, y) \\ & g_2(x, y) \\ & \vdots \\ & g_m(x, y) \end{aligned} \quad \left\{ \begin{array}{l} \text{equality} \\ \text{minimum} \end{array} \right.$$

constrained optimization problem to unconstrained  
problem by adding Lagrangian multipliers and  
penalty terms to the objective function.

eg)  $\min f(x, y) = 2x^2 + 2y^2$   
s.t.  $\begin{cases} g_1(x, y) \\ g_2(x, y) \\ \vdots \\ g_m(x, y) \end{cases}$  } equality

Applying Lagrange multiplier. (to unconstrained)

$$\min L(x, y, \alpha) = f(x, y) - \sum \alpha_i g_i(x, y)$$

eg)  $\min f(x, y) = 2x^2 + 2y^2$   
s.t.  $g(x, y) \geq 0$

$\downarrow$   
 $\min L(x, y, \alpha) = f(x, y) - \alpha g(x, y)$  // can't be converted to unconstrained.  
s.t.  $\alpha \geq 0$

④ Due to m constraint.

$$\begin{aligned} & \min \frac{1}{2} w^T w \\ & \text{s.t. } y_i(w^T x_i + b) \geq 1 \\ & \quad \forall i = 1, 2, \dots, m. \end{aligned}$$

$\downarrow$  Lagrangian multiplier

$$\min L(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^m \alpha_i [y_i(w^T x_i + b) - 1]$$

$$\text{s.t. } \alpha_i \geq 0 \quad i = 1, 2, \dots, m$$



first consider dual then solve, i.e.  
convert primal to dual problem.

e.g) let primal problem be

$$\min_x x^2$$

$$\text{s.t. } x \geq b$$

Ans. Min  $L(x, \alpha) = x^2 - \alpha(x-b)$

$$\text{s.t. } \alpha > 0$$

Corresponding Dual problem

~~Max~~  $\max_{\alpha} d(x, \alpha) =$

$$\max_{\alpha} \min_x [x^2 - \alpha(x-b)]$$

$$\text{s.t. } \alpha \geq 0$$

$$\frac{\partial L}{\partial x} = 2x - \alpha = 0 \Rightarrow x^* = \frac{\alpha}{2}$$

Substituting value of  $x^*$ , we get min

$$x^*, \alpha) \max_{\alpha} \left[ \frac{\alpha^2}{4} - \alpha \left( \frac{\alpha}{2} - b \right) \right]$$

$$\frac{\partial L}{\partial \alpha} = \frac{\partial}{\partial \alpha} \left( \frac{\alpha^2}{4} - \alpha \left( \frac{\alpha}{2} - b \right) \right) = 0$$

$$\Rightarrow b = \frac{\alpha}{2}$$

$$\Rightarrow \alpha = 2b$$

optimal value  
 $\alpha^* = \max(0, 2b)$   
 $x^* = \max(0, b)$

Transfer obj from lower dimensional to higher dimensional space using kernel function.

with the hope that they would become linearly separable in higher dimension

# Converting to dual problem.

To prevent overfitting

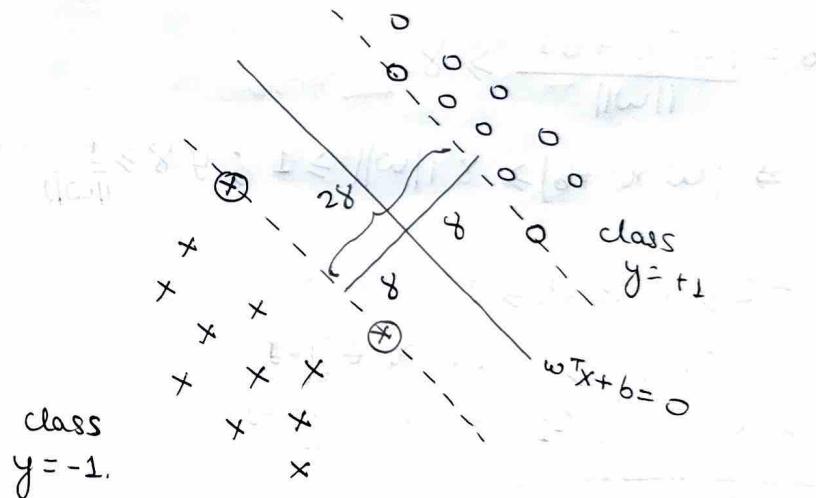
0 Soln  $\in C$ .

Regularisation  
Rectification/Regulation

Trade off bw margin & classification error.

SVM

29/08/2024



The points on the marginal (dotted) line are called Support Vector.

Max 28

$$\omega^T x + b = 0$$

for any  $x_i$  if  $\omega^T x_i + b > 0$   $\omega^T x_i + b < 0$

If  $\omega^T x_i + b \geq 8$   
if  $\omega^T x_i + b \leq -8$

$$\omega^T x + b = 0$$

$$\omega^T x_i$$

Distance of a pt from the plane

$$d = \frac{|\omega^T x_i + b|}{\|\omega\|} \geq 8$$

$$\Rightarrow |\omega^T x_i + b| \geq 8 \|\omega\| \geq 1 ; \text{ if } 8 \geq \frac{1}{\|\omega\|}$$

$$-1 \geq \omega^T x_i + b \geq 1$$

$\omega^T x_i + b \geq 1$  if  $x_i \in y_i = 1$   
 $\omega^T x_i + b \leq -1$  if  $x_i \in y_i = -1$

$$y_i(\omega^T x_i + b) \geq 1$$

← constraint

$$\text{Max } 28 = \frac{2}{\|\omega\|}$$

$$x_i \in (y=1)$$

~~$$x_i \in (y=-1)$$~~

$$x_i \in (y=1)$$

$$x_i \in (y=-1)$$

$$\text{Min } \frac{1}{2} \|\omega\|^2 = \frac{1}{2} \vec{\omega} \cdot \vec{\omega}$$

s.t. constraints

$$y_i(\omega^T x_i + b) \geq 1$$

using Lagrangian multiplier ( $\alpha$ ):

$$L(\omega, b, \alpha) = \text{Min } \left( \frac{1}{2} \omega \cdot \omega - \sum_{i=1}^m \alpha_i [y_i(\omega^T x_i + b) - 1] \right)$$

$$\text{s.t. } \alpha_i \geq 0 ; i = 1, 2, \dots, m$$

Converting to dual problem,

$$\text{Max } \left( \frac{1}{2} \omega \cdot \omega - \sum_{i=1}^m \alpha_i [y_i(\omega^T x_i + b) - 1] \right)$$

$$\frac{\partial L}{\partial \omega} = 0$$

$$\frac{\partial L}{\partial b} = 0$$

$$L(\omega, b, \alpha) = \text{Min } \frac{1}{2} \omega \cdot \omega - \sum \alpha_i y_i \omega^T x_i - \sum \alpha_i y_i b + \sum \alpha_i$$

$$\frac{\partial L}{\partial b} = -\sum \alpha_i y_i = 0$$

$$\sum \alpha_i y_i = 0$$

$$\frac{\partial L}{\partial \omega} = \omega - \sum \alpha_i y_i x_i = 0$$

$$\sum_{i=1}^m \alpha_i y_i x_i = \omega$$

Putting ② in eqn ①,

$$L(\omega, \alpha, b) = \frac{1}{2} \omega \cdot \omega - \sum \alpha_i y_i \omega^T x_i + \sum \alpha_i$$

Putting ③ in eqn ①,

$$L(\omega, \alpha, b) = \cancel{\frac{1}{2} \omega \cdot \omega}$$

$$\frac{1}{2} \sum_{\substack{i=1 \\ j=1}}^m \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

$$- \sum_{\substack{i=1 \\ j=1}}^m \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

$$+ \sum \alpha_i$$

$$\Rightarrow L(\omega, \alpha, b) = \sum \alpha_i - \frac{1}{2} \sum_{\substack{i=1 \\ j=1}}^m \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

converting to dual problem

$$\text{Max} \left\{ \sum \alpha_i - \frac{1}{2} \sum_{\substack{i=1 \\ j=1}}^m \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j) \right\}$$

$$\text{s.t. } \alpha_i \geq 0 \quad ; i = 1, 2, \dots, m.$$

Solving this using Quadratic problem to get  $\alpha_i$ ,

if  $\alpha_i = 0 \Rightarrow x_i$  is inactive constraint

if  $\alpha_i > 0 \Rightarrow x_i$  is active/tight constraints

$x_i$  is support vector

$\vec{x}_i \rightarrow$  training samples

$y_i \rightarrow$  corresponding class

putting value of  $\alpha_i$  as obtained, in.

$\omega = \sum_{i=1}^m \alpha_i y_i \vec{x}_i$ ; we obtain value of  $\omega$ .

let for some  $x_k$ ,  $\alpha_k > 0 \Rightarrow y_k$

$$\omega^T x_k + b \geq 1 \text{ if } (y_k = +1)$$

$$\omega^T x_k + b \leq -1 \text{ if } (y_k = -1)$$

for support vector

$$\omega^T x_k + b = 1 \text{ if } (y_k = +1)$$

$$\omega^T x_k + b = -1 \text{ if } (y_k = -1)$$

$$\Rightarrow y_k = \omega^T x_k + b$$

$$\Rightarrow b = y_k - \omega^T x_k$$

It will give almost same value for every support vector.

$$\text{let } f(x) = \omega^T x + b$$

for any unknown sample  $z$ ,

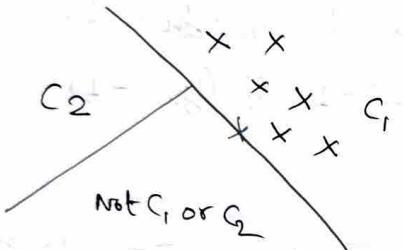
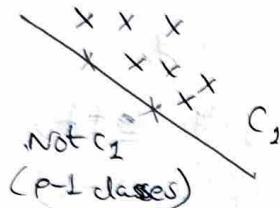
compute  $f(z)$  & check sign of  $f(z)$

if  $f(z) = +ve$  then  $+ve$  class

if  $f(z) = -ve$  then  $-ve$  class.

Training is over upon estimation of parameters successfully.

If no. of classes is more than 2  
Then we need to use multiple binary  
Support vector machines



$$\text{No. of binary SVMs to be constructed} = p-1$$

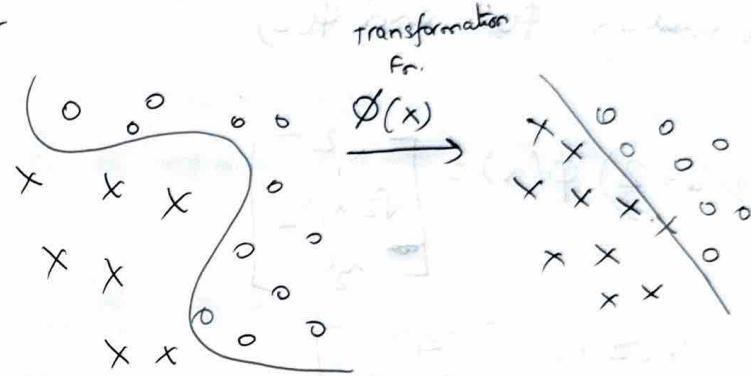
If objects are not linearly separable,  
the objects needs to be transformed  
from one feature space to another  
space (increasing dimension)

The hope is that at higher dimensional  
space, the objects may be linearly  
separable.

Kernels serve this purpose of increasing  
dimensionality of space.

Kernels serve as the black box of a  
transformation function.

Let



Kernel  $k(x, z)$

$$= \underbrace{\phi(x) \cdot \phi(z)}$$

on two training samples  
 $x$  and  $z$ .

Many kernel fns give the same value.  
We do not need to use them directly.

Types of Kernel functions

> Polynomial function

Let  $\phi(x)$  = polynomial function of degree  $d$ .

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad 2 \text{ features each.}$$

If  $d=1$

$$\text{then, } \phi(x) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \phi(z) = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

Usage of Kernel fns  
depends on complexity  
of dataset.

$$k(x, z) = \phi(x) \cdot \phi(z) = x_1 z_1 + x_2 z_2^2 \\ = \vec{x} \cdot \vec{z}$$

$\star k(x, z)$  can directly kernel function instead of individually computing  $\phi(x)$  and  $\phi(z)$

If  $d=2$ )  $\phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$

$$\phi(z) = \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix}$$

$$\begin{aligned} \phi(x) \cdot \phi(z) &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + \\ &\quad x_2^2 z_2^2 \\ &= (\vec{x} \cdot \vec{z})^2 \end{aligned}$$

for dimension  $d$ )  $\phi(x) \cdot \phi(z) = (\vec{x} \cdot \vec{z})^d$

This is considered as Kernel.

2) Polynomial of degree upto  $d$

$$k(x, z) = (\vec{x} \cdot \vec{z} + 1)^d$$

3) Radial / Gaussian Kernel # Generally used.

$$k(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

4) Sigmoid kernel

$$k(x, z) = \tanh(\beta \vec{x} \cdot \vec{z} + \gamma)$$

Linear Kernel: polynomial of degree 1.

$$\text{Max } \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j (\vec{x}_i \cdot \vec{x}_j)$$

$$\text{s.t. } \alpha_i \geq 0 ; i = 1, 2, \dots, m$$

$$\sum \alpha_i y_i = 0 \quad \leftarrow \boxed{\text{from } \frac{\partial L}{\partial b} = 0}$$

$$\Rightarrow \text{Max } \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\text{where, } K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

$$\text{s.t. } \alpha_i \geq 0, \sum \alpha_i y_i = 0$$

$$w = \sum \alpha_i y_i \phi(\vec{x}_i)$$

$$b = y_k - w^T \vec{x}_k$$

$$= y_k - w^T \phi(\vec{x}_k)$$

$$0 \leq \alpha_i < C$$

Regularisation parameter

$C$  helps to control overfitting.  
Value is generally kept very low.

$C \rightarrow$  gives tradeoff b/w

→ marginal space

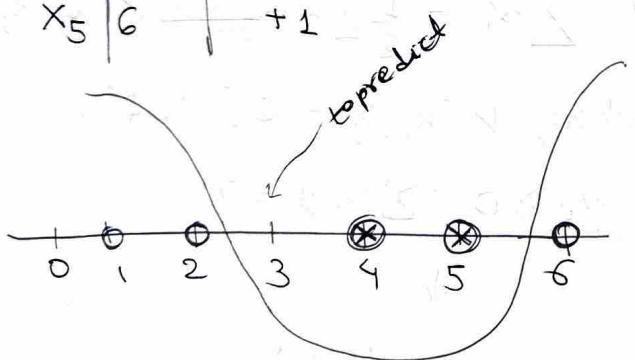
→ classification accuracy.

A high value of  $C$  leads to overfitting.

Problem

Q).

X	X	Y
$x_1$	1	+1
$x_2$	2	+1
$x_3$	4	-1
$x_4$	5	-1
$x_5$	6	+1



Non linear classifier

we need a kernel function,

$$K(X, Y) = (XY + 1)^2$$

let  $C = 50$ , or  $100$

model:

$$\text{Max } \sum_{i=1}^5 \alpha_i - \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

$$\Rightarrow \text{Max } \sum_{i=1}^5 \alpha_i - \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \alpha_i \alpha_j y_i y_j (x_i \cdot x_j + 1)^2$$

$$0 \leq \alpha_i \leq C, \quad \sum_{i=1}^5 \alpha_i y_i = 0$$

After solving, suppose we get the values of  $\alpha$  as follows: —

Lagrangian multipliers.

	X	Y	$\alpha$
$x_1$	1	+1	0
$x_2$	2	+1	2.5
$x_3$	4	-1	0
$x_4$	5	-1	7.3
$x_5$	6	+1	4.8

3 support vectors

Compute.

$$\omega = \sum \alpha_i y_i \Phi(x_i)$$

$$b = y_k - \omega \cdot x_k$$

Then compute  $f(z)$

$$f(z) = \omega^T z + b$$

$$\Rightarrow f(z) = \sum \alpha_i y_i \phi(x_i) \phi(z) + b$$

$$\Rightarrow f(z) = \sum \alpha_i y_i k(x_i, z) + b$$

$$\Rightarrow f(z) = \sum \alpha_i y_i (x_i z + 1)^2$$

$$\begin{aligned} &= (2.5)(1)(2z+1)^2 \\ &+ (7.3)(-1)(5z+1)^2 \\ &+ (4.8)(1)(6z+1)^2 + b \end{aligned}$$

1st & 3rd  
terms are  
zero  
(from table)

$$\Rightarrow f(z) = 0.67 z^2 - 5.33z + b$$

$$b = y_k - \downarrow w \phi(x_k)$$

$$= y_k - \left( \sum \alpha_i y_i \phi(x_i) \right) \phi(x_k)$$

~~$\cancel{\sum \alpha_i y_i}$~~

$$\begin{aligned} &= y_k - [(2.5)(1) \phi(2) + (7.3)(-1) \phi(5) \\ &\quad + (4.8)(1) \phi(6)] \phi(x_k) \end{aligned}$$

$$= y_k - [2.5 \phi(2) - 7.3 \phi(5) + 4.8 \phi(6)] \phi(x_k)$$

consider any one support vector, say  $x_2^k = 2$  ( $k=2$ )  
then

$$b = 1 - [2.5 \phi(2) \phi(2)$$

$$\begin{aligned} &- 7.3 \phi(5) \phi(2) \\ &+ 4.8 \phi(6) \phi(2) \end{aligned}$$

$$= 1 - [2.5 \times 5^2 - 7.3 \times 11^2 + 4.8 \times 13^2]$$

$$\Rightarrow b \approx 9$$

Now, the model is

$$f(z) = 0.67 z^2 - 5.33z + 9$$

Now, find label of  $z=3$

$$f(3) = 0.67 \times 3^2 - 5.33 \times 3 + 9$$

$$= 0.67 \times 9 - 5.33 \times 3 + 9$$

$$= 6.03 - 15.99 + 9$$

$$= 15.03 - 15.99$$

$$= -0.96$$

Since  $f(3) < 0 \Rightarrow z$  is in -ve class

w/ correctly predicted

Exam Problem  $\alpha$  to be given

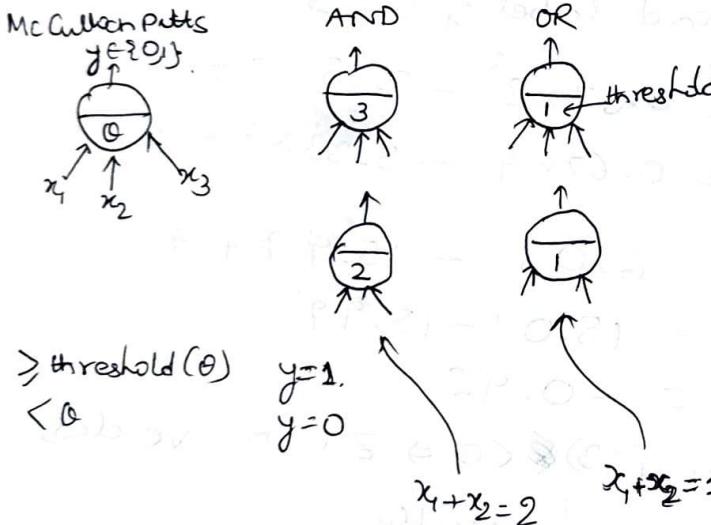
# Artificial Neural Networks

## Comparison of Biological Neuron to ANN

- Dendrites → Inputs → Training Samples
- Nucleus → Node (2 steps)
  - ↳ linear combination of I/Ps with their weights (signal strength)
  - ↳ function to produce off which determines whether neuron will fire or not. (Active  $\rightarrow 1$ ,  $0 \leftarrow$  not fire)

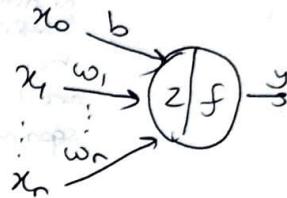
### In ANN

- 1 I/P layer
- multiple hidden layers (for intermediate purposes)
- 1 O/P layer



## Perceptron

Inputs are no longer boolean values like McPitts



$$z \rightarrow \text{linear combination of I/P's}$$

$$= w_0 x_0 + w_1 x_1 + \dots + w_n x_n + b$$

if  $\sum w_i x_i \geq \theta ; y=1$

else.  $\sum w_i x_i < \theta ; y=0$

Stepfunction

$x_0 \rightarrow$  bias term. ;  $x_0=1$ , weight  $= -\theta$  (threshold)

if  $\sum w_i x_i \geq 0 ; y=1$

else  $\sum w_i x_i < 0 ; y=0$

Bias term is used to fit the model well.

There can be multiple lines separating the classes. The weights (parameters) need to be estimated, to minimise the error.

Prof. Prabir Kr Biswas → NPTEL → ANN material

e.g.  $-x_1 - x_2 - 1.5 = 0$

$$x^T w \xrightarrow{f} \text{O/P}$$

$$w_1 = -1$$

$$w_2 = -1$$

$$w_0 = -1.5$$

	<u>AND</u>	<u>OR</u>
$w_0$	-1.5	-0.5
$w_1$	1	1
$w_2$	1	1

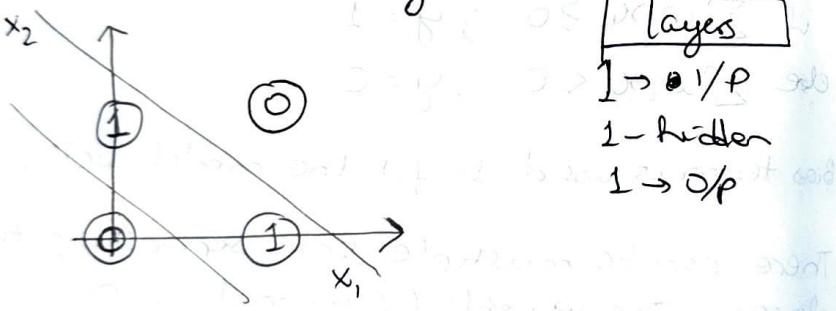
NO hidden layers bcoz their O/Ps are linearly separable

Step function is the Activation function.

$$\begin{cases} \text{step}(x) & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

### XOR

XOR gate requires a hidden layer because it is not linearly separable.



### Activation functions

hidden layer  $\leftarrow$  ReLU / sigmoid / tanh

O/P layer  $\leftarrow$  Sigmoid / Softmax

I/P layer  $\leftarrow$  linear (direct).  
 $\downarrow$  Not applying any activation fn

If a sample has D dimensions there will be D+1 input neurons (1 extra for bias)

Same type of activation fns are used in all hidden layers

coeff.	NAND	AND	OR
$w_0$	1.5	-1.5	-0.5
$w_1$	-1	1	1
$w_2$	-1	1	1

Inputs which are not linearly separable are transformed to another space, where they are linearly separable.

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix} \leftarrow \text{bias}$$

~~combination of binary I/Ps~~

Transformation is done through non-linear function / mapping from existing feature space to higher dimensional feature spaces multiple times such that they can be linearly separable [using linear classifier].