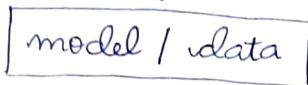


## MACHINE LEARNING

monotonic functions:

If we find  $x$  which minimizes  $y = x$ , it also minimizes  $y = x^2$ . Therefore  $x$  and  $x^2$  are said to be monotonous.

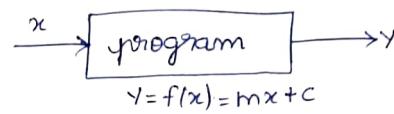
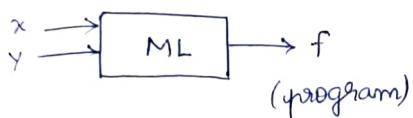
machine learning:

To frame model from given data.

We use regression to formulate the model.

ML

TP (traditional programming)

ML LAB:

- i) UCI  
kaggle }  $\rightarrow$  ML repositories

\* nature of  $x$  and  $y$

data types of  $x$

continuous

discrete

multi-nominal

ordinal / ordered

nominal

data types of  $y$

continuous

discrete

\* classification and regression

- ii) Download IRIS data set and house price prediction data set. Calculate mean, median, mode, standard deviation, covariance, skewness, kurtosis, coefficient of correlation

ML

- \* Machine Learning is the study of algorithms that learn a function from data without being explicitly programmed.

- \* According to Tom Mitchell, ML is represented by the following triplet.

$$\langle T, P, E \rangle$$

task      performance      experience  
measure

\* ML program/algorithm is said to learn a task  $T$  from experience  $E$ , if it is able to improve the performance  $P$  in learning a task  $T$  over a period of time.

#### TERMS OF ML:

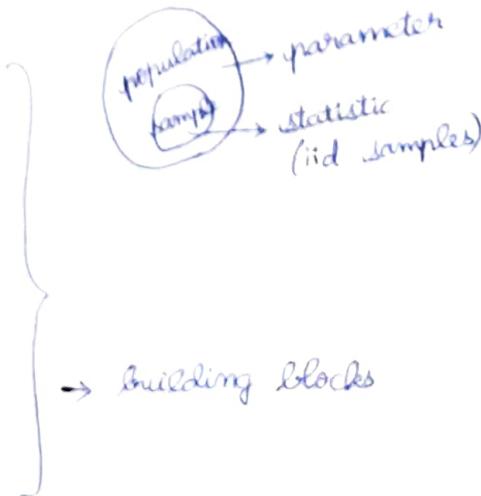
- \* hypothesis / function / model / concept - assumption
- \* hypothesis space  $H = \{h \mid h \text{ is } h(x)\}$  eg:  $H = \{fmx + c \mid m, c \in R\}$
- \* parameters - unknown values ( $m$  &  $c$  are parameters)
- \* hyperparameters - parameters used in training algo  $\rightarrow$   
eg:  $f(x) = w_0 + w_1x + w_2x^2 + \dots + w_dx^d$  d-degree

#### ML ALGORITHMS:

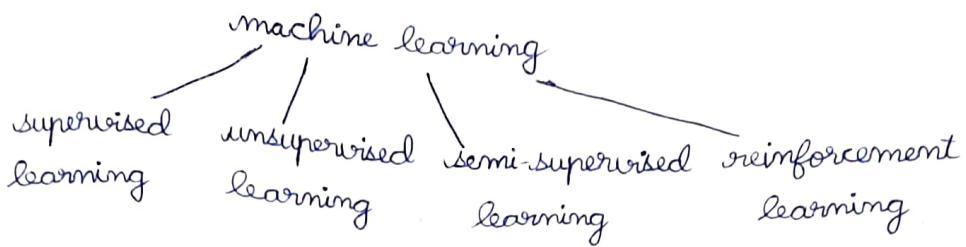
- |                         |                           |
|-------------------------|---------------------------|
| * Linear regression     | * K-nearest neighbours    |
| * polynomial regression | * support vector machines |
| * Naive Bayes           | * decision trees          |



- calculus
- linear algebra
  - eigen vector
  - eigen value
  - projection
  - vector space
  - norms
- matrix factorization
- probability & statistics
- optimization techniques



## CLASSIFICATION BASED ON DATA:

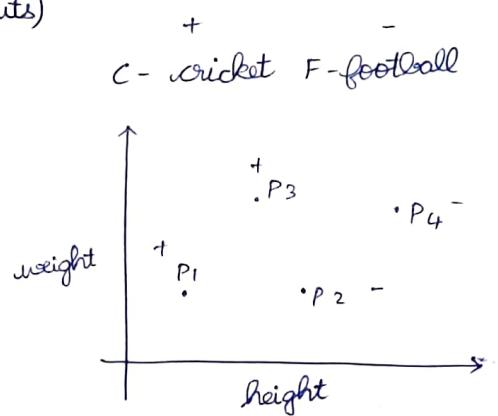


### Dataset :

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^n$$

Dataset  $\mathcal{D}$  consists of  $n$  observations, where  $x_i$  represents input  $i$  and  $y_i$  is its corresponding output.  
eg: 2 dimensional ( $\because$  2 inputs)

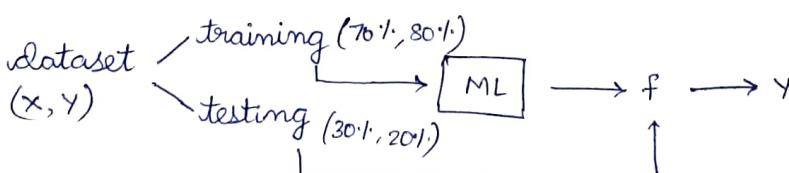
	height	weight	type	
P <sub>1</sub>	5	60	C	+ -
P <sub>2</sub>	6	70	F	-
P <sub>3</sub>	6.1	80	C	+ P <sub>3</sub>
P <sub>4</sub>	5.7	75	F	- P <sub>4</sub>



$x_i$  = input space / feature space / attributes  
 $y$  = output / class label / outcome

↑  
discrete  $x_i \in \mathbb{R}^2$

Dataset is considered as a matrix in which rows are considered as instances, datapoints, observations, examples and columns are features, attributes, characteristics, input vector.



### SUPERVISED LEARNING:

/ \  
classification      regression  
(outcome is      (outcome is  
discrete)      continuous)

## CLASSIFICATION ALGORITHMS:

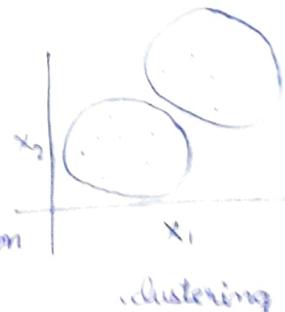
1. Naive Bayes
2. Logistic regression
3. Support vector machine (SVM)
4. k-nearest neighbours (KNN)
5. Linear Discriminant Analysis (LDA)
6. Quadratic Discriminant Analysis (QDA)
7. Decision tree
8. Random forest

## REGRESSION ALGORITHMS:

1. linear regression
2. polynomial regression
3. multiple regression
4. auto regression
5. KNN
6. SVM
7. decision trees

## UNSUPERVISED LEARNING

- only inputs → forms pattern
- no output → no supervision



- \* In clustering, each group share common characteristics. Naming groups using domain knowledge becomes classification
- \* In dimension reduction, d dimensions is reduced to k dimensions  $d \rightarrow k$
- \* Outlier analysis

## ML LAB:

$$\text{variance } \sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \rightarrow \text{standard deviation} = \sqrt{\text{variance}}$$

Boxplot → variations.

- 1. little no. summary (min, max,  $Q_1$ ,  $Q_2$ ,  $Q_3$ )  $Q$ -quartiles
- 2. box-whisker plot
  - median
  - median of 1st half
  - median of 2nd half
- outliers
  - \* mild outliers
  - \* strong outliers

## SEMI-SUPERVISED LEARNING:

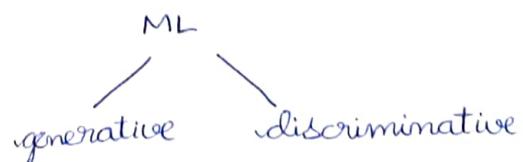
→ few inputs have output and the rest with no outputs.

→ apply supervised algorithms to data with outputs and build a confidence interval. Take data with high accuracy and apply ML algorithm to construct a model.

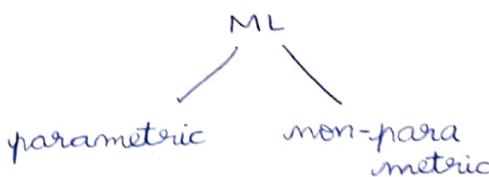
## REINFORCEMENT LEARNING:

- \* Learning from sequence of moves. (penalty, rewards)
- \* No training set
- \* eg: robot learning

## BASED ON MODEL:

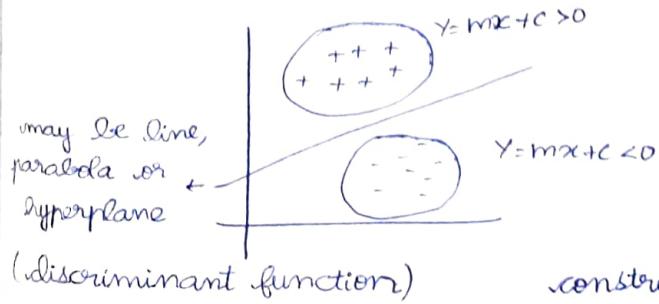


## BASED ON PARAMETERS:



## DISCRIMINATIVE:

Divide data points in feature space according to its feature. (decision boundary)



## ALGORITHMS:

- \* perceptron
- \* SVM
- \* logistic regression

constructs only 1 model

Discriminative finds conditional probability  $P(Y|X)$   
after analysing data points called posterior probability.  
 $Y \rightarrow \text{output}$     $X \rightarrow \text{input}$

### GENERATIVE:

- \* Based on learning data, we generate model.
- \* It finds joint probability.  $P(X, Y)$  \* learns data space
- ALGORITHMS: \* for every class we construct a model

Naive Bayes, Hidden Markov Model (HMM)

### PARAMETRIC:

constructing a model with fixed set of parameters  
irrespective of number of data. Values of parameter  
changes.

### ALGORITHMS:

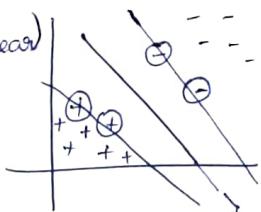
- \* perceptron
- \* SVM \* linear & logistic regression  
(linear)

### NON-PARAMETRIC:

- \* no parameters
- \* number of parameters may vary

eg: SVM (non-linear)

KNN



hyperplane \* support vectors closest to hyperplane

\* circled are called support vectors

\* may vary when number of inputs vary

### LINEAR REGRESSION: (supervised learning algorithm)

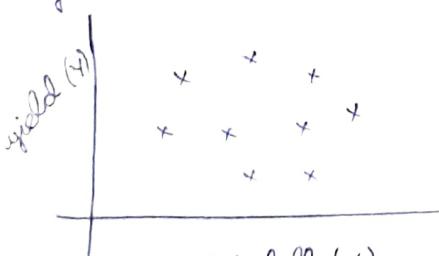
Given  $D = \{(x_i, y_i)\}_{i=1}^n$

$x_i \in \mathbb{R}$     $y_i \in \mathbb{R}$

$f: x_i \rightarrow y_i$

$f \rightarrow \text{true function}$     $\hat{f} \rightarrow \text{estimated function of}$

objective :  $\hat{f}: \mathbb{R} \rightarrow \mathbb{R}$



$x \rightarrow \text{independent}$

$y \rightarrow \text{dependent on } x$

$$y = f(x) = mx + c$$

linear relationship b/w dependent and independent variables

$P(Y/x)$

Probability.

model.

data space  
model

parameters  
inter

m

port

er of

to hyperplane

function of

endent  
es

### DATASET

	x	y
P <sub>1</sub>	5	80
P <sub>2</sub>	6	90
P <sub>3</sub>	10	200
P <sub>4</sub>	20	400
P <sub>5</sub>	4	40
P <sub>6</sub>	7	75

$$f = \beta_0 + \beta_1 x$$

$\beta_0$  and  $\beta_1$  are population parameters

By linear regression,

$$\hat{f} = b_0 + b_1 x_1$$

$b_0 \approx \beta_0$   $b_1 \approx \beta_1$

$b_0$  is an estimator of  $\beta_0$  and

$b_1$  is an estimator of  $\beta_1$

$$\hat{f}: \hat{y} = b_0 + b_1 x$$

$$Y = b_0 + b_1 x + e$$

$$Y = \hat{y} + e$$

$$\therefore \text{error} = Y - \hat{y} \quad \text{squared error } (Y - \hat{y})^2$$

$$\text{sum of squared errors SSE} = \sum_{i=1}^n (Y_i - \hat{y}_i)^2$$

Find  $b_0$  and  $b_1$  which minimizes SSE. (Optimization prob)

$$SSE = \sum_{i=1}^n (Y_i - (b_0 + b_1 x_i))^2$$

$$\frac{\partial}{\partial b_0} SSE = 0$$

$$\frac{\partial}{\partial b_0} \sum_{i=1}^n (Y_i - (b_0 + b_1 x_i))^2 = 0$$

$$\sum_{i=1}^n 2(Y_i - (b_0 + b_1 x_i))(-1) = 0$$

$$-2 \sum_{i=1}^n (Y_i - (b_0 + b_1 x_i)) = 0$$

$$\sum_{i=1}^n Y_i - \sum_{i=1}^n b_0 - \sum_{i=1}^n b_1 x_i = 0$$

$$\sum_{i=1}^n Y_i = \sum_{i=1}^n b_0 + \sum_{i=1}^n b_1 x_i$$

$$\sum_{i=1}^n Y_i = n b_0 + b_1 \sum_{i=1}^n x_i \longrightarrow (1)$$

$$\frac{\partial}{\partial b_1} SSE = 0$$

$$\frac{\partial}{\partial b_1} \sum_{i=1}^n (Y_i - (b_0 + b_1 x_i))^2 = 0$$

$$\sum_{i=1}^n 2(Y_i - (b_0 + b_1 x_i))(-x_i) = 0$$

here!

$$-2 \sum_{i=1}^n (y_i - (b_0 + b_1 x_i)) x_i = 0$$

$$\sum_{i=1}^n y_i x_i - \sum_{i=1}^n b_0 x_i + \sum_{i=1}^n b_1 x_i^2 = 0$$

$$\sum_{i=1}^n y_i x_i = \sum_{i=1}^n b_0 x_i + \sum_{i=1}^n b_1 x_i^2$$

$$\sum_{i=1}^n y_i x_i = b_0 \sum_{i=1}^n x_i + b_1 \sum_{i=1}^n x_i^2 \quad \longrightarrow (2)$$

From (1) and (2)

$$\begin{bmatrix} \sum y_i \\ \sum y_i x_i \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots \\ 1 & x_n \end{bmatrix}_{n \times 2} \quad x^T = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix}_{2 \times n} \quad b = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$x^T x = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}_{2 \times 2} \quad x^T y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}_{2 \times 1}$$

$x b$  gives  $\hat{y}$  values  
 $\therefore$  we add column of 1s

$$x^T y = (x^T x) b$$

$$B = (x^T x)^{-1} x^T y$$

parametric & discriminative

(closed form solution)

### ML LAB:

1. X Soap 4 4.5 5 5.5 6 6.5 7

Y Sud 33 42 45 51 53 61 62

$$x = \begin{bmatrix} 1 & 4 \\ 1 & 4.5 \\ 1 & 5 \\ 1 & 5.5 \\ 1 & 6 \\ 1 & 6.5 \\ 1 & 7 \end{bmatrix} \quad x^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 4 & 4.5 & 5 & 5.5 & 6 & 6.5 & 7 \end{bmatrix}$$

$$y = \begin{bmatrix} 33 \\ 42 \\ 45 \\ 51 \\ 53 \\ 61 \\ 62 \end{bmatrix}$$

2. Tu

in Sh

every

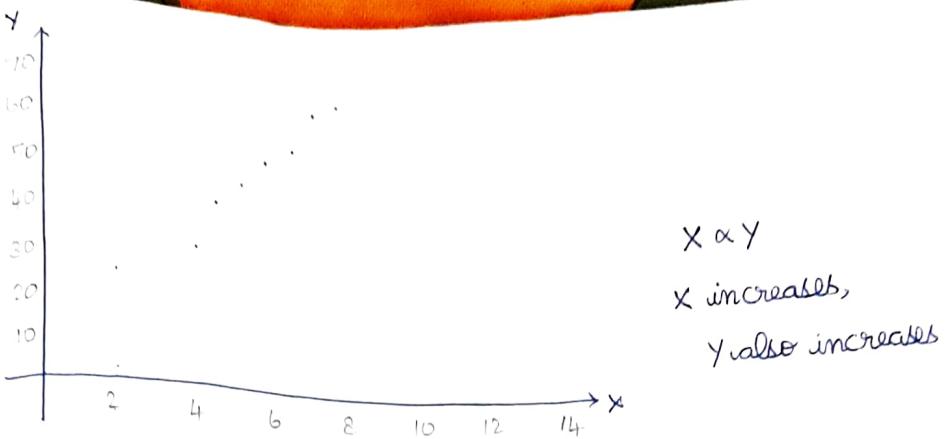
relati

every

diffe

rele

that



$$X^T X = \begin{bmatrix} 7 & 38.5 \\ 38.5 & 218.75 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} 4.454 & -0.785 \\ -0.785 & 0.143 \end{bmatrix} \quad X^T Y = \begin{bmatrix} 347 \\ 1975 \end{bmatrix}$$

$$\begin{bmatrix} b_0 \\ b_1 \end{bmatrix} = (X^T X)^{-1} (X^T Y) = \begin{bmatrix} -1.367 \\ 10.03 \end{bmatrix}$$

Predict std value for soap value = 4.25.

$$\hat{Y} = 37.696$$

$$SSE = 19.96 \quad SE = -2.27$$

### EVALUATION MEASURES :

$$\text{total sum of squares (SST)} = \sum_{i=1}^n (y_i - \bar{Y})^2$$

$$\text{sum of squared error (SSE)} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{regression sum of squares (SSR)} = \sum_{i=1}^n (\hat{y}_i - \bar{Y})^2$$

$$0 \leq SSR \leq SST$$

$$\text{coefficient of determination } r^2 = 1 - \frac{SSE}{SST} \quad (\text{or}) \quad \frac{SSR}{SST}$$

$$SST = SSE + SSR$$

$$0 \leq r^2 \leq 1$$

$$y_i - \bar{Y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{Y})$$

here)

- Identify independent and dependent variables in house price data set. Draw scatter plot for every dep & indep feature and interpret the relationship between model using linear regression for every dep & indep feature and evaluate models using different measures such as SSE, SST, SSR & coeff of determination & find out the most imp feature that influences dep variable. Find outliers if any.

A good model has greater SSR,  $r^2$  and lesser SSE.

Soap Sud:

$$B = (X^T X)^{-1} X^T Y$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 4 & 4.5 & 5 & 5.5 & 6 & 6.5 & 7 \end{bmatrix}_{2 \times 7} \begin{bmatrix} 1 & 4 \\ 1 & 4.5 \\ 1 & 5 \\ 1 & 5.5 \\ 1 & 6 \\ 1 & 6.5 \\ 1 & 7 \end{bmatrix}_{7 \times 2}$$

$$X^T X = \begin{bmatrix} 7 & 38.5 \\ 38.5 & 218.75 \end{bmatrix} \quad X^T Y = \begin{bmatrix} 347 \\ 1975 \end{bmatrix}$$

$$(X^T X)^{-1} = \begin{bmatrix} 4.464 & -0.786 \\ -0.786 & 0.143 \end{bmatrix} \quad \text{※ don't round off}$$

$$B = \begin{bmatrix} 4.464 & -0.785 \\ -0.785 & 0.143 \end{bmatrix} \begin{bmatrix} 347 \\ 1975 \end{bmatrix} = \begin{bmatrix} -3.342 \\ 9.683 \end{bmatrix} \begin{bmatrix} -2.6786 \\ 9.5 \end{bmatrix}$$

$$b_0 = -3.342 \quad b_1 = 9.683$$

$$\hat{Y} = -3.342 + 9.683x$$

$$\bar{Y} = 49.57 \quad \hat{Y} = -2.6786 + 9.5x$$

$$x_1 = 4 \quad \hat{y}_1 = 35.39$$

$$x_2 = 4.5 \quad \hat{y}_2 = 40.2315$$

$$x_3 = 5 \quad \hat{y}_3 = 45.073$$

$$x_4 = 5.5 \quad \hat{y}_4 = 49.9145$$

$$x_5 = 6 \quad \hat{y}_5 = 54.756$$

$$x_6 = 6.5 \quad \hat{y}_6 = 59.5975$$

$$x_7 = 7 \quad \hat{y}_7 = 64.439$$

$$SSE = (33 - 35.39)^2 + (42 - 40.2315)^2 + (45 - 45.073)^2 + (51 - 49.9145)^2 + (53 - 54.756)^2 + (61 - 59.5975)^2 + (62 - 64.439)^2$$

$$= 21.02 \quad 13.914$$

$$SSR = 657.154 \quad 657.75$$

$$SST = 672.5992 \quad 651.7143$$

$P =$

$\sigma_{\epsilon} =$

$P =$

eg: Take

SLR m

petal l

S.N

1

2

3

4

5

6

$X \Rightarrow$

$Y \Rightarrow$

$X^T X$

$X^T Y$

$B =$

$$\gamma^2 = 0.9693$$

$$r = \frac{\text{covar}(x, y)}{\sigma_x \sigma_y} \Rightarrow \text{coefficient of correlation } (r)$$

$$\text{covar}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\sigma_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

$$\sigma_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n-1}}$$

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n-1)\sigma_x \sigma_y}$$

$$-1 \leq r \leq 1$$

$$r = 0.9845$$

Q: Take any 6 rows of iris data set. Construct a SLR model between sepal length and sepal width, petal length and petal width and calculate all measures.

S.NO	SEPAL LENGTH	SEPAL WIDTH	PETAL LENGTH	PETAL WIDTH
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
5	5.0	3.6		
6	5.4	3.9	1.4	0.2
			1.7	0.4

$$X \Rightarrow \text{sepal width} [3.5, 3.0, 3.2, 3.1, 3.6, 3.9]$$

$$Y \Rightarrow \text{sepal length} [5.1, 4.9, 4.7, 4.6, 5.0, 5.4]$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 3.5 & 3.0 & 3.2 & 3.1 & 3.6 & 3.9 \end{bmatrix} \begin{bmatrix} 1 & 3.5 \\ 1 & 3.0 \\ 1 & 3.2 \\ 1 & 3.1 \\ 1 & 3.6 \\ 1 & 3.9 \end{bmatrix}$$

$$= \begin{bmatrix} 6 & 20.3 \\ 20.3 & 69.29 \end{bmatrix}$$

here

$$X^T Y = \begin{bmatrix} 29.7 \\ 100.91 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} 19.6232 & -5.7507 \\ -5.7507 & 1.6997 \end{bmatrix}$$

$$B = \begin{bmatrix} 2.5059 \\ 0.72 \end{bmatrix} \quad \hat{Y} = 2.5059 + 0.72x$$

GRADIENT

\* iter

\* fir

\* um

CONVEX SET

$$\frac{x_1 + (1-x_1)}{x}$$

A set is

Any 2 point

$$\lambda x + (1-\lambda)$$

eg: 

c

intersection

union

CONVEX

$$X \quad 3.5 \quad 30 \quad 32 \quad 31 \quad 36 \quad 39$$

$$Y \quad 51 \quad 49 \quad 47 \quad 46 \quad 50 \quad 54 \quad \bar{Y} = 4.95$$

$$\hat{Y} \quad 5.0259 \quad 4.6659 \quad 4.8089 \quad 4.7379 \quad 5.0979 \quad 5.3139$$

$$SSE = 0.00549 + 0.054 + 0.0118 + 0.01901 + 0.0095 + 0.0074$$

$$SSE = 0.107$$

$$SST = 0.415$$

$$SSR = 0.306$$

$$r^2 = 0.737$$

$$r = 0.858$$

 $X \Rightarrow$  petal width

 $Y \Rightarrow$  petal length

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 0.2 & 0.2 & 0.2 & 0.2 & 0.2 & 0.4 \end{bmatrix} \begin{bmatrix} 1 & 0.2 \\ 1 & 0.2 \\ 1 & 0.2 \\ 1 & 0.2 \\ 1 & 0.2 \\ 1 & 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 6 & 1.4 \\ 1.4 & 0.36 \end{bmatrix}$$

$$X^T Y = \begin{bmatrix} 87 \\ 2.08 \end{bmatrix} \quad (X^T X)^{-1} = \begin{bmatrix} 1.8 & -7 \\ -7 & 30 \end{bmatrix}$$

$$B = \begin{bmatrix} 1.1 \\ 1.5 \end{bmatrix}$$

$$b_0 = 1.1 \quad b_1 = 1.5$$

$$\hat{Y} = 1.1 + 1.5x$$

 $X \quad 0.2 \quad 0.2 \quad 0.2 \quad 0.2 \quad 0.2 \quad 0.4$ 
 $Y \quad 1.4 \quad 1.4 \quad 1.4 \quad 1.5 \quad 1.4 \quad 1.7 \quad \bar{Y} = 1.45$ 
 $\hat{Y} \quad 1.4 \quad 1.4 \quad 1.4 \quad 1.4 \quad 1.4 \quad 1.7$ 

$$SSE = 0.02 \quad SST = 0.095 \quad SSR = 0.075 \quad r^2 = 0.7894 \quad r = 0.8884$$

$$B = (X^T X)^{-1} X^T Y$$

- \* Cannot use if  $X^T X$  is not invertible ie)  $|X^T X| = 0$
- \* when n becomes large, this  $O(n^3)$  algorithm is inefficient (2 matrix multiplication)

To overcome these drawbacks,

## GRADIENT DESCENT METHOD:

\* iterative

\* first order derivative

\* unconstrained convex optimisation problem

## CONVEX SET:



$$0 \leq \lambda \leq 1$$

if  $\lambda = 0$ ,  $y$  | if  $\lambda = 1$ ,  $x$  | if  $\lambda = 0.5$ , midpoint

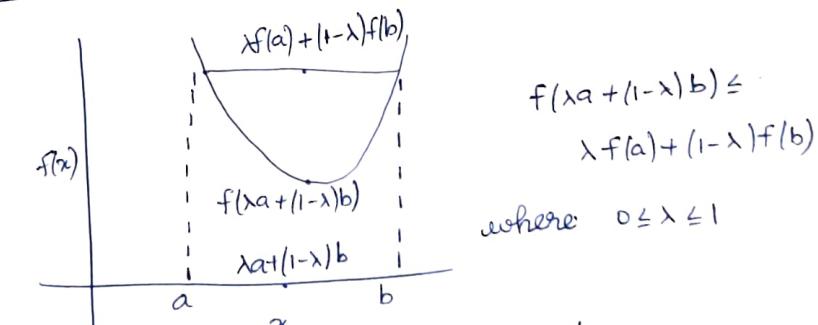
A set is a convex set if all points lie within the set.  
Any 2 points  $x$  and  $y$   $\in$  convex set, then all points of  
 $\lambda x + (1-\lambda)y \in$  convex set, then it is convex set.



line, plane

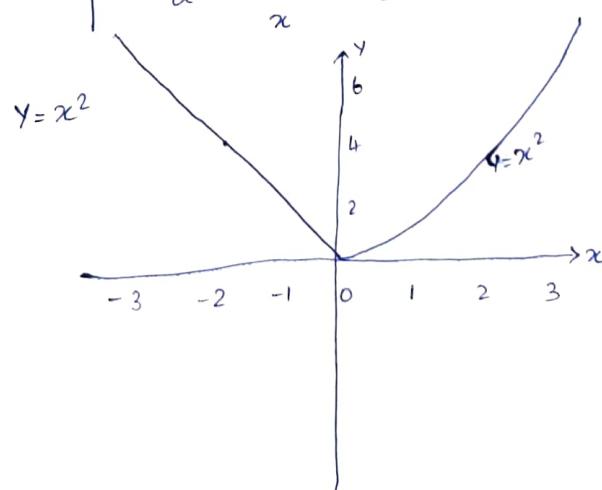
intersection of 2 convex sets is always convex  
union may or may not be convex

## CONVEX FUNCTION:

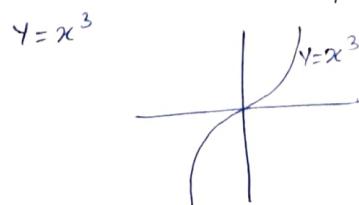


$$\begin{aligned} f(xa + (1-x)b) &\leq \\ xf(a) + (1-x)f(b) & \end{aligned}$$

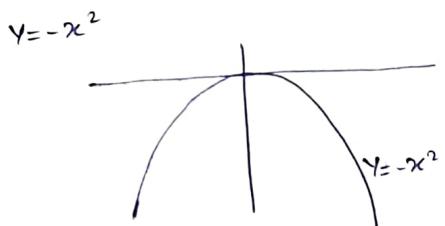
where  $0 \leq x \leq 1$



convex



$y = x^3$  if  $x \geq 0$  then it is convex



not convex function

If  $f(x)$  is convex then  $-f(x)$  is concave and vice versa

All convex functions are  $f''(x) \geq 0 \nabla$  and concave  $f''(x) \leq 0$

Check whether following functions are convex:

- 1)  $x \log x$
- 3)  $e^x \log x$
- 2)  $-x \log x$
- 4)  $\log x$

1.  $f(x) = x \log x$

$$f'(x) = \log x + x \cdot \frac{1}{x} = \log x + 1$$

$$f''(x) = \frac{1}{x} \quad \text{convex for } x \geq 0$$

2.  $f(x) = -x \log x$

$$f'(x) = -\log x + (-x) \cdot \frac{1}{x} = -\log x - 1$$

$$f''(x) = -\frac{1}{x} \quad \text{convex for } x \leq 0$$

3.  $f(x) = e^x \log x$

$$f'(x) = e^x \frac{1}{x} + e^x \log x$$

$$\begin{aligned} f''(x) &= e^x \left( \frac{-1}{x^2} \right) + \frac{1}{x} e^x + e^x \left( \frac{1}{x} \right) + \log x e^x \\ &= e^x \left( \frac{2}{x} - \frac{1}{x^2} + \log x \right) \quad \text{not convex} \end{aligned}$$

4.  $\log(x) \Rightarrow f(x)$

$$f'(x) = \frac{1}{x} \quad f''(x) = -\frac{1}{x^2} \quad \text{not convex}$$

Let  $f(x)$   
To check  
\* for  
\* If

$f$  is con

Jacobian  
mat

All

If

+ve

H/W

f

-

H =

Let  $f(x, y, z)$  is a multivariate function  
To check convexity

\* find Hessian matrix  $H$

\* If  $H$  is positive definite / positive semidefinite, then  $f$  is convex

$$\text{Hessian matrix } J = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \frac{\partial f}{\partial z} \end{bmatrix} = \begin{bmatrix} x & y & z \\ \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial x \partial z} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y \partial y} & \frac{\partial^2 f}{\partial y \partial z} \\ \frac{\partial^2 f}{\partial z \partial x} & \frac{\partial^2 f}{\partial z \partial y} & \frac{\partial^2 f}{\partial z \partial z} \end{bmatrix} \underset{H}{\sim}$$

All eigen values of  $H$  are  $\geq 0 \rightarrow$  +ve semidefinite  
 $> 0 \rightarrow$  +ve definite

If matrix is symmetric, then it is +ve definite /  
+ve semi-definite. check for

$$aHa^T \geq 0 \quad aHa^T > 0$$

H/W Check whether it is convex

$$f(x_1, x_2, x_3) = (x_1 - x_2)^2 + 2x_3^2$$

$$J = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 2(x_1 - x_2)(1) \\ 2(x_1 - x_2)(-1) \\ 4x_3 \end{bmatrix} = \begin{bmatrix} 2x_1 - 2x_2 \\ 2x_2 - 2x_1 \\ 4x_3 \end{bmatrix}$$
here)

$$H = \begin{bmatrix} 2 & -2 & 0 \\ -2 & 2 & 0 \\ 0 & 0 & 4 \end{bmatrix} \quad \begin{aligned} Hx &= \lambda x \\ (H - \lambda I)x &= 0 \\ |H - \lambda I| &= 0 \end{aligned}$$

$H$  is symmetric

$$\left| \begin{array}{ccc} 2-\lambda & -2 & 0 \\ -2 & 2-\lambda & 0 \\ 0 & 0 & 4-\lambda \end{array} \right| = 0$$

$$8 \\ 6 \\ -4 -4$$

$$16 \\ 8 \\ -4 -4$$

$$(2-\lambda) [(2-\lambda)(4-\lambda)] + 2 [-2(4-\lambda)] = 0$$

$$(2-\lambda)[8-4\lambda-2\lambda+\lambda^2] + 2(-8+2\lambda) = 0$$

$$(2-\lambda)(\lambda^2-6\lambda+8) + (-16+4\lambda) = 0$$

$$2\lambda^2 - 12\lambda + 16 - \lambda^3 + 6\lambda^2 - 8\lambda - 16 + 4\lambda = 0$$

$$-\lambda^3 + 8\lambda^2 - 16\lambda = 0$$

$$-\lambda(\lambda^2 - 8\lambda + 16) = 0$$

$$\boxed{\lambda=0}$$

$$\boxed{\lambda=4}$$

$$\boxed{\lambda=4}$$

$\therefore$  all eigen values  $\geq 0$   
+ve semi definite

It is a convex function

### CONVEX OPTIMISATION PROBLEM:

min

$$f(x)$$

→ convex function

s.t.c

$$g(x) \geq 0$$

→ convex set

\* constrained COP

\* unconstrained COP

} types

optimisation problem for LR:

Find  $b_0$  and  $b_1$  which minimises SSE.

$$\hat{y}_i = b_0 + b_1 x_i$$

$$\begin{bmatrix} b_0 & b_1 \end{bmatrix} \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

$$B^T x$$

$$L: \min \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - B^T x_i)^2$$

$$\frac{\partial L}{\partial B} = \sum_{i=1}^n 2(y_i - B^T x_i)(-x_i)$$

$$\frac{\partial^2 L}{\partial B^2} = \sum_{i=1}^n -2x_i(-x_i) = \sum_{i=1}^n 2x_i^2 \geq 0$$

$\therefore$  SSE is a convex function

## LOSS FUNCTIONS: (all are convex)

1. SSE - linear reg, binary classification
2. MSE (mean squared error) - regression, binary classifi
3. negative log likelihood - binary classifi, logistic reg
4. hinge loss - SVM
5. cross entropy - multiple classification

## GRADIENT DESCENT METHOD

- \* unconstrained convex optimization problem
- \* first order partial derivative

Find  $\theta$  which minimizes  $f(\theta, x)$

$\theta$  is a parameter

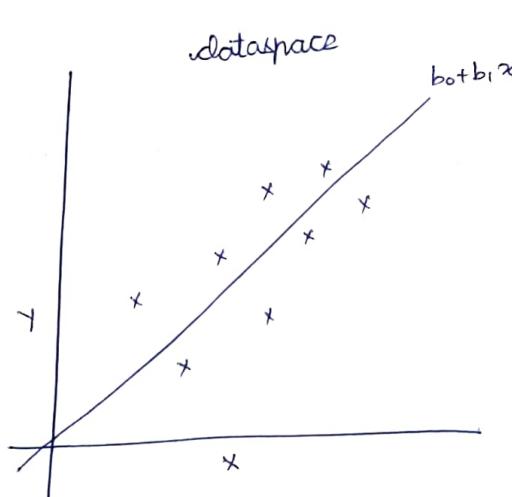
$$\text{gradient} = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \end{bmatrix} \quad (\text{slope})$$

### ALGORITHM:

input  $x, \eta$  (step size, learning rate)  
 output  $\theta^*$   $\theta^* = \underset{\theta}{\operatorname{argmin}} f(x, \theta)$

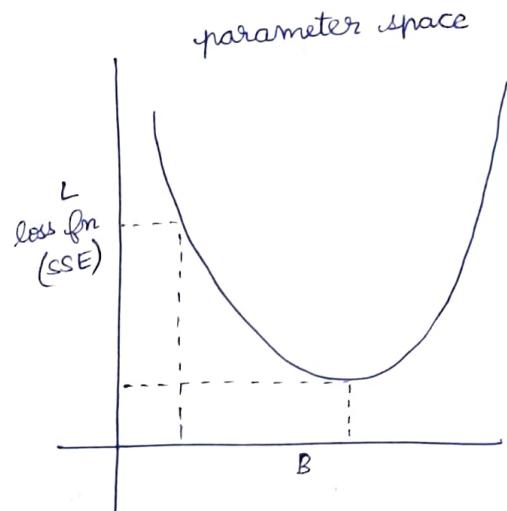
1.  $\theta = \text{initialize}()$
2. while (!convergence)
 
$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \frac{\partial f}{\partial \theta}$$

$$\theta_{\text{old}} = \theta_{\text{new}}$$



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

$$\sum (y_i - B^T x_i)^2$$



$$B^* = \underset{B}{\operatorname{argmin}} \text{SSE}$$

## LINEAR REGRESSION USING GRADIENT DESCENT:

input  $D = \{x_i, y_i\}_{i=1}^n$

output  $\vec{B} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$

1.  $B = \text{initialize}()$

fixed number of epoch

2. while (! convergence)

$\text{grad} = [0, 0]$

$// B_{\text{new}} = B_{\text{old}}$  (or)  $\frac{\partial L}{\partial B} = 0$

a. for each datapoint  $(x_i, y_i)$

$$\hat{y}_i = b_0 + b_1 x_i$$

$$\text{grad} = \text{grad} + (y_i - \hat{y}_i)(-x_i)$$

$$\frac{\partial SSE}{\partial B} = \begin{bmatrix} \frac{\partial SSE}{\partial b_0} \\ \frac{\partial SSE}{\partial b_1} \end{bmatrix}$$

$$b. b_{\text{new}} = b_{\text{old}} - \eta \text{ grad}$$

$$= \begin{bmatrix} S(x_i - \hat{y}_i) / (-1) \\ S(y_i - \hat{y}_i) / (-x_i) \end{bmatrix}$$

$$c. b_{\text{old}} = b_{\text{new}}$$

$$S(x_i - \hat{y}_i) / (-1)$$

### PROBLEM :

$$X = [4, 4.5, 5, 5.5, 6, 6.5, 7] \quad Y = [33, 42, 45, 51, 53, 61, 62]$$

$$\eta = 0.5 \quad B = [0.5, 0.5] \quad \text{grad} = [0, 0]$$

IT-1

$$1. \hat{y}_i = 0.5 + 0.5(4) = 2.5$$

$$\text{grad} = [30.5, -122]$$

$$2. \hat{y}_i = 0.5 + 0.5(4.5) = 2.75$$

$$\text{grad} = [-30.5, -122] + [-39.25, -176.625] = [-69.75, -298.625]$$

$$3. \hat{y}_i = 0.5 + 0.5(5) = 3$$

$$\text{grad} = [-69.75, -298.625] + [-42, -210] = [-111.75, -508.625]$$

$$4. \hat{y}_i = 0.5 + 0.5(5.5) = 3.25$$

$$\text{grad} = [-111.75, -508.625] + [47.75, -262.625] = [-159.5, -771.25]$$

$$5. \hat{y}_i = 0.5 + 0.5(6) = 3.5$$

$$\text{grad} = [-159.5, -771.25] + [-49.5, -297] = [-209, -1068.25]$$

$$6. \hat{y}_i = 0.5 + 0.5(6.5) = 3.75$$

$$\text{grad} = [-209, -1068.25] + [-57.25, -372.125] = [-266.25, -1440.375]$$

$$7. \hat{y}_i = 0.5 + 0.5(7) = 4$$

$$\text{grad} = [-266.25, -1440.375] + [-58, -406] = [-324.25, -1846.375]$$

$$B_{\text{new}} = B_{\text{old}} - \eta \text{ grad}$$

$$= \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} - 0.5 \begin{bmatrix} -324.25 \\ -1846.375 \end{bmatrix}$$

$$B_{\text{new}} = \begin{bmatrix} 162.625 \\ 923.6875 \end{bmatrix}$$

\* one iteration is called epoch

\* This method is called batch gradient descent / vanilla gradient descent method

\* Since this is expensive, we go for stochastic gradient descent where we update weights for each data point

\* Since there will be more variation, we use mini batch gradient descent. If there 1000 inputs we divide it into 10 batches and weights are updated for every batch.

### ASSUMPTIONS OF LINEAR REGRESSION :

1. All observations are independent (iid)
2. Relationship between dependant & independant variables is linear.
3. Y is linear to parameters.
4.  $E(\text{error}) = 0$  (mean)
5. For a  $x_i$ , error follow normal distribution.
6. No auto correlation between errors. (errors are independent)
7. If Errors are heteroscedasticity, use weighted linear regressions.  $\therefore$  Errors are homoscedasticity.  
As  $x$  increases, error becomes constant.

Linear regression is true function:

$$Y = (B_0 + B_1 x) + e \rightarrow \text{irreducible errors}$$

$$\hat{Y} = b_0 + b_1 x + e \rightarrow \text{SSE (reducible errors)}$$

## OUTLIER DETECTION TECHNIQUES :

1. Z Score (for converting data points to follow standard normal dist  $M=0, \sigma=1$ )
- $$\begin{cases} \text{outlier} & \left| \frac{x-M}{\sigma} \right| \geq 3 \\ \text{no} & \text{otherwise} \end{cases}$$

### Box plot

only for data points which follow normal dist  
any points outside lower bound & upper bound

$$M - 3\sigma \Rightarrow LB : Q_1 - 1.5 \text{ IQR}$$

$$M + 3\sigma \Rightarrow UB : Q_1 + 1.5 \text{ IQR}$$

### linear regression

construct model and choose the outliers (data points highly deviated from model). remove outliers and construct model again.

### POLYNOMIAL REGRESSION :

- \* instead of one independant feature  $x$ , we include powers of  $x$ .

- \* when data points itself follow polynomial function



$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_d x^d + \epsilon \quad \hookrightarrow \text{irreducible errors}$$

$$\hat{y} = b_0 + b_1 x + b_2 x^2 + \dots + b_d x^d$$

Given

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^n \quad x_i \in \mathbb{R} \quad y_i \in \mathbb{R}$$

$$f: x_i \rightarrow y_i \quad (\beta_0 + \beta_1 x + \dots + \beta_d x^d) + e$$

objective

$$\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R} \quad b_0 + b_1 x + \dots + b_d x^d$$



## optimization problem

Find  $b = (b_0, b_1, \dots, b_d)$

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

where

$$\hat{y}_i = b_0 + b_1 x + b_2 x^2 + \dots + b_d x^d$$

$$L = \sum_{i=1}^n (y_i - (b_0 + b_1 x + b_2 x^2 + \dots + b_d x^d))^2$$

$$\frac{\partial L}{\partial b_0} = \sum 2(y_i - (b_0 + b_1 x + \dots + b_d x^d))(-1) = 0$$

$$\sum y_i = b_0 n + b_1 \sum x_i + b_2 \sum x_i^2 + \dots + b_d \sum x_i^d$$

$$\frac{\partial L}{\partial b_1} = \sum 2(y_i - (b_0 + b_1 x + \dots + b_d x^d))(-x) = 0$$

$$\sum x_i y_i = b_0 \sum x_i + b_1 \sum x_i^2 + \dots + b_d \sum x_i^{d+1} -$$

$$\frac{\partial L}{\partial b_2} = \sum 2(y_i - (b_0 + b_1 x + \dots + b_d x^d))(-x^2) = 0$$

$$\sum x_i^2 y_i = b_0 \sum x_i^2 + b_1 \sum x_i^3 + \dots + b_d \sum x_i^{d+2}$$

:

$$\frac{\partial L}{\partial b_d} = \sum 2(y_i - (b_0 + b_1 x + \dots + b_d x^d))(-x^d) = 0$$

$$\sum x_i^d y_i = b_0 \sum x_i^d + b_1 \sum x_i^{d+1} + \dots + b_d \sum x_i^{2d}$$

$$Y = X B$$

where

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad X = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ 1 & x_2 & x_2^2 & & x_2^d \\ 1 & x_3 & x_3^2 & & x_3^d \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix}_{(n \times (d+1))} \quad B = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_d \end{bmatrix}_{(d+1) \times 1} \quad (\text{here})$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots & 1 \\ x_1 & x_2 & x_3 & x_4 & \dots & x_n \\ x_1^2 & x_2^2 & x_3^2 & x_4^2 & \dots & x_n^2 \\ \vdots & & & & & \vdots \\ x_1^d & x_2^d & x_3^d & x_4^d & \dots & x_n^d \end{bmatrix}_{(d+1) \times n} \quad X^T X = \begin{bmatrix} n & \sum x_i & \sum x_i^2 & \dots & \sum x_i^d \\ \sum x_i & \sum x_i^2 & \sum x_i^3 & \dots & \sum x_i^{d+1} \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 & \dots & \sum x_i^{d+2} \\ \vdots & & & & \vdots \\ \sum x_i^d & \sum x_i^{d+1} & \sum x_i^{d+2} & \dots & \sum x_i^{2d} \end{bmatrix}_{n \times n}$$

$$X^T Y = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \\ \vdots \\ \sum x_i^d y_i \end{bmatrix} \quad \boxed{B = (X^T X)^{-1} X^T Y}$$

$$\text{eg: } \begin{matrix} x & 1 & 2 & 3 & 4 & 5 & 6 \\ y & 1 & 4 & 9 & 16 & 25 & 36 \end{matrix}$$

$$\hat{y}_i = b_0 + b_1 x + b_2 x^2$$

Construct second degree polynomial and find SSE.  
Compare its SSE with linear regression model.

$$x = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \\ 1 & 5 & 25 \\ 1 & 6 & 36 \end{bmatrix} \quad x^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \\ 1 & 4 & 9 & 16 & 25 & 36 \end{bmatrix}$$

$$x^T x = \begin{bmatrix} 6 & 21 & 91 \\ 21 & 91 & 441 \\ 91 & 441 & 2275 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 4 \\ 9 \\ 16 \\ 25 \\ 36 \end{bmatrix}$$

$$x^T y = \begin{bmatrix} 91 \\ 441 \\ 2275 \end{bmatrix}_{3 \times 1} \quad (x^T x)^{-1} = \begin{bmatrix} 3.2 & -1.95 & 0.25 \\ -1.95 & 1.36964 & -0.1875 \\ 0.25 & -0.1875 & 0.08678 \end{bmatrix}_{3 \times 3}$$

$$B = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \quad \therefore \hat{y}_i = x^2 \quad \boxed{\text{SSE} = 0}$$

By linear regression,

$$\hat{y} = b_0 + b_1 x$$

$$x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \\ 1 & 6 \end{bmatrix} \quad x^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \end{bmatrix}$$

$$x^T x = \begin{bmatrix} 6 & 21 \\ 21 & 91 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 4 \\ 9 \\ 16 \\ 25 \\ 36 \end{bmatrix}$$

$$(x^T x)^{-1} = \begin{bmatrix} 0.8666 & -0.2 \\ -0.2 & 0.05714 \end{bmatrix}$$

$$x^T y = \begin{bmatrix} 91 \\ 441 \end{bmatrix}$$

$$B = (x^T x)^{-1} x^T y = \begin{bmatrix} -9.333 \\ 7 \end{bmatrix}$$

$$\hat{y} = -9.333 + 7x$$

$\hat{y}$	-2.333	4.667	11.667	18.667	25.667	$\dots$
$y$	1	4	9	16	25	$\dots$

$$\text{SSE} = \sum_{i=1}^6 (y_i - \hat{y}_i)^2 = 37333334$$

$$\text{SSE}_{LR} > \text{SSE}_{PR}$$

### ML LAB

Apply polynomial regression to every features and construct an optimum model is via min error (SSE) or max  $r^2$ .

### MULTIPLE REGRESSION :

$$\text{Given } D = \{x_i, y_i\}_{i=1}^n \quad x_i \in$$

$$Y_i = f(x_i)$$

Objective

$$\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R}$$

Hypothesis (model is hypothesis)

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_d x_d$$

$$\hat{Y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_d x_d$$

$$Y = Xw + e$$

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} \quad W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1}$$

$$X^T Y = \begin{bmatrix} 91 \\ 441 \end{bmatrix}$$

$$B = (X^T X)^{-1} X^T Y = \begin{bmatrix} -9.333 \\ 7 \end{bmatrix}$$

$$\hat{Y} = -9.333 + 7x$$

$\hat{Y}$	2.333	4.667	11.667	18.667	25.667	32.667
$Y$	1	4	9	16	25	36

$$SSE = \sum_{i=1}^6 (y_i - \hat{y}_i)^2 = 37333334$$

$$SSE_{LR} > SSE_{PR}$$

### ML LAB

Apply polynomial regression on iris data set between every features and construct an optimal model. An optimum model is a model which has lowest error (SSE) or max r<sup>2</sup>.

### MULTIPLE REGRESSION :

Given  $D = \{x_i, y_i\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$

$$y_i = f(x_i)$$

Objective

$$\hat{f}: \mathbb{R}^d \rightarrow \mathbb{R}$$

Hypothesis (model is hyperplane)

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_d x_d + \epsilon \rightarrow \text{random error}$$

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_d x_d + e \rightarrow SSE$$

$$w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d + e$$

$$Y = XW + e$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} \quad W = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nd} \end{bmatrix} \quad n \times (d+1)$$

## Optimisation problem

Find  $w$

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\min SSE(L) = \sum_{i=1}^n (y_i - (w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}))^2$$

$$\frac{\partial L}{\partial w_0} = 2 \sum_{i=1}^n (y_i - (w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id})) (-1) = 0$$

$$\Rightarrow \sum y_i = n w_0 + w_1 \sum x_{i1} + \dots + w_d \sum x_{id}$$

$$\frac{\partial L}{\partial w_1} \Rightarrow \sum y_i x_{i1} = w_0 \sum x_{i1} + w_1 \sum x_{i1}^2 + \dots + w_d \sum x_{id}^2$$

⋮

$$\frac{\partial L}{\partial w_d} \Rightarrow \sum y_i x_{id} = w_0 \sum x_{id} + w_1 \sum x_{i1} x_{id} + \dots + w_d \sum x_{id}^2$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & x_{31} & \dots & x_{n1} \\ \vdots & \vdots & \vdots & & \vdots \\ x_{1d} & x_{2d} & x_{3d} & \dots & x_{nd} \end{bmatrix} \quad (d+1) \times n$$

$$X^T X = \begin{bmatrix} n & \sum x_{i1} & \sum x_{i2} & \dots & \sum x_{id} \\ \sum x_{i1} & \sum x_{i1}^2 & \sum x_{i1} x_{i2} & \dots & \sum x_{i1} x_{id} \\ \vdots & \vdots & \vdots & & \vdots \\ \sum x_{id} & \sum x_{id} x_{i1} & \sum x_{id} x_{i2} & \dots & \sum x_{id}^2 \end{bmatrix} \quad (d+1) \times (d+1)$$

$$X^T Y = \begin{bmatrix} \sum y_i \\ \sum y_i x_{i1} \\ \vdots \\ \sum y_i x_{id} \end{bmatrix} \quad (d+1) \times 1$$

$$\therefore X^T Y = (X^T X) W$$

$$W = (X^T X)^{-1} X^T Y$$

H/W	CASES OF PRODUCT	7	3	3	4	6	7
DISTANCE (IN FT)	560	220	340	80	150	330	
TIME (IN MINS)	16:68	11:50	12:03	14:88	13:75	18:11	

A soft drink bottler is analysing the vending machine serving routes in his distribution system. He is interested in predicting the time required for by the driver to service the vending machines in an outlet.

$$\hat{Y} = w_0 + w_1 x_1 + w_2 x_2$$

$$Y = \begin{bmatrix} 16.68 \\ 11.50 \\ 12.03 \\ 14.82 \\ 13.75 \\ 18.11 \end{bmatrix}_{6 \times 1}$$

$$X = \begin{bmatrix} 1 & 7 & 560 \\ 1 & 3 & 220 \\ 1 & 3 & 340 \\ 1 & 4 & 80 \\ 1 & 6 & 150 \\ 1 & 7 & 330 \end{bmatrix}_{6 \times 3}$$

$$d=2$$

$$n=6$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 7 & 3 & 3 & 4 & 6 & 7 \\ 560 & 220 & 340 & 80 & 150 & 330 \end{bmatrix}_{3 \times 1}$$

$$W = (X^T X)^{-1} X^T Y$$

$$X^T X = \begin{bmatrix} 6 & 30 & 1680 \\ 30 & 168 & 9130 \\ 1680 & 9130 & 615400 \end{bmatrix}_{3 \times 3}$$

$$X^T Y = \begin{bmatrix} 86.95 \\ 456.14 \\ 25190.2 \end{bmatrix}_{3 \times 1}$$

$$W = \begin{bmatrix} 8.56558 \\ 1.19651 \\ -0.00020 \end{bmatrix}$$

### ASSUMPTION OF MULTIPLE REGRESSION:

no multicollinearity (no correlation between variables)

To eliminate correlated features (check for  $d C_2$  times)

1. find correlation (keep threshold to remove 1 feature)

$$-1 \leq \rho_{AB} = \frac{\text{covar}(A, B)}{\sigma_A \sigma_B} \leq +1$$

2. variance inflation factor (VIF)

$$VIF = \frac{1}{1-r^2} \quad r^2 - \text{coeff of determination}$$

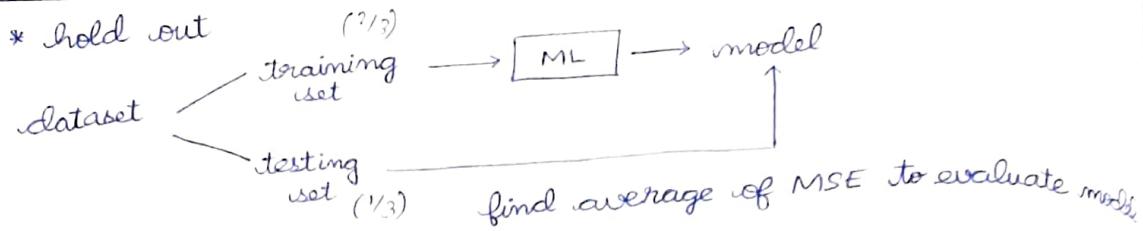
$$0 \leq r^2 \leq 1$$

### MEASURES:

$$\text{adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{(n-d-1)} \quad \begin{array}{l} \text{when } d \text{ increases,} \\ R^2 \text{ decreases} \end{array}$$

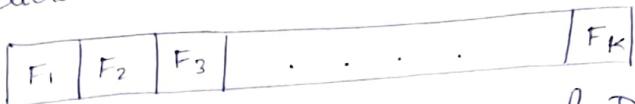
use adjusted  $R^2$  for multiple regression to avoid multicollinearity

## MODEL SELECTION:



- \* k fold cross validation

divide data set into k folds



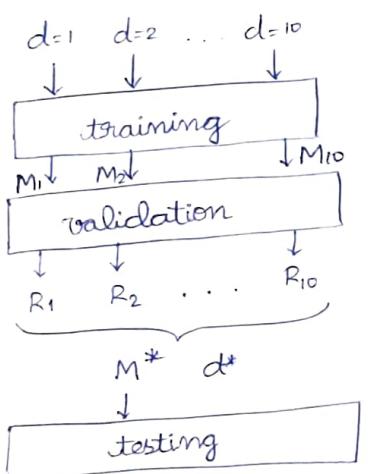
i<sup>th</sup> item used as testing set and  $D - F_i$  as training set

- \* leave one out

n fold cross validation (one fold one datapoint)

## FINDING OPTIMUM VALUE OF HYPERPARAMETER:

(3-way split)



M - model

R - coeff of correlation

d - degree

$$M^* = \underset{k}{\operatorname{argmax}} \{R_{ik}\} \quad 1 \leq k \leq 10$$

SSE is also called training error. If  $\text{SSE}/TE = 0$  then the model passes through all data points.

bias - training error

variance - test error

- \* If bias is low, model fits the training set well / complex model and if bias is high, model is simple / doesn't fit training set well.
- \* If variance is low, then deviation between testing error & training error is less.

low bias & low variance is the best model

	low bias	high bias
low variance	ideal	simple
high variance	complex overfitting	simple underfitting

- \* Overfitting and underfitting occur due to
  - 1) insufficient training set
  - 2) noise / outlier
- \* To increase size of dataset use
  - 1) augmentation (in image analysis, get data from all directions by tilting, rotating)
  - 2) smote analysis (for every data point create another synthetic data)
- \* Overfitting
  - 1) ridge regression / L<sub>2</sub> regularization
  - 2) lasso regression / L<sub>1</sub> regularization
  - 3) elastic net regression (combo of 1 & 2)

} only for parametric problem
- \* regularization (applied only for parametric problems)

Find  $W$

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{regularization term} \rightarrow \|W\|^2$$

$L_1$  norm :  $\text{val}(w_0) + \text{val}(w_1) + \dots + \text{val}(w_d)$

$$L_2 \text{ norm} : \sqrt{w_0^2 + w_1^2 + \dots + w_d^2} \quad (\text{d dimensional circle / sphere})$$

$\|W\|$  is also convex.

$\therefore$  This optimisation problem can be solved by gradient descent method.

#### \* Ridge regularization / L<sub>2</sub> reg

Find  $W$

$$\min \text{SSE} + \lambda \cdot L_2 \text{ norm of } W$$

$\Downarrow$

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|W\|^2$$

$\Downarrow$

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^d w_j$$

## ML LAB:

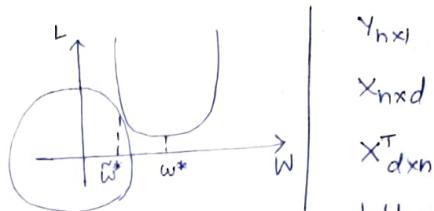
Apply multiple regression on house price prediction by considering continuous features and evaluate the model using hold out method and calculate the average performance measures SSE, SST, SSR, adjusted  $r^2$ . Moreover, before constructing model, eliminate multicollinearity if exists. Apply ridge and lasso regression and calculate measures. Compare the 3 models.

### Ridge regularization:

Find  $w$

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

↓



$y_{n \times 1}$

$x_{n \times d}$

$x^T_{d \times n}$

$w_{d \times 1}$

$w^T_{1 \times d}$

Find  $w$

$$\min \sum_{i=1}^n (y_i - w^T x_i)^2 \quad (x_i \text{ is a col vector})$$

s.t.c

$$w_1^2 + w_2^2 + w_3^2 + \dots + w_d^2 = 1$$

↓

$$\lambda \sum w_j^2 - \lambda$$

$$\text{Find } w \quad \min \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \left( \sum_{j=1}^d w_j^2 - 1 \right)$$

$\lambda \rightarrow$  lagrange multiplier

$$L = (y - xw)^T (y - xw) + \lambda w^T w$$

$L \rightarrow$  scalar  $\otimes$

$(y - xw)_{n \times 1}$

$$L = (y - xw)_{n \times n}^T (y - xw)_{n \times 1} + \lambda w^T w$$

$$= y^T y - \underbrace{y^T x w}_{\text{individually scalars}} - \underbrace{w^T y}_{\text{individually scalars}} + w^T x^T x w + \lambda w^T w$$

$$x^T = x^T x$$

$$L = y^T y - 2 w^T x^T y + w^T x^T x w + \lambda w^T w$$

$$\frac{\partial L}{\partial w} = -2 x^T y + 2 w^T x^T x + 2 \lambda w = 0$$

$$x^T y = w^T (x^T x + \lambda I)$$

$$w = (x^T x + \lambda I)^{-1} x^T y$$

If  $\lambda = 0$ , ridge regularization becomes multiple regression.

If  $\lambda \rightarrow \infty$ ,  $w$  becomes close to 0.  $\therefore$  features are given less importance and collinearity reduces

#### ASSUMPTIONS:

1. No bias ( $w_0 = 0$ )

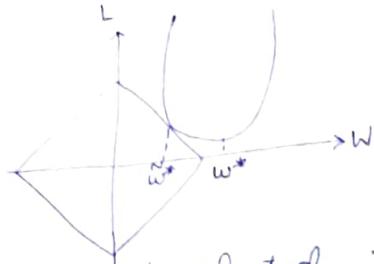
2. Data points are normalized (centered around the origin)

## Lasso Regularization :

Min

SSE + L<sub>1</sub> norm

$$\sum (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^d |w_j|$$

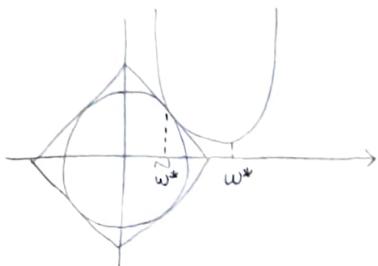


When  $\lambda \rightarrow \infty$ ,  $w$  becomes 0.

If  $w$  is 0 then particular feature is not selected.  $\therefore$  It can be used for feature selection technique.

## Elastic net:

combining Both L<sub>1</sub> and L<sub>2</sub> norm



## CLASSIFICATION :

Given

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^n$$

$x_i \in \mathbb{R}^d$  [discrete / continuous]

$y_i \in \{c_1, c_2\}$  (always discrete)

$$f: x_i \rightarrow y_i$$

objective

$$f: \mathbb{R}^d \rightarrow \{c_1, c_2\}$$

line - linear

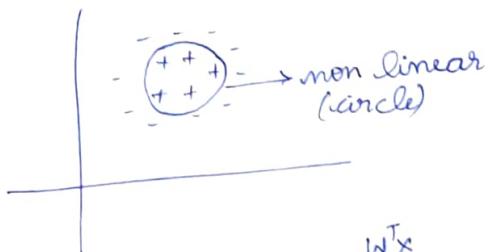
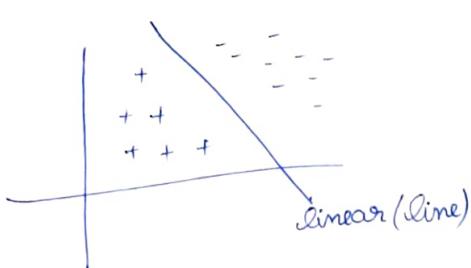
plane - 2D

- nD

hyperplane

## Logistic regression

It is a linear classifier



$\mathcal{D} \Rightarrow$  linear regression  $\Rightarrow w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d \Rightarrow \mathbb{R} (-\infty \text{ to } \infty)$

$$\hat{y} = \begin{cases} c_1 & \text{if } \sigma(w^T x) \geq 0.5 \\ c_2 & \text{o/w} \end{cases}$$

threshold

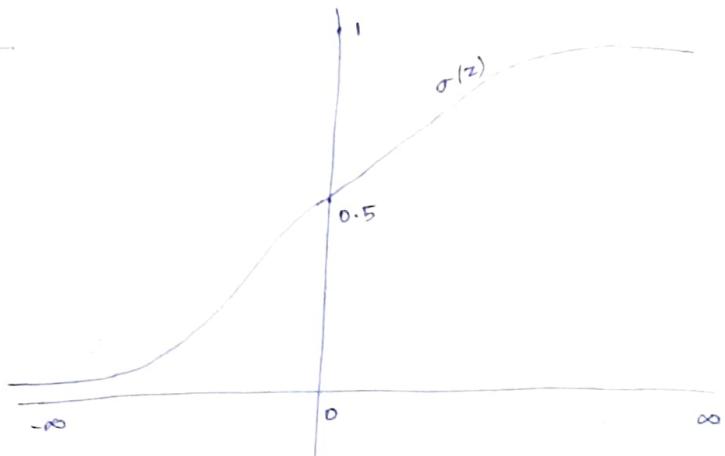
(converts  $\hat{y}$  values from  $-\infty$  to  $\infty$  to 0 to 1 range)

sigmoid function

$0 < \sigma(x) < 1$

### SIGMOID FUNCTION:

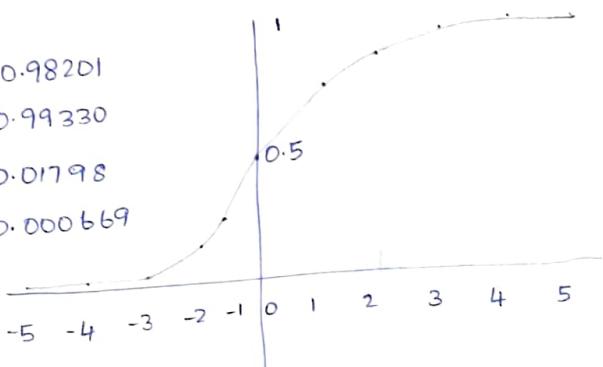
$$\sigma(z) = \frac{1}{1+e^{-z}}$$



Plot  $\sigma(z)$  when  $z$  is in the range of -5 to +5.

$$z = 0 \quad \sigma(z) = 0.5$$

1	0.73105	4	0.98201
2	0.88079	5	0.99330
3	0.95257	-4	0.01798
-1	0.2689	-5	0.000669
-2	0.11920		
-3	0.04742		



### PROPERTIES: S-shaped function

- \*  $0 < \sigma(z) < 1$  for  $z, -\infty \leq z \leq \infty$

- \*  $\sigma(z) = 1/2$  for  $z = 0$

- \*  $z \rightarrow +\infty$  then  $\sigma(z)$  is close to 1

- \*  $z \rightarrow -\infty$  then  $\sigma(z)$  is close to 0

- \* if  $\sigma(z) = \frac{1}{1+e^{-z}}$        $1 - \sigma(z) = \frac{e^{-z}}{1+e^{-z}}$

- \*  $\sigma'(z) = \frac{0+e^{-z}}{(1+e^{-z})^2} = \frac{e^{-z}}{(1+e^{-z})^2} = \sigma(z)(1-\sigma(z))$

Derivative of  $\sigma(z)$  can be written in terms of  $\sigma(z)$  itself.

In classification, SSE gives total number of

misclassification.  $\sum_{i=1}^n (y_i - \hat{y}_i)^2$  (in case of binary classification)

If  $y$  is a Bernoulli random variable,

$$f(y) = P^y (1-P)^{1-y}$$

$$\text{if } y=0 \quad f(y)=P$$

$$\text{if } y=1 \quad f(y)=P$$

## LICISTIC REGRESSION : (discriminative, parametric model)

Given

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^N \quad x_i \in \mathbb{R}^d \quad (x_i \sim N(\mu, \sigma)) \quad y_i \in \{c_1, c_2\}$$

$$f: x_i \rightarrow y_i$$

Objective

$$f: \mathbb{R}^d \rightarrow \{c_1, c_2\}$$

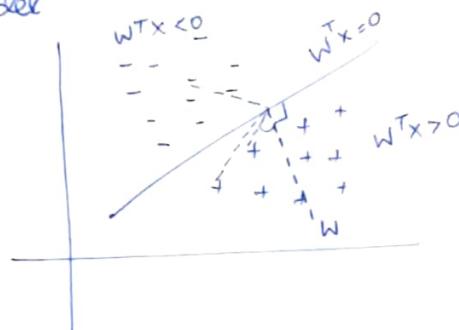
Assumptions:

1.  $x_i$  is continuous  $x_i \sim N(\mu, \sigma)$

2. No or less multi collinearity between features

3. Model  $P(Y=c_1|X) = \frac{1}{1+e^{-w^T x}}$ ;  $P(Y=c_2|X) = 1 - \frac{1}{1+e^{-w^T x}} = \frac{e^{-w^T x}}{1+e^{-w^T x}}$

4. Model



hyperplane  $H = \{x \mid w^T x = 0\}$   
If  $w^T x = 0$  then they are  
lr. (dot product of 2 vectors)  
(or orthogonal)

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \quad w^T x = [w_0 \ w_1 \ \dots \ w_d] \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$= w_0 + w_1 x_1 + \dots + w_d x_d \\ = w_0 + \sum_{j=1}^d w_j x_j$$

$$w \cdot x = \|w\| \|x\| \cos \theta$$

all data points in +ve region have acute angle with w ( $< 90^\circ$ )  
all data points in -ve region have obtuse angle with w ( $> 90^\circ$ )

$$\gamma = \begin{cases} c_1 & w^T x > 0 \\ c_2 & w^T x < 0 \end{cases}$$

$$\hat{y} = \begin{cases} c_1 & \text{if } P(Y=c_1|x) \geq P(Y=c_2|x) \\ c_2 & \text{otherwise} \end{cases}$$

$$= \begin{cases} c_1 & \text{if } \frac{1}{1+e^{-w^T x}} \geq 0.5 \\ c_2 & \text{otherwise} \end{cases}$$

Find  $w$

which correctly classifies all +ve & -ve (or)

$$\min \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{ (or)}$$

$$\min \sum_{i=1}^n (y_i - \sigma(w^T x_i))^2$$

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

$$\frac{\partial L}{\partial w} = -2 \sum_{i=1}^n (y_i - \sigma(w^T x_i)) (\sigma(w^T x_i)(1 - \sigma(w^T x_i))) x_i \quad \left| \begin{array}{l} \hat{y}_i = \sigma(w^T x_i) \\ = -2 \sum (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) x_i \end{array} \right.$$

### GRADIENT DESCENT:

1.  $w_{\text{old}} = \text{initialize}()$

2. while ! convergence {

for each  $(x_i, y_i) \in D$

$$\hat{y}_i = \begin{cases} c_1 & \text{if } \sigma(w^T x_i) \geq 0.5 \\ c_2 & \text{otherwise} \end{cases}$$

if  $y_i \neq \hat{y}_i$

$$w_{\text{new}} = w_{\text{old}} + \eta (y_i - \hat{y}_i) \hat{y}_i (1 - \hat{y}_i) x_i$$

### TEST-1

1.a.i.  $L = y_i \log \hat{y}_i \quad \hat{y}_i = w^T x$

$$L = y_i \log w^T x$$

$$L' = \frac{y_i}{w^T x} (x) = \frac{y_i}{w^T} \quad L'' = \frac{-y_i}{(w^T)^2} \leq 0 \quad \therefore \text{concave}$$

\* epoch - weight updation for every data items after each iteration of data item's grad calculation

2.a.  $SE \text{ (sum of residuals/error)} = 0 = \sum (y_i - \hat{y}_i) \quad \text{for linear regression}$

$x_1 \quad x_2 \quad y$

0 0 0

0 1 1

$$C_1 \Rightarrow 1 \quad C_2 \Rightarrow 0$$

1 0 1

1 1 1

Construct a logistic regression model for the given data using SGD method with  $\eta = 0.5$ ,  $w_0 = 0.5$ ,  $w_1 = -0.5$ ,  $w_2 = 0.5$ .

$x_1$	$x_2$	$w_0$	$w_1$	$w_2$	$y$	$w^T x$	$\sigma(w^T x)$	$\hat{y}$
0	0	0.5	-0.5	0.5	0	0.5	0.622459	$C_1$
0	1	0.42	-0.5	0.5	1	0.92	0.715	$C_1$
1	0	0.42	-0.5	0.5	1	-0.08	0.48	$C_2$
1	1	0.48	-0.44	0.5	1	0.54	0.63	$C_1$
0	0	0.48	-0.44	0.5	0	0.48	0.62	$C_1$
0	1	0.41	-0.44	0.5	1	0.91	0.71	$C_1$
1	0	0.41	-0.44	0.5	1	-0.03	0.49	$C_2$
1	1	0.47	-0.38	0.5	1	0.59	0.64	$C_1$

if  $y$  and  $\hat{y}$  are not equal update  $w$  values

$$\begin{aligned} w_{\text{new}} &= \begin{bmatrix} 0.5 \\ -0.5 \\ 0.5 \end{bmatrix} + (0.5) (-0.62)(0.62)(0.38) \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0.5 \\ -0.5 \\ 0.5 \end{bmatrix} + \begin{bmatrix} -0.073 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.427 \\ -0.5 \\ 0.5 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} w_{\text{new}} &= \begin{bmatrix} 0.42 \\ -0.5 \\ 0.5 \end{bmatrix} + (0.5)(0.52)(0.48)(0.52) \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0.42 \\ -0.5 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0.06 \\ 0.06 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.48 \\ -0.44 \\ 0.5 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} w_{\text{new}} &= \begin{bmatrix} 0.48 \\ -0.44 \\ 0.5 \end{bmatrix} + (0.5) (-0.12)(0.62)(0.38) \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 0.407 \\ -0.44 \\ 0.5 \end{bmatrix} \end{aligned}$$

$$W_{\text{new}} = \begin{bmatrix} 0.41 \\ -0.44 \\ 0.5 \end{bmatrix} + (0.5)(0.51)/(0.49)(0.51) \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 0.41 \\ -0.44 \\ 0.5 \end{bmatrix} + (0.06) \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.47 \\ -0.38 \\ 0.5 \end{bmatrix}$$

PROVE

Logistic regression is a linear model.

PROOF:

$$P(Y=1|X) = \frac{1}{1+e^{-w^T x}}$$

$$P(Y=0|X) = 1 - P(Y=1|X) = \frac{e^{-w^T x}}{1+e^{-w^T x}}$$

$$\text{odds of success} = \frac{P(Y=1|X)}{P(Y=0|X)}$$

$$= \frac{1}{e^{-w^T x}} = e^{w^T x}$$

$$\log \text{odds of success} = \log_e^{w^T x}$$

↓  
(Sigmoid model)

$$= w^T x \text{ is linear}$$

logistic reg → sigmoid function is called logistic function

$$p(B) = 22/38 \quad P(G) = 16/38$$

odds probability: total = 38

# boys = 22    # girls = 16

$$\text{odds of girl} = \frac{16}{22} = \frac{8}{11} / \frac{16}{38}$$

$$\text{odds of boy} = \frac{22}{16} = \frac{11}{8} / \frac{16}{38}$$

$$\text{odd} = \frac{\text{probability (success)}}{\text{probability (failure)}}$$

- \* Find the odds of each category of flower in IRIS dataset:
- \* odd of each class =  $\frac{1}{2}$  ( $\because$  each class 50, total 150,  $P(\text{class}) = \frac{1}{3}$ )

### EVALUATION METRICS / MEASURES:

accuracy

true positive rate (TPR)

precision

false positive rate (FPR)

recall

receiver operating characteristic curve (ROC)

F-measure

area under the ROC curve (AUC)

		predicted	
		+	-
actual	+	true positive false negative	false positive true negative
	-	false positive true negative	

$$P = TP + FN$$

$$N = FP + TN$$

$$\hat{P} = TP + FP \quad \hat{N} = FN + TN$$

1) accuracy =  $\frac{TP+TN}{TP+TN+FP+FN} = \frac{TP+TN}{n}$  n - total no. of datapoints

$$= \frac{TP+TN}{P+N} = \frac{TP+TN}{\hat{P}+\hat{N}}$$

eg:  $n = 100$      $TP = 35$      $TN = 25$   
 $\begin{array}{c} / \\ 50 \\ \backslash \end{array}$      $\begin{array}{c} / \\ 50 \\ \backslash \end{array}$   
balanced dataset

35	15	50
25	25	50
60	40	

$$\text{accuracy} = \frac{60}{100} = 0.60$$

eg:  $n = 100$      $TP = 0$      $TN = 90$   
 $\begin{array}{c} / \\ 90 \\ \backslash \end{array}$      $\begin{array}{c} / \\ 10 \\ \backslash \end{array}$   
skewed dataset / imbalanced

0	10
0	90

$$\text{accuracy} = \frac{90}{100} = 0.90$$

since we are interested in +ve, accuracy is not a good measure for skewed dataset (it didn't predict single +ve correct)

2) precision =  $\frac{TP}{\hat{P}} / \frac{TN}{\hat{N}}$  (no. of +ve predictions w.r.t total +ve predictions)

3) recall =  $\frac{TP}{P} / \frac{TN}{N}$  (no. of +ve predictions w.r.t total +ve in actual dataset)

4) F-measure = harmonic mean of precision and recall

$$= \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2}{\frac{PR}{P+R}} = \frac{2PR}{P+R}$$

(arithmetic mean is not used, since it varies with extreme)

5)  $TPR = \frac{TP}{P}$  (increase TPR)

} both are disjoint

6)  $FPR = \frac{FP}{N}$  (decrease FPR)

7)  $TNR = \frac{TN}{N}$

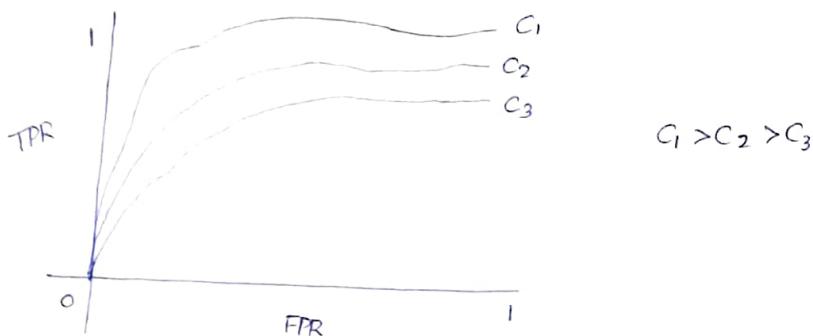
TPR = recall

8)  $FNR = \frac{FN}{P}$

(Calculate all measures)

## RECEIVER OPERATING CHARACTERISTIC CURVE : (ROC CURVE)

TPR Vs FPR



TPR and FPR are chosen because they are disjoint  
TP w.r.t P and FP w.r.t N.

TPR ( $\frac{TP}{P} = \frac{TP}{TP+FN}$ )	FPR ( $\frac{FP}{N} = \frac{FP}{FP+TN}$ )	INFERENCE
0    TP=0, P=FN	0    FP=0, N=TN	all data $\rightarrow$ N Biased to N class
0    TP=0, P=FN	1    FP=1, TN=0, N=FP	all +ve $\rightarrow$ N all -ve $\rightarrow$ P
1    TP=1, P=TP, FN=0	0    FP=0, TN=N	all datapoints correctly classified
1    TP=1, FN=0, P=TP	1    FP=1, TN=0, N=FP	all data $\rightarrow$ P Biased to P class

### 2 KINDS OF CLASSIFIER :

1) discrete  $\rightarrow$  yes or no (perceptron, SVM, DT)

2) probabilistic  $\rightarrow P(\text{yes}/x)$  or  $P(\text{no}/x)$  (logistic regression)

↑

ROC curve only for this  
(change threshold value  $\theta$  to draw ROC curve)

$$\hat{y}_i = \begin{cases} 1 & P(y_i=1/x) \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

arrange in increasing order of x-axis to get smooth curve.

## ALGORITHM

input  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

output : ROC curve

a. for each threshold  $\theta$ , do

a. for each  $(x_i, y_i) \in \mathcal{D}$

$$\hat{y}_i = \begin{cases} 1 & P(y_i = 1 | x_i) \geq \theta \\ 0 & \text{otherwise} \end{cases}$$

$$b. \text{ calculate } TPR_{\theta} = \frac{TP}{P}, FPR_{\theta} = \frac{FP}{N}$$

c. arrange  $(TPR, FPR)$  in increasing order

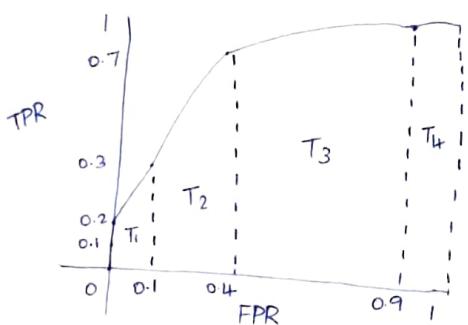
\* Plot the ROC curve with actual class labels and probability scores for classifier.

$y$	$\hat{y}$	$\hat{y}_{0.2}$	$\hat{y}_{0.4}$	$\hat{y}_{0.6}$	$\hat{y}_{0.8}$	$\hat{y}_{0.9}$
P - 0.9	P	P	P	P	P	P
P - 0.8	P	P	P	P	P	P
N - 0.7	P	P	P	P	Z	Z
P - 0.6	P	P	P	Z	Z	Z
P - 0.55	P	P	Z	Z	Z	Z
P - 0.54	P	Z	Z	Z	Z	Z
N - 0.53	P	Z	Z	Z	Z	Z
N - 0.52	P	Z	Z	Z	Z	Z
P - 0.51	P	Z	Z	Z	Z	Z
N - 0.505	P	Z	Z	Z	Z	Z
P - 0.4	P	Z	Z	Z	Z	Z
N - 0.39	P	Z	Z	Z	Z	Z
P - 0.38	P	Z	Z	Z	Z	Z
N - 0.37	P	Z	Z	Z	Z	Z
N - 0.36	P	Z	Z	Z	Z	Z
N - 0.35	P	Z	Z	Z	Z	Z
P - 0.34	P	Z	Z	Z	Z	Z
N - 0.33	P	Z	Z	Z	Z	Z
P - 0.3	P	Z	Z	Z	Z	Z
N - 0.1	Z	Z	Z	Z	Z	Z

$$TPR = \frac{TP}{P} = \frac{0}{10} = 0, \quad \frac{7}{10} = 0.7, \quad \frac{3}{10} = 0.3, \quad \frac{2}{10} = 0.2, \quad \frac{1}{10} = 0.1$$

$$FPR = \frac{FP}{N} = \frac{9}{10} = 0.9, \quad \frac{4}{10} = 0.4, \quad \frac{1}{10} = 0.1, \quad 0, \quad 0$$

FPR	0	0	0.1	0.4	0.9
TPR	0.1	0.2	0.3	0.7	1



$$AUC = \text{area of } (T_1 + T_2 + T_3 + T_4)$$

$$\text{area of } T_1 = \frac{1}{2} \times 0.1 \times 0.3 = 0.025$$

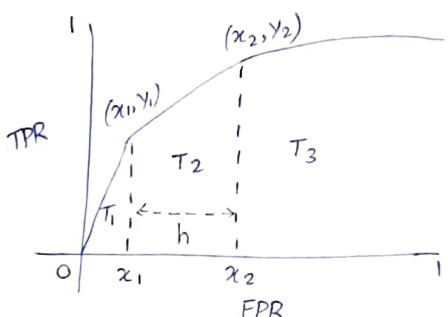
$$\text{area of } T_2 = \frac{1}{2} \times 0.3 \times 1.0 = 0.15$$

$$\text{area of } T_3 = \frac{1}{2} \times 0.5 \times 1.7 = 0.425$$

$$\text{area of } T_4 = \frac{1}{2} \times 0.1 \times 2 = 0.1$$

$$\boxed{AUC = 0.7}$$

### AREA UNDER THE ROC



AUC of ideal model is 1

$$\begin{aligned} \text{area of trapezium} &= \frac{1}{2} \times \text{height} \times (\text{sum of parallel sides}) \\ &= \frac{1}{2} \times (x_2 - x_1) \times (y_1 + y_2) \\ &= \frac{1}{2} \times (FPR_2 - FPR_1) \times (TPR_1 + TPR_2) \end{aligned}$$

DAY	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAY TENNIS
1	sunny	85/hot	85/H	weak	no
2	sunny	80/hot	90/H	strong	no
3	overcast	83/hot	78/H	weak	yes
4	rain	70/mild	96/H	weak	yes
5	rain	68/cool	80/H	weak	yes
6	rain	65/cool	70/N	strong	no
7	overcast	64/cool	65/N	strong	yes
8	sunny	72/mild	95/H	weak	no
9	sunny	69/cool	70/N	weak	yes
10	rain	75/mild	80/N	weak	yes
11	sunny	75/mild	70/N	strong	yes

H - High N - normal

12	overcast	12 / mild	70 / H	strong	yes
13	overcast	21 / hot	75 / N	weak	yes
14	rain	71 / mild	80 / H	strong	no

NBAYE BAYES CLASSIFIER (probabilistic classifier)

Given

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^n$$

$x_i \in$  discrete / continuous

$$y_i \in \{c_1, c_2\}$$

$$f: x_i \rightarrow y_i$$

objective

$$\hat{f}: \mathbb{R}^d \rightarrow \{c_1, c_2\}$$

Hypothesis

maximum a posterior (MAP) probability

$$\hat{y} = \underset{c_k}{\operatorname{argmax}} \{P(c_k | x)\} \quad c_k \in \{c_1, c_2\} \rightarrow \text{map theorem}$$

$$\hat{y} = \begin{cases} c_1 & \text{if } P(c_1 | x) \geq P(c_2 | x) \\ c_2 & \text{otherwise} \end{cases}$$

assumption  $P(F_1 | Y, F_2, F_3 \dots F_d) = P(F_1 | Y)$

all features are conditionally independent given class

\* From play tennis data set,

prior probability: without observing input features and based on class label if we calculate probability then it is called prior probability.

$$P(\text{yes}) = 9/14 \quad P(\text{no}) = 5/14$$

posterior probability: after observing input features we obtain probability

$$P(\text{yes} | x) \rightarrow \text{conditional probability}$$

\* It is a generative model

prior  $P(\text{yes})$   $P(\text{no})$

posterior  $P(\text{yes} | X)$   $P(\text{no} | X)$

joint  $P(X, Y)$

(sum of all joint probabilities)  
denominator of Bayes' theorem ( $P(Y)$ )

marginal / total  $P(X|Y) = \frac{P(X \cap Y)}{P(Y)}$  (or)  $P(X|Y) = \frac{P(X, Y)}{P(Y)} = \frac{P(Y|X)P(X)}{P(Y)}$  (or)  $\#_1 / \#_2$

eg:  $P(O = \text{sunny}, PT = \text{yes}) = \frac{2}{14}$   $P(O = \text{sunny}, PT = \text{no}) = \frac{3}{14}$

$P(O = \text{overcast}, PT = \text{yes}) = \frac{4}{14}$

$P(O = \text{rain}, PT = \text{yes}) = \frac{3}{14}$

$P(O = \text{overcast}, PT = \text{no}) = 0$

$P(O = \text{rain}, PT = \text{no}) = \frac{2}{14}$

H/W Construct joint probability table and conditional probability for all features in dataset.

JPT:

outlook

	yes	no
sunny	$2/14$	$3/14$
overcast	$4/14$	$0$
rain	$3/14$	$2/14$

temperature

	yes	no
hot	$2/14$	$2/14$
mild	$4/14$	$2/14$
cool	$3/14$	$1/14$

humidity

	yes	no
high	$3/14$	$4/14$
normal	$6/14$	$1/14$

wind

	yes	no
weak	$6/14$	$2/14$
strong	$3/14$	$3/14$

CPT:

	yes	no
outlook	$S: 2/9$	$3/5$
	yes	no
O	$4/9$	$0$
R	$3/9$	$2/5$

Temperature H M C

	yes	no
H	$2/9$	$2/5$
M	$4/9$	$2/5$
C	$3/9$	$1/5$

Humidity	H	yes 3/9	no 4/5	wind	S W	3/9 6/9	3/5 2/5
	N	6/9	1/5				

joint probability of 2 RV:

$$P(A, B) = P(A|B)P(B) \text{ (or) } P(B|A)P(A)$$

joint probability of 3 RV:

$$\begin{aligned} P(A, B, C) &= P(A|B, C)P(B, C) \\ &= P(A|B, C)P(B|C)P(C) \end{aligned}$$

joint probability of n RV: (by chain rule)

$$\begin{aligned} P(A_1, A_2, \dots, A_n) &= P(A_1|A_2, A_3, \dots, A_n)P(A_2|A_3, A_4, \dots, A_n)P(A_3|A_4, \dots, A_n) \\ &= P(A_1|A_2, A_3, \dots, A_n)P(A_2|A_3, A_4, \dots, A_n)P(A_3|A_4, \dots, A_n)P(A_4|A_5, \dots, A_n) \\ &= P(A_1|A_2^n)P(A_2|A_3^n) \dots P(A_{n-2}|A_{n-1}^n)P(A_{n-1}|A_n)P(A_n) \\ &= \left[ \prod_{j=1}^{n-1} P(A_j|A_{j+1}^n) \right] P(A_n) \end{aligned}$$

NAIVE BAYES CLASSIFIER:

Given  $X$ , predict  $\hat{Y}$

Probability  $P(\hat{Y}=1|X) = \frac{P(X, \hat{Y}=1)}{P(X)} = \frac{P(X, \hat{Y}=1)}{P(X, \hat{Y}=0) + P(X, \hat{Y}=1)}$

$$P(\hat{Y}=0|X) = \frac{P(X, \hat{Y}=0)}{P(X, \hat{Y}=0) + P(X, \hat{Y}=1)}$$

$$\hat{Y} = \begin{cases} 1 & P(\hat{Y}=1|X) \geq P(\hat{Y}=0|X) \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{Y} = \underset{C_k}{\operatorname{argmax}} P(C_k|X) \quad \text{where } C_k \in \{0, 1\}$$

It is called Naive Bayes classifier, because of the naive conditional independance.

$$P(X_1|X_2, Y) = P(X_1|Y)$$

$$P(X_1|X_2, \dots, X_d, Y) = P(X_1|Y)$$

Numerator of Bayes theorem:

$$P(x|y) P(y) = P(x,y)$$

$$= P(x_1, x_2, x_3, \dots, x_d, y)$$

$$= P(x_1|x_2^d, y) P(x_2|x_3^d, y) \dots P(x_{d-1}|x_d, y) P(x_d|y) P(y)$$

By conditional independence,

$$= P(x_1|y) P(x_2|y) \dots P(x_{d-1}|y) P(x_d|y) P(y)$$

$$\stackrel{d}{=} \left[ \prod_{j=1}^d P(x_j|y) \right] P(y)$$

Denominator of Bayes theorem:

$$P(x) = P(x_1, x_2, \dots, x_d)$$

$$= P(x, y=0) + P(x, y=1)$$

$$= P(x_1, x_2, \dots, x_d, y=0) + P(x_1, x_2, \dots, x_d, y=1)$$

$$= \prod_{j=1}^d P(x_j|y=0) P(y=0) + \prod_{j=1}^d P(x_j|y=1) P(y=1) \quad (\text{from numerator})$$

$$P(y=c_1|x) = \frac{P(x, y=c_1)}{P(x, y=c_1) + P(x, y=c_2)}$$

$$= \frac{P(x_1|y=c_1) P(x_2|y=c_1) \dots P(x_d|y=c_1) P(c_1)}{\prod_{j=1}^d P(x_j|y=c_1) P(y=c_1) + \prod_{j=1}^d P(x_j|y=c_2) P(y=c_2)}$$

$$P(y=c_2|x) = \frac{P(x, y=c_2)}{P(x, y=c_1) + P(x, y=c_2)}$$

$$= \frac{P(x_1|y=c_2) P(x_2|y=c_2) \dots P(x_d|y=c_2) P(c_2)}{\prod_{j=1}^d P(x_j|y=c_1) P(y=c_1) + \prod_{j=1}^d P(x_j|y=c_2) P(y=c_2)}$$

model is conditional & prior probability. It is generative model because it constructs model for each class. It also finds joint probability. It tries to learn data space.

## TEST DATASET:

outlook = rainy    temperature = cool    humidity = normal  
 windy = strong

$$P(PT = \text{yes} | \langle \text{rainy}, \text{cool}, \text{normal}, \text{strong} \rangle)$$

$$= P(O_L = \text{rainy} | \text{yes}) P(T = \text{cool} | \text{yes}) P(H = \text{normal} | \text{yes}) P(W = \text{strong} | \text{yes}) \\ P(\text{yes}) \\ P(\text{rainy} | \text{yes}) P(\text{cool} | \text{yes}) P(\text{normal} | \text{yes}) P(\text{strong} | \text{yes}) P(\text{yes}) + \\ P(\text{rainy} | \text{no}) P(\text{cool} | \text{no}) P(\text{normal} | \text{no}) P(\text{strong} | \text{no}) P(\text{no})$$

$$= \frac{3/9 \times 3/9 \times 6/9 \times 3/9 \times 9/14}{(3/9 \times 3/9 \times 6/9 \times 3/9 \times 9/14) + (2/5 \times 1/5 \times 1/5 \times 3/5 \times 5/14)}$$

$$= \frac{0.01587}{0.01587 + 0.00343} \approx 0.822279793 \quad (\text{Play tennis} = \text{yes})$$

Buys computer dataset: (H/W)

Sunburn dataset:

hair	height	weight	lotion	result
blonde	average	light	no	sunburned
blonde	tall	average	yes	none
brown	short	average	yes	none
blonde	short	average	no	sunburned
red	average	heavy	no	sunburned
brown	tall	heavy	no	none
brown	average	heavy	no	none
blonde	short	light	yes	none

CPT: hair	Blonde	result		height	short	average	tall	result
		2/3	2/5					
Blonde	2/3	2/5	height	short	1/3	2/5		
Brown	0/3	3/5	average	average	2/3	1/5		
Red	1/3	0/5	tall	tall	0/3	2/5		

		result		result			
		light	1/3	1/5			
weight	average		1/3	2/5	lotion		
	heavy		1/3	2/5	yes	0/3	3/5
				no	3/3	2/5	

Predict test data for patient with brown hair, tall, heavy weight & uses lotion.

$$P(\text{yes} | \text{brown, tall, heavy, yes})$$

$$= \frac{P(\text{brown} | \text{yes}) P(\text{tall} | \text{yes}) P(\text{heavy} | \text{yes}) P(\text{yes} | \text{yes})}{P(\text{brown} | \text{yes}) P(\text{tall} | \text{yes}) P(\text{heavy} | \text{yes}) P(\text{yes} | \text{yes}) + P(\text{brown} | \text{no}) P(\text{tall} | \text{no}) P(\text{heavy} | \text{no}) P(\text{yes} | \text{no})}$$

$$= 0 \quad (\text{since insufficient dataset})$$

•  $\times$  disadvantage of Naive Bayes

## 2 CHALLENGES:

- \* zero frequency
- \* continuous features

### zero frequency problem :

$F_1, F_2, \dots, F_d$  rare features

$$F_j = \{v_1, v_2, \dots, v_m\}$$

$$P(F_j = v_m | y_k) = \frac{n_{jmk}}{n_k}$$

where

$n_{jmk}$  = # times the feature  $F_j$  takes value  $v_m$  in class  $y_k$

$n_k$  = # datapoints in class  $y_k$

$$\text{If } n_{jmk} = 0, \quad P(F_j = v_m | y_k) = \frac{n_{jmk} + 1/m}{n_k + 1}$$

### smoothing technique

i) Laplacian (add 1/m)

ii) add one

		yes	no
eg: outlook	S	2/9	3/5 $\frac{3+1}{6}$
	O	4/9	0 $\frac{0+1}{6}$
	R	3/9	2/5 $\frac{2+2}{5+1}$

By Laplacian technique,

		result				result	
		Blonde	$\frac{7}{12}$	$\frac{7}{18}$	short	$\frac{4}{12}$	$\frac{2}{5}$
hair	-	Brown	$\frac{1}{12}$	$\frac{10}{18}$	average	$\frac{7}{12}$	$\frac{1}{5}$
		red	$\frac{4}{12}$	$\frac{1}{18}$	tall	$\frac{1}{12}$	$\frac{2}{5}$
weight -		light	$\frac{1}{3}$	$\frac{1}{5}$	lotion -		$\frac{1}{8}$
		average	$\frac{1}{3}$	$\frac{2}{5}$	yes	$\frac{7}{8}$	$\frac{3}{5}$
		heavy	$\frac{1}{3}$	$\frac{2}{5}$	no	$\frac{1}{8}$	$\frac{2}{5}$

$P(\text{yes} | \langle \text{Brown}, \text{tall}, \text{heavy}, \text{yes} \rangle)$

$$\begin{aligned}
 &= \frac{\frac{1}{12} \times \frac{1}{12} \times \frac{1}{3} \times \frac{1}{8} \times \frac{3}{8}}{\left( \frac{1}{12} \times \frac{1}{12} \times \frac{1}{3} \times \frac{1}{8} \times \frac{3}{8} \right) + \left( \frac{10}{18} \times \frac{2}{5} \times \frac{2}{5} \times \frac{3}{5} \times \frac{5}{8} \right)} \\
 &= 0.0032476 \quad (\text{result} = \text{none}) \quad \frac{0.0001085}{0.0001085 + 0.0323}
 \end{aligned}$$

Buy computer dataset :

age	income	student	credit rating	buy
youth	high	no	fair	no
youth	high	no	excellent	no
middle-age	high	no	fair	yes
senior	medium	no	fair	yes
senior	low	yes	fair	no
senior	low	yes	excellent	yes
middle-age	low	yes	excellent	no
youth	medium	no	fair	yes
youth	low	yes	fair	yes
senior	medium	yes	fair	yes
youth	medium	yes	excellent	yes
middle-age	medium	no	excellent	yes
middle-age	high	yes	fair	yes

CPT:

		yes	no	Buy	add one,
age	youth	2/9	3/4	2/9	10/15
	middle-age	4/9	0/4	4/9	1/15
	senior	3/9	1/4	3/9	4/15
income	high	2/9	2/4		
	low	3/9	1/4		
	medium	4/9	1/4		
student	yes	6/9	1/4		
	no	3/9	3/4		
credit rating	fair	6/9	2/4		
	excellent	3/9	2/4		

$P(\text{no} | \text{medium, senior, no, excellent})$

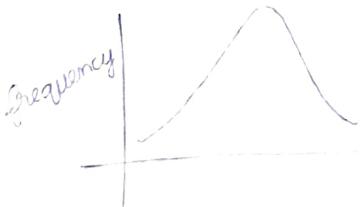
$$= \frac{\frac{1}{4} \times \frac{1}{4} \times \frac{3}{4} \times \frac{2}{4} \times \frac{4}{13}}{\left( \frac{1}{4} \times \frac{1}{4} \times \frac{3}{4} \times \frac{2}{4} \times \frac{4}{13} \right) + \left( \frac{3}{9} \times \frac{4}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{13} \right)} = \frac{0.00721}{0.00721 + 0.011396}$$

$$= 0.3875 \quad (\text{buy computer} = \text{yes})$$

continuous features:

\* convert to discrete values / discretisation

\* assume feature follows gaussian / normal dist



using pdf construct CPT

$$\text{pdf} = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$

known:  $x$

unknown:  $\mu, \sigma$

day	outlook	temperature	humidity	wind	play
1	sunny	85	85	weak	no
2	sunny	80	90	strong	no
3	overcast	83	78	weak	yes
4	rain	70	96	weak	yes
5	rain	68	80	weak	yes
6	rain	65	78	strong	no
7	overcast	64	65	strong	yes
8	sunny	72	95	weak	no
9	sunny	69	70	weak	yes
10	rain	75	80	weak	yes
11	sunny	75	70	strong	yes
12	overcast	72	90	strong	yes
13	overcast	81	75	weak	yes
14	rain	71	80	strong	no

test data: sunny, 40, 70, strong

temperature | play tennis = yes

$$\bar{x} = 73 \quad \sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = 6.164414$$

temperature | play tennis = no

$$\bar{x} = 74.6 \quad \sigma = 7.8930349$$

humidity | play tennis = yes

$$\bar{x} = 78.222 \quad \sigma = 9.88407811$$

humidity | play tennis = no

$$\bar{x} = 84 \quad \sigma = 9.61769203$$

CPT:

	yes	no
temperature	$3.87 \times 10^{-8}$	$2.27 \times 10^{-6}$

humidity	0.028593	0.014271
----------	----------	----------

$$\begin{aligned}
 & \text{yes} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle \\
 & P(\text{yes} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle) = P(\text{yes} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle) P(\text{strong} | \text{yes}) P(\text{yes}) \\
 & + P(\text{no} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle) P(\text{no} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle) P(\text{strong} | \text{no}) P(\text{no}) \\
 & P(\text{yes} | \langle \text{sunny}, \text{uv}, 70, \text{strong} \rangle) = \frac{\frac{2}{9} \times (3.87 \times 10^{-8}) \times (0.028599) \times \frac{3}{9} \times \frac{9}{14}}{\left[ \left( \frac{2}{9} \times (3.87 \times 10^{-8}) \times (0.028599) \times \frac{3}{9} \times \frac{9}{14} \right) + \left( \frac{3}{5} \times (3.37 \times 10^{-6}) \times (0.014371) \times \frac{3}{5} \times \frac{5}{14} \right) \right]} \\
 & = \frac{5.27 \times 10^{-11}}{(5.27 \times 10^{-11}) + (6.2267 \times 10^{-9})} \\
 & = 0.00839 \quad (\text{play tennis} = \text{no})
 \end{aligned}$$

### MULTIVARIATE GAUSSIAN:

Let  $x_1, x_2, \dots, x_k$  be  $k$  features

$$x_1 \sim N(\mu_1, \sigma_1)$$

$$x_2 \sim N(\mu_2, \sigma_2)$$

 $\vdots$ 

$$x_k \sim N(\mu_k, \sigma_k)$$

$$ax_1 + bx_2 \sim N(\mu, \Sigma)$$

where  $\mu$  is mean vector of len k

$\Sigma$  is covar matrix of len  $k \times k$

$$P(x) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

z score gaussian (or)  
Mahalanobi's distance

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1k} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2k} \\ \vdots & \vdots & & \vdots \\ \sigma_{k1} & \sigma_{k2} & \dots & \sigma_{kk} \end{bmatrix}$$

$$|\Sigma| = \sigma_1^2 \sigma_2^2 \dots \sigma_k^2$$



$$\therefore (x - \mu)^T = (x - \mu)(x - \mu)^T \text{ w.r.t } \Sigma \text{ (covar)}$$

If  $\Sigma$  is a unit matrix, then

Mahalanobi's dist = Euclidean dist

$$x_{k \times 1} \quad \mu_{k \times 1} \quad \Sigma_{k \times k}$$

Consider the class means  $\mu$  and  $\Sigma$  for classes  $C_1$  and  $C_2$

$$\mu_1 = 1, 3$$

$$\mu_2 = 5, 5$$

$$\Sigma_1 = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}$$

$$\Sigma_2 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

Classify  $3, 4$  using NB, full NB

$$P(\text{class } C_1) = 1/4 \quad P(\text{class } C_2) = 3/4$$

$$P(\text{class } C_1 | \langle 3, 4 \rangle) = \frac{P(\langle 3, 4 \rangle | C_1) P(C_1)}{P(\langle 3, 4 \rangle)}$$

$$P(\langle 3, 4 \rangle | C_1) = \frac{1}{(2\pi)^{2/2} (10)^{1/2}} e^{-\frac{1}{2} [2, 1] \begin{bmatrix} \frac{1}{10} & 0 \\ 0 & \frac{5}{10} \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix}}$$

$$|C_1| = \begin{vmatrix} 5 & 0 \\ 0 & 2 \end{vmatrix} = 10 \quad (\text{By conditional independence assumption})$$

$$P(\langle 3, 4 \rangle | C_1) = 0.02628 \times 1/4 = 0.0065718 \quad \begin{bmatrix} 4/10 & 5/10 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$P(\langle 3, 4 \rangle | C_2) = \frac{P(\langle 3, 4 \rangle | C_2) P(C_2)}{P(\langle 3, 4 \rangle)} = \frac{1}{(2\pi)^{2/2} (2)^{1/2}} e^{-\frac{1}{2} [-2, -1] \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \end{bmatrix}}$$

$$P(\langle 3, 4 \rangle | C_2) = 0.02512$$

$$P(\langle 3, 4 \rangle | C_2) = 0.02512 \times 3/4 = 0.018843 \quad \begin{bmatrix} -1 & -1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \end{bmatrix} = 3$$

$\therefore$  prediction :  $C_2$

#### \* PRIVACY PRESERVING ML ALGORITHM:

model 1 - model for  $n$  datapoints

model 2 - model for  $n-1$  datapoints

(Build for  $n$  datapoints, remove 1 datapoint randomly)

If model 1 and model 2 are same then it is called

privacy preserving ML.

e.g.: differentially private KNN, decision trees

- \* Add gaussian noise or laplacian noise to X and Y. It is called perturbation

↓  
It is done to the data stored in cloud in order to attain privacy.

algo due to

- \* Security: suppose the importance of wrong data eg: robust kNN added by attacker/hacker to suppress efficiency of algorithm

### EXPLAINABLE ML ALGORITHM / EXPLAINABILITY:

The model should predict the outcome also explain on why it is predicted that way.

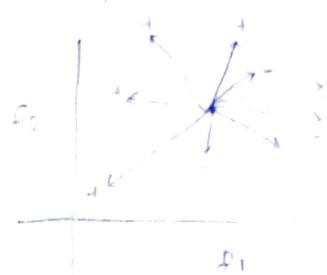
### K-NEAREST NEIGHBOURS (KNN)

- classification
- regression

- \* time consuming as model is constructed in testing phase
  - \* no training phase (memory based model, training set in memory)
  - \* lazy learners (training set req at testing phase) model constructed for each data point. diff model for diff data point (kNN) since it is diff, model is constructed at the time of testing.
  - \* Other algorithms are eager learners - model constructed in training phase (no training set required, only model required at testing)
  - \* instance based
- At the time of testing, find the distance between the test data and other training set data points.

$$X = \{x_1, x_2\} \quad (\text{test data})$$

$$Y = \{y_1, y_2, y_3\} \quad (\text{data points})$$



$$X = \{a, b, c, \dots\}$$

$$X_1 = \{a_1, b_1, c_1, \dots\}$$

$$\text{dis}(x, x_i) = \sqrt{\sum_{j=1}^d (x_j - x_{ij})^2}$$

↓  
can use Euclidean dist / Manhattan dist /  
Minkowski dist / Chebyshev dist / Mahalanobi's dist

K - neighbours, K - hyperparameter

Find optimal value of K by three way split

$$K = 5 \quad C_1 = 2 \quad C_2 = 3$$

test data should be assigned to -ve class  
avoid even values for K incase of tie

ASSUMPTION :

closest data points belong to same class (closeness)

Regression: Don't do voting. Y value of test data  
is the mean of y values of K nearest neighbours.  
can be any measure

$$X_1 = \text{acid durability} \quad X_2 = \text{strength} \quad Y = \text{classification}$$

7	7	bad
7	4	bad
3	4	good
1	4	good

Predict test data (2, 2) & (5, 6) using training set  
with  $\boxed{k=3}$  by applying KNN.

$$* (2,2) \quad \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

<u>training data</u>	<u>distance</u>	<u>class</u>
(7,7)	7.07106781	
(7,4)	5.38516	
(3,4)	2.236	
(1,4)	2.236	

$K$  neighbour  
 $K=3$

$\therefore (2,2)$  is predicted good

$$* (5,6)$$

<u>training data</u>	<u>distance</u>	<u>class</u>
(7,7)	2.236	bad
(7,4)	2.828	bad
(3,4)	2.828	good
(1,4)	4.4721	

$\therefore (5,6)$  is predicted bad

ALGORITHM - KNN: (pre-processing & normalisation) \*

input  $D = \{x_i, y_i\}_{i=1}^n$  output 1 or 0  
(classification)

$k$ , test data  $(x, y)$

1. for each test data  $(x, y) \in \text{test}$

a.  $C_0 = C_1 = 0$

b. for each train data  $(x_i, y_i) \in D$

i.  $\text{dist}_i = \text{distance}(x_i, x)$

$O(n^d)$

c. sort  $\text{dist}_i$

$O(n \log n)$

d.  $N = \{(x_i, y_i) \mid \text{dist}_i \text{ is minimum}\}$  &  $|N| = k \quad O(k)$

e.  $C_0 = \sum_{x_i, y_i \in N} 1(y_i = 0) \quad e. \hat{y} = \frac{\sum_{x_i, y_i \in N} y_i}{k} \quad O(k)$

f.  $C_1 = \sum_{x_i, y_i \in N} 1(y_i = 1) \quad (\text{regression}) \rightarrow \text{mean} \quad O(k)$

f.  $\hat{y} = \begin{cases} 1 & \text{if } C_1 > C_0 \\ 0 & \text{o/w} \end{cases} \quad O(1)$

for one test data,

$$\text{time complexity} = \max(n^d, n \log n)$$

### Normalisation technique:

z-score normalisation

min-max normalisation

Model: K nearest neighbours / closed polygon covers all K neighbours)

### DECISION TREES:

classification & regression

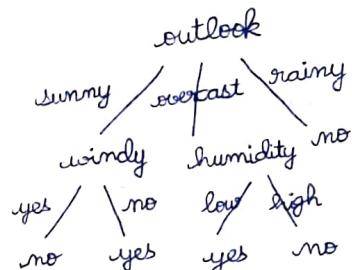
model - tree

pure / homogeneous

If all datapoints in D belong to same class.

objective

Given  $D = \{x_i, y_i\}_{i=1}^n$  & impure



- \* interior nodes - features
- \* branches - possible values
- \* leaf nodes - classes

It aims to convert D into pure by splitting D based on features.

### impurity measures

- entropy
- conditional entropy
- information gain (ID<sub>3</sub>) iterative dicotomizer
- gain ratio (C4.5)
- gini index (CART) classification & regression tree

entropy  $0 \leq H(Y) \leq \log_2 m$  (surprisal)

impurity of data

\* information content in a RV

\* randomness

\* # bits required to represent a RV (from information theory)

(information  $\propto 1/\text{probability}$ )

when  $y$

$H(y)$  entropy (no of bits to represent)

2

1

4

2

:

:

$m$

$\log_2 m$

equally likely

$$H(y) = \log_2 m = -\log_2 (1/m) = -\log P \text{ (unbiased)}$$

If  $y$  is biased,

$$H(y) = -\sum_{i=1}^m p_i \log p_i$$

$$E(y) = \text{mean}$$

Calculate entropy of play tennis, buy computer & iris dataset.

$$* \text{play tennis} = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$* \text{buy computer} = -\frac{8}{13} \log_2 \frac{8}{13} - \frac{5}{13} \log_2 \frac{5}{13} = 0.9611$$

$$* \text{iris dataset} = -\log \frac{1}{3} = 1.58 \text{ (unbiased)}$$

Note:  $\log 0 = 0$

for pure dataset, entropy is 0.

H/W Plot the curve for entropy of Bernoulli RV by finding entropy of probability distribution

P(head)      P(tail)

0

1

0.25

0.75

0.50

0.50

0.75

0.25

1

0

$$\begin{array}{|c|c|} \hline \text{PMF} & p^x (1-p)^{1-x} \\ \hline \text{entropy} & \\ \hline \end{array}$$

$$-0 \log 0 - 1 \log 1 = 0$$

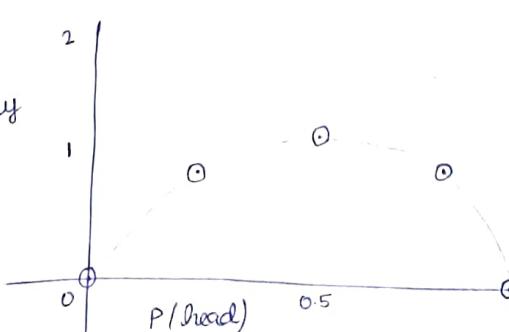
$$-0.25 \log 0.25 - 0.75 \log 0.75 = 0.8113$$

$$-0.50 \log 0.50 - 0.50 \log 0.50 = 1$$

$$-0.75 \log 0.75 - 0.25 \log 0.25 = 0.8113$$

$$-1 \log 0 - 0 \log 1 = 0$$

entropy



concave function

Find using KNN		4	330
H/W	CASES OF PRODUCT	7	3
DISTANCE (IN FT)	560	220	340
TIME (IN MINS)	16.68	11.50	12.03
(7, 330)	(1.32, 0.32)	(x - μ)/σ	18.11 ↑ test data
training data	distance	normalized	
(7, 560)	1.22	μ = 270 σ = 188.41	11.50
(3, 220)	0.442 ✓	(1.32, 1.54)	12.03
(3, 340)	0.1961 ✓	(+0.88, +0.27)	
(4, 80)	1.207	(+0.33, +1.01)	
(6, 150)	0.636 ✓	(0.77, +0.64)	13.75 12.43
check*		mean	

MODE : SD , M+, SHIFT 2  
 $\bar{x}$  - mean    $\sigma_x$  - std

dataset =  $\begin{cases} \text{pure} & H(y) = 0 \\ \text{impure} & \text{otherwise} \end{cases}$

OBJECTIVE: trying to minimise entropy (close to 0)

$\therefore$  entropy of pure dataset = 0

$H(y|x) \rightarrow$  conditional entropy

$H(y|x = \text{value}) \rightarrow$  specific conditional entropy (SCE)

$$H(y|x=v) = - \sum_{i=1}^m P(y_i|x=v) \log P(y_i|x=v)$$

SCE : average number of bits req. to represent y when x takes some specific value.

\* Entropy takes maximum value when it is equally distributed  $1/m$

\* Entropy is minimum when it is pure dataset

e.g. calculate of SCE of play tennis given outlook = sunny, rainy and overcast.

OUTLOOK:

$OL = \text{sunny}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 5 \\ \diagdown \\ \text{no} \end{array}$	<u>scentropy</u> $-\left(\frac{2}{5} \log \frac{2}{5} + \frac{3}{5} \log \frac{3}{5}\right) = 0.97$
--	--

$OL = \text{rainy}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 3 \\ \diagdown \\ \text{no} \end{array}$	$0.97$
--	--------

$OL = \text{overcast}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 4 \\ \diagdown \\ \text{no} \end{array}$	$0$ (pure dataset)
---	--------------------

CONDITIONAL ENTROPY :

$$0 \leq H(Y|X) \leq \text{entropy of } Y \quad H(Y)$$

expectation of SCE

$$H(Y|X) = \sum_{v \in X} P(v) H(Y|X=v)$$

$$H(PT|OL) = P(S) H(PT|S) + P(R) H(PT|R) + P(OC) H(PT|OC)$$

$$= \frac{5}{14} \times 0.97 + \frac{5}{14} \times 0.97 + \frac{4}{14} \times 0$$

$$H(PT|OL) = 0.693$$

$\because$  upper bound is always entropy of  $Y$

SCE

TEMPERATURE :

$T = \text{hot}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 2 \\ \diagdown \\ \text{no} \end{array}$	$-\left(\frac{2}{4} \log \frac{2}{4} + \frac{2}{4} \log \frac{2}{4}\right) = 1$
---	---

$T = \text{mild}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 4 \\ \diagdown \\ \text{no} \end{array}$	$-\left(\frac{4}{6} \log \frac{4}{6} + \frac{2}{6} \log \frac{2}{6}\right) = 0.9183$
--	--

$T = \text{cold}$ $\begin{array}{c} \text{yes} \\ \diagup \\ 3 \\ \diagdown \\ \text{no} \end{array}$	$-\left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) = 0.81125$
--	---

$E: H(PT|T) = P(\text{hot}) H(PT|\text{hot}) + P(\text{mild}) H(PT|\text{mild}) + P(\text{cold}) H(PT|\text{cold})$

$$= \frac{4}{14} \times 1 + \frac{5}{14} \times 0.9183 + \frac{5}{14} \times 0.81125 = 0.91106$$

### Humidity :

$$H = \text{High} \quad \begin{array}{c} \text{yes} \\ \diagup \\ 7 \end{array} \quad \begin{array}{c} 3 \\ \diagdown \\ \text{no} \end{array} \quad - \left( \frac{3}{7} \log \frac{3}{7} + \frac{4}{7} \log \frac{4}{7} \right) = 0.98526$$

$$H = \text{normal} \quad \begin{array}{c} \text{yes} \\ \diagup \\ 7 \end{array} \quad \begin{array}{c} 6 \\ \diagdown \\ \text{no} \end{array} \quad - \left( \frac{6}{7} \log \frac{6}{7} + \frac{1}{7} \log \frac{1}{7} \right) = 0.591674$$

$$\text{CE: } \frac{7}{14} (0.9852) + \frac{7}{14} (0.591674) = 0.788451$$

### Wind:

$$W = \text{weak} \quad \begin{array}{c} \text{yes} \\ \diagup \\ 8 \end{array} \quad \begin{array}{c} 6 \\ \diagdown \\ \text{no} \end{array} \quad - \left( \frac{6}{8} \log \frac{6}{8} + \frac{2}{8} \log \frac{2}{8} \right) = 0.81127$$

$$W = \text{strong} \quad \begin{array}{c} \text{yes} \\ \diagup \\ 6 \end{array} \quad \begin{array}{c} 3 \\ \diagdown \\ \text{no} \end{array} \quad - \left( \frac{3}{6} \log \frac{3}{6} + \frac{3}{6} \log \frac{3}{6} \right) = 1$$

$$\text{CE: } \frac{8}{14} (0.81127) + \frac{6}{14} (1) = 0.8921$$

### NOTE:

- \* When  $H(Y|X)=0$ , purity of dataset increases which means  $X$  highly influences  $Y$ . highly correlated.
- \* When  $H(Y|X)=H(Y)$ ,  $X$  and  $Y$  are independant,  $X$  does not influence  $Y$ .

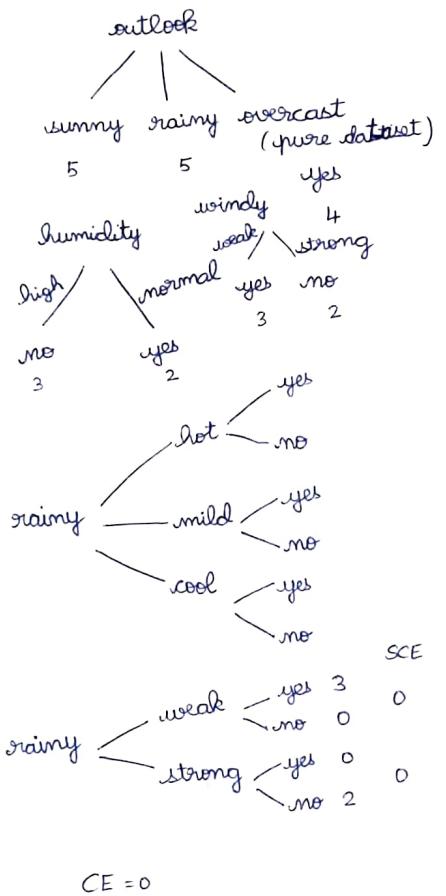
INFORMATION GAIN: (ID3 algorithm)

average # bits saved to represent  $Y$  using  $X$

$$IG(X, Y) = H(Y) - H(Y|X)$$

$\text{best feature} = \begin{cases} \underset{\text{or}}{\text{argmax}} & IG(X, Y) \\ \underset{\text{argmin}}{\text{argmin}} & H(Y X) \end{cases}$	$x^*$
--	-------

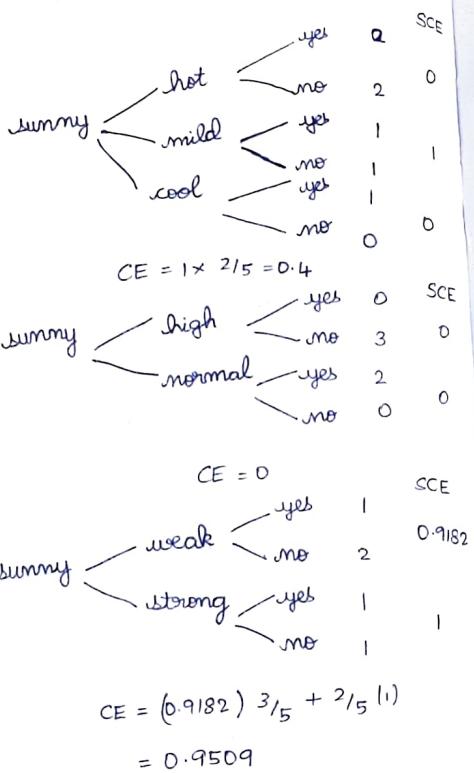
which feature reduces entropy the most will be root node



$$H(PT|T) = 0.4$$

$$H(PT|H) = 0$$

$$H(PT|W) = 0.95$$



CHARACTERISTICS:

- \* non-parametric, discriminative
- \* Decision trees can be converted to sequence of if-then rules. Number of if-then is equal to number of leaf nodes

- \* There are 2 views : data view
    - interior nodes : impure
    - leaf nodes : pure
  - \* feature view
    - features
    - class labels
  - \* training error of decision trees is zero. overfitting  
low bias, high variance (SOLUTION: tree pruning)  
setting leaf (edge)
  - \* feature selection of pre-processing is not needed  
(decision tree automatically selects feature)

ML LAB: Construct decision tree for

Find height and all performance  
of tree

choose 5 random subset of features

Compare with original decision tree.

NOTE:

- \* Inclusion of noise & outlier will  
due to overfitting
- \* To overcome overfitting we can  
since it can be applied only to a
- \* Assume that we include SNo as a  
tenant dataset. Find the information  
what would be the root node of  
in the dataset?

Given SNO, all have SCE as 0.  $\therefore$  SNO will be the root node ( $\therefore$ )  
 Eventually  $IG = H(Y) - H(Y|X)$   
 $IG = H(Y) = 0.94$

- \* IG gives high value which doesn't influence PT, IG
  - \* ID3 can be applied only for SNO.

↓  
limitations.

GAIN RATIO: (C4.5)

$$GR(x, y) = \frac{IG(x, y)}{H(x)} \quad (\text{norm})$$

- \* can be applied for continuous
  - \* non influential due to  $n$

ML LAB: Construct decision tree for Tic Tac Toe dataset. Find height of tree and all performance measures. Then choose 5 random subset of features and construct DT. Compare with original decision tree. Find training err.

NOTE :

- \* Inclusion of noise & outlier will not give accurate DT due to overfitting
- \* To overcome overfitting we cannot use regularization since it can be applied only to parametric model
- \* Assume that we include SNo as a feature in play tennis dataset. Find the information gain of SNo and what would be the root node if SNo is included in the dataset?

In SNo, all have SCE as 0.  $\therefore$  CE of PT given  $SNo = 0$

SNo will be the root node ( $\because \min CE$ )

Eventually  $IG = H(Y) - H(Y|X)$  is maximum

$$IG = H(Y) = 0.94$$

- \* IG gives high value which has more splits. Even though SNo doesn't influence PT, IG choose SNo as best feature.
- \* ID3 can be applied only for discrete features

↓  
limitations.

GAIN RATIO: (C4.5)

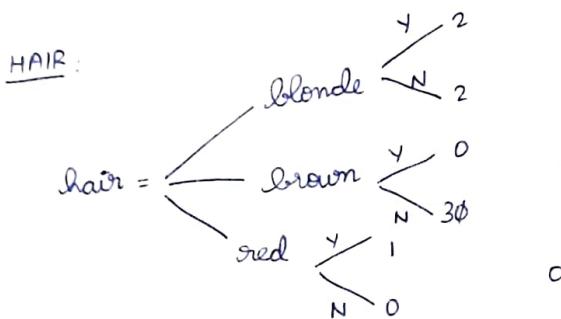
$$GR(X, Y) = \frac{IG(X, Y)}{H(X)} \quad (\text{normalised information gain})$$

- \* can be applied for continuous features.
- \* non influential due to more splits.

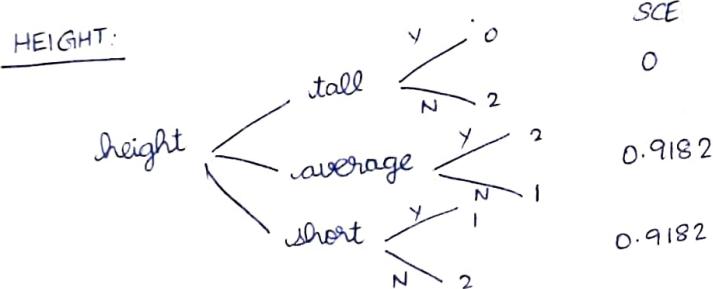
\* Construct DT using C4.5 for sunburn dataset.

$$H(Y) = 0.9544$$

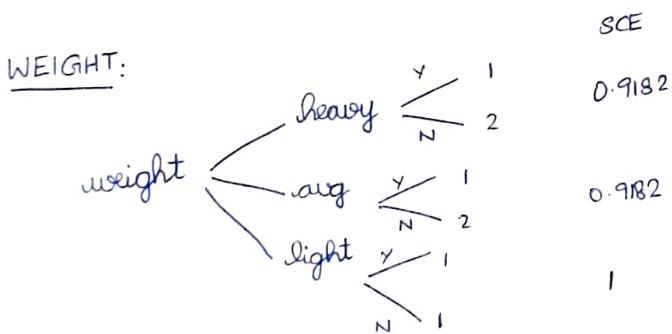
SCE



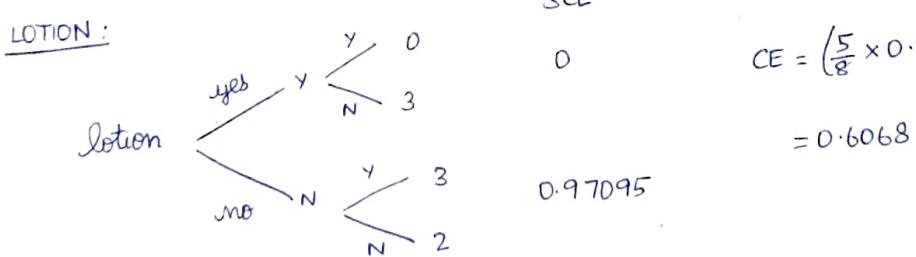
$$CE = \frac{4}{8} \times 1 = \frac{1}{2}$$



$$CE = \left(\frac{3}{8} \times 0.9182\right) + \left(\frac{3}{8} \times 0.9182\right) = 0.68865$$



$$CE = \left(\frac{3}{8} \times 0.9182\right) + \left(\frac{3}{8} \times 0.9182\right) + \left(\frac{2}{8} \times 1\right) = 0.93865$$



$$CE = \left(\frac{5}{8} \times 0.97095\right) = 0.6068$$

grain ratio	Hair	height	weight	lotion
$\frac{IG(x, y)}{H(x)}$	$0.4544$	$0.26575$	$0.01575$	$0.3476$
	$1.4056$	$1.56127$	$1.56127$	$0.9544$
	$0.3232$	$0.17021$	$0.010087$	$0.36420$

✓

root node

## GINI INDEX: (CART algo)

$$GI = 1 - \sum_{i=1}^m p_i^2$$

for pure dataset  $GI = 0$

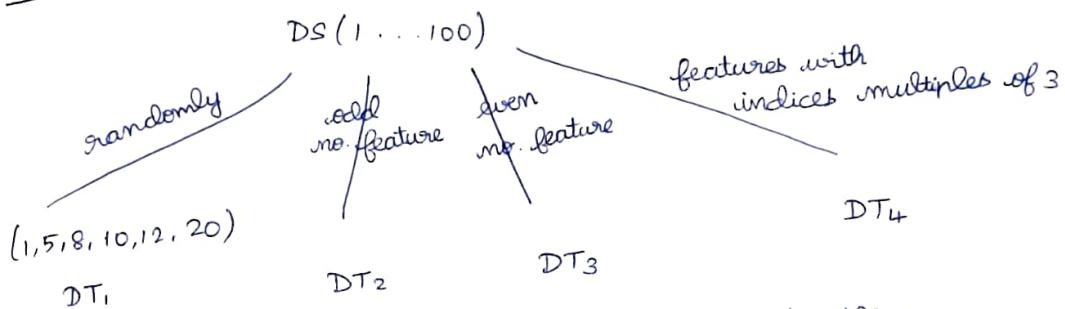
$$x^* = \underset{x}{\operatorname{argmin}} \quad GI(x, y)$$

easy to compute -  $O(n \log n)$

## PRUNING:

cut branches of leaf nodes by applying chi-square test. If it has high confidence values don't cut. (goodness of fit)

## RANDOM FOREST:



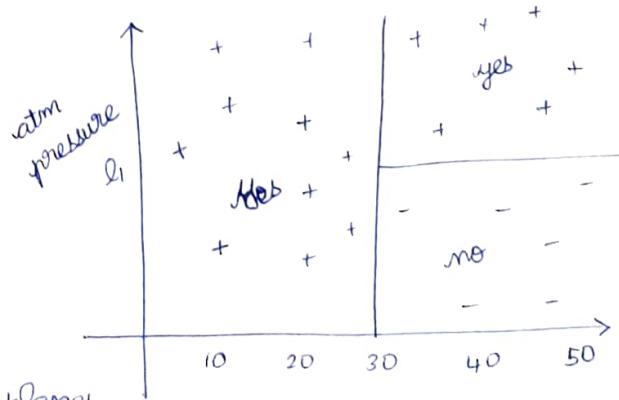
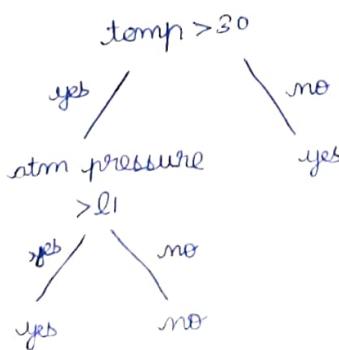
predict using all trees to avoid noise / outliers

- collection of decision trees
- randomly choose subset of features / no. of datapoints

$\underbrace{K \text{ times}}_{D_1 = \text{random}()} \quad / / F_1, F_2, \dots, F_d$   
 construct ID3 ( $D_i$ )

$K = \# \text{ decision trees}$  (hyperparameter)

## GEOMETRIC REPRESENTATION OF DECISION TREES:



combination of parallel axis planes

temp

SUPPORT VECTOR MACHINE: (SVM) (parametric discriminative linear)

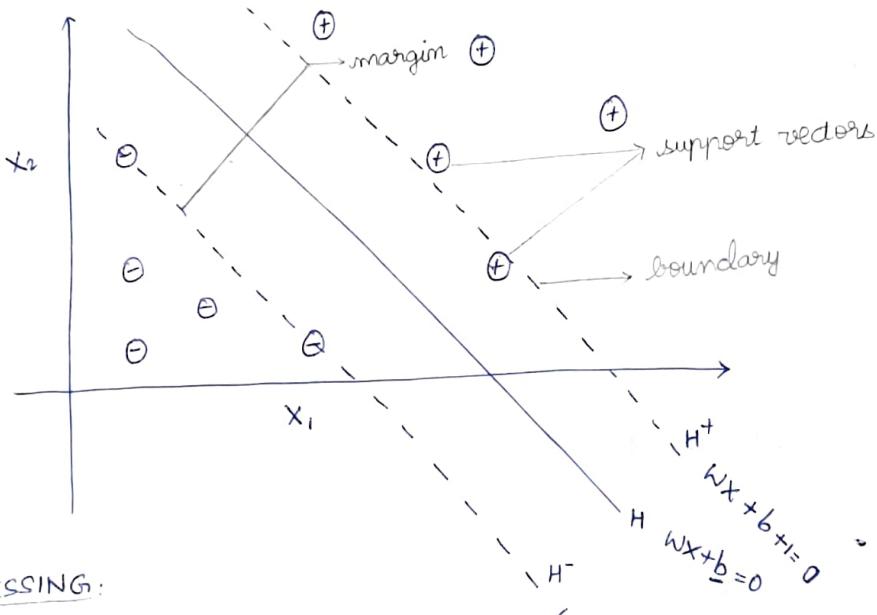
Given  $\mathcal{D} = \{x_i, y_i\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \{+1, -1\}$

$$f: x_i \rightarrow y_i$$

Objective

$\hat{f}: \mathbb{R}^d \rightarrow \{+1, -1\}$  + find an optimal hyperplane equi-distance from +ve & -ve margin

VC dimension: (Vapnik, Chervonenkis) + max distance b/w HP & bounds, to avoid bias  
max no. of datapoints req. to compute measures



PRE PROCESSING:

- \* Z score normalisation,
- \* min-max normalisation,
- \* divide by max val of features (0,1)

slope same

y-intercept diff

$$\hat{y} = \begin{cases} +1 & w x + b \geq 0 \\ -1 & w x + b < 0 \end{cases}$$

$$\hat{y} = \text{sign}(w x + b)$$

MODEL:

<u>dimension</u>	<u>model</u>
1	point
2	line
3	plane
$> 3$	hyperplane

max margin hyperplane

Optimal hyperplane is  $H$  with max margin  
 $\text{max} = \text{distance between boundaries}$   
 Find  $W$  &  $b$  which maximizes margin

$M = \text{distance between } H_+ \text{ & } H_-$

$H_+$  and  $H_-$  are parallel

distance between parallel lines

$$wx + (b-1) = 0$$

$$wx + (b+1) = 0$$

$$M = \frac{|(b-1) - (b+1)|}{\|w\|} = \frac{2}{\|w\|}$$

for simplicity unit distance

$$\text{Max } \frac{2}{\|w\|} \quad \text{s.t.c. } \begin{aligned} wx_i + b &\geq 1 && \text{for } y_i \in +1 \\ wx_i + b &\leq -1 && \text{for } y_i \in -1 \end{aligned} \quad \begin{array}{l} n \text{ constraints} \\ \text{for each datapoint} \end{array}$$

s.t.c

$$y_i(wx_i + b) \geq 1$$

$$\downarrow \quad \text{max margin} = \min \|w\|$$

$$\min \frac{1}{2} \|w\|^2$$

s.t.c

$$y_i(wx_i + b) \geq 1$$

$$\min \frac{1}{2} w^T w$$

s.t.c

$$y_i(wx_i + b) \geq 1$$

constrained convex optimisation problem

(Solve by lagrange multipliers & KKT)  
 Karush, Kuhn, Tucker

$$\downarrow \quad \text{lagrange multiplier:}$$

$$f: \min \frac{1}{2} \underbrace{w^T w}_{w^2} - \sum_{i=1}^n \alpha_i (y_i(wx_i + b) - 1) \quad \text{① } \alpha_i \geq 0$$

PRIMAL

$$\text{② } \frac{\partial f}{\partial w} = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\text{③ } \frac{\partial f}{\partial b} = - \sum \alpha_i y_i = 0$$

$$\sum \alpha_i y_i = 0$$

$$\text{④ } \alpha_i (y_i(wx_i + b) - 1) \geq 0$$

KKT conditions : ①, ②, ③, ④

$$\text{Min}_{w,b} \frac{1}{2} w^T w - \sum_{i=1}^n \alpha_i y_i w^T x_i - \sum_{i=1}^n \alpha_i y_i b + \sum_{i=1}^n \alpha_i$$

$\stackrel{\text{by } \textcircled{2}}{w^T w} - \stackrel{\text{by } \textcircled{3}}{w^T x_i} - 0 + \sum \alpha_i$

$$-\frac{1}{2} w^T w + \sum \alpha_i$$

**DUAL**

$$\text{Min}_{w,b} = \text{Max}_{\alpha} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i y_i x_i^T \alpha_j y_j x_j + \sum \alpha_i \quad (\text{function of } \alpha)$$

$\alpha_i, y_i, y_j$  - scalars  
 $x_i, x_j$  - vectors

s.t.c

$$\alpha_i \geq 0$$

$$\sum \alpha_i y_i = 0$$

$$-\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Note:  $\alpha_i \neq 0$  then data point is said to be support vector  
 $\alpha_i = 0$  then data point is not support vector

$$SV = \{(x_i, y_i) | \alpha_i \neq 0\}$$

$$\textcircled{2} \Rightarrow w = \sum_{\alpha_i \neq 0} \alpha_i y_i x_i$$

$$w = \sum_{\substack{x_i, y_i \in SV}} \alpha_i y_i x_i$$

slope of HP is found from SV

for datapoint  $(x_k, y_k)$

$$y_k = w^T x_k + b$$

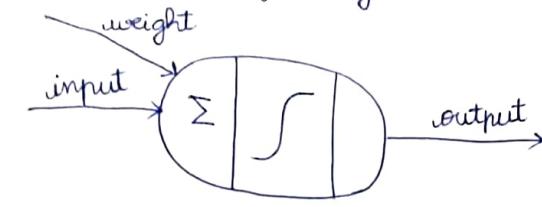
$$b = y_k - w^T x_k$$

"linear SVM is called hard margin"

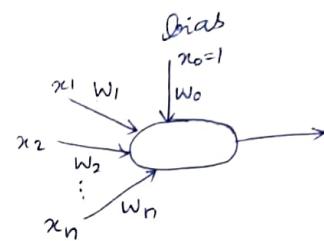
- hard margin

## PERCEPTRON

- artificial neuron
- simulation of single neuron



weighted sum =  $\sum_{j=0}^d x_j w_j$



activation function =  $\begin{cases} 1 & \text{if } \text{es-WS} \geq 0 \\ 0 & \text{otherwise} \end{cases}$   
 (brings non-linearity)

MODEL:  $w^T x$

PARAMETER:  $w$  (for linearly separable)

→  $\begin{cases} \text{Sigmoid function - binary classification} \\ \text{Softmax function - multi-way classification} \end{cases}$

## PERCEPTRON AS CONVEX OPTIMISATION PROBLEM:

parameters  $w = (w_0, w_1, \dots, w_n)$

find  $w$

$$\min \quad SSE = \frac{1}{2} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

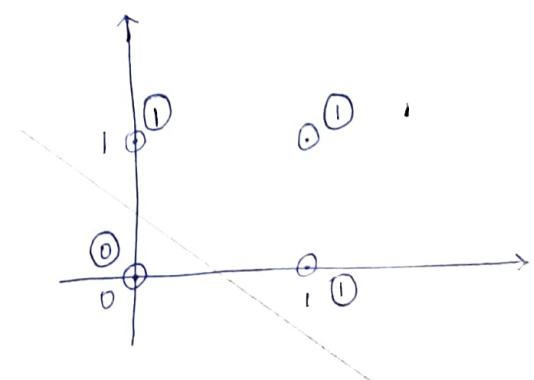
$$\text{gradient } \frac{\partial SSE}{\partial w} = \sum (y_i - w^T x_i) (-x_i)$$

forward pass : predicting  $\hat{y}$       } all deep learning  
 backward pass : updating  $w$       } algo follows

eg: Find  $w$  for initial values  $(0.5, 0.5, 0.5)$ .  $n=1$   
 assume activation function is threshold function.  
 threshold value = 0.5

### OR GATE

0	0	0
0	1	1
1	0	1
1	1	1

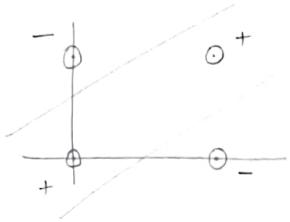


EPOCH-1			$X = [0, 0, 1, 1]$			$Y = [0, 1, 1, 1]$			$B = [0.5 \ 0.5 \ 0.5]$		
$x_0$	$x_1$	$x_2$	$b_0$	$b_1$	$b_2$	$y$	$\hat{y}$	$E_n = E_0 + \eta \text{ grad}$			
1	0	0	0.5	0.5	0.5	0	0	$0.5 >= 0.5 [1, 0, 0]$			
1	0	1	-0.5	0.5	0.5	1	0	$[1, 0, 1]$			
1	1	0	0.5	0.5	1.5	1	1	$[0, 0, 0]$			
1	1	1	0.5	0.5	1.5	1	1	$(0.5 > 0.5)$			

EPOCH-4 converged

XOR is not linearly separable.  $\therefore$  perceptron not applied.

0	0	0
0	1	1
1	0	1
1	1	0



(MLP in PPT, refer gmail)

### SOFTMARGIN SVM :

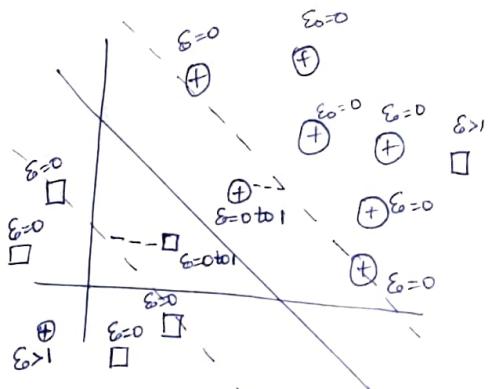
$$\min \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n I(\varepsilon_i > 0)$$

s.t.c  
 $y_i(w^T x_i + b) + \underline{\varepsilon_i} \geq 1$  slack variable

where

C is the penalty

$$I(\varepsilon > 0) = \begin{cases} 1 & \text{if } \varepsilon > 0 \\ 0 & \text{if } \varepsilon \leq 0 \end{cases}$$



\* If C is high, number of errors will be reduced

\* If  $\varepsilon_i = 0$ , datapoints are correctly classified. Either lie on the margin or above.

For a  $\square$  datapoint

$$w x_i + b \leq -1$$



$$w x_i + b + \varepsilon_i \leq -1$$

For a  $\oplus$  datapoint

$$w x_i + b \geq 1$$



$$w x_i + b + \varepsilon_i \geq 1$$

$$\epsilon_i = 0$$

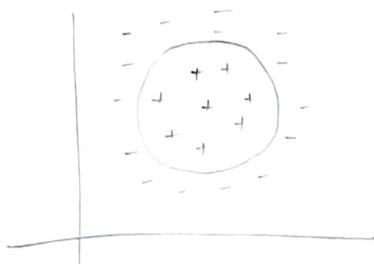
correct side

$$0 \leq \epsilon_i \leq 1$$

within the margin

$$\epsilon_i > 1$$

wrong side

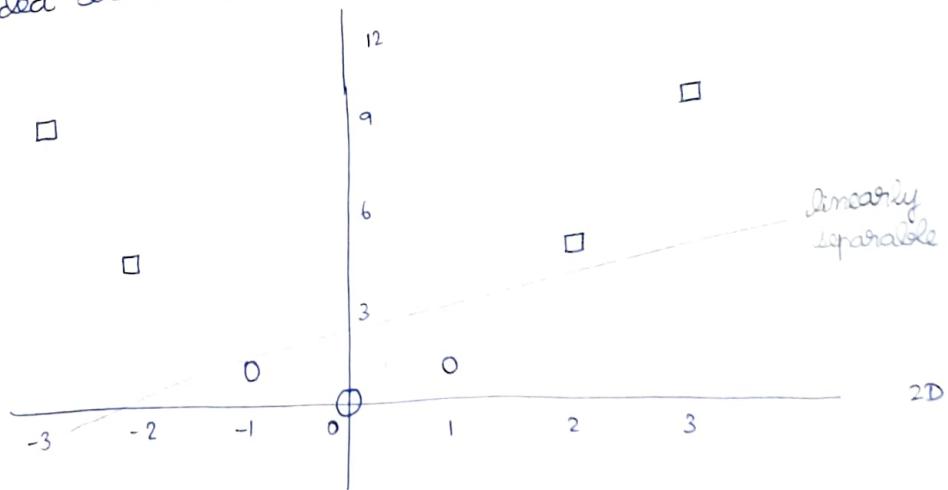


linearly non-separable problem

### SVM WITH KERNEL FUNCTION

A linearly non-separable problem when projected in higher dimension becomes linearly separable.

"Idea behind kernel SVM."



while multiplying  $x_i x_j$  vectors, increase dimension by applying kernel function  $K(x_i) K(x_j)$

### KERNEL FUNCTIONS:

- linear  $(x^T y)^d$  ( $d$  changes according to dim)
- polynomial
- gaussian
- radial basis function (RBF)
- exponential

e.g. for  $\text{dim} = 2$

$$X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad Y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

2 datapoints

$$K(x, y) = (x^T y)^2$$

Kernel fn must be +ve semi-definite

$$x^T y = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = x_1 y_1 + x_2 y_2$$

$$(x^T y)^2 = x_1^2 y_1^2 + x_2^2 y_2^2 + 2 x_1 x_2 y_1 y_2$$

$$\begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2} x_1 x_2 \end{bmatrix} \quad \begin{bmatrix} y_1^2 \\ y_2^2 \\ \sqrt{2} y_1 y_2 \end{bmatrix}$$

## UNSUPERVISED LEARNING

$$\mathcal{D} = \left\{ x_i \right\}_{i=1}^n$$

- clustering
- dimensionality reduction
- outlier analysis

### CLUSTERING:

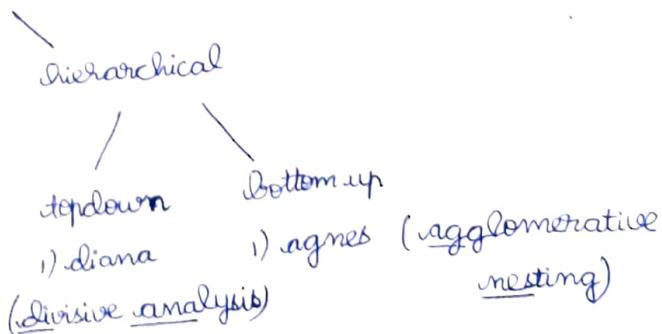
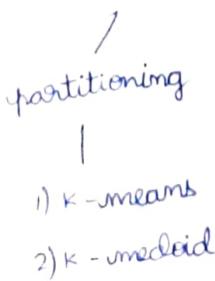
$$\mathcal{D} = \left\{ x_i \right\}_{i=1}^n$$

$$\mathcal{D} = \{ C_1, C_2, C_3, \dots, C_k \}$$

$$C_k = \{ x_i \mid \text{distance}(x_i, \mu_k) = \text{minimum} \}$$

where  $\mu_k$  is called cluster representative / centroid

### CLUSTERING ALGORITHMS

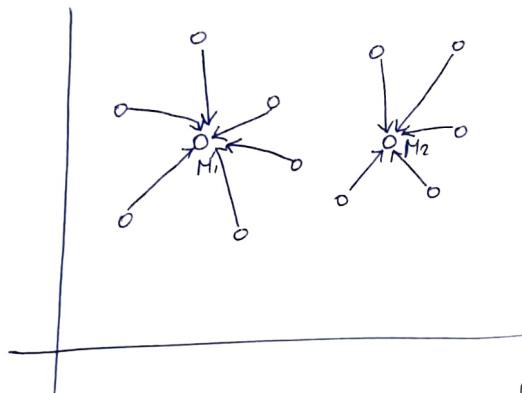


## K-MEANS CLUSTERING:

input  $D = \{x_i\}_{i=1}^n$  &  $k$

If you have domain knowledge, choose centroid.  
else choose any  $k$  centroids randomly.

→ sometimes min & max are chosen



Suppose  $k=2$

{ find distance metric  
{ update centroid  
↓  
for each iteration

- Only in first iteration, a datapoint is centroid.
- Thereafter centroid is mean of all datapoints  $\in$  clusters
- Successive 2 iterations if the centroid is same, optimality obtained.

eg: Apply k-means clustering and formulate 2 clusters.  
use Euclidean distance as the distance metric. Choose centroids randomly.

distance	speed				
71.24	28				
52.53	25				
64.54	27				
55.69	22				
54.58	25				

Euclidean distance

	$M_1$	$M_2$		
(64.54, 27)	12.17	6.7	$C_2$	
(55.69, 22)	4.35	16.667		$C_1$
(54.58, 25)	2.05	16.93		$C_1$

To find optimum value of  $K$ , we use distortion as performance measure.

### DISTORTION:

sum of squares of distances between datapoints & its corresponding centroid

$$\sum_{i=1}^n (x_i - \mu_i)^2$$

intuition → how well it is clustered for minimal  $K$

Minimum the value of distortion, optimum the value of  $K$

### K-MEANS CLUSTERING:

unsupervised learning algorithm

Form  $c_1, c_2, \dots, c_k$

which min distortion / intracenter distance  
(sum of all cluster distance)

$$D = \sum_{x_i \in C_j} \|x_i - \mu_j\|^2$$

ALGORITHM: (time complexity:  $O(nkd)$ )

1. initialize the centroids randomly

2. while ! convergence

a. for each  $x_i \in D$  //cluster formation

$$\mu_i = \underset{\mu_j}{\operatorname{argmin}} \|x_i - \mu_j\|^2 \quad 1 \leq j \leq k$$

b. for each cluster  $j$  //centroid updation

$$\mu_j = \frac{\sum x_i}{n_j}$$

KMC is mixture of gaussian since each cluster follows normal dist. In ND,  $\mu = \sigma$ , it is sphere

### LIMITATIONS:

\* works only for spherical clusters

\* convergence depends on selection of initial centroids

### ASSUMPTION:

\* closeness

## SPECTRAL CLUSTERING:

- based on spectral properties of graph
- given dataset  $D$  is represented as a graph  $G$ 
  - \* nodes - datapoints
  - \* edges - connection between the nodes (affinity)
- $G$  may be directed / undirected
  - \* adjacency matrix  
(complete graph)
  - \* similarity matrix  
(use distance formula)
  - \* KNN neighbour  
(connect only  $K$  nearest neighbours)

eg: k-means clustering dataset used

distance matrix

$$A = \begin{bmatrix} 0 & 18.95 & 6.77 & 16.67 & 16.93 \\ 18.95 & 0 & 12.98 & 4.36 & 2.05 \\ 6.77 & 12.18 & 0 & 10.16 & 11.14 \\ 16.67 & 4.36 & 10.16 & 0 & 3.199 \\ 16.93 & 2.05 & 11.14 & 3.199 & 0 \end{bmatrix}$$

$k=2$   
nearest neighbours

$$\begin{bmatrix} 0 & 0 & 6.77 & 16.67 & 0 \\ 0 & 0 & 0 & 4.36 & 2.05 \\ 6.77 & 0 & 0 & 10.16 & 0 \\ 0 & 4.36 & 0 & 0 & 3.199 \\ 0 & 2.05 & 0 & 3.199 & 0 \end{bmatrix} \quad \text{directed}$$

degree matrix

$$D = \begin{bmatrix} 4 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 4 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{bmatrix}$$

Laplacian matrix

$$L = D - A = \begin{bmatrix} 4 & -18.95 & -6.77 & -16.67 & -16.93 \\ -18.95 & 4 & -12.18 & -4.36 & -2.05 \\ -6.77 & -12.18 & 4 & -10.16 & -11.14 \\ -16.67 & -4.36 & -10.16 & 4 & -3.199 \\ -16.93 & -2.05 & -11.14 & -3.199 & 4 \end{bmatrix}$$

## PROPERTIES OF SPECTRUM:

- \* If there exists  $d \times d$  symmetric matrix, then there will be  $d$  eigen values. Every eigen value has its corresponding eigen vector. All eigen values are real and All eigen vectors are real and orthogonal.
- \* Sum of eigen values of covar matrix will give total variance.

$$\lambda_1 + \lambda_2 + \lambda_3 + \dots + \lambda_d = \sum_{j=1}^d \sigma_j^2$$

- \* Direction of maximum eigen value's eigen vector follows direction of variance.

## ALGORITHM :

input: dataset

output: clusters

1. formulate graph A
2. find degree matrix D as a diagonal matrix in which diagonals are degree
3. find laplacian  $L: D - A$   
in which diagonals  $\rightarrow$  degree  
off diagonals  $\rightarrow$  negative edge weights
4. find eigen values of  $L$   
 $\#$  zero eigen values =  $\#$  connected components / clusters  
first non-zero eigen value = fiedler value
5. find the eigen vector  $V$  for the fiedler value
6. for each value  $x_m$  of eigen vector  $V$

- a. if  $x_m > 0$

$$x_m \rightarrow C_1$$

- b. if  $x_m < 0$

$$x_m \rightarrow C_2$$

fiedler value: 2<sup>nd</sup> smallest eigen value of graph

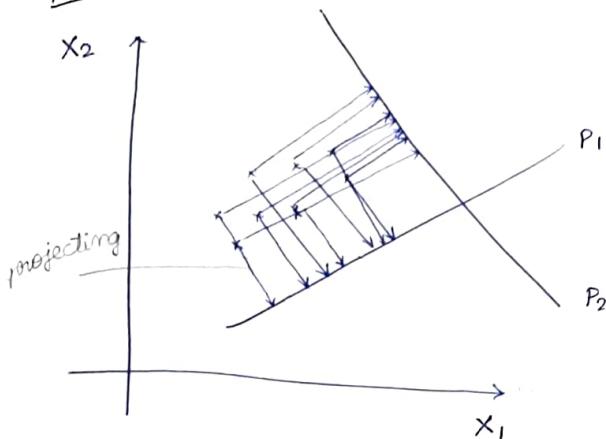
graph cut - fiedler value

no. of edges to be removed to make connected components

## DIMENSIONALITY REDUCTION:

- principal component analysis (PCA)  $\rightarrow$  unsupervised
- linear discriminant analysis (LDA)  $\rightarrow$  supervised

### PCA:



If there exists 2 eigen vectors  $P_1$  and  $P_2$ , choose vector which has more covariance.  
 $\therefore$  data has to be visible.  
(scattered more)

$\boxed{P_1}$

### ALGORITHM:

input :  $X = \{x_1, x_2, \dots, x_n\}$

$$x_i \in \mathbb{R}^d$$

$$z = w^T x$$

$$z \in \mathbb{R}^k \text{ where } k < d$$

1.  $\mu = \{\mu_1, \mu_2, \dots, \mu_d\}$
2. mean subtracted data  $x' = x - \mu$  ( $\mu$  will be 0)
3. calculate covariance matrix / for mean subtracted data  
 $\Sigma_{ij} = \frac{\sum_{i=1}^n (x_i - \mu_i)(x_j - \mu_j)}{n-1} = \frac{\sum_{k=1}^n x_{ki}x_{kj}}{n-1}$
4. eigen values of  $\Sigma$ ,  $\lambda_1, \lambda_2, \dots, \lambda_d$
5. arrange the eigen values in descending order
6. choose the largest eigen value ( $\lambda_1, \lambda_2, \dots, \lambda_d$ )
7. find the corresponding eigen vector EV  
 $w_1, w_2, \dots, w_k$

$$8. z = (w^T x)^T \quad z^T = x^T w$$

$n \times k \quad n \times d \quad d \times k$

def k eigen vector are principal components  
(of covr matrix)

$$\begin{array}{c|cc} X_1 & X_2 \\ \hline P_1 & 2 & 6 \\ P_2 & 1 & 7 \end{array}$$

$$\begin{array}{c|cc} X_1 & X_2 \\ \hline P_1 & 0.5 & -0.5 \\ P_2 & -0.5 & 0.5 \end{array}$$

$$\Sigma_{\text{cov}} = \begin{bmatrix} X_1 & X_2 \\ X_2 & X_2 \end{bmatrix} = \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 0.5 \end{bmatrix}$$

$$|\Sigma - \lambda I| = 0$$

$$\begin{aligned} \Sigma - \lambda I &= \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 0.5 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \\ &= \begin{bmatrix} 0.5 - \lambda & -0.5 \\ -0.5 & 0.5 - \lambda \end{bmatrix} \end{aligned}$$

$$|\Sigma - \lambda I| = 0$$

$$(0.5 - \lambda)(0.5 - \lambda) - (-0.5)(-0.5) = 0$$

$$0.25 + \lambda^2 - \lambda + 0.25 = 0$$

$$\lambda(\lambda - 1) = 0$$

$$\boxed{\lambda = 0} \quad \boxed{\lambda = 1}$$

eigen values

funding eigen vector for  $\lambda = 1$  (max)

$$\Sigma - \lambda I = \begin{bmatrix} 0.5 - 1 & -0.5 \\ -0.5 & 0.5 - 1 \end{bmatrix} = \begin{bmatrix} -0.5 & -0.5 \\ -0.5 & -0.5 \end{bmatrix}$$

to form characteristic equation

$$\begin{bmatrix} -0.5 & -0.5 \\ -0.5 & -0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$AX = 0$$

eigen vectors determine?

linear

$$ax + by = c$$

$$a_1x_1 + a_2x_2 + b_1y_1 + b_2y_2 = c_1$$

$$\frac{a_1}{a_2} = \frac{b_1}{b_2}$$

$$\begin{cases} -0.5x_1 - 0.5x_2 = 0 \\ -0.5x_1 - 0.5x_2 = 0 \end{cases}$$

and the 2nd eqn is linear

$$\Rightarrow -0.5x_1 = 0.5x_2$$

$$-x_1 = x_2$$

$$\boxed{x_1 = 1}$$

$$\boxed{x_2 = -1}$$

to convert into unit vector divide by norm

$$\text{norm} = \sqrt{2} \quad (\sqrt{1^2 + (-1)^2})$$

$$w = \left( \frac{1}{\sqrt{2}}, \frac{-1}{\sqrt{2}} \right)$$

unit vectors act like axis

$$z = w^T x$$

$$= \begin{bmatrix} \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 0.5 & -0.5 \\ -0.5 & 0.5 \end{bmatrix}$$

$$z = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

+ *unit vector divide by norm*

20PH13

**PSG COLLEGE OF TECHNOLOGY, COIMBATORE - 641 004**  
**Department of Applied Mathematics and Computational Sciences**  
**M.Sc. (SS) – V Semester**

**TEST - I**  
**20XW53 - MACHINE LEARNING**

**Time: 1 Hour 15 min.**

**Maximum Marks: 40**

**INSTRUCTIONS:**

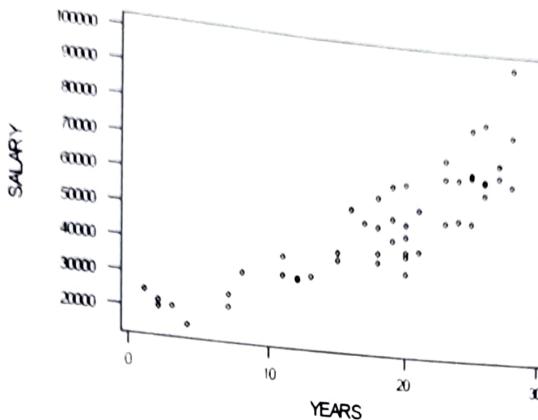
1. Answer **ALL** questions. Each question carries 20 Marks.
2. Subdivision (a) has 2 questions, each carries 2 marks, subdivision (b) carries 6 marks each and subdivision (c) carries 10 marks each.

1. a) i. Check the convexity of cross entropy loss  $L = y_i \log \hat{y}_i$  where  $\hat{y}_i = W^T X$ . (2)  
ii. Differentiate generative and discriminative models. (2)  
  
b) Define parameters and hyperparameters of a model. State the parameters and hyperparameters (if any) for the following ML algorithms? i. Polynomial regression ii. Multiple regression iii. Ridge regularization. Elaborate a method on finding the optimal value of hyperparameters? *3 way split*  
  
c) What are the characteristics of gradient descent method? Design an algorithm for finding a solution of an optimization problem. Explain its variants with procedure. What are the possible criteria for convergence in Gradient descent (GD) method? Justify the convergence mathematically. Given a data set with 2000 instances and 10 features, how many times parameter will be updated in an epoch when the dataset is applied on GD and its variants?
2. a) i. The scatter plot shows the employees' years of experience vs their salary for a sample of 50 managers. The regression model obtained from data is  
$$\text{salary} = 11369 + 2141 * \text{years}$$
 and what is the sum of residuals (errors) of the model?

$$SE = \sum (y_i - \hat{y}_i) = 0 \text{ for linear regression} \quad (2)$$

no convert into unit vector divide by norm

$$\rightarrow \begin{pmatrix} 1 \\ 1, 2, \dots, n \end{pmatrix}$$



- ii) Write an unconstrained optimization problem for Ridge regularization and a solution to the problem.
- $$\sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \sum_{j=1}^d (w_j^2)$$
- $$w = (x^T x + \lambda I)^{-1} x^T y$$

- b) What are the inferences about the model (M) and data (D) when M overfits D. Elaborate the impact of following factors in Multiple regression. How do you alleviate these concerns in an ML model?

- i. Multicollinearity
- ii. Outliers

Assume that a regression data set D contains 100 independent features and a dependant variable. A ML engineer who has no knowledge about the relevance and importance of these features influencing dependant variable, applied Multiple Regression on data set D. He encountered an overfitting issue. Which technique do you suggest him to resolve the overfitting? Why? Give reasons.

Lasso regression

- c) Given a data set with a continuous independent variable X, a dependant variable Y, construct a quadratic polynomial regression model and answer the following:

Infer the relationship between dependant and independent variables from the parameters. Obtain the irreducible error of the model and find the proportion of variance explained by X to determine Y. SSR

X	1	2	3	4	4	5	5	6	7	8	9	10
Y	7	8	9	8	9	11	10	13	14	13		

\*\*\*\*\*

**PSG COLLEGE OF TECHNOLOGY, COIMBATORE - 641 004**  
**Department of Applied Mathematics and Computational Sciences**

M.Sc. (SS) – V Semester

**TEST - II**

**20XW53 - MACHINE LEARNING**

**Time: 1 Hour 15 min.**

**Maximum Marks: 40**

**INSTRUCTIONS:**

1. Answer **ALL** questions. Each question carries 20 Marks.
2. Subdivision (a) has 2 questions, each carries 2 marks, subdivision (b) carries 6 marks each and subdivision (c) carries 10 marks each.

1. a) i. Given a dataset D with  $n$  observations and each observation consists of an input vector  $x_i \in R^d$  and an outcome  $y_i \in \{c_1, c_2, \dots, c_m\}$ , infer the information, purity, surprisal of Y when the Entropy  $H(Y)$  is  $\log_2 m$ ? What is bias in decision tree? (2)  
ii. Differentiate lazy learners with eager learners. (2)
- b) Prove that logistic regression is a linear model. Let us assume that a binomial logistic regression is applied on a dataset with an input entrance marks, and a class variable with values pass or fail. The model parameters obtained are  $w_0 = 1$ ,  $w_1 = 8$ . Compute the prediction for a student with entrance marks is 70.
- c) Apply Naïve Bayes classifier for the below data and construct a model. Assume that taxable income is a continuous feature and normalize it by dividing the value by 1000. Do smoothing in case of zero frequency problem. Also predict test data with values <No, married, 120K>

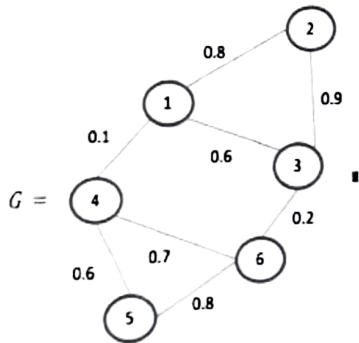
Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes ✓
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes ✓
9	No	Married	75K	No
10	No	Single	90K	Yes ✓

$$P(Y) =$$

to convert into unit vector divide by norm

$$\text{norm} = \sqrt{1 + (-1)^2}$$

2. a) i. Is it possible to apply k-means clustering for clusters of non-convex sets? Justify your answer.
- ii. Suppose we have three cluster centroids  $C_1 = (1, 2)$ ,  $C_2 = (-3, 0)$ ,  $C_3 = (4, 2)$ . Furthermore, we have a training example  $X = (-1, 2)$ . Based on k-means clustering with Euclidean distance,  $X$  will be assigned to which cluster?
- b) Consider a hyperplane function of Linear SVM for two variables is  $w_1x_1 + w_2x_2 + b$ . If two hyperplanes considered in a hypotheses space are  $2x_1 + 3x_2 + 4$  and  $10x_1 + 3x_2 + 1$ . Find the cost or objective function for each hyperplane and pick an optimum hyperplane. Write the optimization problem for soft margin SVM. What is the impact of hyper-parameters in it?
- c) How is spectral clustering advantageous to k-means clustering? List the steps in Spectral clustering. Below graph  $G$  gives the representation of a dataset with six data points.



Obtain the Laplacian matrix  $L$  from the above weighted undirected graph. List the characteristics of Laplacian matrices. Infer the possible Eigen values. What would be the possible Fiedler value for  $L$ ? Find the sum of all Eigen values. Assume that Fiedler vector is  $[ -0.4 \ -0.3 \ -0.5 \ 0.6 \ 0.7 \ 0.8 ]$ . How do you form clusters? Justify your answer wherever required.

\*\*\*\*\*