

Ex-1 Time series Data (forecasting the stock index)

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X
	Price of Petrol	Price of Gold	Conflict b/w neighbour countries	Trap (Y/No)	Previous Day's value	Mid capy	Fixed asset	\mathbb{R}
x_1	x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_1
x_2								
x_3								
x_i	x_{i1}	x_{i2}	x_{i3}	x_{i4}	x_{i5}	x_{i6}	x_{i7}	y_i
x_N								
x_{new}	$x_{new,1}$	$x_{new,2}$	$x_{new,3}$				$x_{new,7}$	y_{new}

$y \in \mathbb{R}$

Ex-2: (Time Series) : Predicting Rainfall in a Particular locality. (Regression)

(amount of Rainfall in 3rd week)

X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	Y_{new}
2001								y_1
2002								y_2
2003								y
⋮								⋮
2024								y_{20}

$y_{25} = ?$
 $y_{20} = ?$

Mathematical formulation of this problem

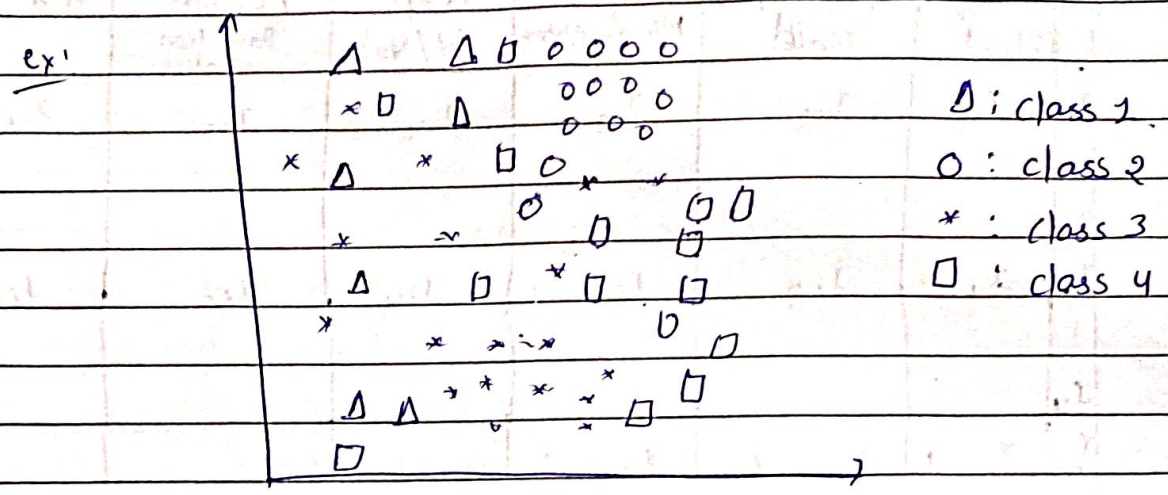
Problem statement:

Given $D = \{x_i, y_i\}_{i=1}^{N=24}$, $x_i \in \mathbb{R}^{p=8}$, $y_i \in \mathbb{R}$

Predict y_{25} for x_{25} where $x_{25} \neq x_i, \forall i = 1$ to 24

KNN Classification

↳ K Nearest Neighbours.



$$d_1 = \sqrt{\sum_{j=1}^N (x_1 - x_{1j})^2}$$

$$d_2 = \sqrt{\sum_{j=1}^N (x_2 - x_{2j})^2}$$

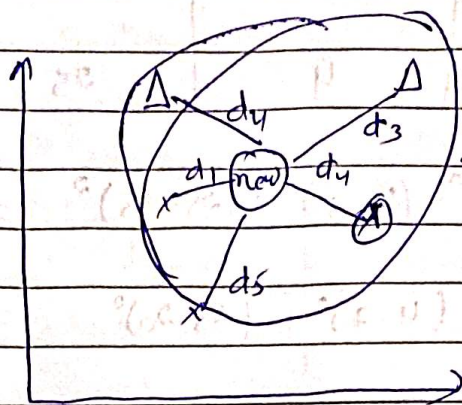
K odd for proper distribution
 1, 3, 5, 7, 9
 (3, 5, 7, 9)

$$d_N = \sqrt{\sum_{j=1}^N (x_N - x_{Nj})^2}$$

	x_1	x_2	x_3	...	x_{10}	y
x_1	x_{11}	x_{12}	x_{13}		$x_{1,10}$	y_1
x_2						y_2
x_3						y_3
⋮						⋮
x_{100}	$x_{100,1}$	$x_{100,2}$	$x_{100,3}$		$x_{100,10}$	y_{100}
new sample x_{new}	$x_{new,1}$	$x_{new,2}$	$x_{new,3}$...	$x_{new,10}$	$y_{new} = ?$

$$d_1 = \sqrt{(x_{new,1} - x_{11})^2 + (x_{new,2} - x_{12})^2 + (x_{new,3} - x_{13})^2 + \dots + (x_{new,10} - x_{1,10})^2}$$

$N = 100$
 $P = 10$



$K=5$

N_1

New Sample \in Class of x

a

The following table represents the training samples. a factory produces a new paper that pass laboratory with acid durability. Strength is 4 kg per sq m \rightarrow Cost is 25 Rs per kg. apply K-NN classification of with $K=5$ to predict the quality of the tissue paper from the given data set.

x_1 Acid durability (hr)	x_2 Strength (kg/m ²)	x_3 Cost ₹ Per kg	y opinion from Peopls.
$x_1 \rightarrow 8$	6	15	Good
$x_2 \rightarrow 9$	7	20	bad
$x_3 \rightarrow 2$	5	18	bad
$x_4 \rightarrow 1$	3	9	bad
$x_5 \rightarrow 6$	7	8	Good
$x_6 \rightarrow 4$	5	10	bad
$x_7 \rightarrow 5$	5	6	Good
$x_8 \rightarrow 3$	6	10	Good

X_{new}	5	4	25	?
-----------	---	---	----	---

$$d_1 = \sqrt{(5-8)^2 + (4-6)^2 + (25-15)^2} = \sqrt{9+4+100} = \sqrt{113}$$

$$d_2 = \sqrt{(5-9)^2 + (4-7)^2 + (25-20)^2} = \sqrt{16+9+25} = \sqrt{50}$$

$$d_3 = \sqrt{(5-2)^2 + (4-5)^2 + (25-18)^2} = \sqrt{9+1+49} = \sqrt{59}$$

$$d_4 = \sqrt{(5-1)^2 + (4-3)^2 + (25-9)^2} = \sqrt{16+1+256} = \sqrt{273}$$

$$d_5 = \sqrt{(5-6)^2 + (4-7)^2 + (25-8)^2} = \sqrt{1+9+225} = \sqrt{235}$$

$$d_6 = \sqrt{(5-4)^2 + (4-5)^2 + (25-10)^2} = \sqrt{1+1+225} = \sqrt{227}$$

$$d_7 = \sqrt{(5-5)^2 + (4-5)^2 + (25-6)^2} = \sqrt{0+1+361} = \sqrt{362}$$

$$d_8 = \sqrt{(5-3)^2 + (4-6)^2 + (25-10)^2} = \sqrt{4+4+225} = \sqrt{233}$$

$$d_9 = \sqrt{(5-2)^2 + (4-7)^2 + (25-8)^2} = \sqrt{9+9+225} = \sqrt{307}$$

step-2 ^{after} arrange the distance in non-decreasing order
 $d_2 < d_3 < d_1 < d_6 < d_8$

∴ The neighbour are = x_2, x_3, x_1, x_6, x_8
! ! ! ! !
bad bad Good bad Good

3 → bad
 2 → Good

as no. of bad quality = 3 > No. of Good quality
 ∴ New item is of Bad Quality \square

Q. 2. A Mobile manufacturing Company has manufactured a new mobile asset. The Company want to push adds to potential users to buy new mobile. The Company has data set that contains multiple through the social network profile the dataset contains lots of information about the users but the approx. salary, age, and amount of time spent with the mobile for internet are the 3-key predictors to determine the potential candidate. The target is the purchase item earlier of similar type. apply K-NN - to predict whether a new user is a potential candidate for the add or not.

Q. Write a program in Python to cleanse the data set, determine the principle components, plot the graph & Predict the class for any new items.

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x
	User ID	MC ID	age	Gender	Salary	PIN code	time spent 1 day with mobile with internet	Purchase item of similar type earlier
x_1	1509784	204.192.1.32	32	M	50K	751024	2.5	Yes.
x_2	1508788	1.1.1.32	45	F	60K	751021	1	No
x_3	201569	20.20.20.3	19	F	5K	751001	5	No
x_4	192036	192.168.1.3	21	M	30K	621569	3	No
x_5	2016167	11.20.32.69	27	M	45K	321892	2	Yes.
x_6	172962	192.168.15	25	F	30K	552651	5	Yes
x_7	281969	10.6.42.10	17	F	50K	456521	4	No
x_{new}	242552	92.10.3.46	22	M	70K	66102	3.4	(?)

$K=3$

- determination of a potential candidate to purchase a new mobile depends upon, age, salary & time spent

$$d_1 = \sqrt{(20-32)^2 + (70-50)^2 + (3.4-2.5)^2} =$$

$$d_2 = \sqrt{(20-45)^2 + (70-60)^2 + (3.4-1)^2} =$$

$$d_3 = \sqrt{(20-19)^2 + (70-5)^2 + (3.4-5)^2} =$$

$$d_4 = \sqrt{(20-28)^2 + (70-30)^2 + (3.4-3)^2} =$$

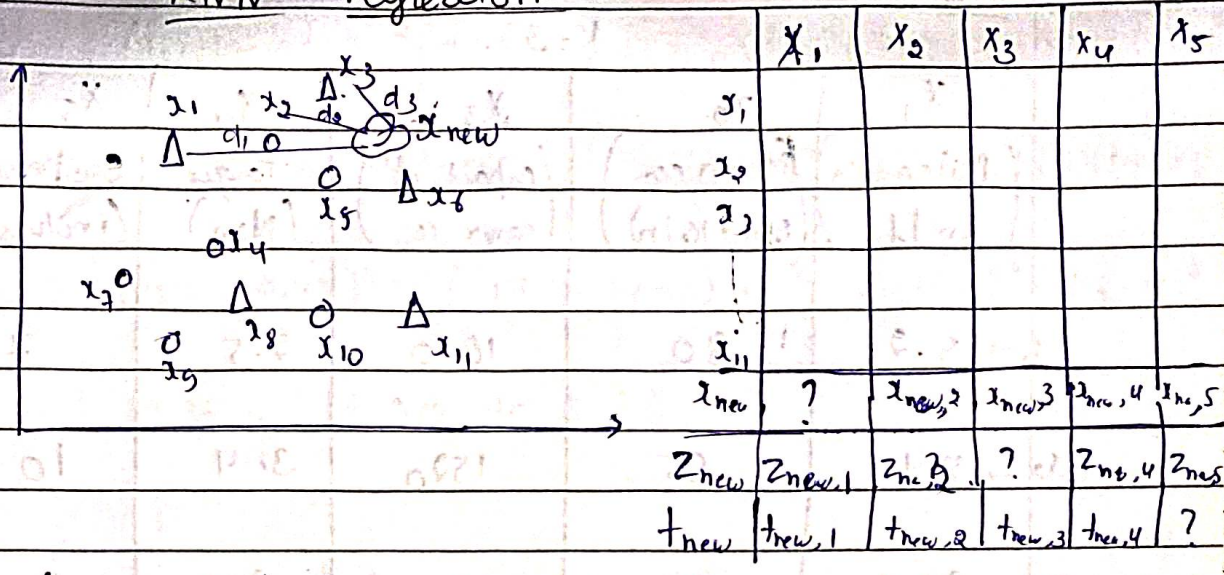
$$d_5 = \sqrt{(20-27)^2 + (70-45)^2 + (3.4-2)^2} =$$

$$d_6 = \sqrt{(20-25)^2 + (70-30)^2 + (3.4-5)^2} =$$

$$d_7 = \sqrt{(20-17)^2 + (70-50)^2 + (3.4-4)^2} =$$

15-Jan

KNN Regression:



i) Compute the distances from new-sample to the given samples.

Let the distances are: $d_1, d_2, d_3, \dots, d_{11}$

ii) Arrange the distances in non-decreasing order for N -neighbours.

Let $k=3$ (given)

$$d_5 < d_6 < d_3$$

iii) The 3 nearest neighbours are: x_5, x_6 & x_3

iv) \therefore target of $x_{new} = \frac{x_{51} + x_{61} + x_{31}}{3}$

z_{new}

Problem:

$K=3$

X_1 (Mileage in km/l.)	X_2 (Mean eco-speed (km/h))	X_3 (Capacity of engine cc)	X_4 (Torque (Nm))	X_5 (Suspension (inches))
$x_1 \rightarrow 15.2$	80	1600	2.5	7
$x_2 \rightarrow 18.4$	95	1520	3.4	10
$x_3 \rightarrow 20.4$	98	1400	3.0	8
$x_4 \rightarrow 25.2$	84	1360	2.8	6.5
$x_5 \rightarrow 22.8$	76	1440	4.8	9
Q1. 18.4		1480	7.2	8
Q2. 19.5	85	1600	7	9.2

$$d_1 = \sqrt{(18.4 - 15.2)^2 + (1480 - 1600)^2 + (7.2 - 2.5)^2 + (8 - 7)^2}$$

$$= \sqrt{10.24 + 14,400 + 22.09 + 1} \Rightarrow \sqrt{14433.33}$$

$$d_2 = \sqrt{(18.4 - 18.4)^2 + (1480 - 1520)^2 + (7.2 - 3.4)^2 + (10 - 8)^2}$$

$$= \sqrt{0 + 1600 + 14.44 + 4} = \sqrt{1618.44}$$

$$d_3 = \sqrt{(18.4 - 20.4)^2 + (1480 - 1400)^2 + (7.2 - 3.0)^2 + (8 - 8)^2}$$

$$= \sqrt{4 + 6400 + 17.64 + 0} \Rightarrow \sqrt{6421.64}$$

$$d_4 = \sqrt{(18.4 - 25.2)^2 + (1480 - 1360)^2 + (7.2 - 2.8)^2 + (8 - 6.5)^2}$$

$$= \sqrt{46.24 + 14,400 + 19.36 + 2.25} = \sqrt{14467.85}$$

$$d_5 = \sqrt{(18.4 - 22.8)^2 + (1480 - 1440)^2 + (7.2 - 4.8)^2 + (8 - 9)^2}$$

$$= \sqrt{19.36 + 1600 + 8.76 + 1} = \sqrt{1626.12}$$

arranging the distance in non-decreasing order

$$d_2 < d_5 < d_3$$

∴ The nearest neighbour are: X_2, X_5, X_3

$$X_{new} = \frac{95 + 98 + 76}{3} \approx 89.6$$

2. Q $d_1 = \sqrt{(19.5 - 15.2)^2 + (85 - 80)^2 + (1600 - 1600)^2 + (9.2 - 7)^2}$

$$= \sqrt{48.33}$$

$$d_2 = \sqrt{(19.5 - 18.4)^2 + (85 - 95)^2 + (1600 - 1520)^2 + (10 - 9.2)^2}$$

$$= \sqrt{1.21 + 100 + 6400 + 0.64} = \sqrt{6501.85}$$

$$d_3 = \sqrt{(19.5 - 20.4)^2 + (85 - 98)^2 + (1600 - 1400)^2 + (9.2 - 8)^2}$$

$$= \sqrt{1.21 + 169 + 40,000 + 1.44} = \sqrt{40,171.85}$$

$$d_4 = \sqrt{(19.5 - 25.9)^2 + (85 - 84)^2 + (1600 - 1360)^2 + (9.2 - 6.5)^2}$$

$$= \sqrt{32.49 + 1 + 57,600 + 7.29} = \sqrt{57,640.78}$$

$$d_5 = \sqrt{(19.5 - 22.8)^2 + (85 - 76)^2 + (1600 - 1440)^2 + (9.2 - 9)}$$

$$\sqrt{10.89 + 81 + 25600 + 0.04} \Rightarrow \sqrt{25691.93}$$

after array the distance in row - decreases and.

$$d_1 < d_2 < d_5$$

The three nearest neighbors are x_1, x_2, x_5

$$\Rightarrow Y = \frac{2.5 + 3.4 + 4.8}{3} \Rightarrow \frac{x_1 + x_2 + x_5}{3}$$

$$\Rightarrow \frac{10.7}{3} = 3.566$$

18-Jan

Simple Linear Regression

Problem statement:

Given a dataset $D = \{x_i, y_i\}_{i=1}^N$, $x_i \in \mathbb{R}$, $y_i \in \mathbb{R}$

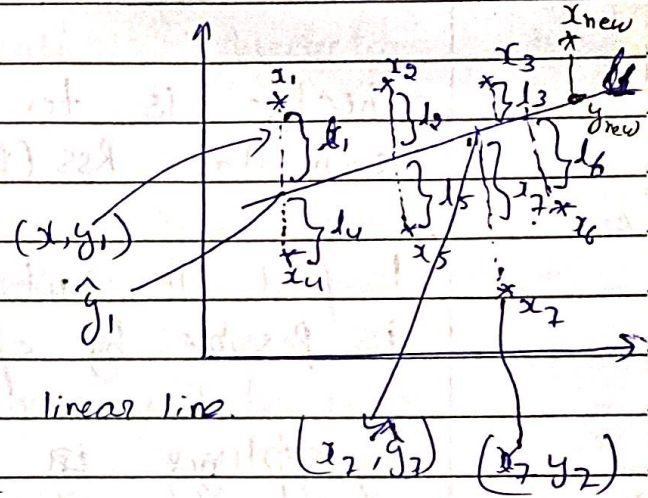
Then for any new sample $x_{new} \neq x_i \forall i=1$ to N
 Predict the target $y_{new} = ?$

Explanation:

Now as it is a linear regression model, we can represent it as:

$$y_i = \beta_0 + \beta_1 x_i \quad \text{--- (1)}$$

where:
 β_1 \rightarrow Slope of the linear line.
 β_0 \rightarrow intercept



As we have to predict:

$$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \text{--- (2)}$$

where $\hat{\beta}_0, \hat{\beta}_1$ are the estimate of β_0 & β_1 , & \hat{y}_i is the estimate of y_i

- i.e. we have to estimate the two unknown $\hat{\beta}_0$ & $\hat{\beta}_1$.
 ie to determine the suitable $\hat{\beta}_0$ & $\hat{\beta}_1$, there are several methods. However, we may use the RSS (Residual sum square) error method to estimate $\hat{\beta}_0$ & $\hat{\beta}_1$.

i.e. $RSS(\beta) = e_1^2 + e_2^2 + e_3^2 + \dots + e_N^2$

$$e_i^2 = (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + \dots + (y_i - \hat{y}_i)^2 + \dots + (y_N - \hat{y}_N)^2$$

$$= (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 + \dots + (y_N - \hat{\beta}_0 - \hat{\beta}_1 x_N)^2$$

$$\boxed{RSS(\beta) = \sum_{i=1}^N (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2} \quad (3)$$

objective is to determine $\hat{\beta}_0$ & $\hat{\beta}_1$ such that $RSS(\beta)$ will be minimum

to find the suitable value $RSS(\beta)$ will be minimum is possible by applying application of derivatives.

by applying in the calculus of derivative w.r.t to β we have $\hat{\beta}_1$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^N (x_i - \bar{x})^2} \quad (4a)$$

and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (4b)$

Now, by using the data set we can estimate $\hat{\beta}_1$

D:

	Predictor	Target
	x_1	y
observation	x_{11}	y_1
	x_{21}	y_2
	x_{31}	y_3
	\vdots	
	x_{N1}	y_N
	x_{new}	?

Then substituting $\hat{\beta}_1$ in eq - (4b) we can estimate $\hat{\beta}_0$.

⇒ finally we can estimate $\hat{y}_{new} = \hat{\beta}_0 + \hat{\beta}_1 x_{new}$ Ans

- limitation of simple linear equation module:

1) for any sample with multiple predictors this method is not suitable.

2) this method assumes that there exist a linear relationship b/w the predictor & the target.

→ However if not then error will be more

• Solution: to overcome ~~the~~ difficulty of the limitation 1 multiple linear equation module may be use

→ to overcome limitation 2. any Non linear ^{Regression} module is used. for ex: SVM (support vector Machine)

W...

Get back to cells amount of on item form
to sign date)

amount paid	amount of cells / to material amount (\$)
-------------	--

6850	205068
20582	6088922
12522	55006382
5875	2222924
22305	406049
42225	?

$$RSS(\beta) = (Y - X\beta)^T (Y - X\beta) \quad \text{--- (1)} \quad RSS(\beta) \in \mathbb{R}$$

We have to find 1st order derivative
 - we have to differentiate eq-(1) w.r.t. β

let $z = Y - X\beta$

$z^T \Rightarrow (Y - X\beta)^T$

$\therefore RSS(\beta) = z^T z$

$$\Rightarrow \frac{\partial}{\partial \beta} (RSS(\beta)) = \frac{\partial}{\partial \beta} (z^T z) \quad \text{--- (1A)}$$

for References:

Case: 1:

$\frac{\partial Y}{\partial \beta}$	$\frac{\partial Y}{\partial \beta_0}$
$\frac{\partial Y}{\partial \beta_1}$	$\frac{\partial Y}{\partial \beta_1}$
$\frac{\partial Y}{\partial \beta_j}$	$\frac{\partial Y}{\partial \beta_j}$
$\frac{\partial Y}{\partial \beta_p}$	$\frac{\partial Y}{\partial \beta_p}$

Y : scalar
 β : vector

Case: 2

$\frac{\partial Y}{\partial \beta_0}$	$\frac{\partial Y_1}{\partial \beta_0}$	$\frac{\partial Y_2}{\partial \beta_0}$	$\frac{\partial Y_3}{\partial \beta_0}$	\dots	$\frac{\partial Y_N}{\partial \beta_0}$
$\frac{\partial Y}{\partial \beta_1}$	$\frac{\partial Y_1}{\partial \beta_1}$	$\frac{\partial Y_2}{\partial \beta_1}$	$\frac{\partial Y_3}{\partial \beta_1}$	\dots	$\frac{\partial Y_N}{\partial \beta_1}$
$\frac{\partial Y}{\partial \beta_j}$	$\frac{\partial Y_1}{\partial \beta_j}$	$\frac{\partial Y_2}{\partial \beta_j}$	$\frac{\partial Y_3}{\partial \beta_j}$	\dots	$\frac{\partial Y_N}{\partial \beta_j}$
$\frac{\partial Y}{\partial \beta_p}$	$\frac{\partial Y_1}{\partial \beta_p}$	$\frac{\partial Y_2}{\partial \beta_p}$	$\frac{\partial Y_3}{\partial \beta_p}$	\dots	$\frac{\partial Y_N}{\partial \beta_p}$

Y : vector $\in \mathbb{R}^N$
 β : vector $\in \mathbb{R}^p$

Case: 3

if A : matrix y $\frac{\partial y}{\partial x}$
 x : vector Ax A^T (5A)

$x^T A$ A (5B)

$x^T x$ $2x$ (5C)

$x^T Ax$ $Ax + A^T x$ (5D)

Further if $(y = f(x); u = g(x)) - \frac{dy}{dx} = \frac{\partial y}{\partial u} \cdot \frac{du}{dx}$

Chain Rule

if case of matrix/vector, we have

$z = f(B)$, $RSS = \phi(z)$

$\therefore \frac{\partial (RSS)}{\partial B} = \frac{\partial z}{\partial B} \cdot \frac{\partial (RSS)}{\partial z}$

Now from eq (4A) we have

$\frac{\partial RSS(B)}{\partial B} = \frac{\partial (z^T z)}{\partial B}$

$= \frac{\partial z}{\partial B} \cdot \frac{\partial (z^T z)}{\partial z}$

$= \frac{\partial (y - xB)}{\partial B} \cdot 2z$ (\because by eq 5C)

$$= \frac{\partial Y}{\partial \beta} \cdot 2Z - \frac{\partial}{\partial \beta} (X\beta) \cdot 2Z$$

$$= 0$$

$$= 0 - X^T \cdot 2(Y - X\beta) \quad [\because \text{by eq-5A}]$$

$$\Rightarrow \frac{\partial}{\partial \beta} (RSS(\beta)) = -2X^T Y + 2X^T X \beta = 0 \quad [\because \text{for minimization}]$$

$$\Rightarrow 2X^T X \beta = 2X^T Y$$

$$\Rightarrow X^T X \beta = X^T Y$$

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

[\because by Multiplying $(X^T X)^{-1}$ on B.S]

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

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

Algorithm (Multiple Linear Regression)

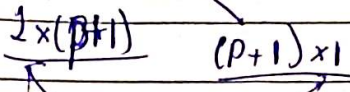
1. Input : Dataset = $\{x_i, y_i\}_{i=1}^N$ $x_i \in \mathbb{R}^P$, $y_i \in \mathbb{R}$

and new sample $x_{\text{new}} \neq x_i, \forall i = 1 \text{ to } N$

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

3. formulate $\tilde{x} = \begin{bmatrix} 1 \\ x_{\text{new}} \end{bmatrix} = \begin{bmatrix} 1 \\ x_{\text{new},1} \\ \vdots \\ x_{\text{new},P} \end{bmatrix}$

4) estimate $y_{new} = \hat{\beta}^T \tilde{x}$



$1 \times 1 = \text{Real No.}$

y_{new} is a no.

$y_{new} \in \mathbb{R}$

Q

$X_1 = \text{GPA}$

$X_2 = \text{IQ}$

$X_3 = \text{level}$

$X_4 = \text{interaction b/w GPA \& IQ}$

$X_5 = \text{interaction b/w GPA \& level}$

$\hat{\beta}_0 = 50, \hat{\beta}_1 = 20, \hat{\beta}_2 = 0.07, \hat{\beta}_3 = 35$

$\hat{\beta}_4 = 0.01, \hat{\beta}_5 = -10$

(Di)

X_1 (GPA)	X_2 IQ	X_3 (level)	X_4 1: Coll 0: HS	X_5 GPA x IQ	X_6 GPA x level	Y score
$\beta_0 = 50$						
$\beta_1 = 20$						
$\beta_2 = 0.07$						
$\beta_3 = 35$						
$\beta_4 = 0.01$						
$\beta_5 = -10$						

a) $Y_{\text{new-college}} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_3 + \hat{\beta}_4 X_4 + \hat{\beta}_5 X_5$

GPA: $G = 20 \Rightarrow 50 + 20G + 0.07 \cdot Q + 35 \times 1 + 0.01 \cdot GQ$

~~IQ: $Q = 20$~~ $+ (-10) \cdot G \times 1$

$\Rightarrow 50 + 20G + 0.07Q + 35 + 0.01GQ - 10G$

$\Rightarrow 85 + 10G + 0.07Q + 0.01GQ$

$Y_{\text{new-high school}} \Rightarrow 50 + 20G + 0.07Q + 35 \times 0 + 0.01 \cdot GQ - 10G$

$= 50 + 20G + 0.07Q + 0 + 0.01GQ - 0$

$= 50 + 20G + 0.07Q + 0.01GQ$

\Rightarrow College value is more.

College: $85 + 10G$

High School: $50 + 20G$

$\Rightarrow 85 + 10G = 50 + 20G$

$85 - 50 = 20G - 10G$

$35 = 10G$

$G = 3.5$

b

$50 + 20G + 0.07Q + 35 + 0.01GQ$

$50 + 20 \times 4 + 0.07 \times 110 + 35 + 0.01 \times 4 \times 110$

$50 + 80 + 7.7 + 35 + 0.01 \times 4 \times 110$

$\Rightarrow G=4 \ \& \ Q=110$

$85 + 10G + 0.07Q + 0.01GQ$

$85 + 10 \times 4 + 0.07 \times 110 + 0.01 \times 110 \times 4$

$\Rightarrow 85 + 40 + 7.7 + 4.4 = 137.1$

$\Rightarrow 137.1$

Line 1 = $y = a + bx$

Line 2 = $z = a + bx + cy$

Line 3 = $w = a + bx + cy + dz$

$y_0 = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5$

$y_{P+1} = \beta_1 x_{P+1} + \beta_2 x_{P+2} + \dots + \beta_P x_{P+P}$

dimensi (P+1) - variabel P

Limitation of least square regression

- We the predictors & the target are not linear related so application of least square regression

$$\text{least square model} = \hat{y}_{new} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 \dots + \hat{\beta}_p x_p$$

- The accuracy of model is not as desired hence accuracy.

The result obtained by those models are not accepted. To improve the performance of the model we may regularise the parameters (B-values) obtained.

- in this situation, there may exist some unnecessary predictors which affect the performance of the model. So by shrinking the impact of those predictors, the performance can be improved.

- there are several methods for this..

- i) Ridge regression / Ridge regularisation
- ii) LASSO Regression / Regularisation

RIDGE Regression (shrinkage Method)

Analytics: To reduce the impact of unnecessary predictors Ridge regression is used.
- it is one of the shrinkage method

Given: $D = \left\{ x_i, y_i \right\}_{i=1}^N$, $x_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$

Predict y_{new} for $x_{new} \neq x_i$; $i = 1$ to N
with desired accuracy.

Now,

step 1: Given matrix: $X =$

x_{11}	x_{12}	x_{13}	...	x_{1p}
x_{21}	x_{22}	x_{23}	...	x_{2p}
\vdots	\vdots	\vdots	\vdots	\vdots
x_{i1}	x_{i2}	x_{i3}	...	x_{ip}
...
x_{N1}	x_{N2}	x_{N3}	...	x_{Np}

$y =$

y_1
y_2
y_3
\vdots
\vdots
y_N

→ Step 2: Compute $\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_p$

$$\text{by } \bar{x}_j = \frac{\sum_{i=1}^N x_{ij}}{N}$$

and subtract \bar{x}_j from corresponding element of that column.

$$\therefore \text{New } X = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & & x_{2p} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{i1} - \bar{x}_1 & x_{i2} - \bar{x}_2 & & x_{ip} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & & x_{Np} - \bar{x}_p \end{bmatrix}$$

→ Step 3: Compute.

$$\beta_0 = \frac{\sum_{j=1}^N y_j}{N}$$

→ Step 4: Then the updated $y = \begin{bmatrix} y_1 - \beta_0 \\ y_2 - \beta_0 \\ \vdots \\ y_N - \beta_0 \end{bmatrix} = \begin{bmatrix} y_1 - \beta_0 & y_2 - \beta_0 \\ & y_3 - \beta_0 \\ & & \ddots \\ & & & y_N - \beta_0 \end{bmatrix}$

→ Step 5: Now we have to determine β by minimizing the below function.

$$\text{Ridge RSS}(\beta) = \sum_{j=1}^N \left[y_j - \left[\beta_0 + \sum_{i=1}^p \beta_i x_{ij} \right] \right]^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Multiple Linear Regression.

Regularization is to be selected judiciously.

(K fold cross validation is used to estimate λ)

So by differential eq (*) for minimization
Ridge

$$\Rightarrow \frac{\partial \text{RSS}(\beta)}{\partial \beta} = -2X^T(Y - X\beta) + \lambda 2\beta$$

$$\Rightarrow \frac{\partial \text{RSS}(\beta)}{\partial \beta} = -2X^T(Y - X\beta) + \lambda 2\beta = 0 \quad (\because \text{Minimize})$$

$$\Rightarrow -2X^T(Y - X\beta) + 2\lambda\beta = 0$$

$$\Rightarrow -2X^TY + 2X^TX\beta + 2\lambda\beta = 0$$

$$\Rightarrow 2(X^TX + \lambda I)\beta = 2X^TY$$

$$\Rightarrow (X^TX + \lambda I)\beta = X^TY$$

$$\Rightarrow (X^TX + \lambda I)^{-1} (X^TX + \lambda I)\beta = (X^TX + \lambda I)^{-1} X^TY$$

$$\Rightarrow I\beta = (X^TX + \lambda I)^{-1} X^TY$$

$$\Rightarrow \hat{\beta} = (X^TX + \lambda I)^{-1} X^TY$$

extra only. compare to multiple dim β .

$(X^T X)$ is not a full rank matrix

2. the cross validation method may be used to determine λ

3. Ridge regression shrinks the coefficients of sum of the predictors which have less significant effect on response target.

5/02/24

Date: / /

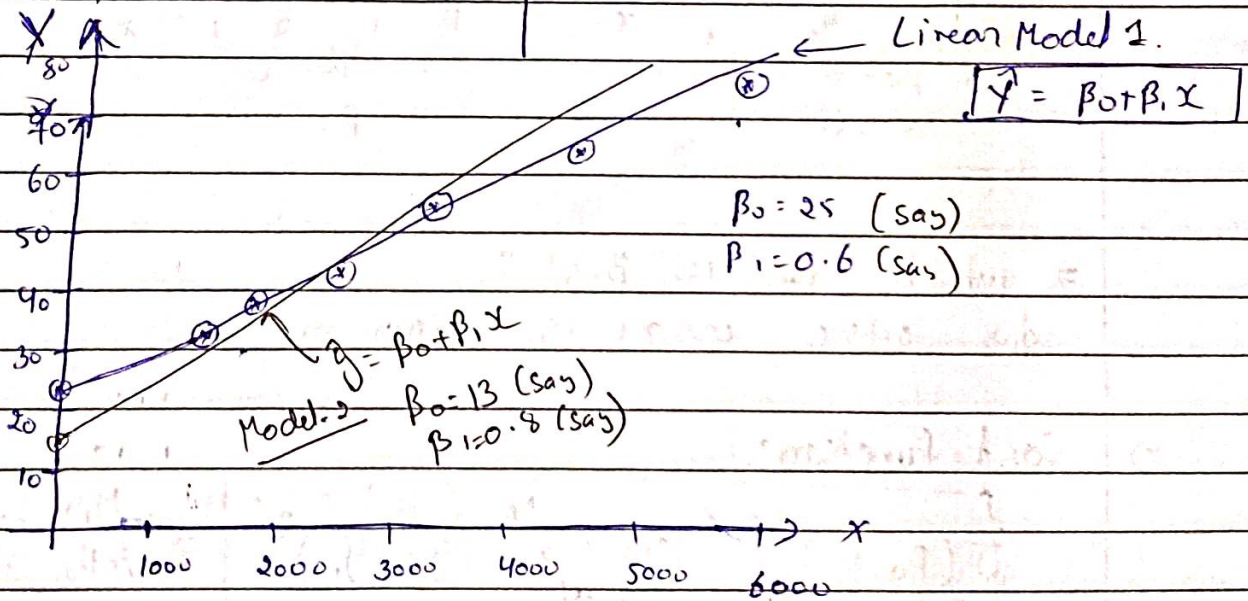
SIMPLE Linear Regression with one variable

Model Representation:

- It is supervised learning alg.
- Given any new sample, it is expected to estimate accurately the response / target.

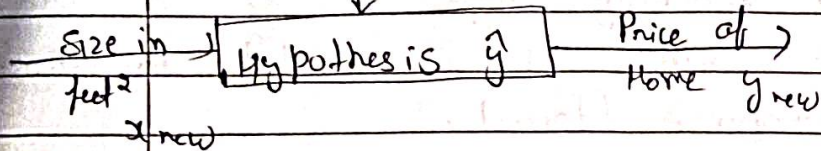
- ex: Estimating Price of a House?

Size of flat (x)	Price of House in Lakh (₹)
$x_1 \rightarrow 1474$	34
$x_2 \rightarrow 1632$	39
$x_3 \rightarrow 3407$	56
$x_4 \rightarrow 4408$	67
$x_5 \rightarrow 5569$	79
$x_6 \rightarrow 2696$	42
$x_7 \rightarrow$	



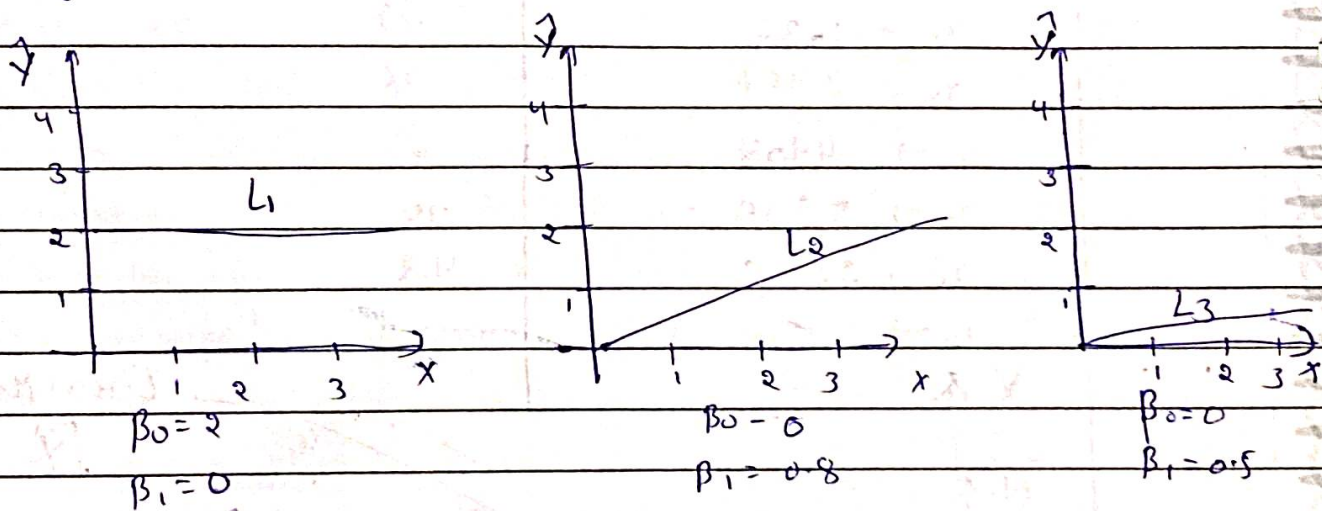
⇒ Training Set:

↓
Learning Algo.



→ How to Represent \hat{y}

$$\hat{y} = \beta_0 + \beta_1 x$$



⇒ which one is Best?

Ans: where error is Minimum

⇒ Cost function:

$$J(\beta_0, \beta_1) = \text{Min}_{\beta_0, \beta_1} \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

RSS(β)

$$\Rightarrow \text{Min}_{\beta_0, \beta_1} \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

β is vector of two param.

W 20-1-1

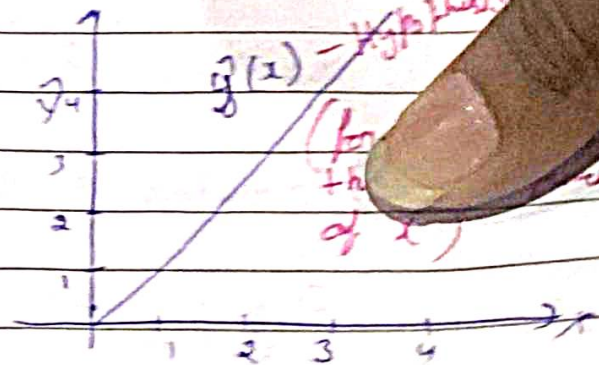
Cost function intuition-1

$$\hat{y} = \beta_0 + \beta_1 x$$

Parameters: β_0, β_1

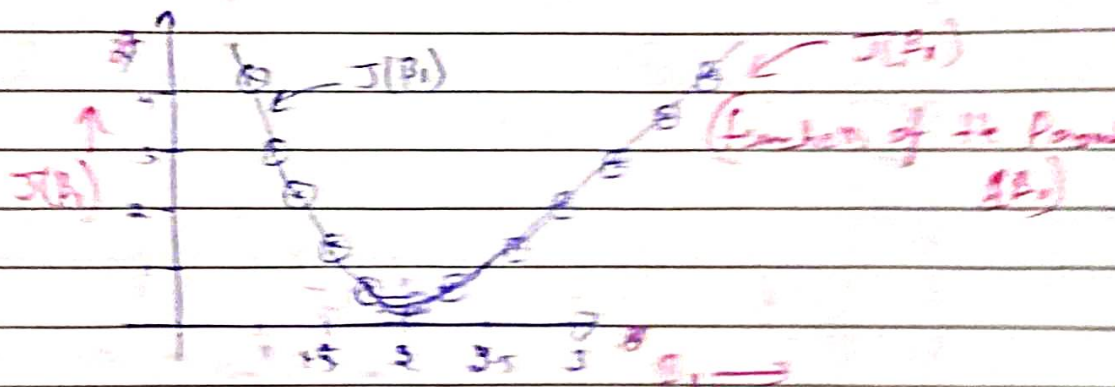
Cost function:

$$J(\beta_0, \beta_1) = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$



Goal: Min $J(\beta_0, \beta_1)$

Simplified version:



$$\hat{y} = \beta_1 x$$
$$J(\beta_1) = \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Goal: Min $J(\beta_1)$

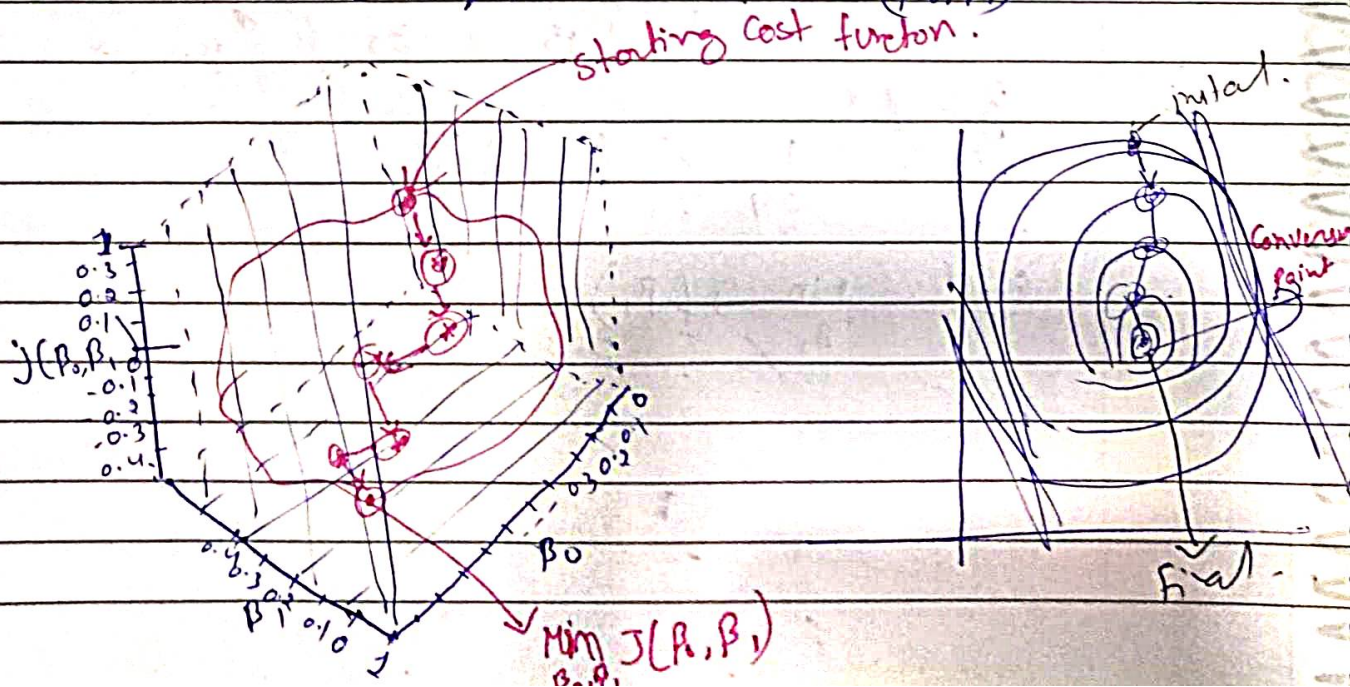
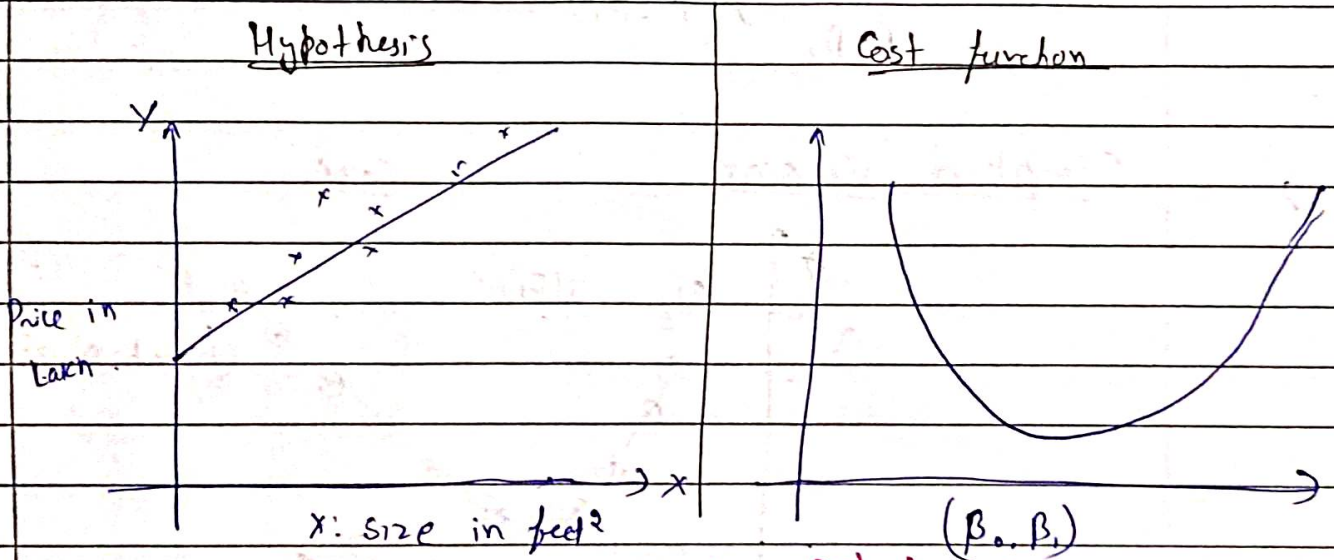
→ Cost function: Introduction-2

Hypothesis: $\hat{y} = \beta_0 + \beta_1 x$

Parameters: β_0, β_1

Cost function = $J(\beta_0, \beta_1) = \frac{1}{2N} \sum_{j=1}^N (y_j - \hat{y}_j)^2$

Goal: $\text{Min}_{\beta_0, \beta_1} J(\beta_0, \beta_1)$



7/02/24

Date _/ _/ _

#

GRADIENT DESCENT

algo 1. start with the parameter

$$\beta_0 = 0$$

$$\beta_1 = 0$$

2. Go on changing the values of β_0 & β_1 s.t the cost function $J(\beta_0, \beta_1)$ get reduced until we reach a minimum.

Repeat until convergence {

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} J(\beta_0, \beta_1)$$

for simulation purpose:

$$\text{temp}_0 = \beta_0 - \alpha \frac{\partial}{\partial \beta_0} J(\beta_0, \beta_1)$$

$$\text{temp}_1 = \beta_1 - \alpha \frac{\partial}{\partial \beta_1} J(\beta_0, \beta_1)$$

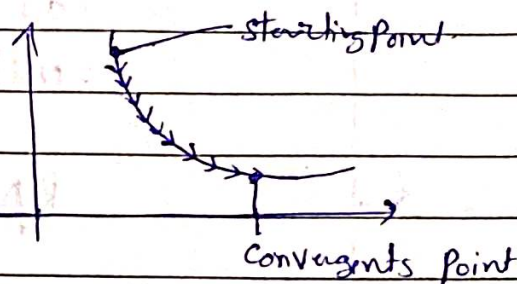
$$\beta_0 = \text{temp}_0$$

$$\beta_1 = \text{temp}_1$$

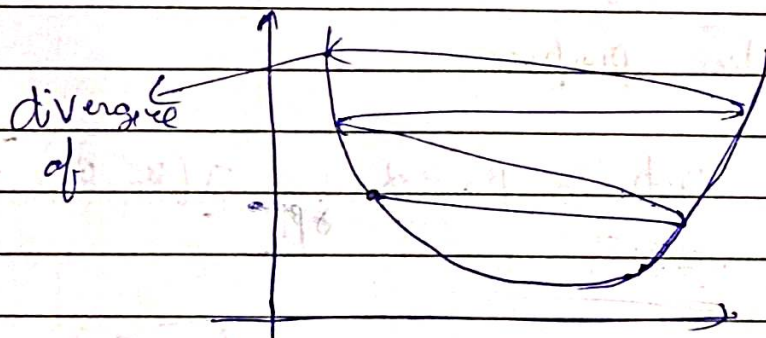
- Intuition of

1. it is a iterative optimization process to minimize objective function.

2. if α (is too small) \Rightarrow gradient descent is slow. (GD)



3. if α (is too large) \Rightarrow gradient descent may not converge. it may diverge.



4. Gradient descent may converge to a local minimum even with learning rate α is fixed.

5. when we approach a local minimum G-D will automatically takes smaller steps

$$\beta_j = \beta_j - \alpha \frac{\partial}{\partial \beta_j} (Cost)$$

$$\frac{1}{2N} \sum_{j=1}^N (y_i - \hat{y}_i)^2$$

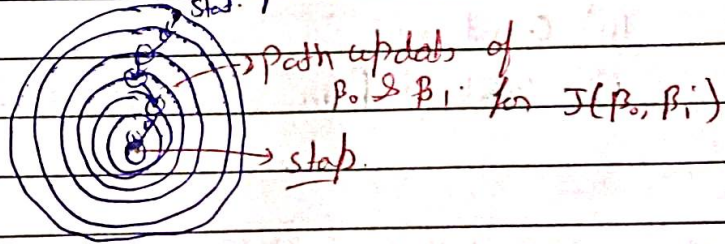
$$\frac{1}{2N} \sum_{i=1}^N e_i^2$$

types of Gradient Descent Method

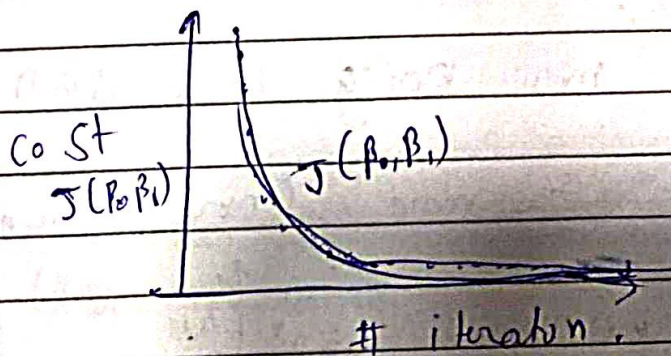
1. Batch Gradient Descent (BGD)
2. Stochastic GD : (SGD)
3. Mini-Batch GD : (MBGD)

1# Batch Gradient Descent: (BGD):

- i) in case of BGD the whole training sample is used to compute the parameter updates.
- ii) that is the parameters β_0, β_1 are updated after computing the gradient of the error w.r.t to entire training set.
- iii) time complexity is more
- iv) it makes smooth updates on model parameter that is



- v) when no. of iteration increases



2# Stochastic GD (SGD)

Pseudocode (SGD):

1. Initialization $\beta_0 \leftarrow \text{random}_0$
 $\beta_1 \leftarrow \text{random}_1$
2. Define K : No. of max iteration
 α : Learning Rate: (0 to 1)
 ϵ : error tolerance: > 0
 $\text{cx} = 0.0001$
 $\text{Count} = 0$
3. do
 { Repeated until convergence or reaches max. iteration }
 - i) $X \leftarrow X$ shuffle
 - ii) $\beta_{\text{new}} = \beta_{\text{old}} - \alpha \frac{\partial}{\partial \beta_j} (P_0, P_1)$
 - iii) $\text{Count}++$
 - iv) while $(\| \beta_{\text{new}} - \beta_{\text{old}} \| \geq \epsilon \text{ or } \text{Count} \leq K)$

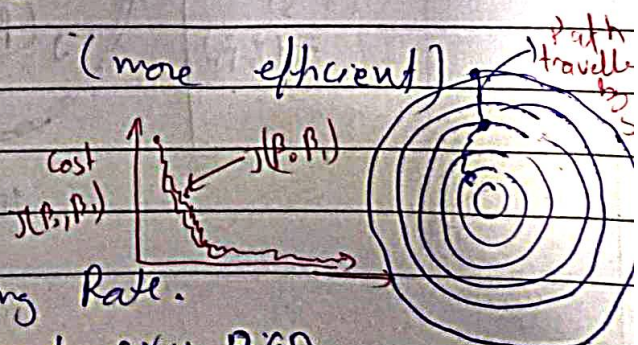
1) it is computational efficient than BGD

2) as it selects the training sample randomly that's why it is Stochastic

3) generally it is more noisy than BGD

4) No. of iteration are more (more efficient)

5) it avoid local minima.



6) However it is sensitive to Learning Rate.

Due to less expensive it is Preferred over BGD.

02/02/24
Date _/ _/ _

#3. Mini-Batch Gradient Descent Alg (MBGD)

i) The parameters are updated after computing the Gradient of the error w.r.t the subset of training set.

ii) The updates are faster in comparison to BGD & slower in comparison to SGD (if it uses less no. of training sample)

iii) Depending upon the BATCH-SIZE the updates can be of less noisy.

-i.e Greater is the Batch size lesser is the Noise.

Algo: MBGD

initialization:

$$\beta_0 = \beta_{init}$$

$$\beta_1 = \beta_{init}$$

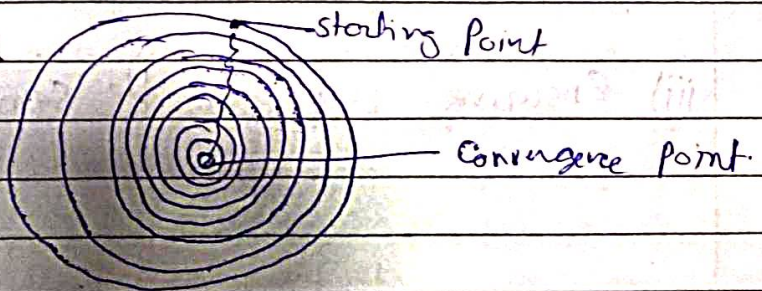
$$K = \text{Max. iteration (No. of epochs)}$$

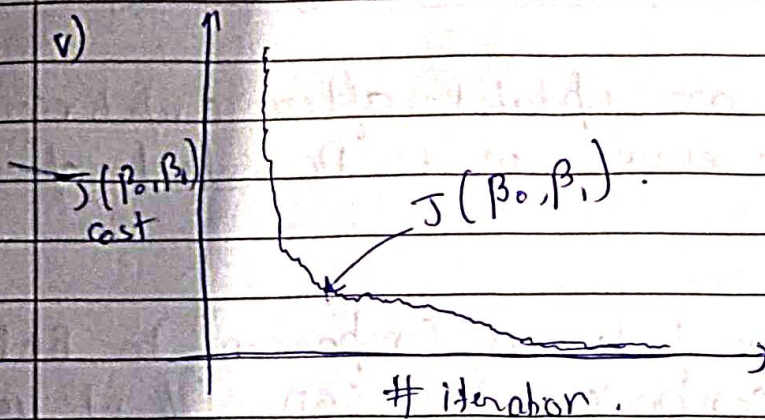
for $i=1$ to K

for mini-Batch (x_{min}, y_{min}).

$$\beta_{new} = \beta_{old} - \alpha \frac{\partial}{\partial \beta_j} J(\beta_0, \beta_1)$$

iv) The Convergence





Feature scaling

- i) Normalization
- ii) Standardization

- Feature scaling is important for building accurate and effective ml. model by using either Normalization or standardization processes.

- that is feature scaling is the process of transforming the data to specific ranges. so as to reduce the impact of outliers and to reduce dominance of huge data of one feature.

Advantages:

- i) improve performance of the model
- ii) Reduces the impact of outlier
- iii) Ensuring data in the same scale.

for ex:

Student	CGPA
A	3.5 in the scale of 0 to 5
B	5.2 " " " 0 to 10
C	75%

→ Performance of the students

- Here we need to scale to performance of the students in same standardized method (scale)

Student	% of final exam X_1	Salary in the company (lacs) X_2
A	3.5	60
B	3.2	70
C	4.0	55
D	4.8	62

range [3, 5] range [50, 100]

The two feature X_1 & X_2 are not in respective scale (Range)

→ Why should we need feature scaling & when? (models)

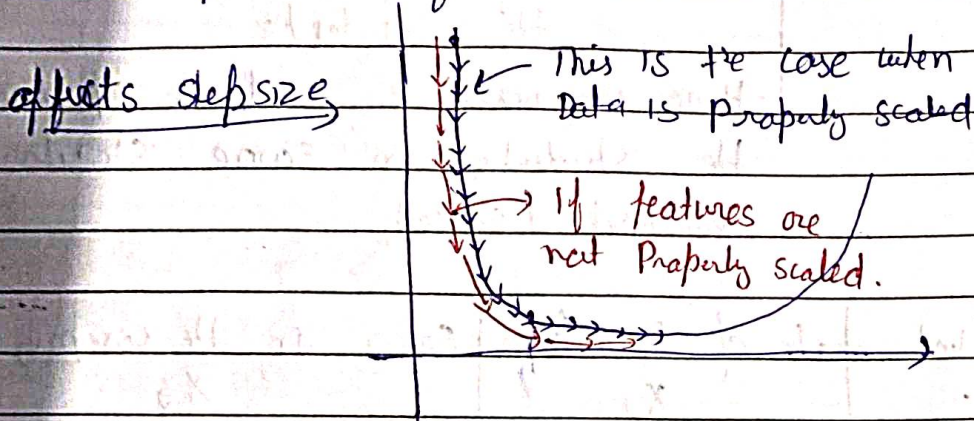
[∵ Some ml Algo are sensitive to feature scaling where as some are insensitive to feature scaling]

#1 Gradient Descent Algo (BGD, MBGD, SGD)

- The ml model such as LR, Logistic regression, Neural NW, PCA (Principal Component Analysis) that use Gradient Descent Algo as an optimization technique. require data to be scale.

$$\beta_j = \beta_j - \alpha \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i) x_j^{(i)}$$

- The presence of feature X that is $x_j^{(i)}$ will effect the step size of Gradient descent.



- to ensure that GD moves smoothly towards the minimum & the steps of GD are updated at the same time & Rate for all the feature. We have to scale the data before fitting to the model feeding (taking input)

#2. Distance based Algo:

The ML models KNN, K means, SVM. are effected by the range of the features. So scaling is important for those algo. also.

activity

for ex:	student	x_1 CGPA	x_2 Salary (lakh)
Scale (0,1)	A	3.5	66
	B	3.2	70
	C	4.0	55
	D	4.8	62

Range [3.5] Range [50,100] The two feature x_1 & x_2 are respective scale (Range)

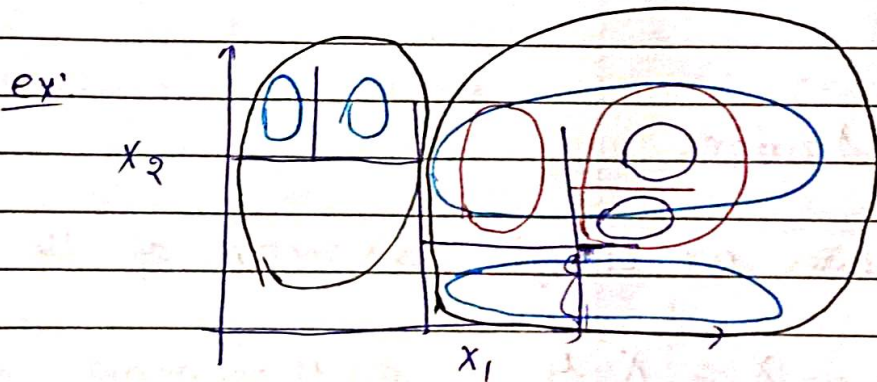
#3. tree - Based Algo.

ML models like decision tree, Random are fairly insensitive to the scale of features

- the decision tree splits a node by selecting the correct feature & then putting splitting a node at a suitable splitting point.

- in this case by splitting it increases the homogeneity of the node.

- other features don't influence this split of feature.



12/02/24 W 2024 --- ---

1. Normalization:

Normalization for feature scaling

In this method the features are rescaled to the range $[0, 1]$. It is also known as min-max feature scaling method.

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

i.e. initially the range of feature was $[X_{\min}, X_{\max}]$ and after Normalization the range becomes $[0, 1]$

2. Standardization:

The features are closed to mean of the corresponding feature.

$$X' = \frac{X - \mu}{\sigma}, \quad \begin{array}{l} \mu = \text{mean} \\ \sigma = \text{Standard Deviation.} \end{array}$$

$$= \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

Note: Variance = σ^2

$$= E[(X - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

$$= E[X^2] - (E[X])^2$$

for 3 marks

Q. When we should use Normalization and when to use standardization as feature scaling.

<u>Normalization</u>	<u>Standardization</u>
1) The features are scaled to the range $[0,1]$.	1) The features are scaled within range $[\mu-\sigma, \mu+\sigma]$
2) when the distribution of data set is unknown & not Gaussian.	- Distribution & Gaussian unknown & Gaussian
3) Retains the shape of the distribution	- original distribution is not retained
4) Sensitive to outliers	- less sensitive to outlier
5) $X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$	$x' = \frac{x - \mu}{\sigma}$
6) The relationship b/w the data point may not be preserved	- relationship b/w data point is preserved.

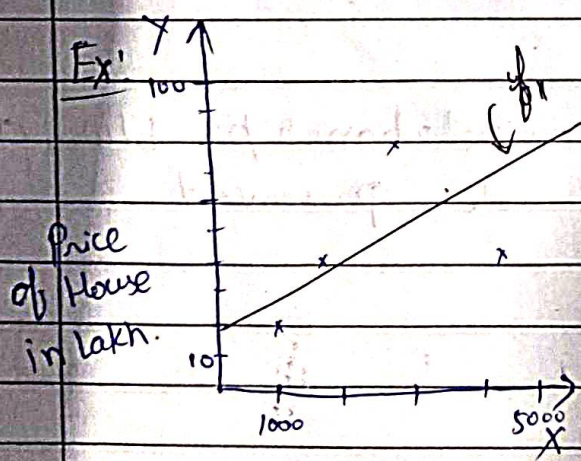
~~Define~~ overfitting & Underfitting.

underfitting: Vs overfitting

underfitting: when a complex model / Problem is modeled by using a simple method then the model may not fit to the data set & may not achieve the desired accuracy.

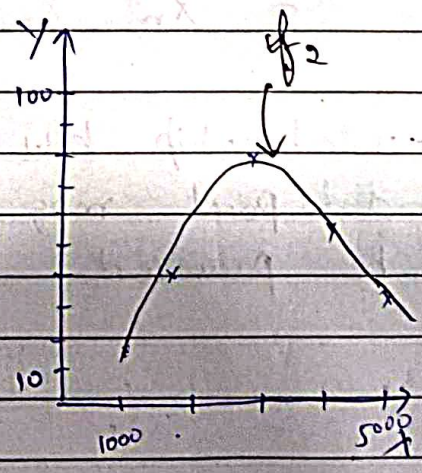
overfitting: when No. of features in data set is large the learned Hypothesis may fit the dataset very well (i.e. $J(\beta) = \frac{1}{N}$

(i.e. $J(\beta) = \frac{1}{2N} \sum_{i=1}^n (y_i - \hat{y}(x_i))^2 \approx 0$) but it couldn't approximate any new data accurately.



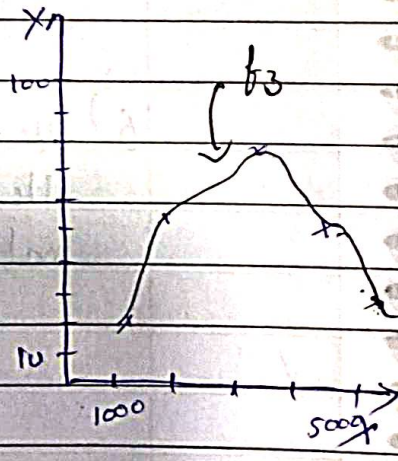
$f_1 = \beta_0 + \beta_1 x_1$
degree - 1

- Under fit



$f_2 = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2$
degree - 2

Just fit



$f_3 = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \beta_3 x_3^4 + \beta_4 x_4^9$
degree - 14

Overfit

Underfit	Just fit	Overfit
<ul style="list-style-type: none"> - Training error <u>large</u> - Test error <u>large</u> 	<ul style="list-style-type: none"> - Training & Test error at <u>Par.</u> (both accepted) 	<ul style="list-style-type: none"> - Training error is <u>minimum</u> - Test error is <u>large</u>

How to address overfit? Method

1) reduce no. of features.

- Manually select the suitable features to keep
- apply model selection Algo.

2) Regularization

- Keeping all the features & reducing the magnitudinal value of the parameters β_j

The methods are:- i) Ridge Regularization ^{of Parameters}
 ii) LASSO Regularization ^{or regularization}

- Scale ~~or~~ Re-scale the features so that all the features can contribute uniformly i.e.,
- $X_1, X_2, X_3, \dots, X_p$ should contribute to Y_{new} uniformly (∵ By regulating the ranges of features)
 → features are regulated)

BIAS & Variance Trade-off

Objective: To decide for a any given set of Data which method produces the best ~~not~~ result i.e. selection of a suitable ML model for a given set of data is the most challenging task for a data Scientist w.r.t ^{stati}Statistical learning in practice.

How to?

- 1) We have to measure the quality of fit of the ML model (least error).

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}(x_i))^2$$

Training Mean Square error

estimated.

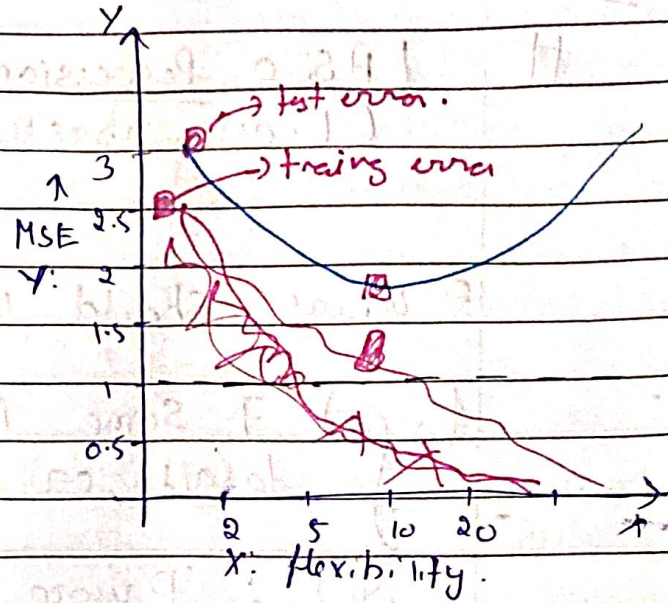
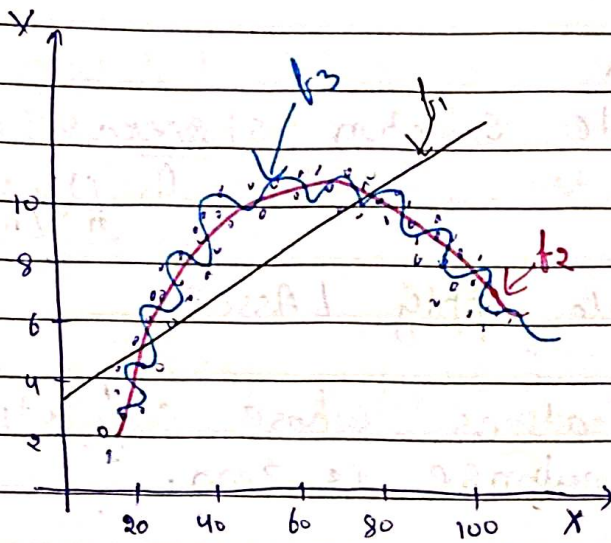
Analysis: Here we are not interested where or not the model accurately predicts the training Data set.

We are interested in accurately predicting on the new sample.

So, for given $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$, $x_i \in \mathbb{R}^p$, $y_i \in \{0, 1\}$

we obtain the estimate \hat{f}

- We can then compute $\hat{f}(x_1), \hat{f}(x_2), \hat{f}(x_3) \dots \hat{f}(x_n), \hat{f}(x_{new})$



- f_1 → linear Regression
- f_2 → approximated spline.
- f_3 → approximated spline.
- training error.

19/02/24

LASSO Regression: (Least absolute selection shrinkage operator Algo)

- i) Ridge (non zero)
- ii) LASSO (zero)

Q Where should we apply LASSO?

A. i) \exists Some features whose contribution towards target or response is zero.

ii) \therefore Presence of those features in the ML model incurs

- a) Unnecessary complexity
- b) affect the performance of the model
- c) Desired accuracy is not achieved

These unnecessary features should be made zero

Note: LASSO is also considered as Features selection.

x_1	x_2	x_3	x_4	x_5	...	x_p	y
	0	0		0			

K - features are selected out of p features:
 $K < p$.

How To?

In case of LASSO regression the parameters β_j ($j=1$ to p) are estimated by minimizing the function as:

$$RSS(\beta) = \sum_{i=1}^N \underbrace{\left[y_i - \left(\beta_0 + \sum_{j=1}^p \beta_j x_{ij} \right) \right]^2}_{\text{error term}} + \lambda \sum_{j=1}^p |\beta_j|$$

$\lambda > 0$
 Regularization term.
 L_1 -Regularization

$\Rightarrow L_2$ -Regularization
 : Ridge Regression.

on. Minimizing:

$$RSS(\beta) = \sum_{i=1}^N \left[y_i - \left(\beta_0 + \sum_{j=1}^p \beta_j x_{ij} \right) \right]^2$$

Subject to $\sum_{j=1}^p |\beta_j| \leq t$, where $t > 0$,
 t is constant

lasso is better than Ridge:

- it discards unnecessary feature.
- Performance is better.
- Accuracy is also good

LOGISTIC Regression (classification)

* it is a supervised ML model for classification.

* it is a statistical model.

Goal: To Predict the class label of any new data X_{new} . by estimating probab $(Y=1 | X=X_{new})=?$

if $(P(Y=1 | X=X_{new}) > 0.5) \Rightarrow X_{new} \in C_1$
(class 1)

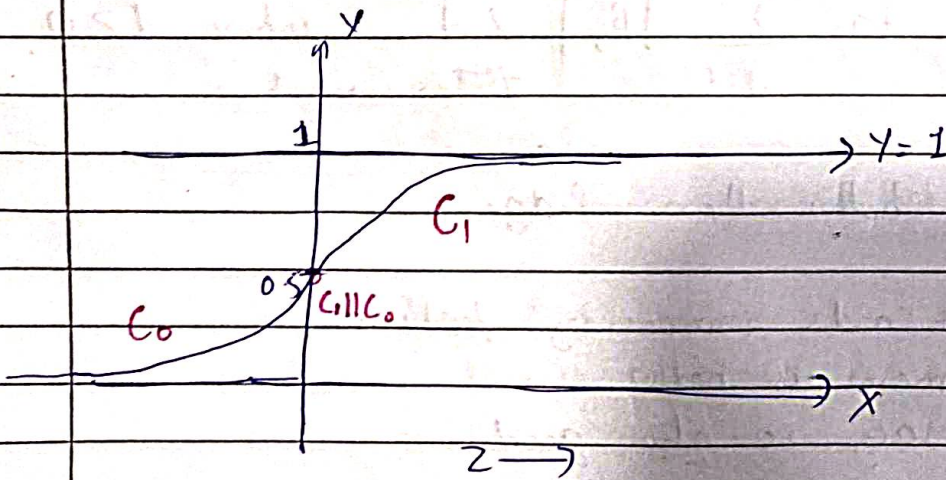
else if $(P(Y=1 | X=X_{new}) < 0.5) \Rightarrow X_{new} \in C_0$
(class 0)

else $(P(Y=1 | X=X_{new}) = 0.5) \Rightarrow X_{new} \in C_1 \cup C_0$

To estimate the Probability SIGMOIDAL function.

is used.

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



if $z=0$
 $g(z)=0.5$

if $z \rightarrow \infty$
 $g(z) \rightarrow 1$

if $z \rightarrow -\infty$
 $g(z) \rightarrow 0$

types of logistic Regression

i) Binomial: 2-class classification problem (0 or 1)

ii) Multinomial: No. of classes are 3 or more & not ordered
eg: Tiger, Lion, Cat, Dog.

iii) ordinal: 3 or more number of ordered classes.
eg: low, Medium, High
: F, D, C, B, A, E, O (grades)

Conditions for implementing logistic Regression (Assumption)

i) Features should be independent. (independent features)

ii) Binary dependent variables.

iii) \nexists any outliers.

iv) Large Data size

v) \exists linear relationship b/w independent variables & log odds (class label)

→ Pseudocode (logistic regression).

Given the dataset $D = \{x_i, y_i\}_{i=1}^N$, $x_i \in \mathbb{R}^p$, $y_i \in \{0, 1\}$

Predict the class label of x_{new} ?

1. initialize β^{old} , choose some small $\epsilon > 0$
(Eg: $\epsilon = 0.0001$), $\alpha \in (0, 1)$

2. Compute Y, X

3. for $i = 1$ to N

$$p_i = \frac{1}{1 + e^{-\beta^{old T} x_i}}, \text{ where } \tilde{x}_i = \begin{bmatrix} 1 \\ x_i \end{bmatrix}$$

4. Compute p : $p_1, p_2, p_3, \dots, p_N$, where $p_i = P(x_i, \beta^{old})$

5. compute $\beta^{new} = \beta^{old} + \alpha X^T (Y - p)$ \rightarrow (*gradient ascent*)

6. if $\|\beta^{new} - \beta^{old}\| < \epsilon \Rightarrow$ Return β . \rightarrow /* it is done */

7) else $\beta^{old} = \beta^{new}$ and goto step-3

$$8) P(Y=1 | X=x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p)}}$$

9) if $(P(Y=1 | X=x)) > 0.5 \Rightarrow x$ is in class 1.

10) else if $(P(Y=1 | X=x)) \leq 0.5 \Rightarrow x$ is in class 0

explanation:

$$\| \beta^{new} - \beta^{old} \|$$

$$\Rightarrow (\beta_1^{new} - \beta_1^{old})^2 + (\beta_2^{new} - \beta_2^{old})^2 + \dots + (\beta_p^{new} - \beta_p^{old})^2$$

Now if $(\dots) \leq 0.0001$

$$\Rightarrow \text{Return } \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}$$

β is vector.

Error Analysis

Why?

We need to develop a

- i) reliable
- ii) robust
- iii) least bias
- iv) least variance

So ~~we~~ generalized ML model for real-life Problem.

To achieve such a model the following parameters are considered for model's performance evaluation.


The parameters are:

- i) accuracy of the models
- ii) Precision " " "
- iii) Recall " " "
- iv) F1-score " " "

1. error analysis in detail

It is the process of isolating, observing, & diagnosing the erroneous ML predictions so as to determine the high & low performing ML model.

Case-Study - I: Image Recognition ML model for Dog detection.

 Benchmark	In day light accuracy: 100%	In Night 70%	In Rainy night 65%	in a water pool 80%

Alg ₁	— (86%)	(62%)	(60%)	(76%)
Alg ₂	— (85%)	(65%)	(63%)	(74%)
Alg ₃	— (85%)	(68%)	(58%)	(41%)

- → accepted
- → May be accepted
- → Rejected

2. Error identification & diagnosis

Case study - ~~data~~ (features are not apparent)

When the types are text, graphic, audio, image.

ML model: Cat classification.

Suppose we train a model on: 5000 images (those images include including cat in there different scenario)

No. of test sample: 1000

% of accuracy: (85%)

850 images are correctly classified

150 images are miss classified.

Confusion matrix:

Total = P+N = 1000		Actual.	
		Cat (P)	Not Cat (N)
Predicted	Cat (PP)	700 (TP)	(100) (FP) → Type I error
	No cat (PN)	(50) (FN) → Type-2 error	150 (TN)

Here, the confusion matrix doesn't provide clarity when and what is being miss-classified.

- Consider separate situation & analyse the % of error.

Solⁿ

Case-Study II. (when data is more apparent)

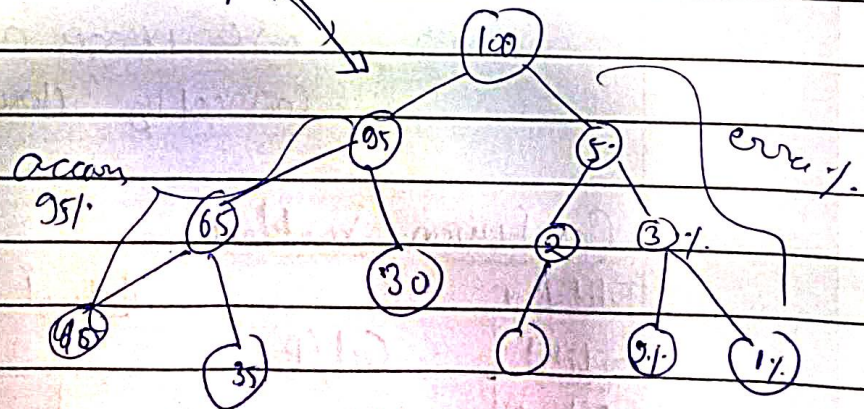
- eg: Predicting the Mileage of a car
- Estimating the Credit Card limit of a person
- Recommending a product

The features should be clearly categorized

→ Segregation of features: can help identifying & further diagnose the error distribution across features & certain values of those features.

methods are: i) Heat maps
ii) Tree maps

x_1	x_2	x_3	x_4	...	x_p
8%	10%	5%	21%		9%



skd 9

#3: How to resolve the error:

after identifying & analysing the errors we need to resolve those errors

Step 1: Collect ~~more~~ more data where errors are happening

Step 2: Augment the data (By transformation)

↓
(i) translation, (ii) scaling,
(iii) rotation, (iv) shearing)

Step 3: Use other ML model or tune the Hyper Parameters differently.

Problem solving (Using Confusion Matrix)

i) Type-1 error (FP): over estimation (False alarm)

ii) Type-2 error (FN): under estimation. (Miss)

iii) Precision: $\frac{TP}{PP} \Rightarrow \frac{TP}{TP + FP}$

iv) Recall: $\frac{TP}{P} \Rightarrow \frac{TP}{TP + FN}$ (Sensitivity)

v) accuracy: $\frac{TP + TN}{P + N}$

vi) F-1-score = $\frac{2 * Precision * Recall}{Precision + Recall}$