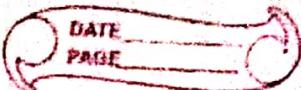


09/01/2022  
Tuesday



# MACHINE LEARNING

→ searching for  
the best hypothesis  
(stmts)

Intelligence (Human intelligence)

Learn new things, creative, adapt to environment,  
ability to learn.

Artificial Intelligence

System should replicate human  
intelligence

Expert System

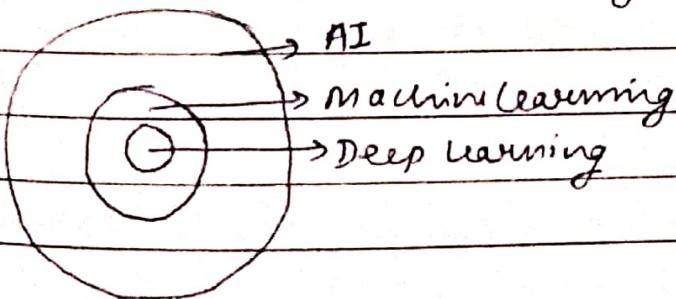
PROLOG and LISP → <sup>AI</sup> languages

Artificial Intelligence :-

- (1) Learning :- It focuses on acquiring data & creating rules for how to turn it into actionable information.
- (2) Reasoning :- It focuses on choosing the right rule to reach out come.
- (3) Self-correction :- It focuses to continuously fine tune rules and ensure they provide the most accurate results possible.
- (4) creativity :- Generate new images, new text, new music.

Generating A.I.

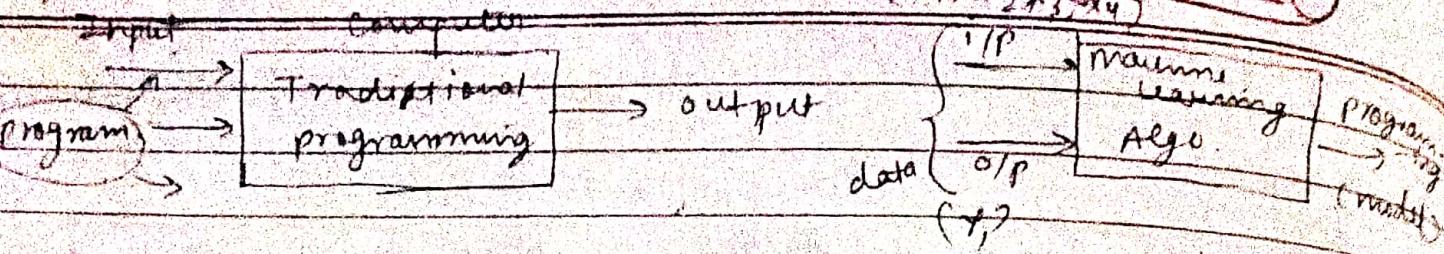
creativity := Deep learning



Machine learning :-

Arthur Samuel (1959) "The field of study that gives computers the ability to learn without being explicitly programmed."

10/01/2024 8:37



Machine learning - Tom M. Mitchell

"A machine is said to learn if it is able to take experience &

where it such that its performance improves upon similar experiences in the future."

A computer program is said to learn from experience E with respect to some task T & performance measure P; if its performance at task T, improves with experience E.

Task :- Decision Making, Prediction, Classification

Experience (Data) :- Collection of different data items related to task.

Performance :- Accuracy, precision, error.

Information related to data :-

Attributes / Features / Dimensions

A	B	C	D	E	T
x1	x2	x3	x4	x5	y1
x6	x7	x8	x9	x10	y2
x11	x12	x13	x14	x15	y3

O/p Attribute = Target attribute  
or O/p label

or ground truth

I/P      O/P

Different O/P

Values are called

Sample Space or Instances or all possible O/P space attributes

Take an example of learning problem  
explore each process in details



### Example:-

#### ① Handwriting recognition problem

Task:- Recognition & classification of handwritten words within images

Performance:- % of words correctly classified

Data:- Dataset of handwritten words with given classification

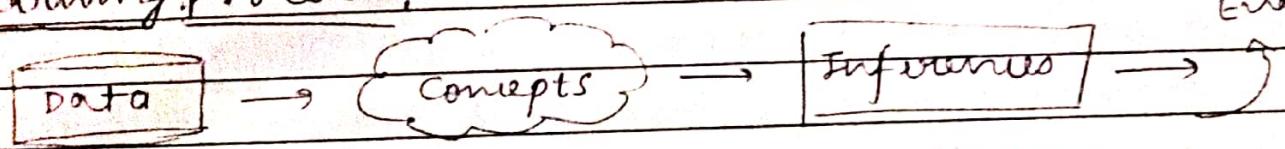
#### ② A robot driving problem.

Task:- Driving on road using vision sensors

Performance:- Average distance travelled before an error

Data:- sequence of images & steering command recorded while observing a human driver.

### Learning process :-



Data Storage      Abstraction      Generalization

① Data :- It utilizes observation, memory storage & provide a factual basis for the further reasoning.

② Abstraction:- (i) ~~Extracting~~ knowledge about the stored data.

③ Inference (ii) creating general concepts about the data as a whole

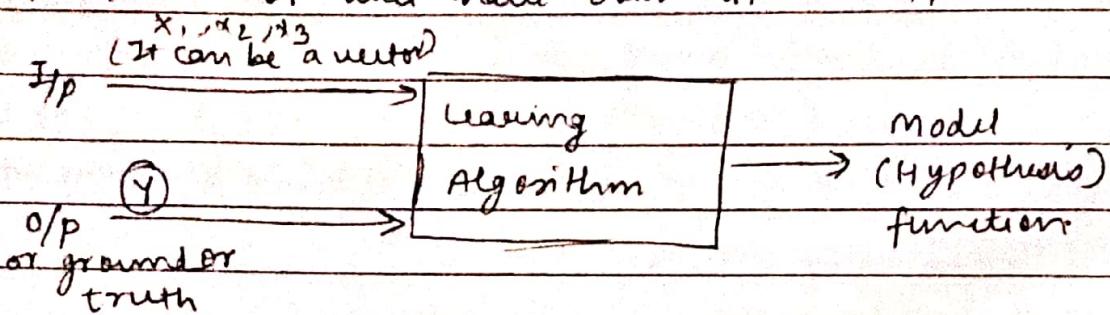
④ Generalization:- Process of turning the knowledge about stored data a form that can be utilized for future actions

⑤ Evaluation:- It is the process of giving feedback to the user to measure the utility of the learned knowledge.

## # Types of Machine Learning

- ① supervised learning
- ② unsupervised learning
- ③ semi-supervised learning
- ④ Reinforcement learning

① Supervised learning: In SL, a model is trained on a "labelled dataset" → It will have both I/P and O/P.



In SL, a training set of examples with correct responses is provided and based on this training set, the algorithm generalizes to respond correctly to all possible I/P.

### Types of SL

will

(i) Classification:- It deals with prediction of historical target variables, which represent discrete classes or labels.

Ex: (i) Classification of image as spam or not spam.

(ii) Predicting whether a patient has high risk of heart disease.

~~Ex :-~~ a) Decision Tree

c) SVM

b) Logistic Regression

d) Naive Bayes

(ii) Regression:- predicting continuous target variables which represents numerical values.

~~Ex :-~~ Regression Algo. learn to map the I/P features to a continuous numerical O/P values.



2024  
Page

format.

Ex. (a) Predicting the price of houses based on its size, location, amenities.

linear regression CGPA

(b) Forecasting the sales of a product.

Ex. a) Linear Regression

c) Decision Tree

b) Polynomial Regression

## ② Unsupervised Learning :-

Data is not labelled.

In unsupervised learning correct responses are not provided but instead the algorithm tries to identify similarities within the inputs so that inputs that have something in common are categorized together.

a) Clustering

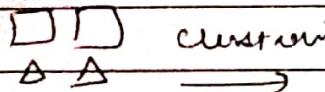
b) Association

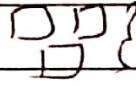
Grouping the data on basis of some similarity  
→ K-means clustering

c) Market Basket Analysis

Bread & Butter & Milk.

a) It is the process of grouping data points into clusters based on their similarities.

 clustering



a. 1) K-means algorithm

a. 2) Mean-shift algorithm



b) It is a technique for discovering relationship between items in a dataset.

Ex. Market basket analysis :- items that people buy together separately frequently. ex., Bread, milk, butter.

b. 1) Apriori Algorithm

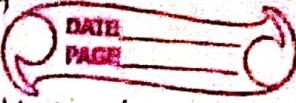
b. 2) FD-growth Algorithm

Scanning we do not get labelled dataset ?

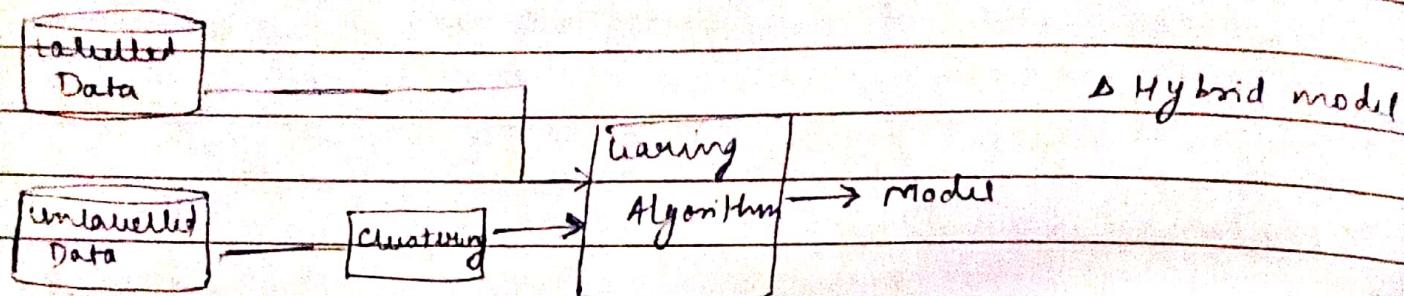
③ Semi-supervised learning Algorithm :-

There are huge collection of unlabelled data & some labelled data. Labelling is a costly process & difficult to perform by the humans.

TPU : matrix multiplication  
in 1 cycle. ( $4 \times 4$ )  
Neural N/W →  
hpc



The main aim of the <sup>semi-</sup> supervised learning is to efficiently use all the available data.

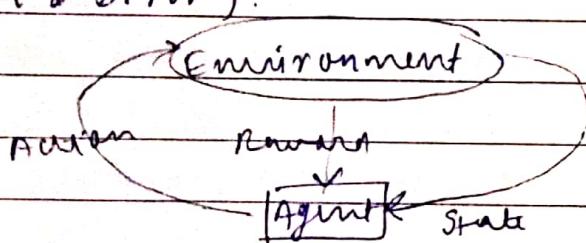


Ex:- Text classification

④ Reinforcement Learning :-

Say driving cars

In it there is no I/P data as we find in supervised or unsupervised learning algo. It is about the learning the optimal behaviour or action in an environment to either obtain maximum reward. A learner is not told what action. but learner will discover which action yields the most reward by trying them. (trial & error).



Algorithm

Ex: Q-learning

- i) Train computer to play chess.
- ii) Train the dog on new tricks

two challenges of Machine Learning :-

① Problems :- all for algorithm

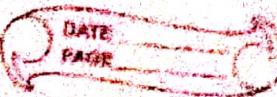
b Task + Performance / Parameter / Categories / names

② Ill posed problem o/p

$\text{cp}(x_1, x_2)$	target	$y = x_1 + x_2$
(1, 1)	1	$y = x_1 \cdot x_2$
(2, 1)	2	$y = \frac{x_1}{x_2}$
(3, 1)	7	$y = \frac{x_1}{x_2} - x_2$
(4, 1)	5	$y = x_1 - x_2$
(5, 1)	3	

Many models on  
same data set.

17.01.2024  
Wednesday



② Huge data - primary requirement

↳ quantity, pure

no missing data, incorrect data concept.

③ High performance coding

Big Data      Machines using GPU + TPU

④ Complexity of approximate Algorithm :-

to design

select

evaluate

⑤ Overfitting & underfitting.

A model that too fits the training data correctly but fails to fit test data.

# Application of Machine Learning :-

Movie Recommendation

# Designing a Learning System

to Unknown target function  
 $f: x \rightarrow y$

Training Examples  
 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$h$  is subset of  $H$

Hypothesis set  
 $H$

Learning Algorithm

Final Hypothesis  
 $h$

$h \in H$

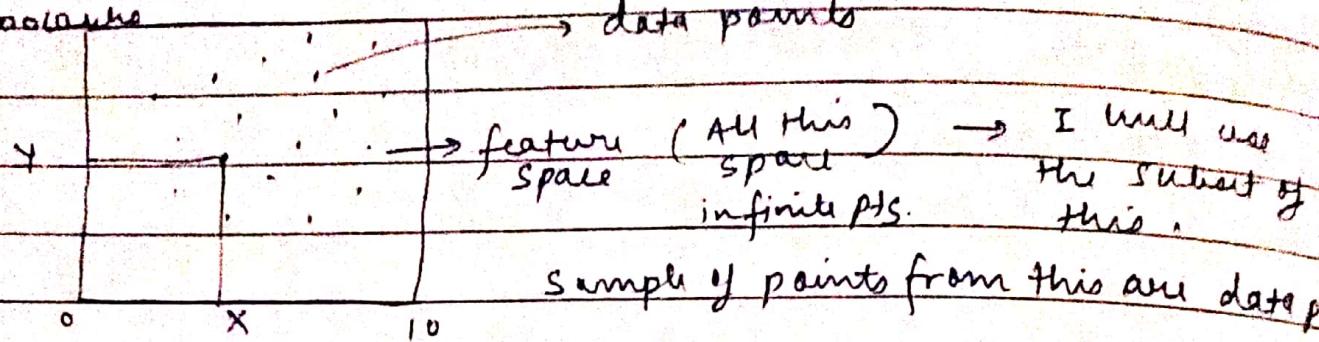
linear

Class, form : User decide the form  $\rightarrow$  Polynomial

Function has {  
2 things      Attributes ; m & c value for linear Regression  
(after op we decide)

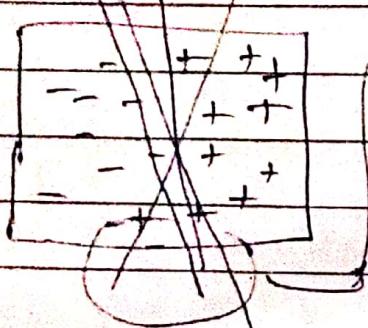
Feature vector ;  $(x_1, x_2, x_3, x_4)$

variables  $\rightarrow$  data points



sample of points from this are data points

Classification problem

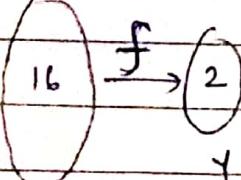


Hypothesis space : Set of all possible hypothesis  
 Why hypothesis?  $\rightarrow$  Because it is not proved yet

Calculation of size of hypothesis space :-

Suppose  $x_1, x_2, x_3, x_4$  are 4 features & the value of these features is boolean.  $x_1, x_2, x_3, x_4 \in \{0, 1\}$

$$\text{Set 1} \quad \text{Set 2}$$



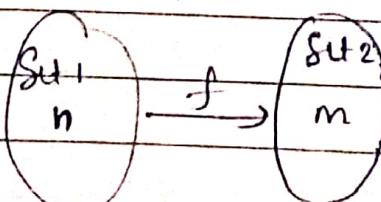
$$\left\{ \begin{array}{l} \text{Possible instances} = 2^4 = 16 \text{ for } x \\ \text{Possible instances} = 2 \end{array} \right. \text{for } y$$

$$2^{16} \text{ functions}$$

are possible

$=$  Hypothesis space  $\} \rightarrow$  Set 1 after applying bias.

$m \& n$  are no. of instances



$$m^n = \text{functions possible}$$

One is 4 instances & other 3 instances &  $y$  has 3 instances

$$3 \times 3 \times 3 \times 4 \times 2 \times \text{instances} \leq 108$$

$$3 = 7 \text{ instances}$$

$$3^{108}$$

~~if f is linear then~~

3.01.2021  
Tuesday



Two types of Bias :- (that can be added to Hypothesis)

- (1) constraints
- (2) polynomial regression for hypothesis

(2) polynomials

→ I am going to select lower order polynomial.

$$y = w_0 + w_1 n + w_2 n^2 + w_3 n^3 + \dots$$

Only this

## LINEAR REGRESSION

① Regression Analysis :- It is the used to model the relationship b/w one or more independent variables  $x_i$  and independent variable  $y_i$

Eg. prices of real estate

area (independent variable) | price (dependent variable)

(reality : city → can be other parameters that defines price.)

	$A_1$	$A_2$	$A_3$	... $A_n$	OP	$y$
1						
2						
3						
$m$						

- (i) No. of features are :  $n$
- (ii) No. of datapoints :  $m$
- (iii)  $i^{th}$  data point  $x_i$   $1 \leq i \leq m$
- (iv) OP associated with  $i^{th}$  datapoint :  $y_i$

(v) predicted OP :  $\hat{y}_i$

(vi) Hypothesis function, weights  
 $h_w(x) = w_0 + w_1 x_i$   
(or  $y = mx + c$ )

In Linear Regression:-

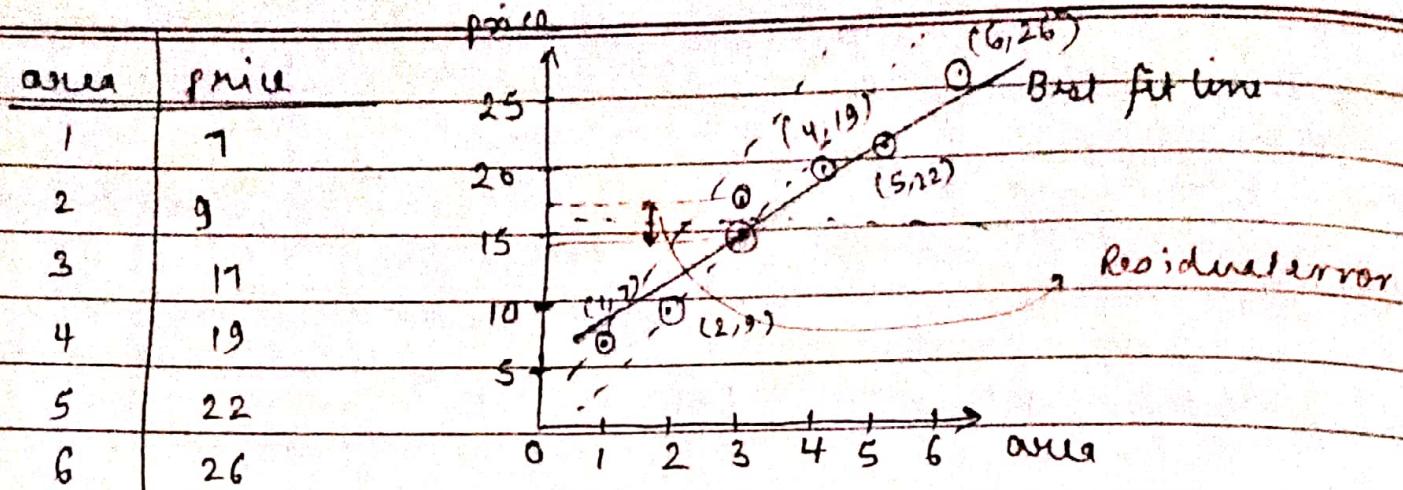
A line is fitted upon given training data for finding

the linear regression relationship

b/w independent & dependent variable.

# Gradient Descent

## L3 Basis of Neural Net & dep. learning models.



Residual error :- Difference in actual value and predicted value  
of  $y = y_i - \hat{y}_i$

Actual error by given line  $= \sum \text{Residual error} = \sum (y_i - \hat{y}_i)$

Linear Regression uses mean square error as loss function / cost function.

loss function / Cost function  
(most used)

It's differentiation is done mostly, we put extra  $\frac{1}{2}$  in formula to make calculation easier.

$$\text{Let } J = \frac{1}{2m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

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

Mean square error

Evaluation metrics

MSE

MAE (Mean Absolute error)

R<sup>2</sup> error score

RMSE Root mean square error

$$h_w(x) = w_0 + w_1 x$$

$$\text{Let } w_0 = 0$$

$$h_w(x) = w_1 x = \hat{y}$$

$$\begin{matrix} x \\ 1 \end{matrix}$$

$$\therefore J = \frac{1}{2m} \sum_{i=1}^m (y_i - w_1 x_i)^2$$

$$\begin{matrix} x \\ 2 \end{matrix}$$

$$w_1 = 0.5$$

$$\begin{matrix} x \\ 3 \end{matrix}$$

$$w_1 = 1$$

$$\begin{matrix} x \\ 4 \end{matrix}$$

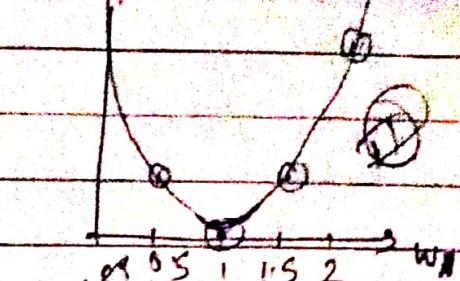
$$w_1 = 1.5$$

Training data  
Normal

$$w_1 = 2$$

(D) Non convex function can have multiple minima but ~~convex~~ convex parity has only one minimum

	$w_0 = 0.5$	$w_1 = 1$	$w_2 = 1.5$	$w_3 = 2$	$J$
1	1	0.5	0	-0.5	-1
2	2	1	0	-1	-2
3	3	1.5	0	-1.5	-3
4	4	2	0	-2	-4
$J$		$0.25$	$\overline{0}$	$0.9$	$3.75$
$J = \frac{1}{2} \ w - w^*\ ^2$		$2.25$	$7.50$		$\frac{30}{8} = 3.75$
$w^*$		$7.50$	$8$	$0.9$	$\frac{16}{30}$ Sorry !



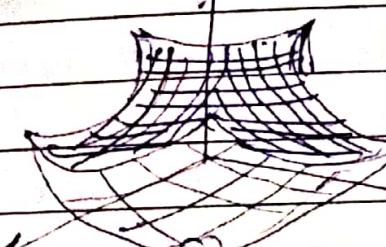
Convex function

[Same can be drawn for no.]

$$w_{new} = w_{old} - \eta \frac{dJ}{dw}$$

{ When you are away  
then  $dJ$  is big }

{ as you come nearer  $dJ$  becomes small. }  $\frac{dJ}{dw}$



3D Convex function

for which  $w_0$  &  $w_1$ ,  $J$  value is minimum

→ Aim

give the most optimal set for Cost function

Decide Stopping Criteria → min. error permissible

$$\text{ideal } \frac{dJ}{dw} = 0$$

→ number of iterations

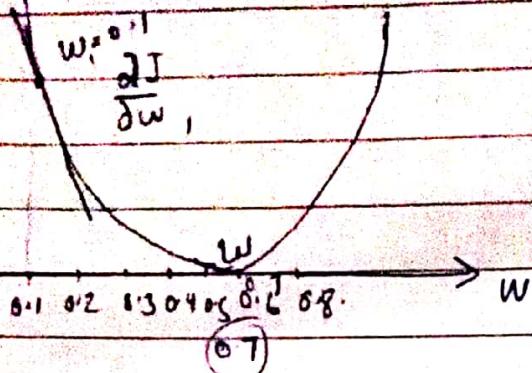
Convex vs Non convex f → (D)  
Global minima vs Local minima

## Linear Regression using Gradient Descent.

- It is an optimization technique (Error minimization)
- It is an iterative algorithm for finding a local minimum of a differentiable f.

$$w(\text{new}) = w(\text{old}) - \alpha \left( \frac{\partial J}{\partial w_i} \right)$$

for  $w_i \quad \frac{\partial J}{\partial w_i} < 0$



$$w_i(n) = w_i(0) + \alpha \left| \frac{\partial J}{\partial w_i} \right|$$

(67)

### Gradient Descent Algorithms (GDA):-

Step.1: choose hyperparameters ( $\alpha, E$ ) values

$\alpha$  = Learning rate       $E$  = stopping criteria

Step.2: Initialize  $w$ 's values.

Step.3:  $w(\text{new}) = w(\text{old}) - \alpha \frac{\partial J}{\partial w}$

Step.4: check stopping criterion ' $E$ '. If satisfied stop else repeat from step. (3)

$$J = \frac{1}{2m} \sum_{i=1}^m (y_i - h_w(x_i))^2 \Rightarrow J = \frac{1}{2m} \sum_{i=1}^m (y_i - (w_0 + w_1 x_i))^2$$

$$\frac{\partial J}{\partial w_0} = \frac{1}{2m} \sum_{i=1}^m 2 \times (y_i - (w_0 + w_1 x_i))^2 (-1)$$

$$\Rightarrow \frac{\partial J}{\partial w_0} = \frac{1}{m} \sum_{i=1}^m (w_0 + w_1 x_i - y_i) \quad \text{or} \quad \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)$$

$$\frac{\partial J}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m (y_i - (w_0 + w_1 x_i)) \times (-w_1) \quad \text{or} \quad \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) w_1$$

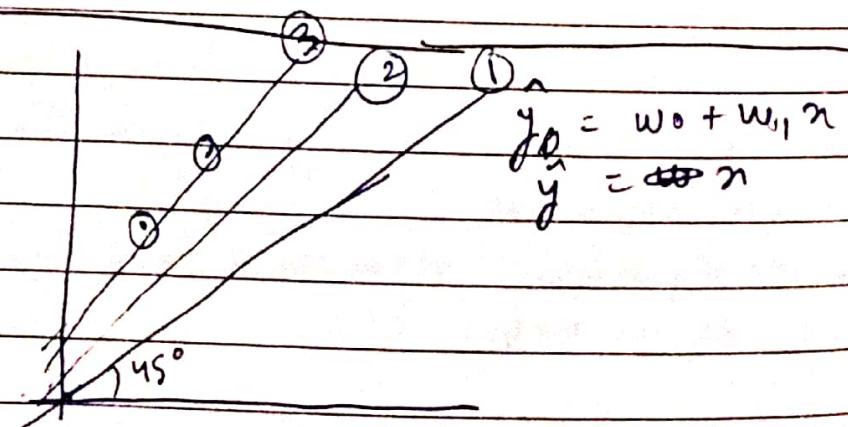
x	y
2	5
3	8

So:-

GDA for 2 iterations

$$w_0 = 0, w_1 = 1$$

$$\alpha = 0.01$$



$$\hat{y}_0 = w_0 + w_1 x$$

$$\hat{y} = \text{constant}$$

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

Iteration - 1

$\hat{y}_i = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) = \frac{1}{2} \{ (2 - 5) + (3 - 8) \}$   
 $= \frac{1}{2} \{ (-3) + (-5) \} = -4$

$\frac{\partial J}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m (-4 \times 1) = -4 \times \frac{1}{2} (-3 \times 2 + (-5 \times 3)) = -21 =$

$(w_0)_n = w_0 - \alpha \frac{\partial J}{\partial w_0} = 0.04 - \frac{2}{-10.5} = -10.5$

$(w_1)_n = w_1 - \alpha \times \frac{(-10.5)}{2} = 0.04 - 0.01 \times (-10.5) = 0.04 + 0.105 = 0.105$

Iteration - 2

$\hat{y}_i = (w_0)_n + (w_1)_n x_i \quad \rightarrow (2)$

$\hat{y}_0 = 0.04 + 0.105 \times 2$

x	y	$\hat{y}$
2	5	2.25
3	8	3.355

$\frac{\partial J}{\partial w_0} = \frac{1}{2} \{ (2.25 - 5) + (3.355 - 8) \}$   
 $= \frac{1}{2} (-2.75 + -4.645) = -3.6975$

$\frac{\partial J}{\partial w_1} = \frac{1}{2} \{ (-2.75) \times 2 + (-4.645) \times 3 \} = -9.7175$

$w_0 = 0.04 - 0.01 \times (-3.6975) = 0.0769$

$w_1 = 0.04 - 0.01 \times (-9.7175) = 1.137$

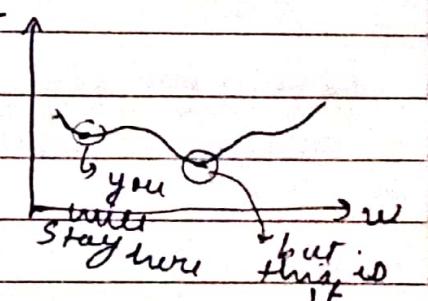
$\hat{y} = 0.0769 + 1.137 x \quad \rightarrow$  Line is becoming better fit becoming

(3)

### Issues in GDA :-

#### ① Effect of loss function on GDA

In case of non-convex  $f^n$  GDA can stuck on local minima



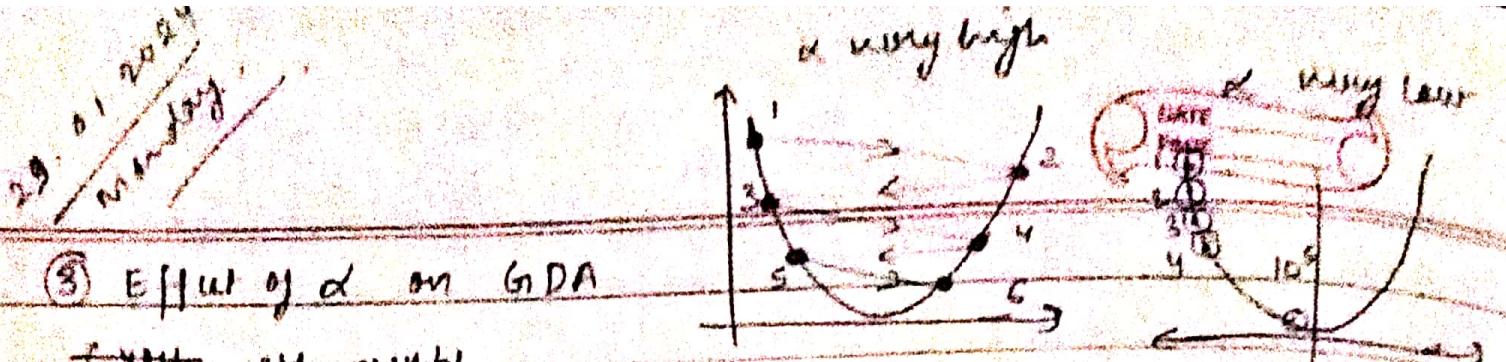
#### ② Effect of data on GDA

A1	A2	A3	Y
0.1 0.5	1 1.05	1 1.09	

Range of attributes is different

GDA will take more time to converge

$\therefore$  Scaling of values is done.



### ③ Effect of $\alpha$ on GDA

If  $\alpha$  is very small,

$\alpha$  is high, you will not reach minima.

$\alpha$  is very low, you will take much more time.

(Take various  $\alpha$  values, generally  $\alpha = 0.1$  or  $0.01$ )  
Hit & trial

## MULTIPLE LINEAR REGRESSION

We have multiple independent variables.

The basic assumptions of MLR:-

① Independent variables are not highly correlated. & hence multicollinearity problem does not exist.

ex: Height & weight in database.  $\rightarrow \sum |y_i - \hat{y}_i|$

② It is assumed that residuals are normally distributed

$$x_1, x_2, \dots, Y \quad Y = f(x_1, x_2)$$

$$Y = w_0 + w_1 x_1 + w_2 x_2$$

For  $n$  independent variable

$$Y = f(x_1, x_2, x_3, \dots, x_n)$$

$$y_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_n x_{in}$$

$\hookrightarrow$  (value of  $y_i$  at  $i$ th row)

## Linear Regression in matrix form

Ordinary least square approach. (OLS)

$$\hat{Y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1} \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix}_{n \times n}$$

$\boxed{C_{m \times n}}$   $m \times (n+1)$

residual =  $y_i - \hat{y}_i$



$$Y = XW$$

$$\begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad (n+1) \times 1$$

### RULES OF ALGEBRA

- (1)  $A \cdot A^T = a^2 + b^2 + c^2 + \dots$
- (2)  $(AB)^T = B^T A^T$
- (3)  $A^T B = B^T A$
- (4)  $d X^T X = 2X$

$$\begin{bmatrix} a & b & c \end{bmatrix} \begin{bmatrix} a & b & c \end{bmatrix}^T = a^2 + b^2 + c^2$$

$$J = \sum_{i=1}^m (y_i - \hat{y}_i)^2 \quad \hat{y}_i: \text{predicted}$$

$$= \sum_{i=1}^m (Y - XW)(Y - XW)^T \quad [\text{From (4)}]$$

$$= Y(Y - XW)^T - (XW)^T(Y - XW)$$

$$= Y^T Y - Y^T X W - \cancel{X^T} (XW)^T Y - (XW)^T X W$$

$$\therefore J = Y^T Y - 2Y^T X W - W^T X^T X W \quad [A^T B = B^T A^T]$$

$$\therefore \frac{dJ}{dW} = 0 - 2Y^T X - 2X^T X W = 0 \quad [\text{Using (4)}]$$

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

~~④~~ Multiply both ~~④~~  $(X^T X)^{-1}$

$$\therefore W = (X^T X)^{-1} (Y^T X) \quad \text{④} \quad \text{⑤}$$

or

$$\boxed{W = (X^T X)^{-1} (X^T Y)} \quad \text{④} \quad \text{⑤} \quad \text{⑥}$$

$$[A^T B = B^T A]$$

- Q) Find Linear Regression of the data of week 4 product sales (in thousands given below)

Chaitin  
Information  
and Communication

DATE  
PAGE

X	y
1	1
2	3
3	4
4	8

Use LR  
in matrix  
form to  
find the  
Hypothesis

$$\text{② } n_w(x_i) = w_0 + w_1 x_i$$

Sol:

$$y = \begin{bmatrix} 1 \\ 3 \\ 4 \\ 8 \end{bmatrix} \quad x = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}$$

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

$$x^T x = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

$$\text{adj}(x^T x) = \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix} \quad (x^T x)^{-1} = \frac{1}{120 - 100} = 20$$

$$(x^T x)^{-1} = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix}$$

$$(x^T x^{-1}) x^T = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 3 & 4 \end{bmatrix}$$

$$x^T y = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 4 \\ 8 \end{bmatrix} =$$

OLS is cost  
computationally  
costly operation

$\therefore$  Gradient  
Descent is  
better.

$$x^T y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 4 \\ 8 \end{bmatrix} = \begin{bmatrix} 16 \\ 51 \end{bmatrix}$$

$$(x^T x)^{-1} x^T y = \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} \begin{bmatrix} 16 \\ 51 \end{bmatrix} = \begin{bmatrix} -1.5 \\ 2.2 \end{bmatrix} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} = w$$

## Gradient Descent Method -

(i) Batch Gradient Descent      (ii) mini Batch GD      (iii) Stochastic GD

randomly  
only one data  
point

## Regression Metrics:-

(i) mean absolute error (MAE) =  $\frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$  otherwise sum error + w  
+ we  $\rightarrow$  center

### Advantage of MAE:-

① MAE unit and  $y$  unit is same.

② It is robust to outliers  $\rightarrow$  extreme data points

### Disadvantage of MAE:-

① It is not differentiable at all points

(ii) Mean Squared Error (MSE) =  $\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$

It is the sum of square of residuals.

### Advantage:-



### MSE Graph

① It is differentiable at all the points

### Disadvantage:-

① Unit of MSE and  $y$  are not same.

② Not robust to outliers.

(iii) Root mean Square Error (RMSE) =  $\sqrt{\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2}$

### Advantage:-

① RMSE unit is same  $y$  unit

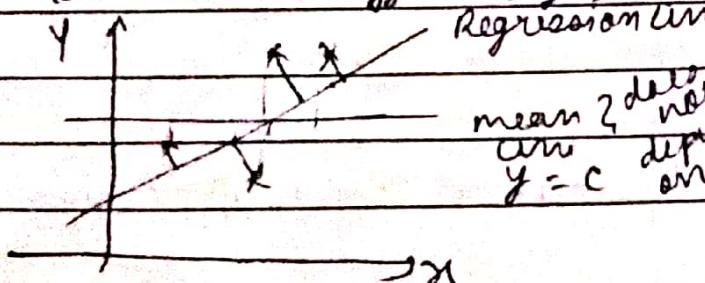
### Disadvantage:-

② not much robust to outliers.

∴ Matrix depends on content of problem ?

∴ We use R2-Score

(iv) R2-Score :- coefficient of determination

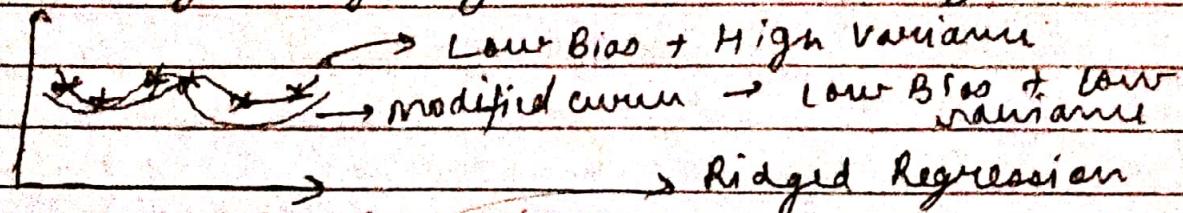


$$R^2 \text{ score} = \frac{RSS}{TSS}$$

RSS = sum of square of residual

TSS = Total sum of squares  
(It is sum of all squares of differences between actual value & overall mean  $\bar{y}$ .)

**Regularizations:** Regularization technique used to reduce the errors by fitting the function approximately the given training set and avoiding overfitting. → appropriately



3 techniques for regularization :- → square / power of 2

① Ridge Regularization (L2 Regularization)  $w \neq 0$

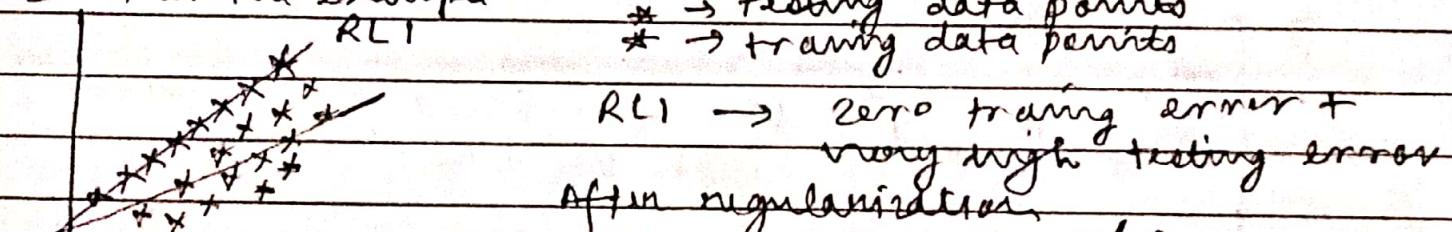
② Lasso Regularization (L1 Regularization)  $w = 0$

③ Elastic Net Regularization (L1 & L2 Regularization)

$$\hookrightarrow = \text{Ridge} + \text{Lasso}$$

① Ridge Regularization (L2 Regularization)

Consider the example -



RL1 → zero training error + very high testing error

After regularization,

Modified cost function (in OLS)

$$J = \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda \sum_{i=1}^m w_i^2$$

↑ can be done gradient descent?

↑ penalty of Ridge Regularization + bias

$$J = \sum_{i=1}^m (y_i - \hat{y}_i)^2 + \lambda w_i^2$$

$$\hat{y}_i = w_0 + w_i x_i$$

$$w_0 = \bar{y} - w_1 \bar{x}, \quad w_i = \frac{1}{m} \sum_{i=1}^m (y_i - \bar{y})(x_i - \bar{x})$$

$$\cancel{\star} \cancel{\star} \cancel{\star} \sum_{i=1}^m (x_i - \bar{x})^2 + \lambda$$

$$\boxed{J = 0 \Rightarrow \text{ordinary least square}} \quad \lambda \rightarrow 0 \text{ but } w_i \neq 0 \quad [\lambda = 0 + \infty]$$

$\lambda \rightarrow \infty \Rightarrow w_i \downarrow \infty$  ( $\rightarrow 0$ ) but never becomes 0.

Reducing importance of  $w_i$ . (smoothing of curve, reducing sharpness from curve)

for  $\lambda = 0 \rightarrow N$ : tuning  $\lambda$  to get better results?

Reduce complexity & removing overfitting

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

$R^2 < 0 \Rightarrow$  mean line  
near zero  
near zero  
Regression line

## LASSO regularization (Least Absolute shrinkage & selection operator)

In ridge weights will never  $\equiv$  equals 0

while in LASSO it can be zero. (absolute function)

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

$w_1 = 0$

$L^1$

No effect of  $x_1$  on  $y$

This new  $\Rightarrow$  Regularization.

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

For very high value of  $\lambda$ , some weights will become zero.

How some weights will become zero in Lasso?

$$\hat{y}_i = w_0 + w_i x_i$$

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

$$w_0 = \bar{y} - w_i \bar{x}$$

$$w_i = \frac{\sum_{i=1}^m (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

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

→ Not differentiable

When  $w_i > 0$

$$w_i = \frac{\sum_{i=1}^m (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$\left\{ \begin{array}{l} \text{When } \lambda = \sum_{i=1}^m (y_i - \bar{y}) \\ N=0 \\ \therefore w_i = 0 \end{array} \right.$$

$w_i < 0$

$$w_i = \frac{\sum_{i=1}^m (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$\left\{ \begin{array}{l} \text{Same } \lambda - w_i \\ \text{for } w_i < 0 \end{array} \right.$$

Shrinking at Normally  $|w_i|$  +

use Lasso

(L1 better than L2 in this case)

Dataset is too large in terms of no. of attributes.  
~~If there is multicollinearity (dependency)~~

## ELASTIC NET

Elastic Net is a hybrid method of combining both Ridge and Lasso regression

Method:

$$J = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \underbrace{\lambda_1 \sum_{i=1}^m w_i^2}_{\text{Ridge}} + \underbrace{\lambda_2 \sum_{i=1}^m |w_i|}_{\text{Lasso}}$$

When  $\lambda_1$  &  $\lambda_2 = 0$  both then it will work as OLS

But  $\lambda_1 = 0.5$  &  $\lambda_2 = 0.5$  then turn it

- If  $\lambda_1 = 0$  Then O. it is Lasso.

- If  $\lambda_2 = 0$  Then it is Ridge.