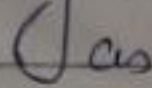


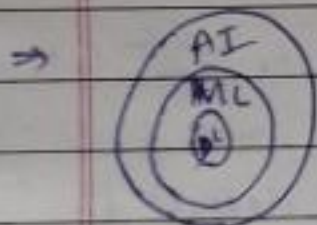
Q. Form

Machine Learning :- It is a boom in computer world. 

Machine Learning is a set of data which make a machine to perform some specific task by using best suited algorithms.

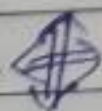
In today's world, we found machine with limited task having.

ML is AI which provide machines the ability to learn the action.



Machine Learning :- "The field of study that gives computer the ability to learn without being explicitly programmed."

A computer program is said to learn from experience E, with respect to some task T, and some performance measure P, if its performance on T, as measured by P, improves with experience E.

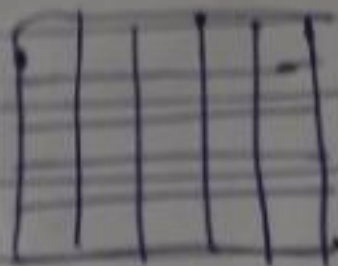


Machine Learning :-

- Well-defined Learning Problem :-
- The class of tasks,
- The measure of performance to be improved
- The source of experience.

Example: 1. checker's hearing problem:

Find: 3-
T → To move ^{/plan} the checker

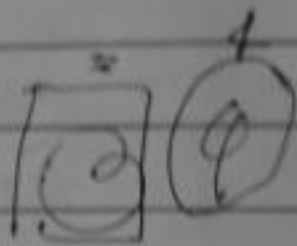


E → where the move happen previously

P → win or loose

Find:

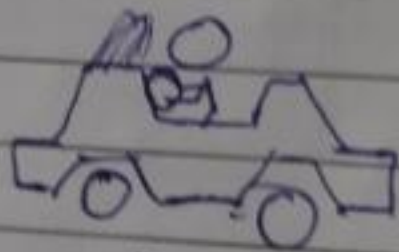
T → to detect the digit



E → How ~~to~~ previously digits ^{lost} are detected.
→ a database of hand written words with given class

P → ^{percent} How much correct digit detected ^{fraction}

Find

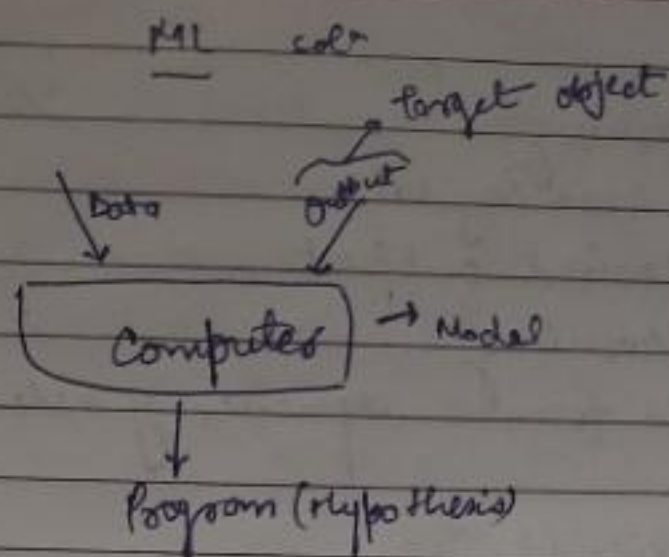


T → ~~How~~ robot drive the car "driving a car"

E → from dataset of how to drive a car or
how many time robot drive previously.

P → percent of driving vehical successfully.

A sequence of image & steering commands recorded while observing a human driver.



→ Features to identify objects :-

- * Shape
- * Color
- * texture

* Supervised → when we have a target value.

* Unsupervised → when the result is uncertain.
i.e. no fix target.

⇒ Where ML use :-

- when human expertise does not exist.
- when we work with big data.
- Model must be customized.

Supervised (inductive) learning

- Given training data + desired output.
- desired outputs (labels)

Unsupervised learning

- Given training data
- without desired output.

Semi supervised learning

- Given: training data
- few desired output.

Reinforcement learning

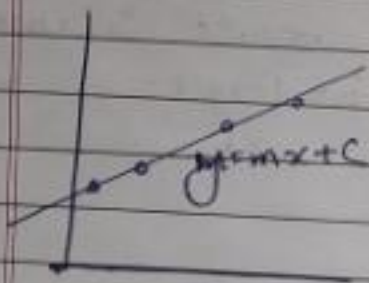
- Rewards from sequence of actions.

~~X~~ → REGRESSION {Supervised}

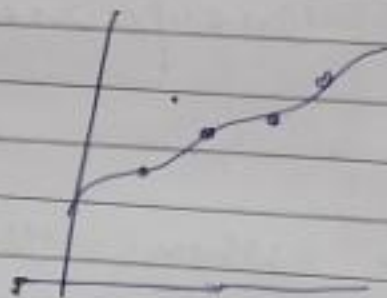
for x & y continuous values
to get prediction
→ Given $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
Learn a function $f(x)$ to predict y given x
 $x-y$ is real valued \therefore regression.

$$y = f(x)$$

to find without error.



Linear Regression



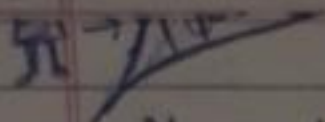
Polynomial regression

→ Regression gives continuous value as a target value.

Classification

for discrete value
to get prediction
→ Learn a function $f(x)$ to predict y given x
 $x-y$ is categorical \therefore classification
→ Regression gives discrete value as a target value.

Input $x \rightarrow$ []
Output $y \rightarrow$ []
Output 2 \rightarrow []
Output 3 \rightarrow []
Output 4 \rightarrow []

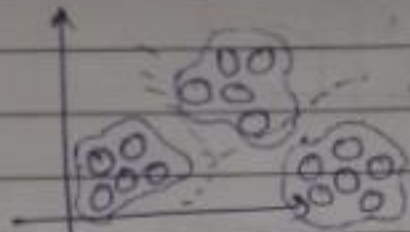


It depend on target value.

* learnbesttech.blogspot.com

Unsupervised learning

- Given x_1, x_2, \dots, x_n (without labels)
- Output hidden structure behind the x_i 's
 - Clustering (unsupervised)

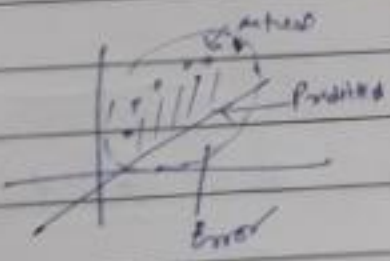


$x_i \rightarrow$ Autonomous data

Reinforcement learning

- Agent do action on rewards
- Agent also got penalty for wrong work.

\Rightarrow Bias :- Difference b/w predicted value & actual value i.e. Error.

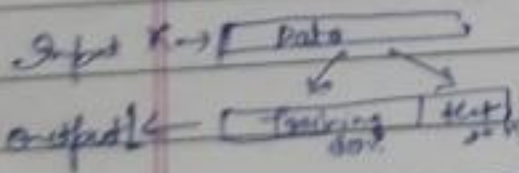


High bias because error is underfit higher.



Model Underfit \rightarrow Sum of Error is high...

Model Overfit \rightarrow Sum of Error is low...



Output 1: Train 40%, Test 60%

Output 2: Train 60%, Test 40%

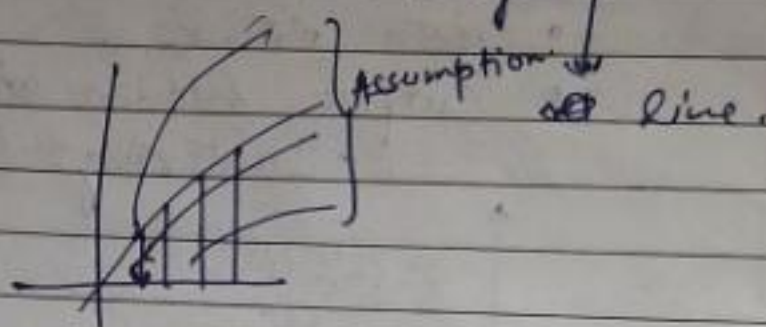
Output 3: Train 80%, Test 20%

Output 4: Train 90%, Test 10%

babba

Variance → difference in results.

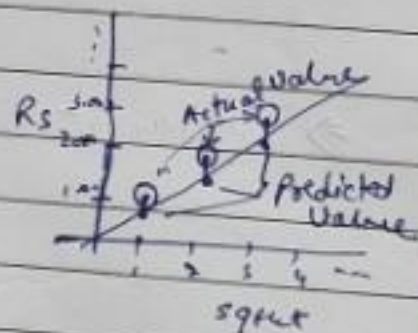
Inductive bias → Set of Assumption to predict the target.



Choose a model (Assumptions)

x	y	Model assumption
2	3.4	1- $y = bx + a$
3	5.9	2- $y = e^{-bx}$
5	7.8	3- $y = \sin(bx)$
7.8	6.5	4- $y = bx^2$
9.2	11.7	5- $y = \sqrt{bx}$
10.4	15.3	
11.8	17.6	

Let plot be of size $1m^2 = 1000R$
 $2m^2 = 2000R$



Inductive bias also called learning bias

$$\text{Inductive bias} = a + bx$$

↓
infinite

Best fit ⇒ Minimum error.

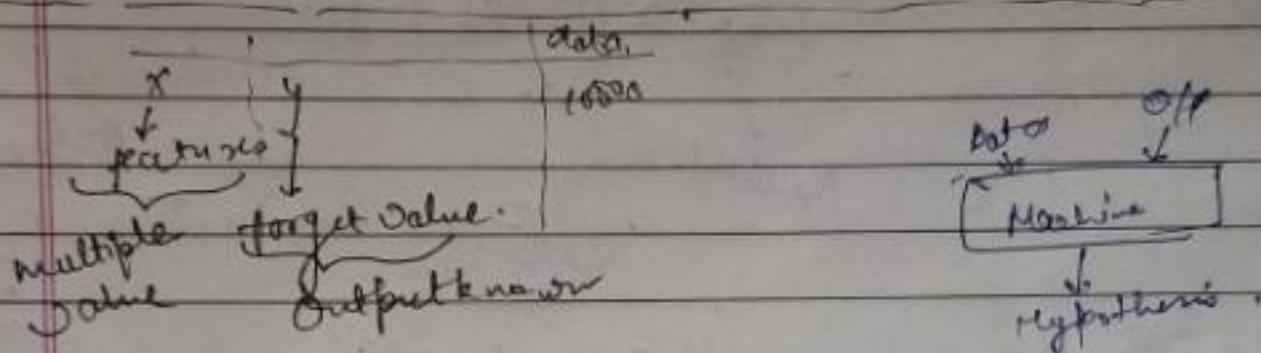
Hypothesis space → when there are no. of spaces set of possible legal hypothesis.



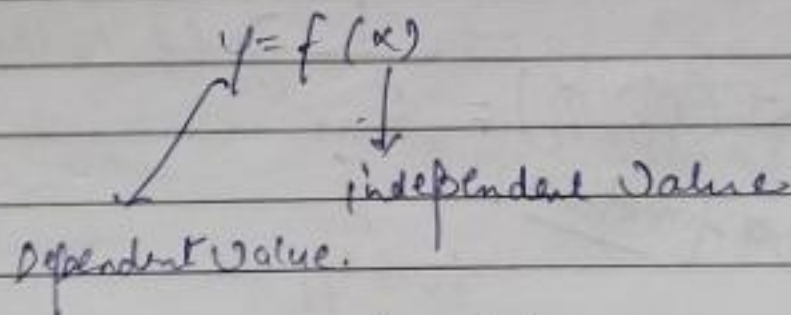
Best solution = hypothesis.

- Bias - error in training data
- Variance - error in test data.

For Best optimal result we should take care for low Bias & low Variance.



LINEAR REGRESSION With ONE VARIABLE →
[finding relationship b/w dependent & independent variable]
best fit line



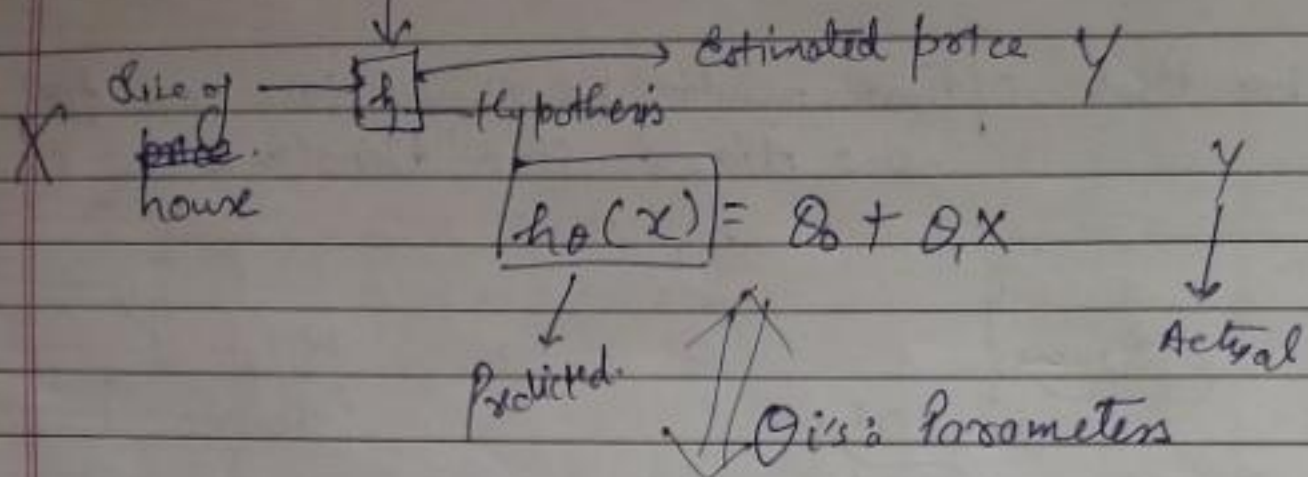
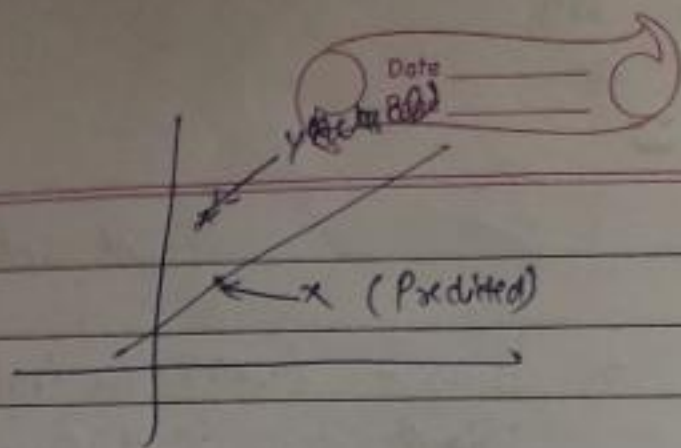
Best fit → Minimum error.

$$\text{ERROR} = \text{Given data (Actual data)} - \text{predicted data.}$$
$$e = y(\text{Actual}) - \hat{y}(\text{predicted})$$

HYPOTHESIS

- m = no. of training examples
- x 's = "input" variable/features
- y 's = "Output variable/targets".

Training set
↓
Learning Algo



$$y = c + mx$$

Cost function $J(\theta_0, \theta_1) \Rightarrow$ Actual - predicted.

or called squared error function.

$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=0}^m (h_0(x^{(i)}) - y^{(i)})^2$

↑ training set

Square make absolute value

Predicted Actual

$$h_0(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

We set as data whose cost size is minimum to get best fit hypothesis.

Single variable

Hypothesis $\rightarrow h_0(m) = \theta_0 + \theta_1 x$

Parameters $\rightarrow \theta_0, \theta_1$

Cost function $\rightarrow J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0 x^{(i)} - y^{(i)})^2$

Goal \rightarrow minimize $J(\theta_0, \theta_1)$

ML Log :-

Training set / Dataset contain features

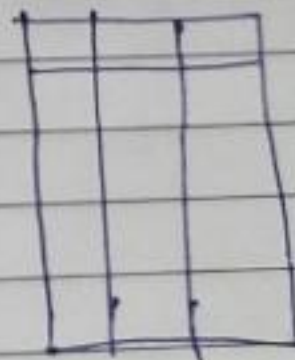
Dataset $\begin{cases} \text{features} \\ \text{Output / label} \end{cases}$

Training set has both features & Output
Testing set has only features & we get outputs
we check accuracy.

\Rightarrow Numpy \rightarrow Library having inbuilt function for array related...

\Rightarrow $\begin{matrix} \text{Library Name} \\ \text{import numpy as np} \end{matrix}$ \rightarrow Object name
 $\begin{matrix} \text{importing library in program} \\ \text{array} \end{matrix}$ $\begin{matrix} \text{Variable Name} \\ \text{a} \end{matrix}$ $\begin{matrix} \text{Type of data structure} \\ \text{array} \end{matrix}$ $\begin{matrix} \text{Data} \\ [1, 2, 3, 4] \end{matrix}$
 $\begin{matrix} \text{array} \\ \text{print(a)} \end{matrix}$ \rightarrow Output command.

ML
 $\begin{cases} \text{Supervised} \\ \text{Unsupervised} \end{cases}$
 \rightarrow KNN
 \rightarrow Random Forest



KNN - { K-Nearest Neighbours }

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

calculate the higher frequency value & return predicted value

→ Kaggle → Use to get dataset :-

→ to read file - `pd.read_csv("C:\\...")`

```
ex> dataset = pd.read_csv("address")
dataset
// Display the data
```

```
ds = dataset.columns → to select column
print(ds)
Index(['PM2.5', ...], dtype='object')
```

outcome

Rain

Sunny

drizzle

fog

0

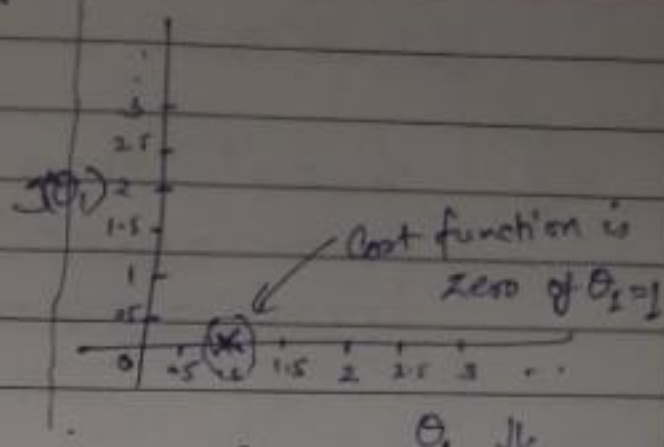
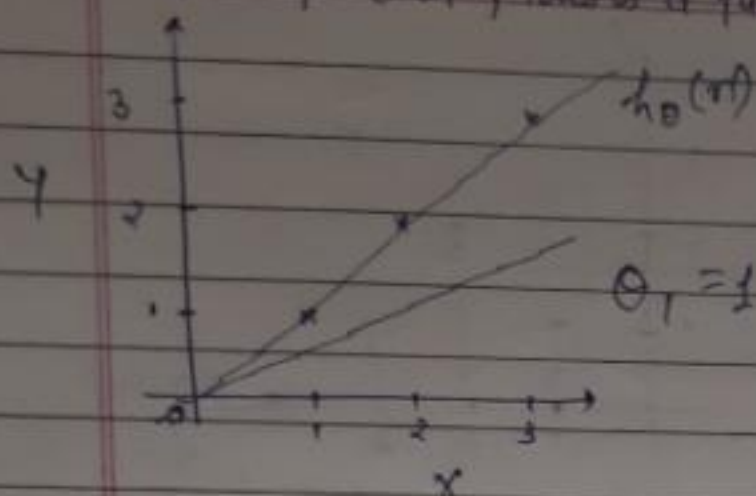
1

2

3

$h_0(x)$
for fixed θ_1 , this is a function of x

J function chart



Cost function: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^i) - y^i)^2$

$\theta_1 \downarrow$
 $J(\theta) = 0$
 $\theta_1 = 1$

$= \frac{1}{2 \times 3} [(1-1)^2 + (2-2)^2 + (3-3)^2]$

$J(\theta) = \frac{1}{6} \times 0 = 0$

Cost function = 0 means best fit

LINEAR

\Rightarrow hypothesis

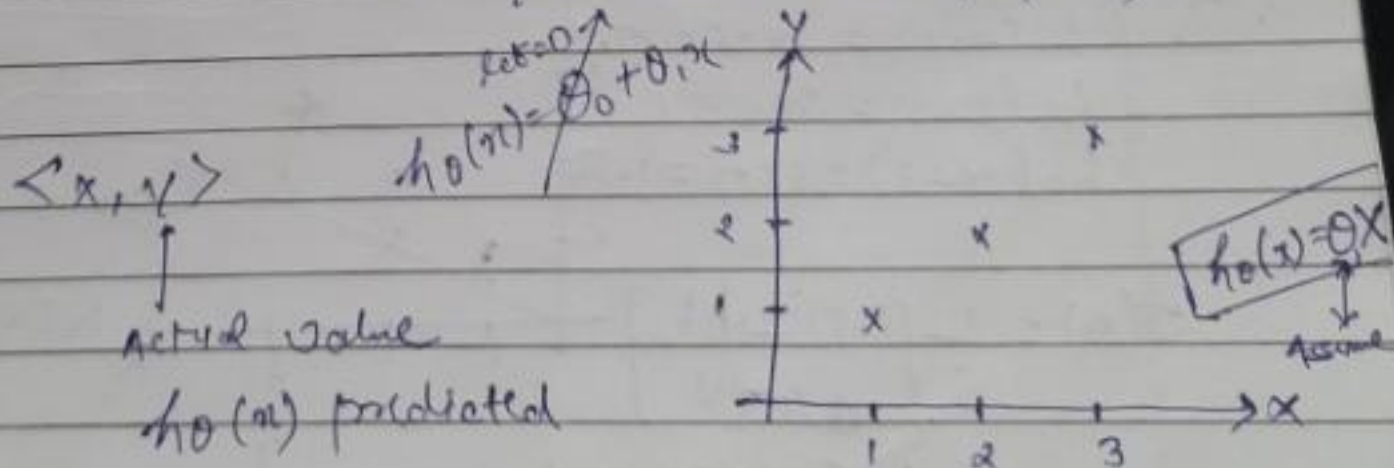
① $h_0(x) = \theta_0 + \theta_1 x$

② θ_i : parameter $i \in \{0, 1\}$

③ Cost function

$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^i) - y^i)^2$

④ Goal Minimize $J(\theta_0, \theta_1)$



* Cost function \rightarrow Error b/w whole data set.

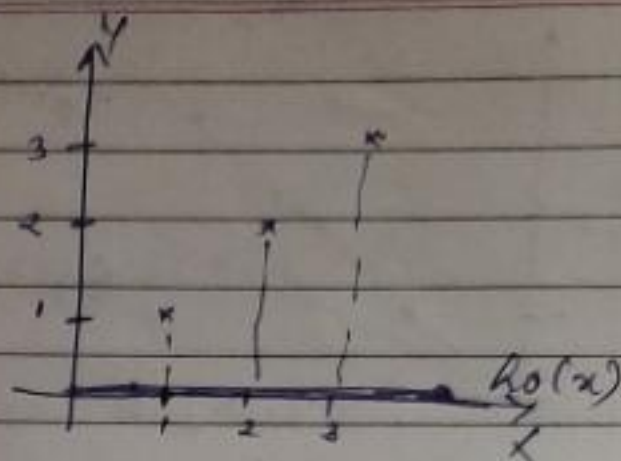
* Loss function \rightarrow Error b/w single data point

Hypothesis 1 for $\theta_1 = 0$

$$h_0(x=1) = 0 \times 1 = 0$$

$$h_0(x=2) = 0 \times 2 = 0$$

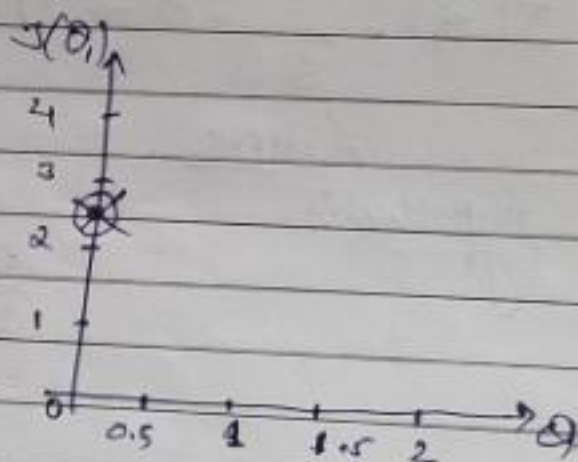
$$h_0(x=3) = 0 \times 3 = 0$$



$$J(\theta_1) = \frac{1}{2 \times 3} \left[(1-0)^2 + (2-0)^2 + (3-0)^2 \right] \quad h_0(x) = \theta_1 x$$

$$J(\theta_1) = \frac{1}{6} \times 14 \Rightarrow \frac{7}{3} = 2.\bar{3}$$

$J(\theta_1)$ vs θ_1



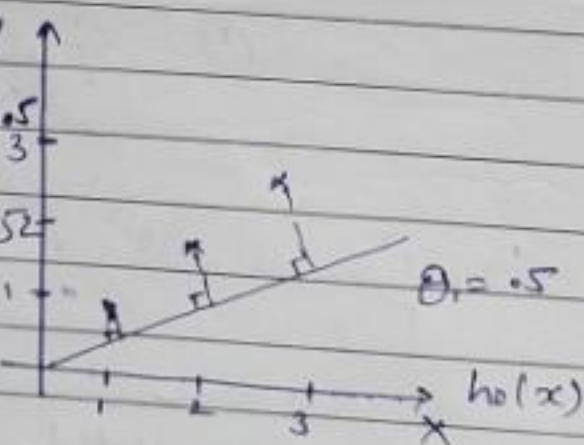
Hypothesis 2

for $\theta_1 = 0.5$

$$h_0(x=1) = 1 \times 0.5 = 0.5$$

$$h_0(x=2) = 2 \times 0.5 = 1$$

$$h_0(x=3) = 3 \times 0.5 = 1.5$$

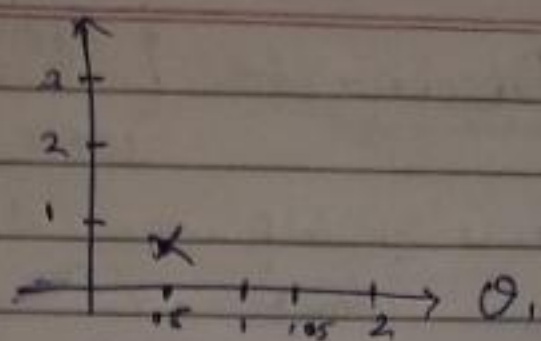


$$J(\theta_1) = \frac{1}{6} \left[(0.5)^2 + (1)^2 + (1.5)^2 \right]$$

$$= \frac{1}{6} [0.25 + 1 + 2.25]$$

$$\Rightarrow \frac{1}{6} \times 3.5 \Rightarrow \frac{3.5}{6} = 0.58\bar{3}$$

$J(\theta_1)$

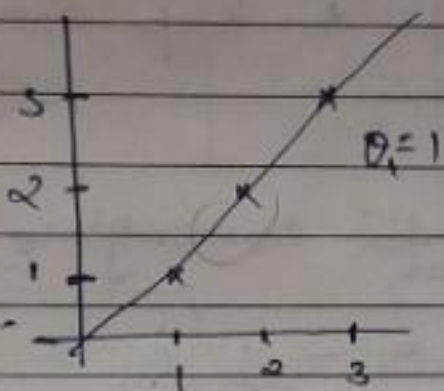


Date _____
Page _____

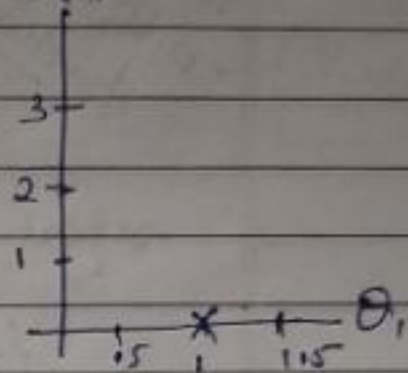
Hypothesis - 3

for $\theta_1 = 1$

$J(\theta_1) = 1$

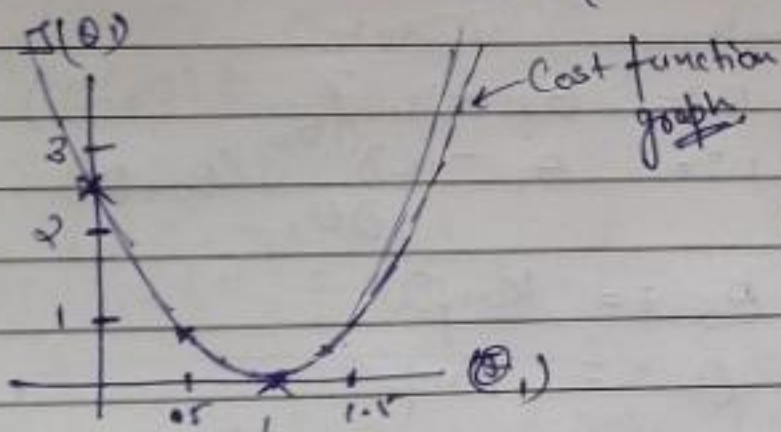


$J(\theta_1)$



Now

Combine Cost function graph:-



→ * the best fit hypothesis is of $\theta_1 = 1$ because it is the most bottom point (min) of θ_1 vs $J(\theta_1)$ graph.

* Also the error is zero.

GRADIENT DESCENT :-

Optimization technique to find out best parameters w.r.t minimum cost function?

- Have some function $J(\theta_0, \theta_1)$
Want min $J(\theta_0, \theta_1)$
 θ_0, θ_1
- Start with θ_0, θ_1
• Keep changing θ_0, θ_1

or for finding best parameters with min. cost function.

down to the slope to reach bottom.

repeat until convergence

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \text{ for } j=0 \text{ to } j=1$$

Learning Rate

correct: Simultaneous update

$$\text{temp } \theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp } \theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp } \theta_0$$

$$\theta_1 := \text{temp } \theta_1$$

- Learning Rate is generally optimal i.e. neither big value nor small value.
- Every step to find bottom line.

eqn of line Hypothesis $h(x) = \theta_0 + \theta_1 x$
for $j=0,1$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

for 3D it is a plane
for 5D it is a plane

LINEAR REGRESSION WITH MULTIPLE VARIABLE →

ND > 3D with multiple plane :- hyper plane

Parameter matrix $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{n-1} \\ \theta_n \end{bmatrix}$ $X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix}$ for all $x_0 = 1$

↑
features

for multiplication change order :-
 $h_0(x) = \theta^T X$

$$h_0(x) = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

X

$h_0(x) = \theta^T \cdot X$

for one features $h_0(x) = [\theta_0 \ \theta_1] \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$

$h_0(x) = \theta_0 + \theta_1 x_1$

$$h_0(x) = \theta^T \cdot X$$

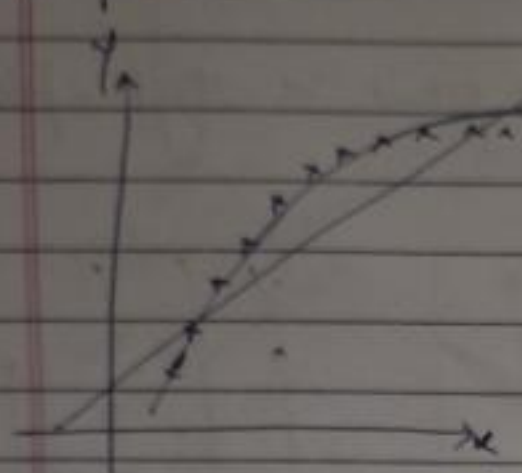
↓
 $\theta^{\text{transpose}}$

$$h_0(x) = \theta_0 + \sum_{n=1}^n (\theta_n x_n)$$

Simple / Single Linear Regression - { Univariate }
Multiple / Multi Linear Regression - { Multivariate }

power more than one ...

POLYNOMIAL REGRESSION :-



$$h_0(x) = \theta_0 + \theta_1 x$$

$$\text{Polynomial} = h_0(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots$$

Way to get polynomial hypothesis.

at least one
→ in hypothesis have power greater than 1.

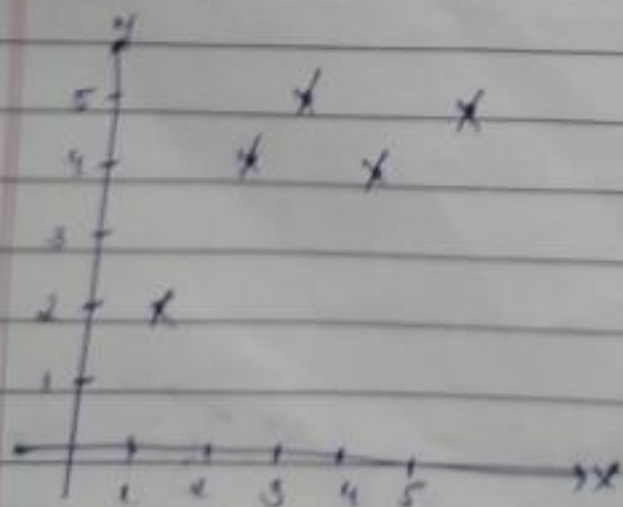
Hypothesis for Polynomial regression.

$$h_0(x) = \theta_0 + \theta_1 x^2 + \dots$$

Ex

For the given point find
best fit. Given $x=6$

x	y
1	2
2	4
3	5
4	4
5	5



Least square Method :-

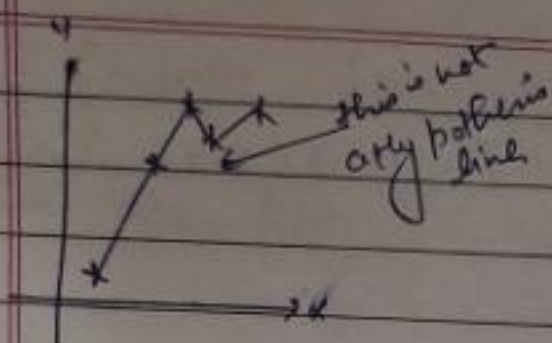
$$h_0(x) = \theta_0 + \theta_1 x$$

$$y = mx + c$$

$$m = \theta_1, c = \theta_0$$

$$m = \frac{y_2 - y_1}{x_2 - x_1}$$

$$c = y_1 - m x_1$$



Now for single line.

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

$$\bar{x} = \frac{\sum x}{n} \quad \bar{y} = \frac{\sum y}{n}$$

$$C = \bar{y} - m \cdot \bar{x}$$

Now calculating
 \bar{x}, \bar{y} & m & C

$$\bar{x} = 3 \quad \bar{y} = 4$$

~~$$m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$~~

$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
-2	-2	4	4
-1	0	0	1
0	1	0	0
1	0	0	1
2	1	2	4
$\sum 0$	$\sum 0$	$\sum 6$	$\sum 10$

$$m = \frac{6}{10} \Rightarrow 0.6$$

$$C = 4 - 0.6 \times 3$$

$$= 4 - 1.8$$

$$C \Rightarrow 2.2$$

$$\boxed{h_0(x) = 0.6x + 2.2} \text{ Hypothesis.}$$

$$h_0(6) = 0.6 \times 6 + 2.2$$

$$\Rightarrow 3.6 + 2.2 = 5.8$$

Linear regression is continuous

Date _____
Page _____

Hours spent in Sub X	Grade Y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
6	82	0.37			
10	88	4.37			
2	56	-3.63			
4	64	-1.63			
6	77	0.37			
7	92	1.37			
0	23	-5.63			
10	41	5.63 -4.63			
8	80	2.37			
5	59	0.63			
3	44	-2.63			

What is the line for the given data & predict the grade for those who spent 9.4 hrs.

Sol Here $\bar{x} = \frac{52}{11}$ $\bar{y} = \frac{709}{11}$
 ≈ 4.72 $= 64.45$

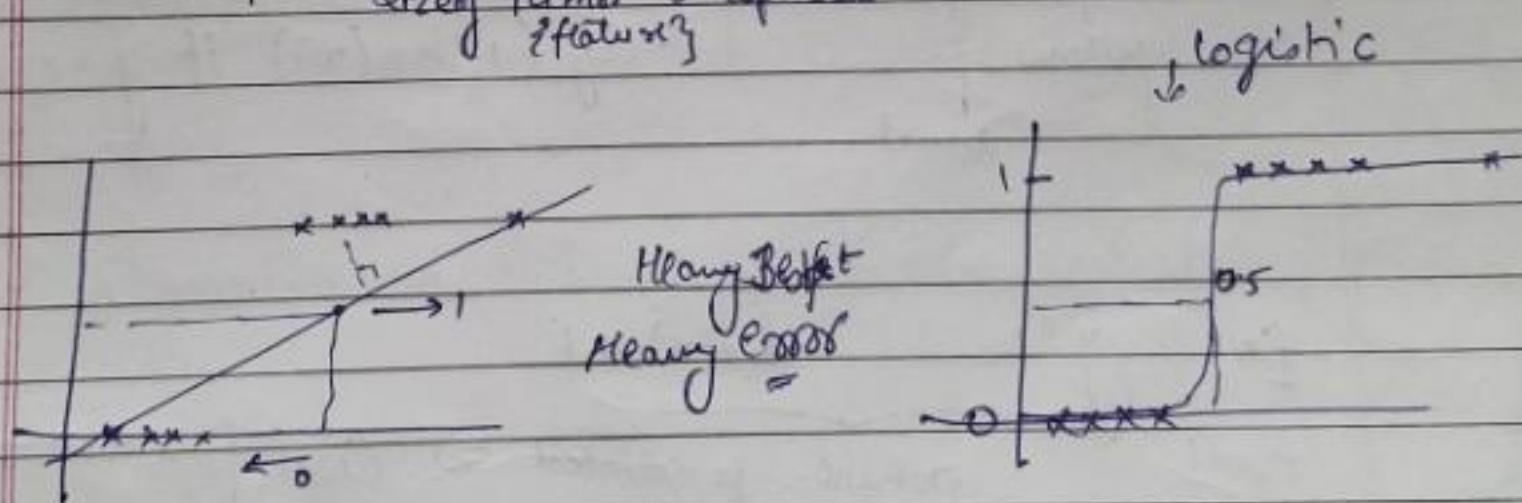
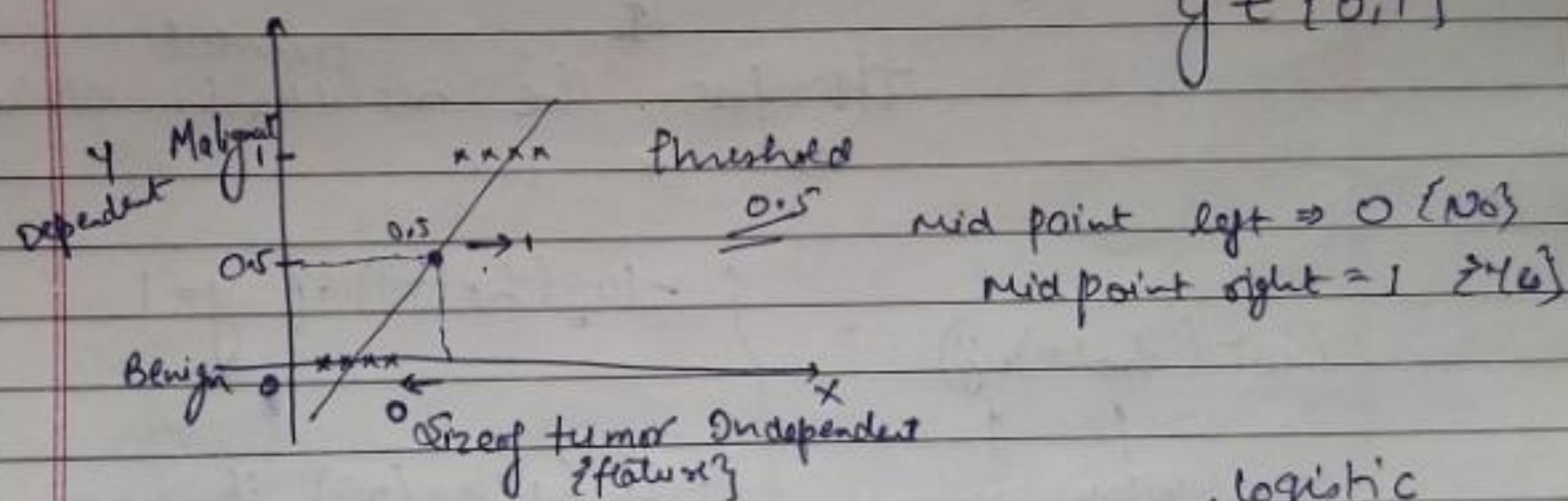
$y = mx + c$
 $y = 6.49x + 33.82$
 $y = 6.49 \times 9 + 33.82$
 $y = 58.41 + 33.82$
 $y = 92.23$

Swing

Date _____
Page _____

Logistic Regression → {Classification technique}

Classification	Spam filtering	Spam/Not spam	yes/no 0/1 Binary
	Fraud detection	Fraud/Not fraud	Binary
	Cancer	Malignant/Benign	Binary $y \in \{0,1\}$



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

Hypothesis of logistic regression

$$h_0(x) = g(\theta^T x)$$

$$1 \geq h_0(x) \geq 0$$

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

Cost function

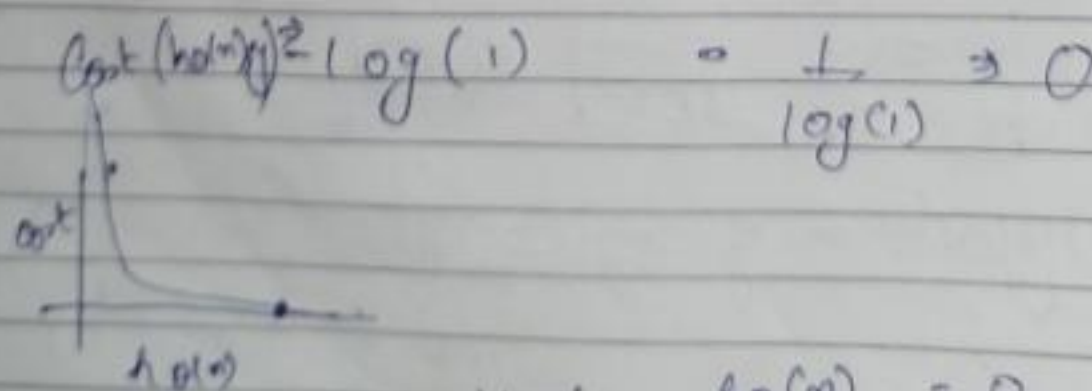
Here for result there is only two choice either 0 or 1

therefore the ^{predicted} result is either 0 or 1.

$$\text{Cost}(h_{\theta}(n), y) = \begin{cases} -\log(h_{\theta}(n)) & \text{if } y=1 \\ -\log(1-h_{\theta}(n)) & \text{if } y=0. \end{cases}$$

\uparrow Predicted \uparrow Actual

for $y=1$ $h_{\theta}(n)=1$
 $\text{Cost} = -\log(1) \Rightarrow \text{Cost function} = 0$
 actual = predicted \Rightarrow Cost function = 0



for $y=0$ $h_{\theta}(n)=0$
 $\text{Cost} = -\log(0) = \infty$

Cost

$h_{\theta}(n)$

Date _____
Page _____

$$\text{Cost}(h_0(x), y) = -y \log(h_0(x)) - (1-y) \log(1-h_0(x))$$

for $y=1$ putting $y=1$ in

$$\text{Cost}(h_0(x)) = -\log(h_0(x)) - 0$$

for $y=0$ putting $y=0$ in

$$\text{Cost}(h_0(x)) = 0 - (1) \log(1-h_0(x))$$

$$\text{Cost}(h_0(x)) = -\log(1-h_0(x))$$

⇒ error is high when we get opposite result.

Logistic Regression → Because we use the concept of line $h_0(x)$
 Classification → Because result is either 0 or 1.

GRADIENT DESCENT

$$J(\theta) = \frac{1}{m} \left[\sum_{i=1}^m y_i \log(h_0(x^i)) - (1-y_i) \log(1-h_0(x^i)) \right]$$

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

where
 $J(\theta) = \text{Cost function}$

use to evaluate
Confusion Matrix \rightarrow for classification model
Size of matrix \propto depend on the no. of classes

ex: $n=2$

dealing with 2 class
Binary classification

above 2 \rightarrow multiclass

\rightarrow class is 2 means matrix size 2×2 .

T \rightarrow true
F \rightarrow false
N \rightarrow -ve
P \rightarrow +ve

to evaluate
classification
model.
we have to
make
confusion
matrix

Predicted	Actual	
	1 \rightarrow P	0 \rightarrow F
+ve	TP	FP
-ve	FN	TN

\rightarrow False positive
also called
type I error

type II error

① Accuracy $\rightarrow \frac{TP + TN}{TP + TN + FP + FN}$

② Precision \rightarrow Actual positive from total true.
$$\frac{TP}{TP + FP}$$

③ Recall \rightarrow Actual predicted from total
$$\frac{TP}{TP + FN}$$

④ F1 score → Harmonic mean of precision & recall.

$$= \frac{1}{\left(\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}\right)} \Rightarrow \frac{2 \times \text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

Ex) For the given values of cell in matrix, calculate all 4 evaluating parameters.

	0	1	2
0	30	10	0
1	0	30	0
2	0	0	930

Accuracy = $\frac{30+30+930}{1000} = 0.96$

TP	FP
30	30
FN	TN
10	930

There is a huge variance b/w Accuracy & precision because of unbalanced dataset.

30 → 930

$$A = \frac{960}{1000} \Rightarrow 0.96 \Rightarrow 96\%$$

$$P = \frac{30}{60} \Rightarrow 0.5 \Rightarrow 50\%$$

$$R = \frac{30}{40} \Rightarrow 0.75$$

Precision → when we check false positive (High FP value)

$$F1 = \frac{2 \times 0.5 \times 0.75}{0.5 + 0.75} = \frac{0.75}{1.25} = 0.6$$

Recall → when we check FN (High FN value)

→ For the medical purpose we check recall value...

KNN - { K-Nearest Neighbors Algorithm }

• supervised classification
• similarity or distance function

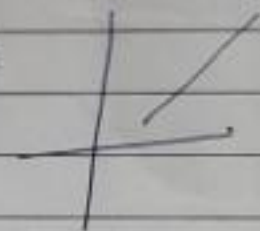
→ stores all available cases & classifies new cases based on a similarity measure (distance function).

→ Generally used for statistical estimation & pattern recognition.

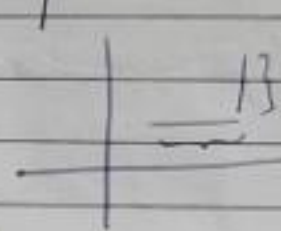
⇒ if $K=1$ then assigned to neighbor with minimum distance.

D
ST
AN
CE
FUNCTION


Euclidean function $\Rightarrow \sqrt{\sum_{i=1}^K (x_i - y_i)^2}$



Manhattan $\Rightarrow \sum_{i=1}^K |x_i - y_i|$



Minkowski $\left(\sum_{i=1}^K (|x_i - y_i|)^q \right)^{1/q}$

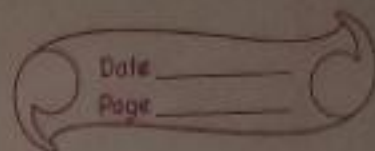


for c points are maxi
in some distance
B go to c.

- Day . . .
- Calculate distance
 - Sort
 - Get top K
 - Get most frequent
 - Return predicted class.



$$\frac{140^2 + 160^2}{160 + 140}$$



Q1

for feature-1

Gender	Height (cm)	Distance	Rank	Sorted
M	178	3	2	M-1
M	174	1	1	M-3
F	163	12	8	F-4
F	168	7	6	5
M	181	6	5	6
F	170	5	4	7
M	184	9	7	9
F	171	4	3	12

$k=3$

Predict class of Gender / person whose height 175 where $k=3$

Distance matrix used is Euclidean & find the distance with other

for $k=3$ the sorted one is

M - 174 - 1

M - 178 - 3

F - 171 - 4

Since Male is in higher frequency therefore

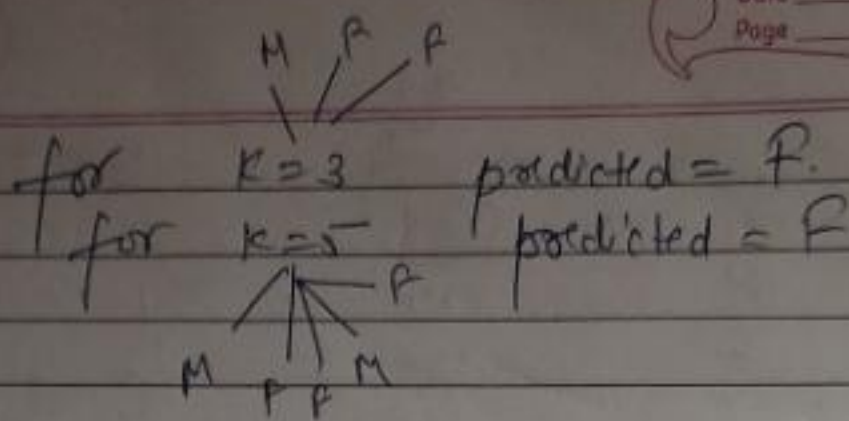
the Gender of person is Male 'M'.

Q2

for feature-2

R	Gender	Height (cm)	Weight (kg)	Distance (Q24)
1	M	178	72	$\sqrt{(178-170)^2 + (72-70)^2} = \sqrt{64+4} = \sqrt{68}$
4	M	174	81	$\sqrt{157} \Rightarrow 11.70$
7	F	163	55	$\sqrt{274} \Rightarrow 16.55$
5	F	168	58	$\sqrt{148} \Rightarrow 12.16$
8	M	181	98	$\sqrt{905} \Rightarrow 30.08$
2	F	170	70	$\sqrt{100} \Rightarrow 10$
6	M	184	78	$\sqrt{260} \Rightarrow 16.12$
3	F	171	59	$\sqrt{122} \Rightarrow 11.04$

Predict for 170 cm & 70 kg for $k=3$ & $k=5$



How to Choose K :-

- * K value must be odd.
- * By Elbow Method.

Feature Scaling → we want to normalize the features.

↓
max ka value
min ka value

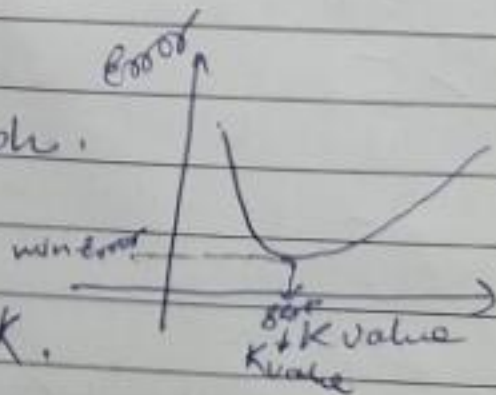
divide the max value
with min value.

- By default euclidean distance.
- confusion matrix → create matrix
- Classification - report → Precision, Recall, F1 score, Support.

Elbow Method

→ Based on the graph.

Lowest point of elbow is
best value for K.



NAIVE BAYES ALGORITHM

Probability

CLASSIFIER

→ Classification

- * It assumes that the presence of a particular feature in a class is unrelated to the presence of any other features.
- * It is very useful for large data set.

Experiment → Planned activity under certain condition

Sample space → Total outcome.

Event → Set of outcome

Exhaustive → at least one event exist

Independent → occurrence of one doesn't affect other

$$P(A \cap B) = P(A \cup B) = P(A) * P(B)$$

Conditional probability → for A, B already occur.

$$P(A|B) = P(A \cap B) / P(B)$$

BAYES Theorem →

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

We process new feature and refine our hypothesis at each step...

$P(A)$ is called Prior probability & $P(B)$ is called Evidence.

$P(B|A)$ is called Likelihood and $P(A|B)$ is called posterior probability.

Date _____
Page _____

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

↑ Posterior
↑ Likelihood
↑ Prior

\Downarrow
 Probability of A when B is already
 occur.

Evidence

⇒ All the features are independent to each other? &

$$P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

$$X = [X_1 X_2 \dots X_n] \rightarrow \text{Features}$$

$$P(X) = P(X_1, X_2, X_3, \dots, X_n) \rightarrow \text{Independent features.}$$

$$P(X) = P(X_1) * P(X_2) * P(X_3) * \dots * P(X_n)$$

Based on
 Conditional
 Independence

$$P(Y|X) = \frac{P(X_1|Y) P(X_2|Y) P(X_3|Y) \dots P(X_n|Y)}{P(X_1) P(X_2) \dots P(X_n)} \cdot P(Y)$$

finally the Y for which $P(Y|X)$ is maximum
is our predicted class.

- The application of Bayes theorem is Naive Bayes classifiers.

PCA - { PRINCIPAL COMPONENT ANALYSIS }

Feature Selection

→ Imp features {Property}

↓

When we take subset of feature from whole set.

Features Extraction

→ Creation of new feature from the previous/original features.

- Algo work with features subset.

- Algo work for new features.

→ Information gain, Randomforest.

- PCA - { A technique for features extraction }
- SVD - { Singular Value Decomposition }

(PCA) As we increase the no. of features, the dimension got increased.

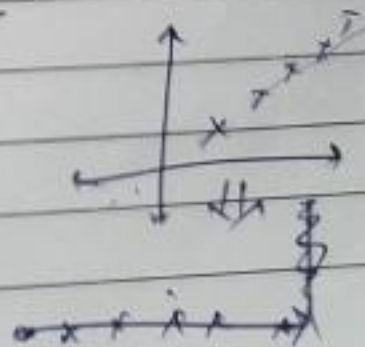
⇒ Dimensionality reduction $r < d$.

Remove inconsistencies

Redundant data

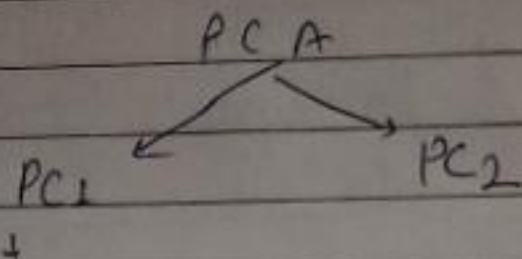
Highly correlated data

$$\begin{aligned} x^w &\in \mathbb{R}^2 (2-d) \\ x^w &\in \mathbb{R} (1-d) \end{aligned}$$



$$\begin{aligned} x^1 &\rightarrow z^1 \\ x^2 &\rightarrow z^2 \\ x^m &\rightarrow z^m \end{aligned}$$

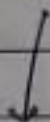
$$x^{(i)} \in \mathbb{R}^3 (3-d)$$



we get all features
without having
too much loss.

##

• Standardization of data



• Computing the covariance matrix



• Calculating the eigenvectors & eigenvalues



• Computing the principal components

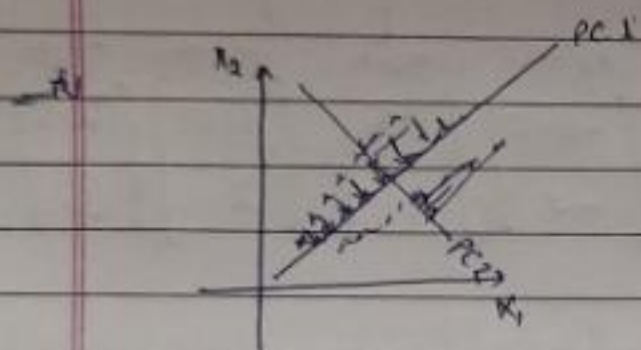
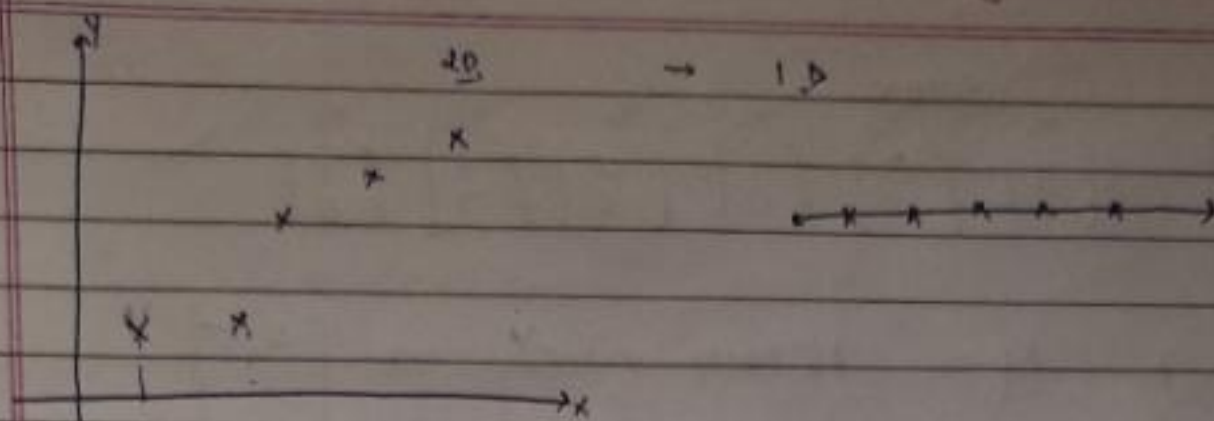


• Reducing the dimensions of the data

• $p \times p$ (Symmetric Matrix) f_1, f_2, \dots, f_p
Dimension

• Principal Components $p = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix}$

+ve - correlated
-ve - inversely
correlated



PC1
PC2
 $n \text{ Dimension} = n \text{ PC}$

Max. Information - PC1

Ex Given data

$\{2, 3, 4, 5, 6, 7; 1, 5, 3, 6, 7, 8\}$ compute the PC using PCA algorithm.

\Rightarrow Steps involved in PC Algorithm

- ① Get data
- ② Compute the mean vector (μ)
- ③ Subtract mean from the given data
- ④ Calculate the co-variance matrix
- ⑤ Calculate the eigen vectors & eigen values of the covariance matrix
- ⑥ Choose components & forming a feature vector.
- ⑦ Deriving the new data set.

Data \rightarrow 2, 3, 4, 5, 6, 7
1, 5, 3, 6, 7, 8.

①

	x_1	x_2	x_3	x_4	x_5	x_6
Vector	$\begin{bmatrix} 2 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 3 \\ -5 \end{bmatrix}$	$\begin{bmatrix} 4 \\ 3 \end{bmatrix}$	$\begin{bmatrix} 5 \\ 6 \end{bmatrix}$	$\begin{bmatrix} 6 \\ -7 \end{bmatrix}$	$\begin{bmatrix} 7 \\ 8 \end{bmatrix}$

② $M \Rightarrow H(\text{Mean vector}) \Rightarrow \begin{bmatrix} 4.5 \\ 5 \end{bmatrix}$

③ ~~vector~~

	$x_1 - H$	$x_2 - H$	$x_3 - H$	$x_4 - H$	$x_5 - H$	$x_6 - H$
	$\begin{bmatrix} -2.5 \\ -4 \end{bmatrix}$	$\begin{bmatrix} -1.5 \\ 0 \end{bmatrix}$	$\begin{bmatrix} -0.5 \\ -2 \end{bmatrix}$	$\begin{bmatrix} 0.5 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 1.5 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 2.5 \\ 3 \end{bmatrix}$

④ formula Covariance

$$\text{Covariance}_{\text{with } H} = \frac{\sum (x_i - H)(x_i - H)^T}{n}$$

$$m_1 = \begin{bmatrix} -2.5 \\ -4 \end{bmatrix} \begin{bmatrix} -2.5 & -4 \end{bmatrix} = \begin{bmatrix} 6.25 & 10 \\ 10 & 16 \end{bmatrix}$$

$$m_2 = \begin{bmatrix} -1.5 \\ 0 \end{bmatrix} \begin{bmatrix} -1.5 & 0 \end{bmatrix} = \begin{bmatrix} 2.25 & 0 \\ 0 & 0 \end{bmatrix}$$

$$m_3 = \begin{bmatrix} -0.5 \\ -2 \end{bmatrix} \begin{bmatrix} -0.5 & -2 \end{bmatrix} = \begin{bmatrix} 0.25 & 1 \\ 1 & 4 \end{bmatrix}$$

$$m_4 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} \begin{bmatrix} 0.5 & 1 \end{bmatrix} = \begin{bmatrix} 0.25 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$m_5 = \begin{bmatrix} 1.5 \\ 2 \end{bmatrix} \begin{bmatrix} 1.5 & 2 \end{bmatrix} = \begin{bmatrix} 2.25 & 3 \\ 3 & 4 \end{bmatrix}$$

$$m_6 = \begin{bmatrix} 2.5 \\ 3 \end{bmatrix} \begin{bmatrix} 2.5 & 3 \end{bmatrix} = \begin{bmatrix} 6.25 & 7.5 \\ 7.5 & 9 \end{bmatrix}$$

$$m = \begin{bmatrix} 17.5 & 22 \\ 28 & 34 \end{bmatrix}$$

Covariance matrix = $m/n = n=6$

$$M = \begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix}$$

⑤ Eigen value = λ

$$|M - \lambda I| = 0 \Rightarrow \begin{vmatrix} 2.92 - \lambda & 3.67 \\ 3.67 & 5.67 - \lambda \end{vmatrix} = \begin{vmatrix} \lambda & 0 \\ 0 & \lambda \end{vmatrix}$$

$$\begin{vmatrix} 2.92 - \lambda & 3.67 \\ 3.67 & 5.67 - \lambda \end{vmatrix} = 0$$

$$8.67, 0.82 = \lambda$$

$$(2.92)(5.67 - \lambda) - (3.67)(3.67) = 0$$

$$16.55 - 2.92\lambda - 5.67\lambda + \lambda^2 - 13.46 = 0$$

$$\lambda^2 - 8.59\lambda + 3.09 = 0$$

$$\lambda = \frac{+ 8.59 \pm \sqrt{(8.59)^2 - 4 \times 3.09}}{2} = \frac{8.59 \pm 7.03}{2}$$

$$\lambda_1 = 8.22, \lambda_2 = 0.38$$

Clearly λ_2 is very small so we neglect λ_2 .
So, Eigen vector ~~for~~ for λ_1 , i.e. for 8.22 value.

features
↓
eigen vector $MX = \lambda X$

$$\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 8.22 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

After solving $x_1 = 2.55$
 $x_2 = 3.67$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix} \quad \text{PC1}$$

8

Data	
1	
2	
3	
5	
8	
9	
11	

5.2

14/10

5

(X) - R

Yes - 1
No - 0
Single -
Married -
Divorced -

Data (10, n)

x	y	$\sqrt{y+36}$	$10 \cdot 01$	A
1	1	$\sqrt{1+36}$	10.01	A
2	3	$\sqrt{3+36}$	8.94	A
3	4	$\sqrt{4+36}$	8.54	A
5	3	$\sqrt{3+36}$	6.4	A
8	6	$\sqrt{6+36}$	2.14	B
8	8	$\sqrt{8+36}$	2.14	B
9	6	$\sqrt{6+36}$	1.41	B
11	7	$\sqrt{7+36}$	1	B

8/14/18

$$(10, 17) \rightarrow B'$$

(X) - R-NO 0
 H-M 2
 T-120K 120

} ^{2wade} NO

Yes	No	Yes-1	1	2	code	Points	
0	1	0	1	125	0	57.5	(2)
0	0	0	2	100	0	20	(3)
0	0	0	1	70	0	50	(4)
0	0	0	2	120	0	1	(1)
0	0	0	3	95	1	85.2	(5)
0	0	0	2	60	0	60	(6)
0	0	0	3	220	0	100	(10)
0	0	0	1	85	1	35	(5)
0	0	0	2	75	0	45	(6)
0	0	0	1	90	1	30	(9)