

ROC curve helps us in deciding the threshold.
 → under the ROC curve tells us the accuracy.

Name _____ Std _____ Sec _____

Roll No. _____ Subject _____ School/College _____

Sl. No.	Date	Title	Page No.	Teacher Sign/ Remarks
---------	------	-------	----------	-----------------------

AUC Curve and ROC curve



(True positive rate) :- No. of positive cases correctly classified as positive by total no. of positive cases at that threshold.

(False positive rate) :- No. of negative cases classified as positive by total no. of negative cases.

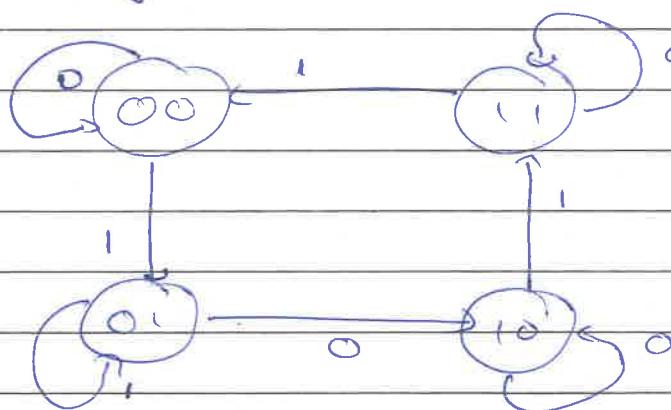
$$\begin{aligned} \text{True positive rate} &= \text{Sensitivity} \\ &= \frac{TP}{TP + FN} \end{aligned}$$

$$\begin{aligned} \text{False positive rate} &= 1 - \text{specificity} \\ &= 1 - \left(\frac{(TN)}{TN + FP} \right) = \frac{FP}{TN + FP} \\ &\rightarrow \text{No. of -Ve cases correctly classified as -Ve by total number of -Ve cases} \end{aligned}$$

The greater is the AUC for a ROC curve, better is the model.

DLD

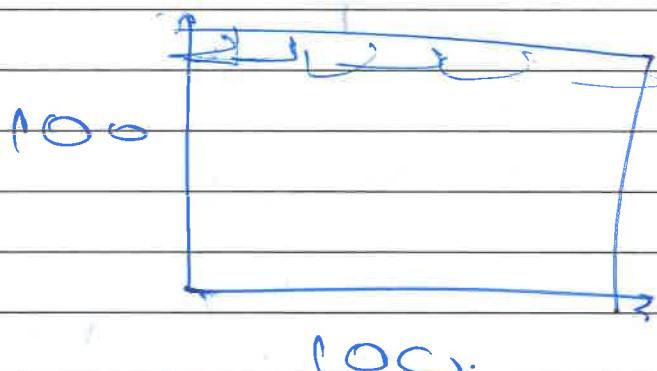
Design procedure for sequential circuits.



Q) we wish to design a synchronous & sequential circuit where state diagram is shown. The type of flip flop to be used is **JK**.

Q_A	Q_B	X	Q_A^+	Q_B^+	T_A	T_B
0	0	0	0	0	0	0
0	0	1	0	1	0	1
0	1	0	1	0	1	1
0	1	1	0	1	0	0
1	0	0	1	0	0	0
1	0	1	1	1	0	1
1	1	0	1	1	0	0
1	1	1	0	0	1	1

A	B	T
0	0	0
0	1	1
1	0	1
1	1	0



$t_1 \quad P_2 \quad P_3 \quad P_4 \quad \dots \quad P_{100}$

$P_1 \quad P_2 \quad P_3 \quad P_4 \quad \dots \quad P_{100}$

10000

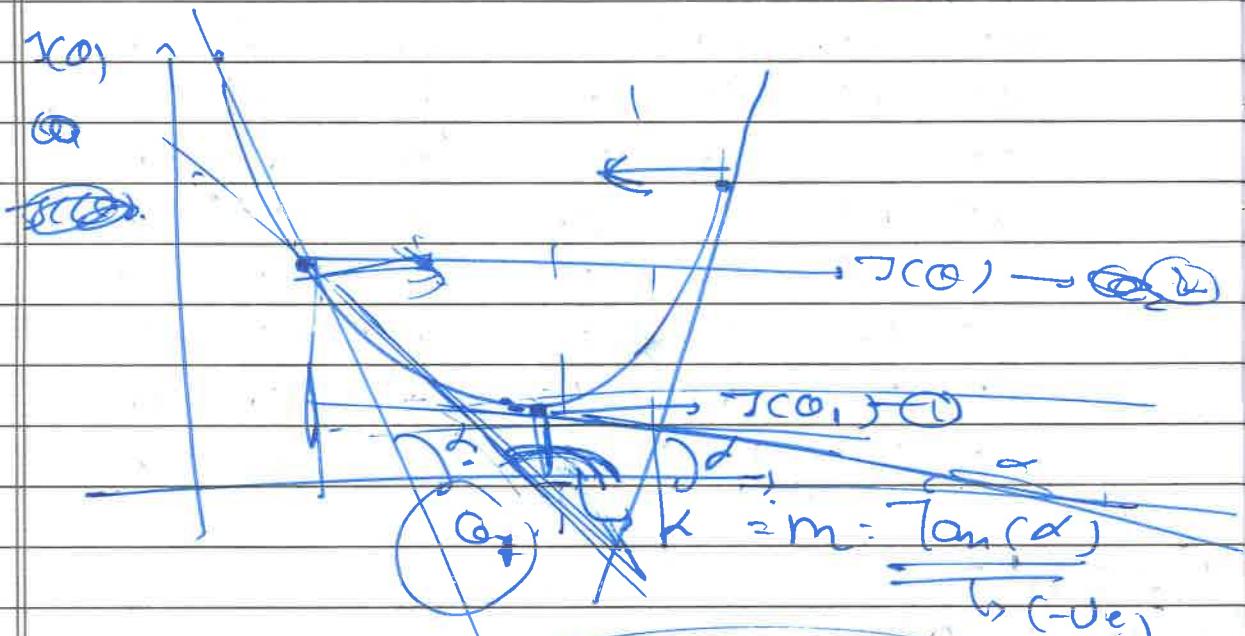
$$y = F \cdot O_1 + S \cdot O_2$$

papergrid

Floor	size	Predicted	Actual	Date:
(F)	(S)	(y)	(A)	
(1)	F_1	S_1	y_1	A_1
(2)	F_2	S_2	y_2	A_2
(3)	F_3	S_3	y_3	A_3
...				
(n)	F_n	S_n	y_n	A_n

$$\text{Cost} = \frac{(y_1 - A_1)^2}{2} + \frac{(y_2 - A_2)^2}{2} + \frac{(y_n - A_n)^2}{2}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (y - A_i)^2$$

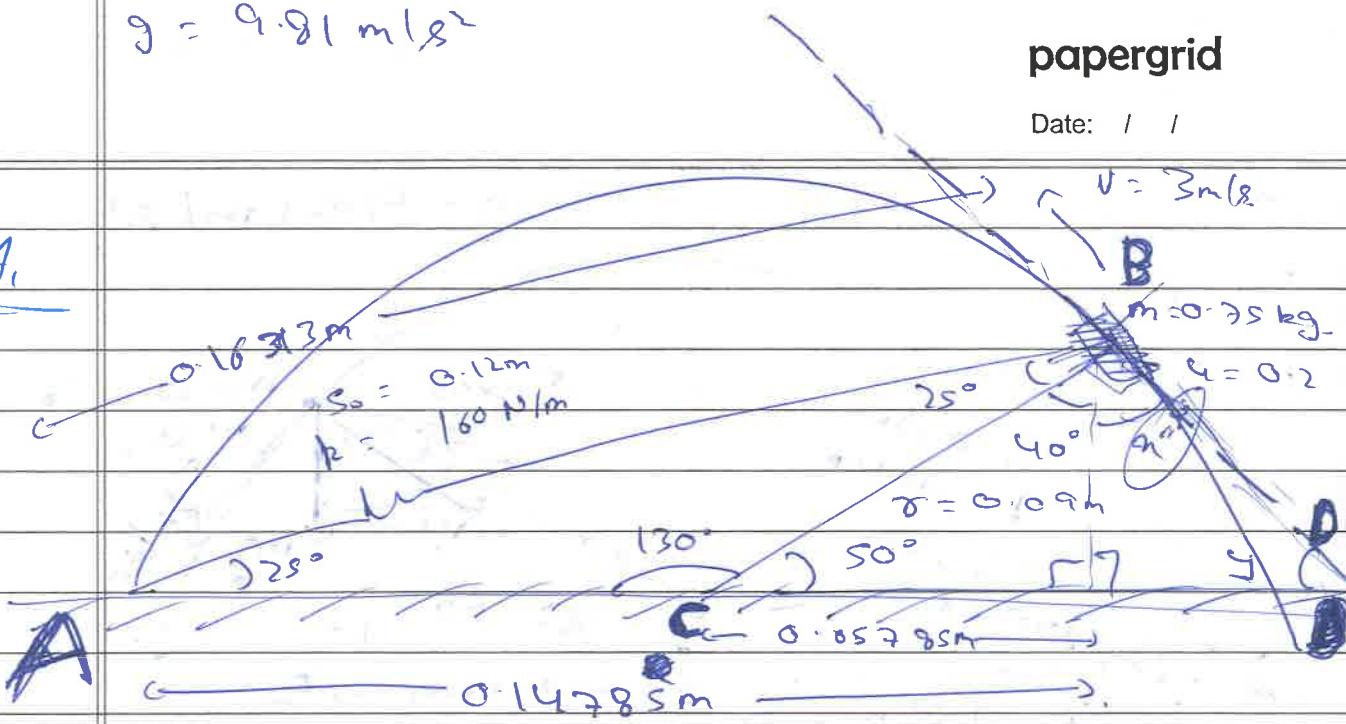


$$J(\theta) = J(\theta_0) - \left(\frac{\partial J(\theta_0)}{\partial \theta} \right) \cdot \Delta \theta$$

$$\Delta \theta =$$

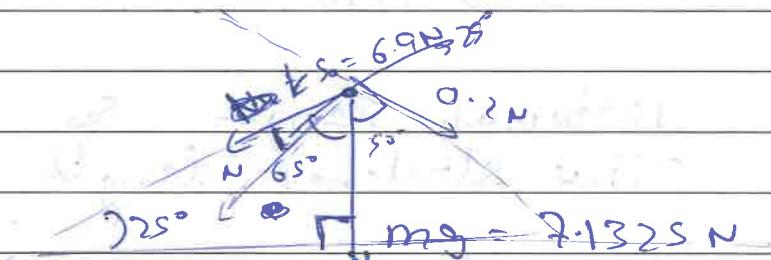
$$y = F^2 O_1 + F O_2 + S \cdot \dots + F_{n+2}$$

$$g = 9.81 \text{ m/s}^2$$



$$\text{Spring force} = (0.043135) \times (160) \\ = 6.9 \text{ N.}$$

F.B.D :-



~~Diagram~~ $\angle ABO = y = 90^\circ + 30^\circ = 120^\circ$

$$90^\circ + 65^\circ + x + y = 180^\circ$$

$$\Rightarrow 90^\circ + x + y = 180^\circ$$

$$\Rightarrow x + y = 90^\circ \quad \text{--- (1)}$$

$$\frac{d_2 + 1}{d_2} = 0$$

$$d_2 = -\frac{1}{2}$$

~~Diagram~~ $\angle ABD = 90^\circ + 30^\circ = 120^\circ$

$$N = 6.9 \cdot \cos 25^\circ$$

$$N = 6.2535$$

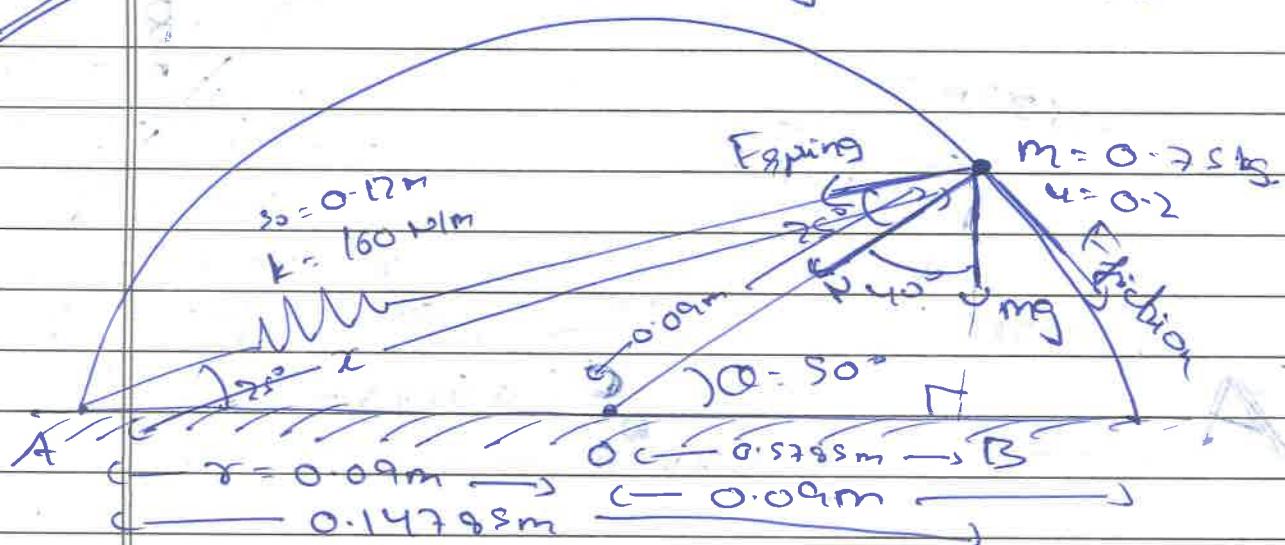
$$F = 1.2507$$

$$a = \frac{V^2}{r} = \frac{(3)^2}{0.09} = 100 \text{ m/s}^2$$

$$y = \left(\frac{-x}{2}\right)^2 + 2 - \frac{1}{2}$$

Q11.

$$g = 9.81 \text{ m/s}^2$$



$$x = (AB) \cdot \sec(25^\circ) = \text{current length}$$

$$\Rightarrow x = (0.14798) \cdot \sec(25^\circ) \text{ of spring}$$

$$\Rightarrow x = 0.16313 \text{ m}$$

$$\text{natural length} = s_0 = 0.12 \text{ m}$$

$$\text{Thus stretched length} = (0.16313 - 0.12) \text{ m}$$

$$= 0.04313 \text{ m}$$

$$\text{thus spring force} = (0.04313) \cdot k$$

$$= (0.04313) \cdot (160)$$

$$= \underline{\underline{6.9 \text{ N}}}$$

$$mg \text{ (gravitational force)} = (0.75)(9.81)$$

$$= \underline{\underline{7.1325 \text{ N}}}$$

$$\text{Normal force} = (\text{Spring force}) \cdot (\cos 25^\circ)$$

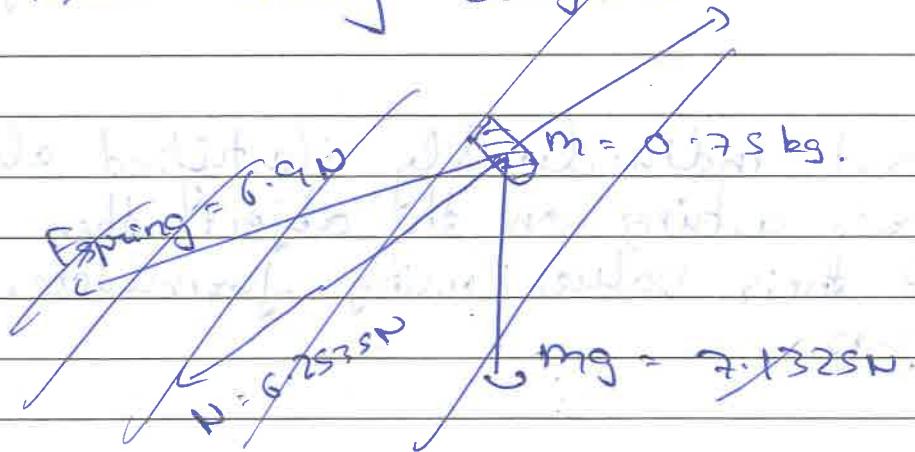
$$= \underline{\underline{6.2535 \text{ N}}}$$

$$\text{Frictional force} = (\text{Normal}) \cdot (\text{friction coefficient})$$

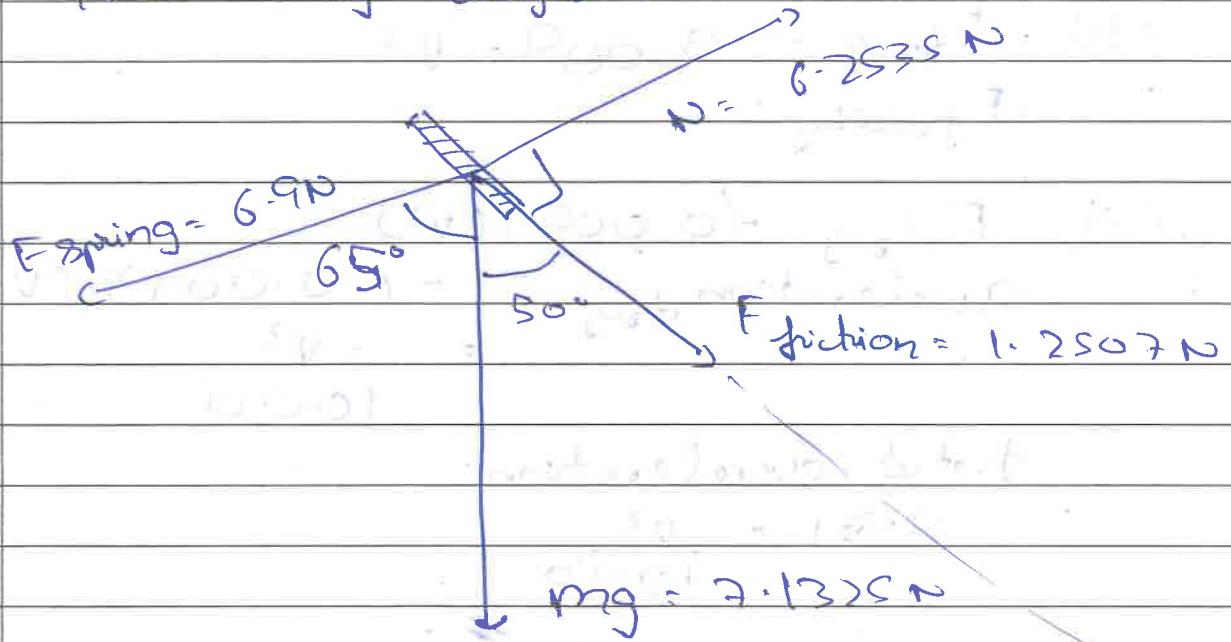
$$= (6.2535) \cdot (0.2)$$

$$= \underline{\underline{1.2507 \text{ N}}}$$

1a) ~~Free body diagram:-~~

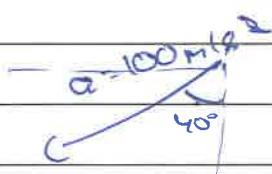


1b) ~~Free body diagram~~



1b) $N = 6.2535 \text{ N}$

1c) acceleration = $\frac{v^2}{r} = \frac{(30)^2}{0.09} = \frac{9}{0.09}$
 $= 100 \text{ m/s}^2$



~~$\vec{F}_{\text{net}} = (100 \sin 40) \hat{i} - 100 \cos 40 \hat{j}$~~

$$\vec{F} = (64.278761\hat{i} - 76.6\hat{j}) \text{ N} \text{ } (2)$$

~~(1)~~ First individually identified all the forces acting on the object then calculated their values using formulae and F.B.D.

~~(2)~~

$$\text{Friction} = 0.009 \cdot v^2$$

~~Friction~~ =

$$\text{Fdrag} = (0.009)(v^2)$$

$$\begin{aligned} \text{Acceleration drag} &= - (0.001) \cdot (v^2) \\ &= \frac{-v^2}{1000} \end{aligned}$$

Total acceleration:-

$$9.81 - \frac{v^2}{1000}$$

Acceleration =

The bird will keep accelerating until its total acceleration becomes 0.

Then until:-

$$9.81 = \frac{v^2}{1000}$$

$$\Rightarrow v^2 = 9810$$

$$\Rightarrow v = 99.045 \text{ m/s}$$

\hookrightarrow $\underline{\text{99.045}}$ (maximum velocity)

182

~~Q2a)~~ First found out the drag force using the final velocity, acceleration due to gravity. Then found out the work done by multiplying drag force with distance travelled. Then to calculate maximum velocity, found the velocity at which the total downward acceleration would become 0.

~~Q3a)~~

$$V_d = R \cdot \omega \\ = (1.5) \cdot (2.5) \\ \Rightarrow V_d = 3.75 \text{ m/s}$$

$$\alpha_{dd} = \omega^2 R$$

A	T	I	S	H	A	T
7	7	7	7	7	7	7

$T = O(n)$ \rightarrow Linear Search

0	1	2	3	4	5	6
10	20	30	40	50	60	70

bottom = 0

60

mid = $\text{int}(\text{len}(a)/2)$ $\text{top} = 6$

60 > 40 \rightarrow 60 > mid

bottom = mid + 1

bottom = 50

60 > 60 \rightarrow mid = $\frac{\text{bottom} + \text{top}}{2} = 55$

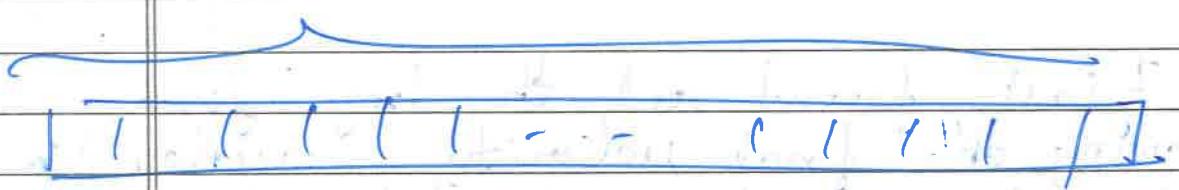
60 > 60

$O(\log(n))$

drop = bottom $\rightarrow 0$.

papergrid

Date: 1/1



len = $n \rightarrow 0^{\text{th}}$

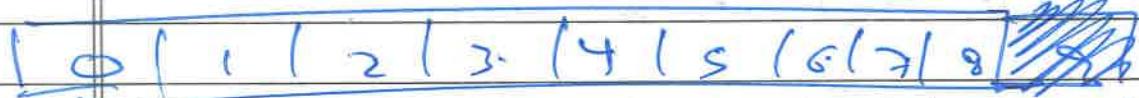
len = $\frac{n}{2} \rightarrow 1^{\text{st}}$

len = $\frac{n}{2^2} \rightarrow 2^{\text{nd}}$

$$\frac{n}{2^k} = 1$$

$$\frac{n}{2^k} = n$$

$$k = \log_2(n)$$



} bottom = 0 ; drop = 8 ; mid = 4
bottom = 5 ; drop = 8 ; mid = 6
bottom = 7 ; drop = 8 ; mid = 7
bottom = 9 ; drop = 8

$$\log_2 8 = 3$$

Emp Nr.	Name	Salary
1	A	10
2	B	20
3	C	30
4	D	40
5	E	50

→ 8 8

??

Search

→ Search Salary of Emp Nr = 8 1 1

Emp Nr	Name	Salary
5	E	50
7	G	70

Index

8

1

3

9

number of Searches = n

$n \log(n) =$

key = k	Value
5	50
7	70
8	80
1	10
3	30
9	90
2	20

key → Hash Table (function) → Value
 (7) → 70

Key	Value
5	50
7	70
8	80
1	10
3	30
9	90
2	20

(i) linear probing → index
 $f(2k+3)$ $f((2k+3) \mod \text{len})$
 ~~$20 \mod (10+3) \mod 7$~~
~~6~~

hash(Say)

5000 → New address
10000 (5000)

when say 5000 → 5000

Columnar db

~~Row DBMS~~

1, 3, 7

f(Say) → malloc(address)

1 → 10 (maddr 1L)
2 → 20 (maddr 2L)
5 → 50 (maddr 5L)

→ f(Say) → 5000 → (1, 5, 15) 5L

maddr 10 5L

n samples.

papergrid

Date: / /

1

$$\frac{(n)(n+1)}{2}$$

2

3

1

3

multiple features (Variables)

Size (feet ²) (x_1)	Number of bed rooms (x_2)	No. of floors (x_3)	Age of home (years) (x_4)	Price (in \$1000) (y)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	35	178
-	-	-	-	-
-	-	-	-	-

notations:-

 n = number of features $x^{(i)}$ = input (features) of i^{th} training example $x_j^{(i)}$ = value of feature j in i^{th} training example

$$x^{(i)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

hypothesis :-

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$\text{Eg. } h_0(x) = 80 + 0.1(x_1) + 0.01(x_2) + 3(x_3) - 2(x_4) \dots$$

$$\Rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \\ = \theta^T x$$

$$= [\theta_0 \ \theta_1 \ \dots \ \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Cost function :-

$$J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \underbrace{\frac{1}{2m}}_{J(\theta)} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Gradient descent:- (for linear regression)
Repeat {

$$\theta_j := \theta_j - \alpha \underbrace{\frac{\partial}{\partial \theta_j} [J(\theta_0, \theta_1, \dots, \theta_n)]}_{J(\theta)}$$

} (simultaneously update for every $j = 0, 1, 2, 3, \dots, n$)

Gradient descent:- (for multiple regression)..

{ p. 7. 0 }

repeat until convergence of

$$\theta_j := \theta_j - \alpha \left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right)$$

(simultaneously update for $j = 0, \dots, n$)

}

$$\theta_0 := \theta_0 - \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

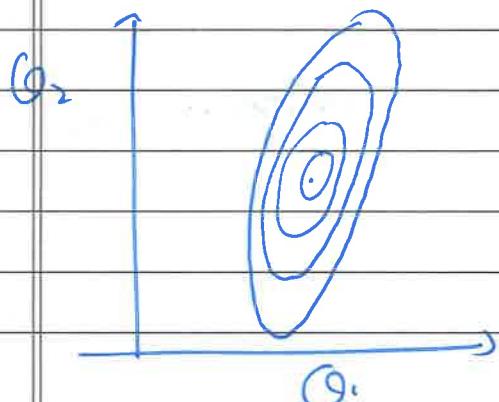
$$\theta_2 := \theta_2 - \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

* Feature scaling

idea: make sure features are on similar scale

Eg. x_1 = size (0 - 2000 feet²)

x_2 = number of bedrooms (1-5)



$$\rightarrow x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$\rightarrow x_2 = \frac{\text{number of bedrooms}}{5}$$

Get every feature into approximately
 $-1 \leq x_i \leq 1$ range

if we get :-

$$-0.0001 \leq x_n \leq 0.0001$$

$\underbrace{\hspace{10em}}$ This is too small.

mean normalization :-

replace x_i with $x_i - \bar{u}_i$ to make
 features have approximately 0 mean.

$$\text{Eg. } x_1 = \frac{\text{size} - 1000}{2000}$$

$$x_2 = \frac{\text{bedroom} - 2}{5}$$

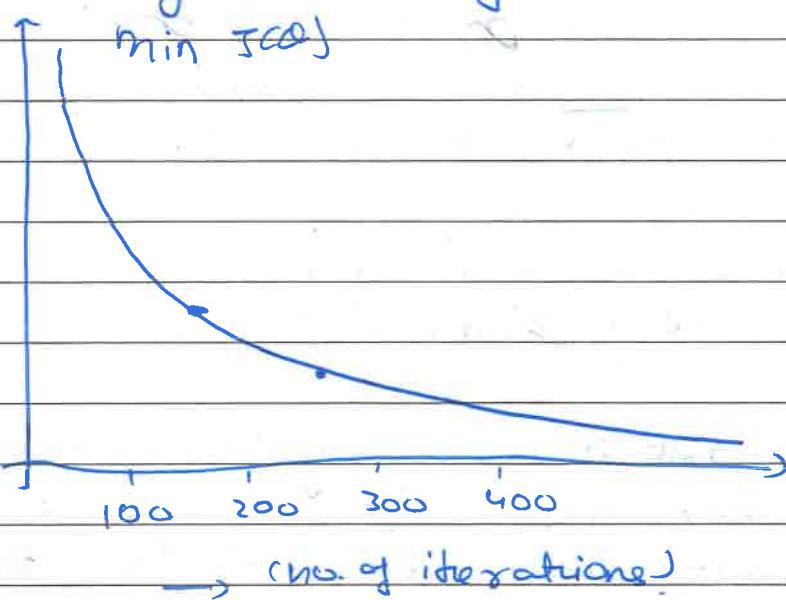
thus we can replace

$$x_i \text{ with } \frac{x_i - \bar{u}_i}{s_i}$$

avg. value
(of x_i in
training set)

(range $\rightarrow \text{max-min}$)

making sure gradient descent is working correctly :-

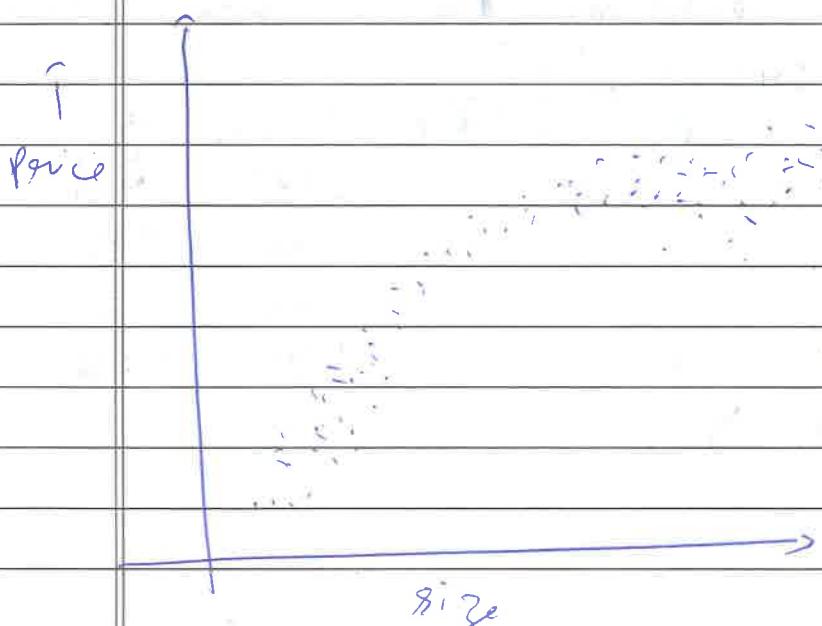


Q1 Polynomial regression :-

let y be the price of the house and it depends upon length & breadth. Thus

$$y = \theta_0 + \theta_1(\text{length}) + \theta_2(\text{breadth})$$

So we could say that in a way the price depends upon the size (area) we are given :-



so we could in a way ~~in~~ formulate

$$h_0(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_2}$$

so let $\sqrt{x_2}$ be x_2

so we have :-

$$h_0(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

Here feature scaling is important
so if we have the range of x_2
as 1000 then as $\sqrt{1000} \approx 33.33$
we could do :-

$$x_1 = \frac{x_1}{1000} \quad \text{and} \quad x_2 = \frac{x_2}{33.33}$$

~~Q1~~ Normal equation :-

$$\theta = (X^T X)^{-1} X^T y$$

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	952	2	1	30	178

$$m = 4$$

$$X = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(3)} \\ x^{(4)} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 952 & 2 & 1 & 36 \end{bmatrix} \text{ (m} \times (n+1)\text{)}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \text{ (m-dimensional vector)} \\ (m \times 1)$$

normal equation	gradient descent
→ No	
→ Need to choose alpha	→ Need to choose alpha
→ No need of iterations	→ Needs many iterations
→ $O(n^3)$	$O(kn^2)$
→ slow when n is large	works well when n is large

When implementing normal equation in Octave we want to use (pinv) rather than (inv)

If $X^T X$ is non-invertible, common causes may be:-

- (i) Two features may be very closely related (linearly dependent)
- (ii) Too many features ($m \leq n$). In this case delete some features, or use regularization.

$$\mathbf{\Theta} = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_n \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$$

Unvectorized implementation:-

$$\text{prediction} = \mathbf{\Theta} \cdot \mathbf{x};$$

for $j = 1:n+1$,

$$\text{prediction} = \text{prediction} + \Theta_j \cdot x_j$$

end;

Vectorized implementation:-

$$\text{prediction} = \mathbf{\Theta} \cdot \mathbf{x}$$

It is a classification model.

Date: / /

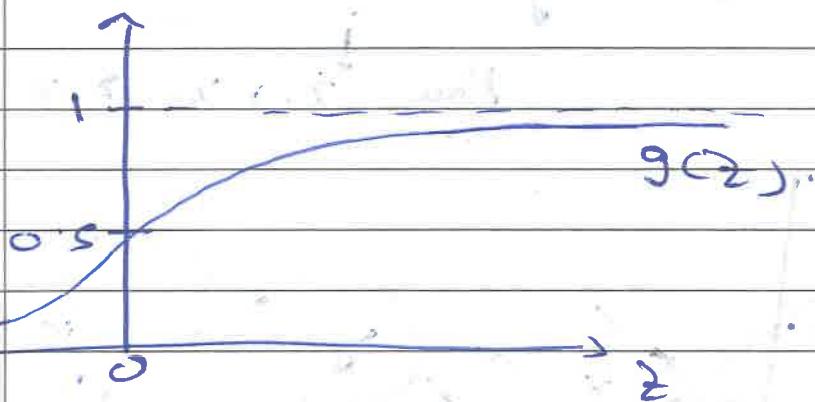
~~At~~ Logistic regression:-

$$\text{Want } 0 \leq h_{\theta}(x) \leq 1$$

$$h_{\theta}(x) = g \cdot (x^T \theta)$$

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

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



Interpretation of hypothesis output

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x .

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$

$$h_{\theta}(x) = 0.7$$

Patient has 70% chance of malignant tumor

$$\text{If } h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$\text{If } h_{\theta}(x) < 0.5 \rightarrow y=0$$

→ we notice that $g(z) > 0.5$ when $z \geq 0$ then $h_{\theta}(x) > 0.5$ when

$$\theta^T x \geq 0$$

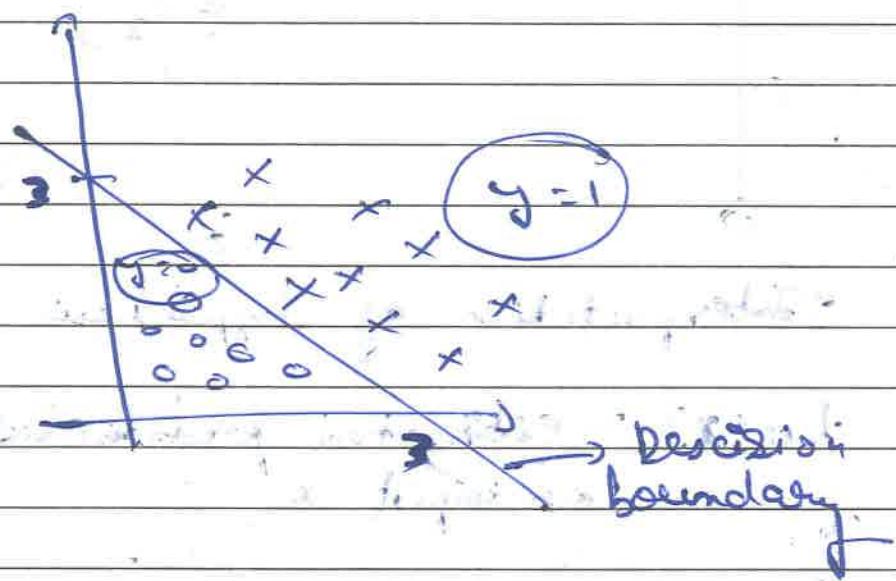
$$\text{let } h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\text{let } \theta_0 = -3; \theta_1 = 1; \theta_2 = 1$$

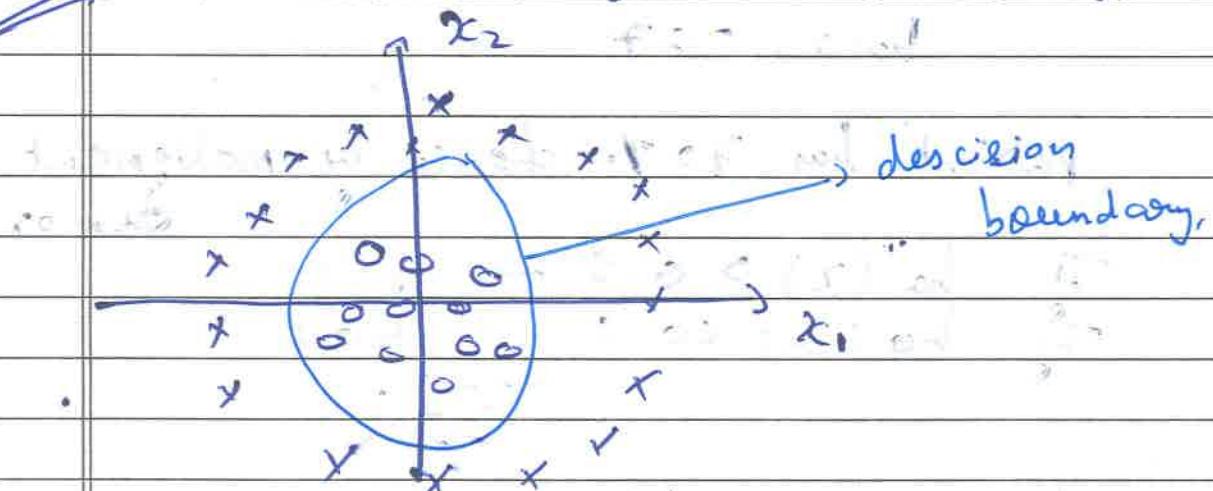
then

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

then predict "y=1" if $\underline{\underline{\theta_0 + x_1 + x_2 \geq 0}}$
then $x_1 + x_2 \geq 3$



Q1. Non-linear decision boundaries.



$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

predict " $y = 1$ " if $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 \geq 1$$

$$\alpha = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$



Training set:-

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

m examples:- $x \in \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^3$ $x_0 = 1$ $y \in \{0, 1\}$

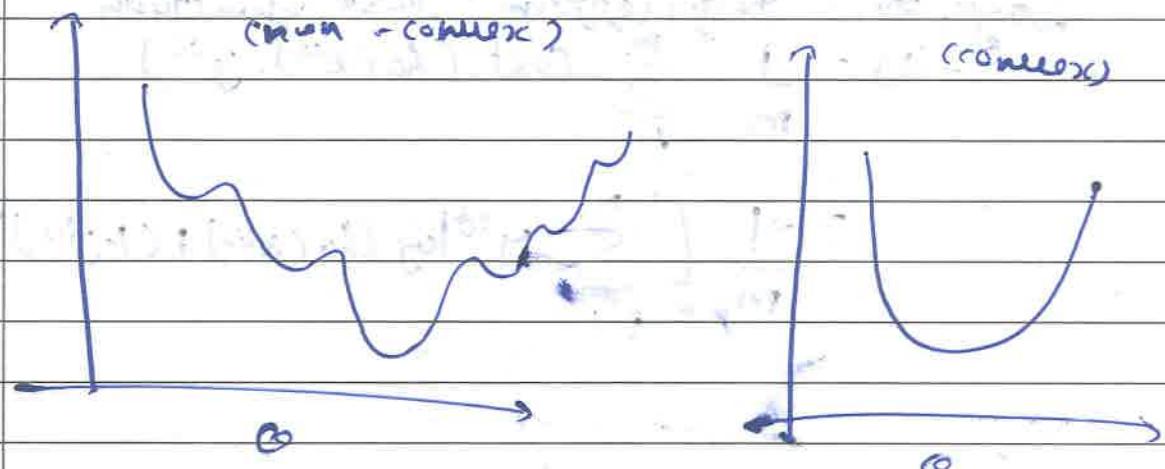
$$h_\alpha(x) = \frac{1}{1 + e^{-\alpha^T x}}$$

now we have to decide parameter α .

cost function:-

$$\text{linear regression: } J(\alpha) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\alpha(x^{(i)}) - y^{(i)})^2$$

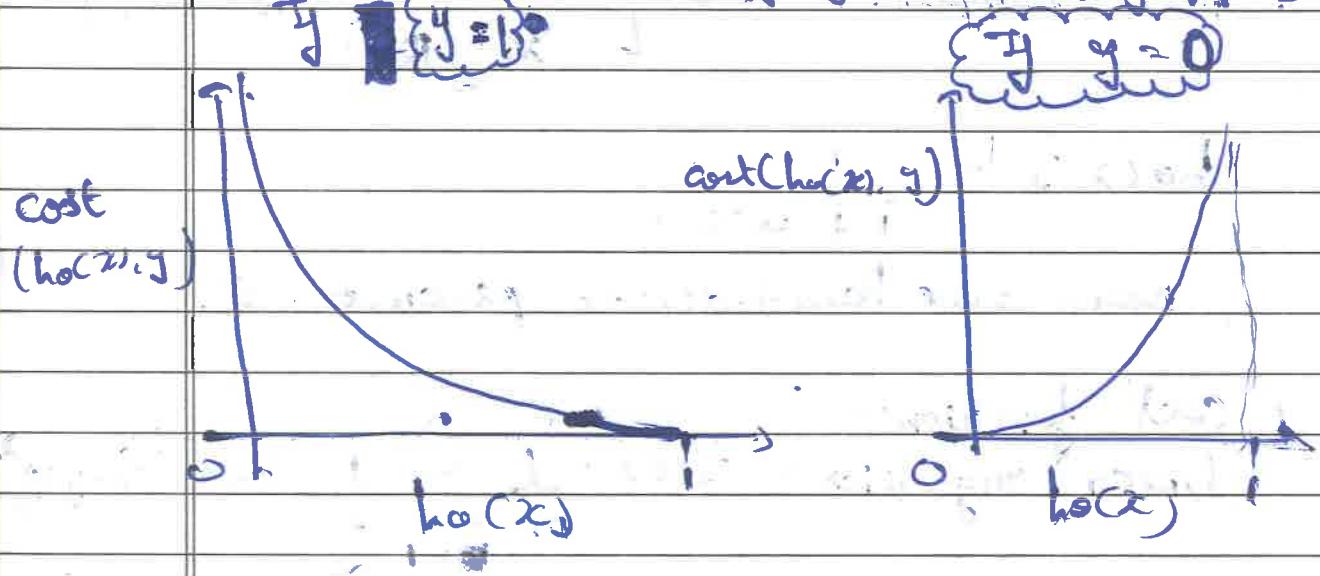
$$\text{cost}(h_\alpha(x), y) = \frac{1}{2} (h_\alpha(x) - y)^2$$



If we have a non-convex function then we won't be able to get the global minimum through the gradient descent algorithm. And we get a non-convex graph if we choose $h(x) = \frac{1}{1+e^{-\theta^T x}}$ for our cost function. Then we change our cost function.

Logistic regression cost functions-

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



$$\text{cost}(h_\theta(x), y) = (1+y) \log(h_\theta(x)) + (1-y) \log(1-h_\theta(x))$$

Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x^{(i)}), y^{(i)})$$

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

To fit parameters θ :-

$$\min J(\theta)$$

$$\theta$$

To make a prediction given by new x :-

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

We want to minimize $J(\theta)$.

Repeat {

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

(Simultaneously
update for all θ_j).

On solving :-

Repeat {

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(Simultaneously update all θ_j)

Algorithm looks identical to linear regression.

Optimization algorithms:-

→ Gradient descent

→ Conjugate gradient

→ BFGS

→ L-BFGS

Advantages:-

- No need to

manually pick α

- Often faster than

gradient descent

Disadvantages:-

- More complex.

Example :-

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_0 - 5)^2 + (\theta_1 - 5)^2$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = 2(\theta_0 - 5)$$

$$\frac{\partial \theta_0}{\partial \theta_1}$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = 2(\theta_1 - 5)$$

$$\frac{\partial \theta_1}{\partial \theta_2}$$

function [Jval, gradient] =

costFunction(theta)

$$Jval = \theta_0(1) - 5)^2 + \theta_1(2) - 5)^2$$

~~gradient =~~ costFunction(theta)

function [Jval, gradient] =

costFunction(theta)

$$Jval = (\theta_0(1) - 5)^2 + (\theta_1(2) - 5)^2$$

$$gradient = 2 \times \theta_0(2, 1);$$

$$gradient(1) = 2 * (\theta_0(1) - 5);$$

$$gradient(2) = 2 * (\theta_1(2) - 5);$$

~~gradient = optimset('GradObj', 'on', 'MaxIter', 100);~~

~~initialTheta = zeros(2, 1);~~

options = optimset('GradObj', 'on', 'MaxIter', 100);

initialTheta = zeros(2, 1);

[optTheta, functionVal, exitFlag]

= fminunc(@(costFunction, initialTheta, options))

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

function [J; val, gradient] = costFunction(theta)

; Val = [code to compute $J(\theta)$];

gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

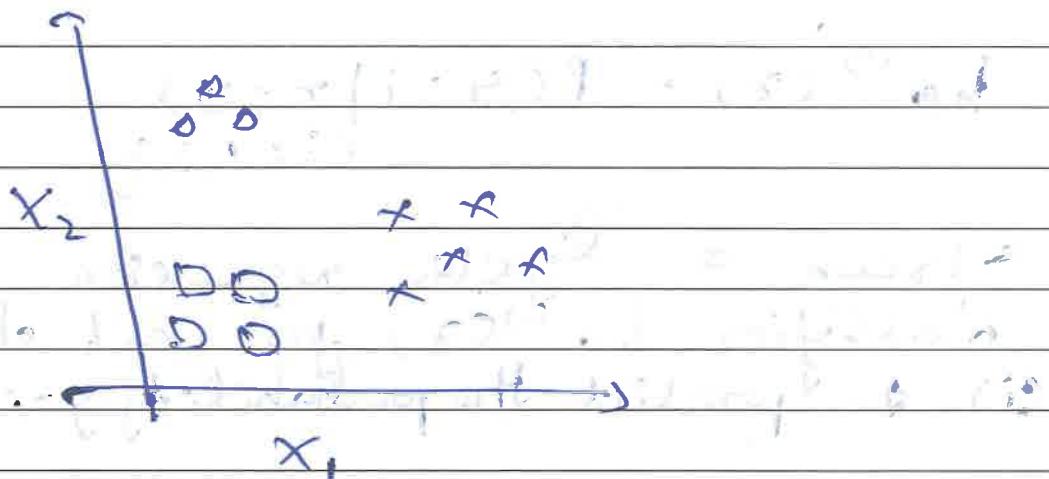
gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

gradient(n) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

Multi-class classification.

Eg) Email filtering / tagging: work, friends, family, hobby

Eg) medical diagnosis:- cold, flu, not ill.



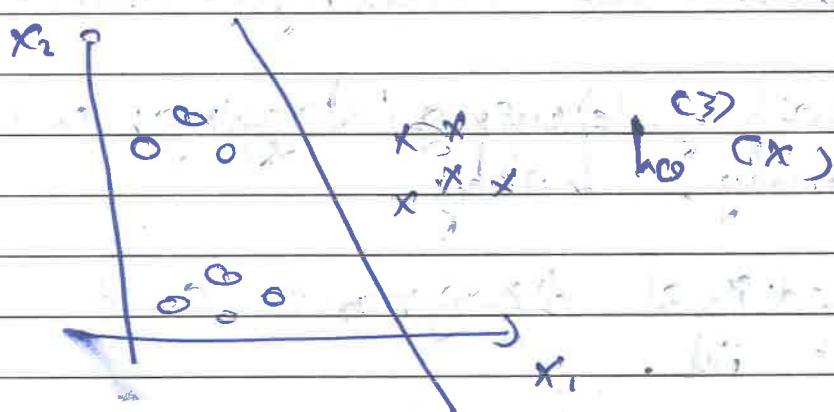
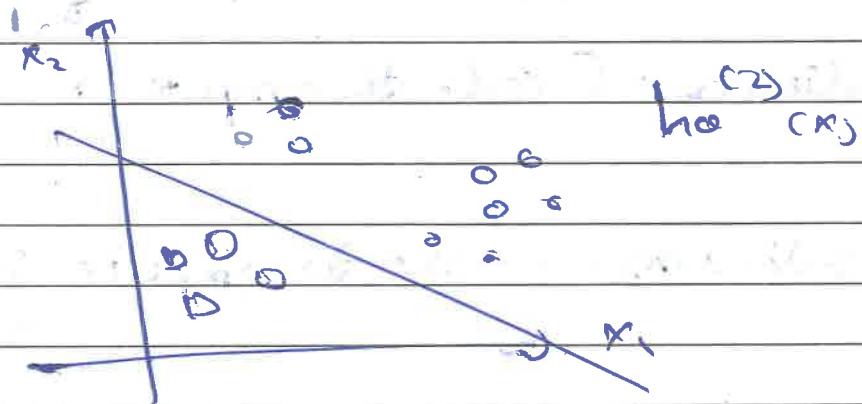
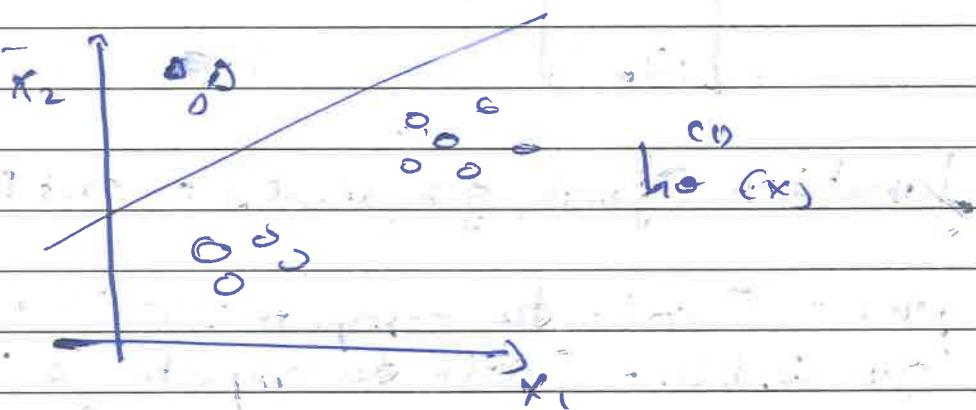
class 1 \rightarrow OOO
class 2 \rightarrow OOO
class 3 \rightarrow XXX

papergrid

Date: / /

We will solve this by making
3 binary classification.

Thus -



$$h_0^{(i)}(x) = P(y=i|x; \theta) \quad (i=1,2,3)$$

Train a logistic regression
classifies $h_0^{(i)}(x)$ for each class
(i) to predict the probability $y=i$



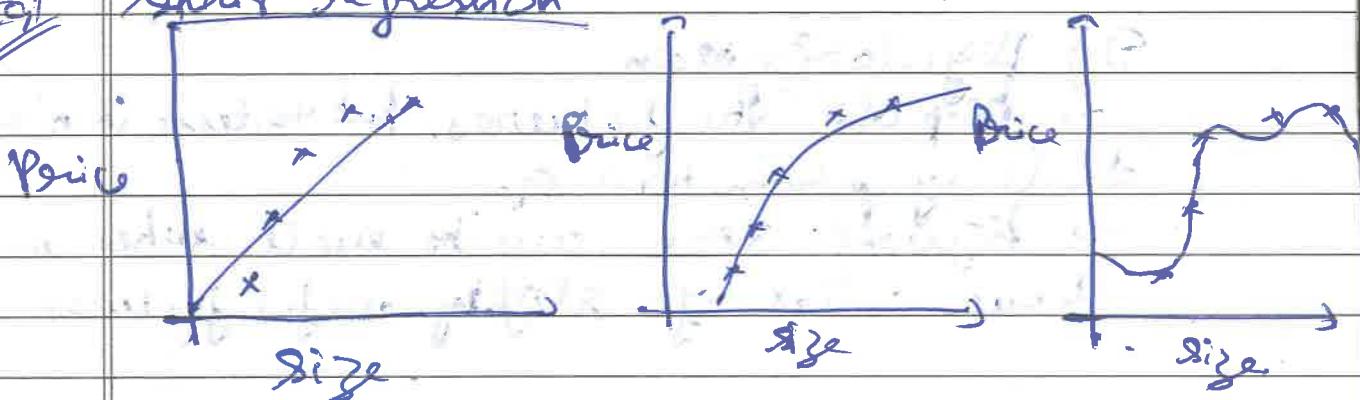
Regularization!

The problem of overfitting.

Eg:

Sometimes logistic & linear regression work poorly because of overfitting.

linear regression



$$\rightarrow \theta_0 + \theta_1 x$$

"Underfit"

"High bias"

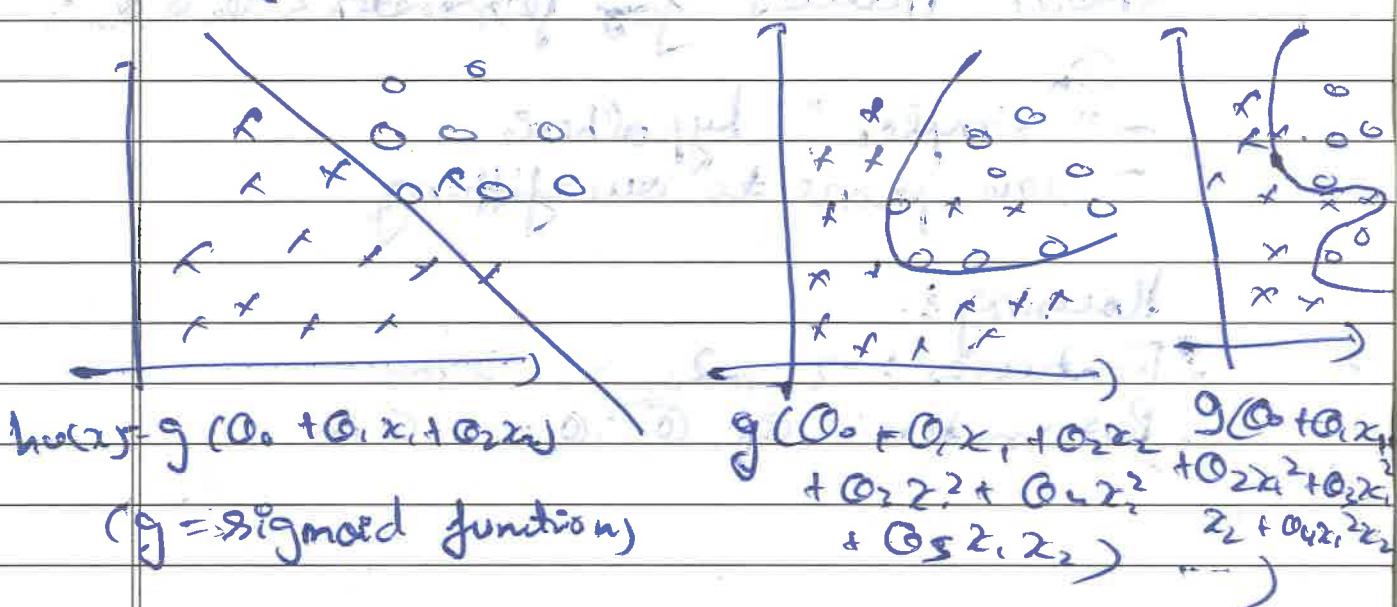
$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

"Just right"

"Overfitting"

Overfitting: - If we have too many features, the learned hypothesis may fit the training set very well but fail to generalize to new examples.

Eg: logistic regression



Addressing overfitting

1.) Reduce number of features.

→ Manually select features to keep

→ Model selection algorithm

2.) Regularization

→ keep all the features, but decrease the magnitude of parameters θ_1

→ regularization works well when we have a lot of slightly useful features.

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4$$

Suppose we penalize and make θ_3, θ_4 very small.

$$\min_{\theta} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

$$\theta_3 \approx 0 \quad \theta_4 \approx 0$$

Small values for parameters $\theta_0, \theta_1, \theta_2, \theta_3, \theta_4$

θ_0
- "simple" hypothesis

- less prone to overfitting.

Housing :-

- Features: x_0, x_1, \dots, x_{100}

- Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

Notes: we do not penalize θ_0 (the constant term). we could penalize though, that won't make much difference, but we don't.

Date: 1/11/2023

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^n \theta_j^2 \right]$$

regularization parameter

This term should decrease the values of θ .

Regulated linear regression.

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^n \theta_j^2 \right]$$

gradient descent:

~~$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right] + \frac{\lambda}{m} \theta_j$$~~

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$$\Rightarrow \theta_j := \theta_j \left(1 - \frac{\alpha \lambda}{m} \right) - \alpha \cdot \frac{1}{m} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

3 P.T.O

usually, $1 - \alpha \left(\frac{1}{m}\right) < 1$ then

$$\Theta_j \left(1 - \alpha \frac{1}{m}\right) < \Theta_j$$

as we know:-

$$\Theta_j := \Theta_j \left(1 - \alpha \frac{1}{m}\right) - \frac{\alpha}{m} \sum_{i=1}^m (h_\Theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

This term
 would become a little
 less than Θ_j .
 And this term
 is the term
 earlier

Normal equation:-

$$X = \begin{bmatrix} (x^{(1)}) \rightarrow \\ (x^{(2)}) \rightarrow \\ (x^{(3)}) \rightarrow \\ \vdots \\ (x^{(m)}) \rightarrow \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\Theta = (X^T X + \lambda \cdot I)^{-1} X^T y$$

where $I = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$

excluding
 entries λ on the diagonal

where n is the number of parameters.

Regularized Logistic regression.

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \dots)$$

cost function:-

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log h_{\theta}(x^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))] \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

gradient descent:-

Repeat {

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$+ \frac{\alpha \lambda}{m} \cdot \theta_0$$

Advanced optimization:-

function [Jval, gradient] = costFunction(theta)

jVal = [code to compute J(theta)];

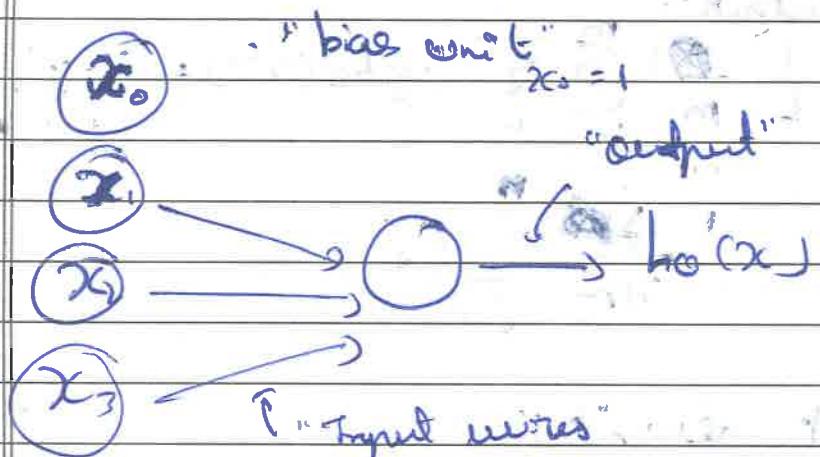
gradient(1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

gradient(2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

gradient(n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];

~~Neural networks~~

- Origins :- Algorithms that try to mimic the brain.
- State of the art technique for many applications.



$$h_0(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Sigmoid (Logistic) activation function

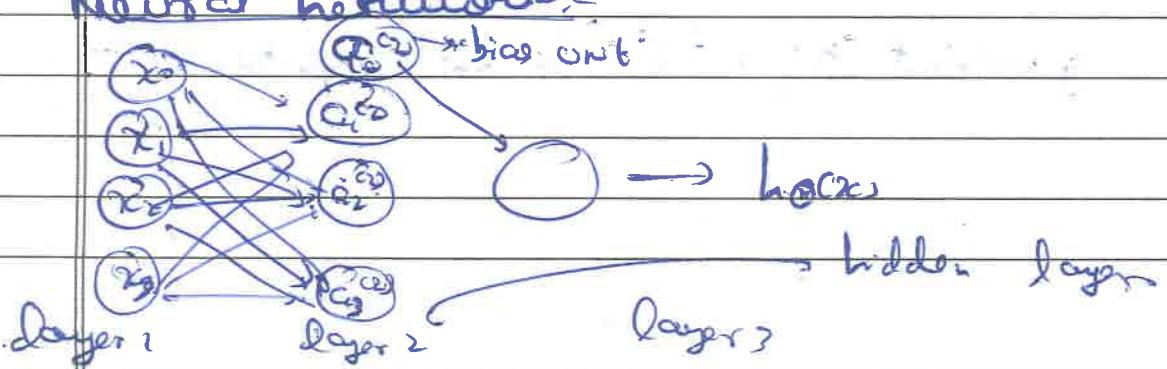
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

↳ these are called as weights = parameters.

Neural networks



$a_i^{(j)}$ = "activation" of unit i in layer j

$\mathbb{H}^{(j+1)}$ = matrix of weights controlling function mapping from layer j to layer $j+1$

$$z_i^{(2)} \rightarrow$$

$$a_1^{(2)} = g \left(\mathbb{H}_{10}^{(2)} x_0 + \mathbb{H}_{11}^{(2)} x_1 + \mathbb{H}_{12}^{(2)} x_2 + \mathbb{H}_{13}^{(2)} x_3 \right)$$

$$a_2^{(2)} = g \left(\mathbb{H}_{20}^{(2)} x_0 + \mathbb{H}_{21}^{(2)} x_1 + \mathbb{H}_{22}^{(2)} x_2 + \mathbb{H}_{23}^{(2)} x_3 \right)$$

$$a_3^{(2)} = g \left(\mathbb{H}_{30}^{(2)} x_0 + \mathbb{H}_{31}^{(2)} x_1 + \mathbb{H}_{32}^{(2)} x_2 + \mathbb{H}_{33}^{(2)} x_3 \right)$$

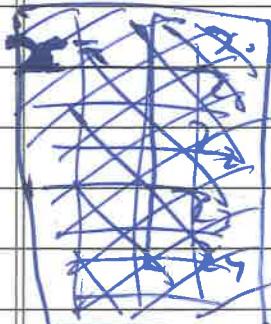
$$h_0(x) = a_1^{(3)} = g \left(\mathbb{H}_{10}^{(2)} a_1^{(2)} + \mathbb{H}_{11}^{(2)} a_2^{(2)} + \mathbb{H}_{12}^{(2)} a_3^{(2)} \right)$$

If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\mathbb{H}^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$

$$a_1^{(2)} = g(z_1^{(2)})$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)})$$



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \mathbb{H}^{(2)} x \Rightarrow \text{is equal to } a^{(2)}$$

$$a^{(2)} = g(z^{(2)})$$

$$\text{now add } Q_0^{(2)} = 1$$

$$z^{(3)} = (H)^{(2)} Q^{(2)}$$

$$h_{(2)}(x) = Q^{(3)} = g(z^{(3)})$$

This process is called forward propagation -
vectorized implementation.

so at the end we get:

~~$$h_{(2)}(x) = g((H)^{(2)} Q^{(2)} + (A)^{(2)})$$~~

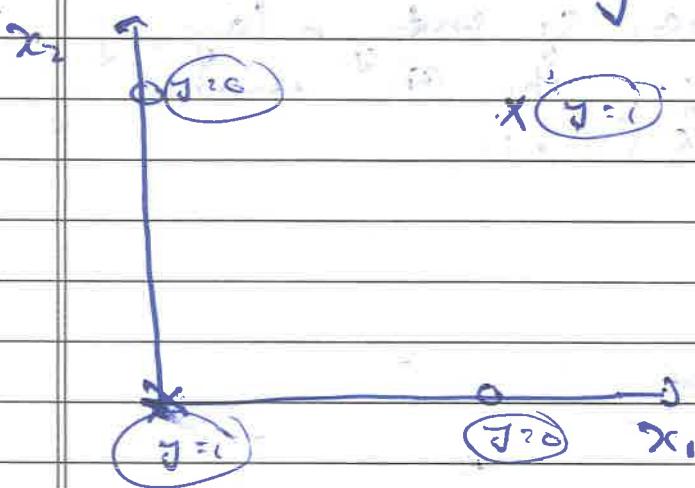
so at the end we get:

$$h_{(2)}(x) = g((H)_{1,0}^{(2)} Q_0^{(2)} + (H)_{1,1}^{(2)} Q_1^{(2)} + (H)_{1,2}^{(2)} Q_2^{(2)} + (H)_{1,3}^{(2)} Q_3^{(2)})$$

~~Non-linear classification example:~~

XOR

$\rightarrow x_1, x_2$ are binary (0 or 1)



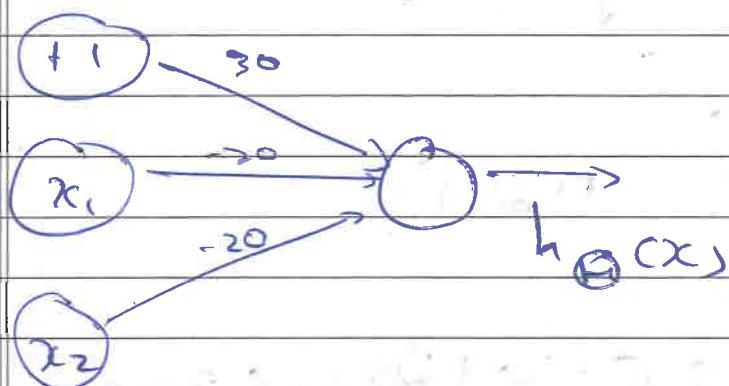
Simple example: AND

$x_1, x_2 \in \{0, 1\}$

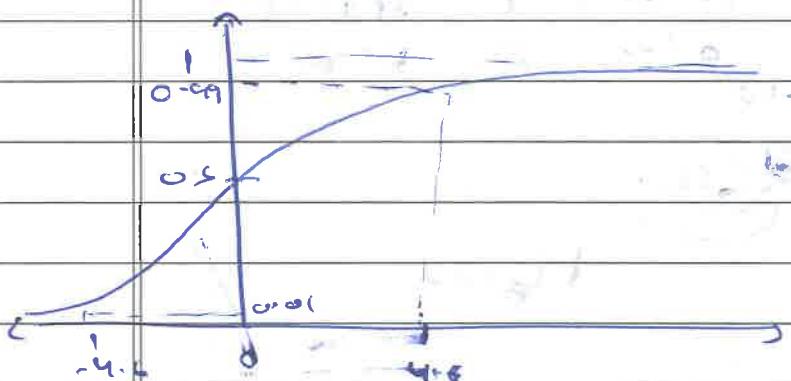
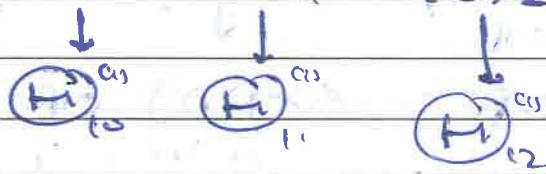
$y = x_1 \text{ AND } x_2$



→ Simple examples:- AND



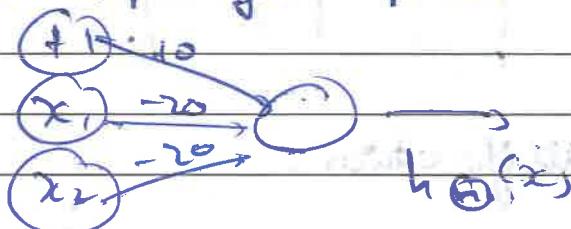
$$h \otimes (x) = g(30 + 20x_1 + 20x_2)$$



Truth table:-

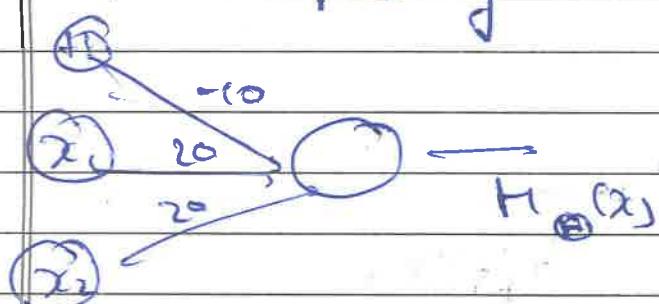
\$x_1\$	\$x_2\$	\$x_1 \wedge x_2 = h \otimes (x)\$
0	0	0
0	1	0
1	0	0
1	1	1

→ Example for NOR



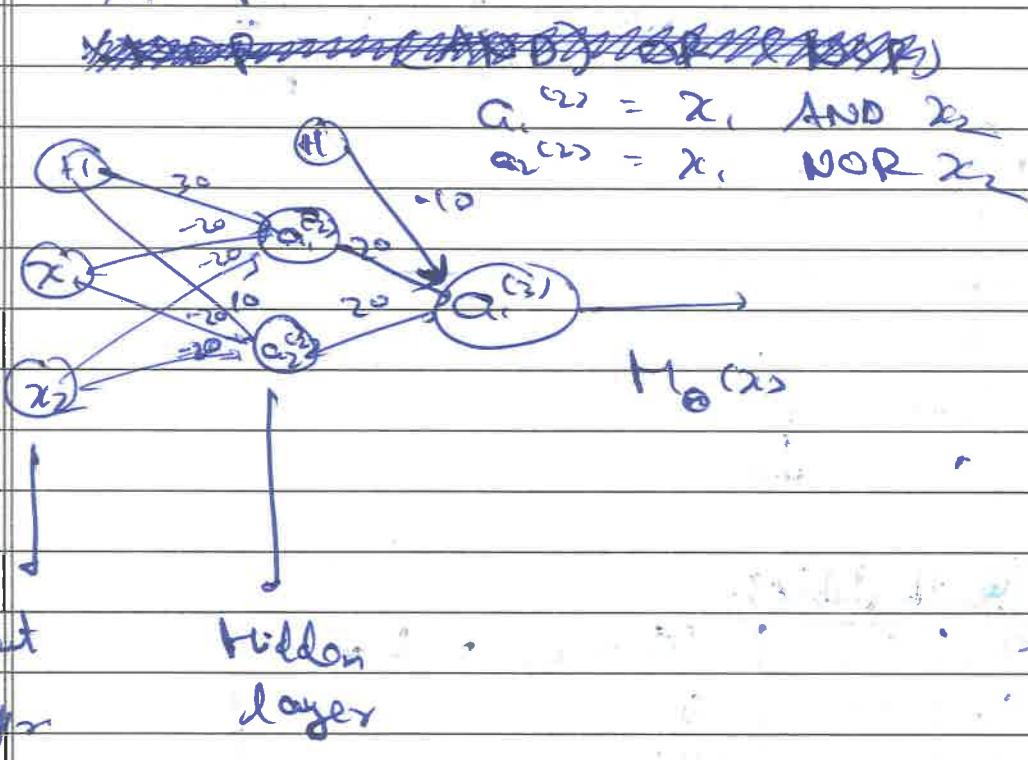
$$h \otimes (x) = g(10 - 20x_1 - 20x_2)$$

→ Example of OR



$$h_0(x) = g(-10 + 20x_1 + 20x_2)$$

→ Now putting it all together for XNOR

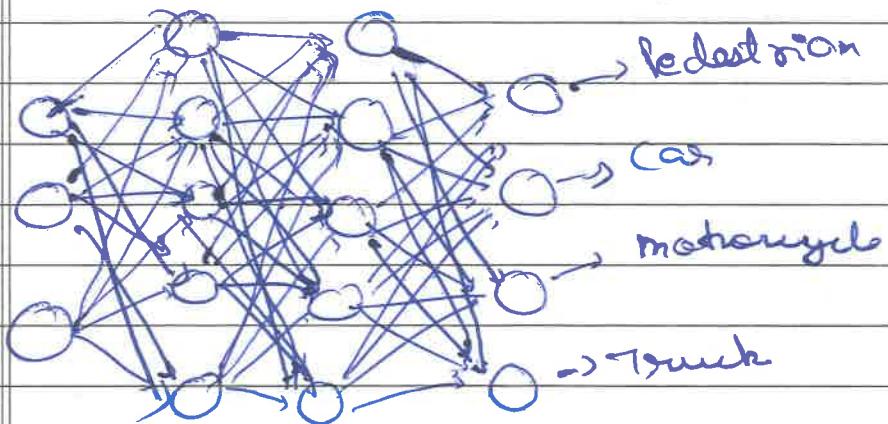


Input layer Hidden layer

x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_0(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

→ Multiclass classification
one-vs-all.

(Pedestrian) - (Car) - (Motorcycle) - (Truck)



etc.

Neural network (Classification)
 $(x^1, y^1), (x^2, y^2), (x^3, y^3), (x^4, y^4) \dots (x^m, y^m)$

④ Training Set

l = Total no. of layers in network

s_i = no. of units (not counting bias unit) in layer i .

7 There are k output units
then $S_L = k$

Cost function for logistic regression

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

$$+ \frac{\partial}{\partial m} \sum_{j=1}^n \theta_j^2$$

- For neural network as we have multiple output nodes:-

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[(y_k^{(i)}) \log \left(h_{\Theta}(x^{(i)})_k \right) + (1 - y_k^{(i)}) \log \left(1 - h_{\Theta}(x^{(i)})_k \right) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{s_l} (\Theta_{j,i}^{(l)})^2$$

\hookrightarrow Here i does not refer to the drawing example

As we need Θ for $\min J(\Theta)$
we need to compute

$$\frac{\partial}{\partial \Theta_j} (J(\Theta))$$

"Backpropagation" is a neural network terminology for minimizing our cost function, just like we were doing with gradient descent in logistic and linear regression.

~~$\delta_j^{(l)}$ is the error for $g^{(l+1)}$ of $\delta_j^{(l)}$ layer l and is calculated as:

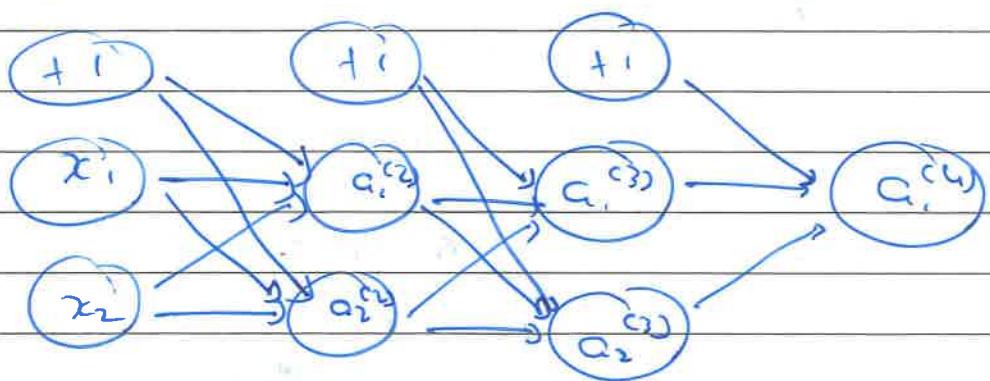
$$\delta_j^{(l)} = a_j^{(l+1)} - y_j^{(l)}$$

\hookrightarrow Predicted value~~

2 p. T. O)

Intuitively, $\delta_j^{(l)}$ is the "error" for $a_j^{(l)}$ (unit j) in layer l .

Eg. 1



$$\delta_1^{(4)} = -y^{(4)} + a_1^{(4)}$$

So first we would do forward propagation :-

$$a^{(0)} = x$$

$$z^{(1)} = \Theta^{(0)} a^{(0)}$$

$$a^{(1)} = g(z^{(1)}) \quad (\text{add } a_0^{(0)})$$

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(1)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)})$$

Now we compute $\delta^{(l)} = a^{(l)} - y$

Now we would do back propagation:-

$$\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) \cdot a^{(l)} \cdot (1 - a^{(l)})$$

The delta values of layer l are calculated by multiplying the delta values in the next layer with the theta matrix of layer l .

We then element wise multiply that with a function called g' or g -prime, which is the derivative of the activation function g evaluated with the input values given by $z^{(0)}$.

The g -prime derivative forms can also be written out as :-

$$g'(z^{(0)}) = a^{(0)} \cdot (1 - a^{(0)})$$

Then we calculate $\delta_{ij}^{(0)}$

First set $\delta_{ij}^{(0)} = 0$ (for all i, j)

$$\delta_{i,j}^{(0)} = \delta_{i,j}^{(0)} + a_j^{(0)} \delta_{i,i}^{(0)}$$

Hence we update our new Δ matrix.

$$\Delta_{i,j}^{(0)} := \frac{1}{m} (\Delta_{i,j}^{(0)} + \lambda \cdot \Theta_{i,j}^{(0)}), \quad j \neq 0$$

$$\Delta_{i,j}^{(0)} := \frac{1}{m} \Delta_{i,j}^{(0)} \quad j = 0$$

The capital - delta matrix Δ is used as an "accumulation" to add up our values as we go along and eventually compute our partial derivative.

$$\text{Then we get } \frac{\partial}{\partial \Theta_{i,j}^{(0)}} J(\Theta) = \Delta_{i,j}^{(0)}$$

~~31~~ "Unroll" into vectors for use in octave

{P.T.O}

function [J, grad] = costFunction(theta)

~~optTheta~~

optTheta = fminunc(@costFunction,
initialTheta, options)

Neural Network (L=4)

→ $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices ($\Theta_1, \Theta_2, \Theta_3$)

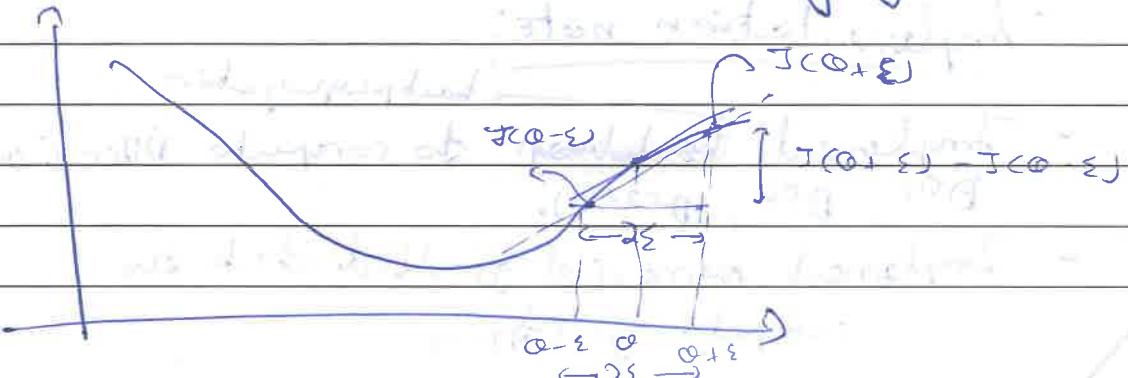
→ $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D_1, D_2, D_3)
gradient

Example

$S_1 = 10, S_2 = 10, S_3 = 1$
 $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$
 $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$

thetaVec = [theta1(:); theta2(:); theta3(:)];
DVec = [D1(:); D2(:); D3(:)];
theta1 = reshape(thetaVec(1:110), 10, 11);
theta2 = reshape(thetaVec(111:220), 10, 11);
theta3 = reshape(thetaVec(221:231), 1, 11);

Numerical estimation of gradient



$$\frac{\partial J(\Theta)}{\partial \Theta} = \frac{J(\Theta + \Sigma) - J(\Theta - \Sigma)}{2\Sigma}$$

Parameter vector Θ

$\Theta \in \mathbb{R}^n$ (E.g. Θ is "unrolled" version of $(\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)})$)

$$\Theta = [\Theta_1, \Theta_2, \Theta_3, \dots, \Theta_n]$$

$$\frac{\partial J(\Theta_1)}{\partial \Theta_1} = \frac{J(\Theta_1 + \Sigma, \Theta_2, \Theta_3, \dots, \Theta_n) - J(\Theta_1 - \Sigma, \Theta_2, \Theta_3, \dots, \Theta_n)}{2\Sigma}$$

$$\frac{\partial J(\Theta)}{\partial \Theta_2} = \frac{J(\Theta_1, \Theta_2 + \Sigma, \Theta_3, \dots, \Theta_n) - J(\Theta_1, \Theta_2 - \Sigma, \Theta_3, \dots, \Theta_n)}{2\Sigma}$$

$$\frac{\partial J(\Theta)}{\partial \Theta_n} = \frac{J(\Theta_1, \Theta_2, \dots, \Theta_{n-1} + \Sigma) - J(\Theta_1, \Theta_2, \dots, \Theta_{n-1} - \Sigma)}{2\Sigma}$$

for $i = 1:n$

$$\text{thetaplus} = \text{theta};$$

$$\text{thetaplus}(i) = \text{thetaplus}(i) + \text{EPSILON};$$

$$\text{thetaminus} = \text{theta};$$

$$\text{thetaminus} = \text{thetaminus}(i) - \text{EPSILON};$$

$$\text{grad Approx}(i) = \frac{(J(\text{thetaplus}) - J(\text{thetaminus}))}{(2 * \text{EPSILON})}$$

end i

check that $\text{grad Approx} \approx \text{DVec}$

Implementation notes:-

- Implement ~~backward pass~~ to compute DVec (curls $D^{(1)}, D^{(2)}, D^{(3)}$).
- Implement numerical gradient check to compute grad Approx

- make sure they give similar values
- Turn off gradient checking. Using backprop code for learning.

↳ make sure to disable gradient checking code before training your classifier because it is very slow

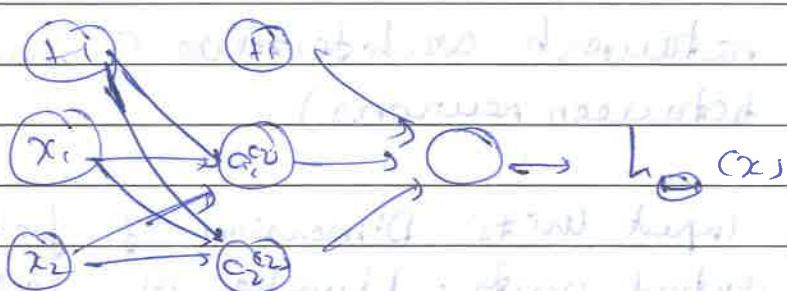
~~1~~ 1 Random initialization

Initial value of Θ

For gradient descent and advanced optimization methods, need initial value of Θ .

Zero initialization would not work
~~1~~ unlike logistic regression.

because:



after each update we will get
 $a_1^{(2)} = a_2^{(2)}$, also $\beta_1^{(2)} = \beta_2^{(2)}$

After each update, parameters corresponding to inputs going into each of two hidden units are identical

{P=7.07}

To solve this we do:
 Random initialization & (symmetry breaking)

→ Initialize each $\Theta_{ij}^{(0)}$ to random value in $(-\alpha, \alpha)$
 $(i.e. -\alpha \leq \Theta_{ij}^{(0)} \leq \alpha)$

Eg. \rightarrow random (10x11) matrix between $(0, \alpha)$

$\Theta_{\text{theta1}} = \text{rand}(10, 11) * C2^* \text{INIT-ALPHA}$
 $\Theta_{\text{theta2}} = \text{rand}(1, 11) * C2^* \text{INIT-ALPHA}$

$\hookrightarrow [-\alpha, \alpha]$

$\Theta_{\text{theta2}} = \text{rand}(1, 11) * C2^* \text{INIT-ALPHA}$
 $\Theta_{\text{theta1}} = \text{rand}(10, 11) * C2^* \text{INIT-ALPHA}$

~~→~~ Putting it all together.

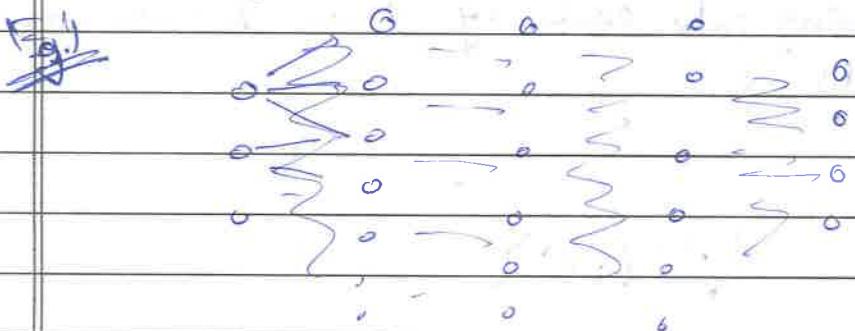
(i) Pick a network architecture (connectivity pattern between neurons)

No. of input units: Dimension of features ~~or~~

No. of output units: Number of classes.

Reasonable default :-

1 hidden layer, or if > 1 hidden layer,
 have same no. of hidden units in every
 layer (usually the more the better)



- (ii) Randomly initialize weights.
- (iii) Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$.
- (iv) Implement code to compute cost function $J(\Theta)$.
- (v) Implement backpropagation to compute partial derivatives $\frac{\partial}{\partial \Theta^{(0)}} J(\Theta)$
 $\frac{\partial}{\partial \Theta^{(1)}} J(\Theta)$

for $i = 1:m$

perform forward propagation & backpropagation using example $(x^{(i)}, y^{(i)})$

get activations $a^{(l)}$ and delta terms

$\delta^{(l)}$ for $l=2, \dots, L$.

$$\delta^{(l)} := \delta^{(l-1)} + \delta^{(l+1)} (a^{(l)})^T$$

\vdots

compute $\frac{\partial}{\partial \Theta^{(0)}} J(\Theta)$

$$\frac{\partial}{\partial \Theta^{(0)}_{i,k}}$$

- (vi) Use gradient checking to compare $\frac{\partial}{\partial \Theta^{(0)}} J(\Theta)$

computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.

→ Then disable gradient checking code.

- (vii) Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ .

by the way $J(\Theta)$ → non convex

(This one can end up in a local minimum instead).
 (But mostly here that is not a problem.)

3. P.T.O

Diagnosing can take time to information out, but detecting so can be a way good we of your time performance.

A test that you can run to gain through what is the cause of your illness a diagnosis algorithm, and gain quickness as do the test to improve it.

Machine learning diagnosis:-

Adding polynomial features. (x^2, x^3, \dots)
Getting additional features.
Smaller sets of features
by joining more descriptive ontology.
What could we do next:-
comes down to in if model doesn't support, we have to tell the user.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

So we can see regression the model of learning process.

Suppose you have implemented sequential
Dealing a learning algorithm.

Date: / /

papergrid

10. L-4

(207.) $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5$
 (208.) $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5$
 (209.) $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5$
 (210.) $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5$

Without any x^5 the boat would sink

1. $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4$
2. $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5$
3. $h_0(x) = 0. + 0.1x + 0.2x^2 + 0.3x^3 + 0.4x^4 + 0.5x^5 + 0.6x^6$

Model Selection:-

$$\text{Test Error} = \frac{1}{m} \sum_{i=1}^m \text{error}(h_0(x_i))$$

The average test error is the measure of fit.

$$\text{error}(h_0(x, y)) = \left\{ \begin{array}{l} 0 \text{ otherwise} \\ 1 \text{ if } h_0(x) \geq 0.5 \text{ and } y = 1 \\ 0 \text{ if } h_0(x) < 0.5 \text{ and } y = 0 \end{array} \right\}$$

Q11) Which of the following is correct?

→ For classification - misclassification error

→ Mean percentage error from learning data

→ If data in the train set and test set

Evaluation using ~~hyperfine~~ ~~hyperfine~~

$I_a(\theta) = I_{\text{max}} \cos(\theta)$

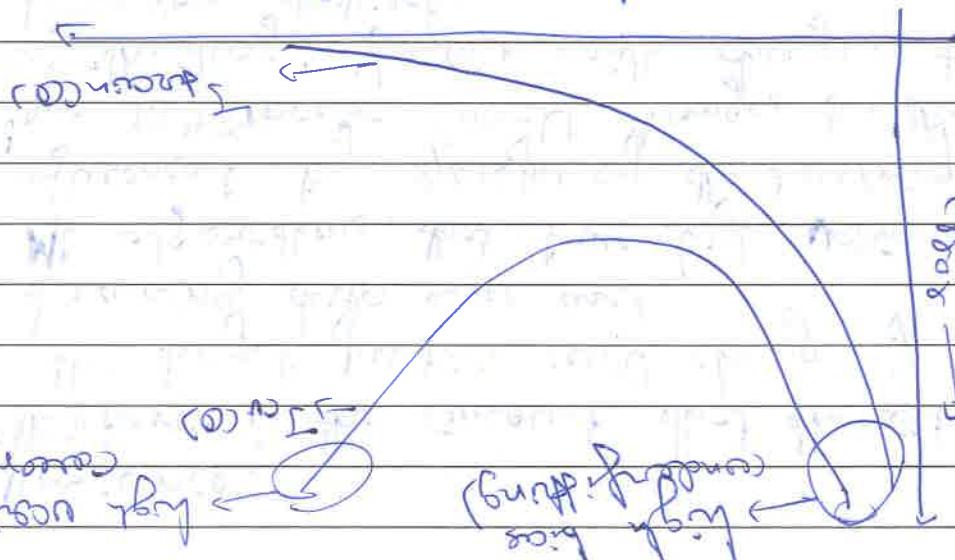
Max. intensity: $I_a(\theta) = I_{\text{max}}$ when $\theta = 0$

$I_a(\theta) \approx I_{\text{max}}(0)$

Intensity variation: $I_a(\theta) = I_{\text{max}} - \text{undiffracted}$

(p)

→ degree of polarization



Zero reflection case: $I_a(\theta) = I_{\text{max}}(0) - I_{\text{max}}(0) = 0$

Max. $I_a(\theta)$

Max. reflection case: $I_a(\theta) = I_{\text{max}}(0) - I_{\text{max}}(0) = I_{\text{max}}(0)$

Max. $I_a(\theta)$

- Find the reflected monochromatic degree of polarization $I_a(\theta)$ when reflection using the formula $I_a(\theta) = I_{\text{max}} \cos(\theta)$ (p = 100% from papergrid)
- Max. reflection using the reflection reflection case using the formula $I_a(\theta) = I_{\text{max}} - I_{\text{max}} = 0$ (p = 100% from papergrid)
- Max. reflection using the reflection reflection case using the formula $I_a(\theta) = I_{\text{max}} - I_{\text{max}} = I_{\text{max}}$ (p = 100% from papergrid)

Max. reflection using the reflection reflection case using the formula $I_a(\theta) = I_{\text{max}} - I_{\text{max}} = I_{\text{max}}$ (p = 100% from papergrid)

- Max. reflection using the reflection reflection case using the formula $I_a(\theta) = I_{\text{max}} - I_{\text{max}} = I_{\text{max}}$ (p = 100% from papergrid)

• Higher the Bias less flexibility

• Because even the simplifying assumption made by a model do not do justice to the situation as per the lesson.

• Log of base: - Design True K-means

• Log of K-means and K-means without modification.

• Log of K-means without modification, logistic regression

Q. How does - Suggest your assumptions about the form of the deposit function.
A. Hgcl. Br₃ :- Suggest more assumptions about the form of the deposit function.

-: Co_3O_4 is magnetic.

10. L. d)

separable

$$\int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx = \int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx$$

$$\int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx = \int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx$$

$$\int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx = \int_{\text{left}}^{\text{right}} (h_0(x) - g_0(x)) \, dx$$

when \rightarrow small \rightarrow error of fitting

when \rightarrow large \rightarrow underfitting

$$h_0(x) = Q_0 + Q_1 x + Q_2 x^2 + Q_3 x^3 + Q_4 x^4$$

linear regression with sigmoid function

as low bias

non-linear ml algorithm:

high bias

linear ml algorithm - high bias &

should return 0 with 0.583

basic could increase then again

vergence

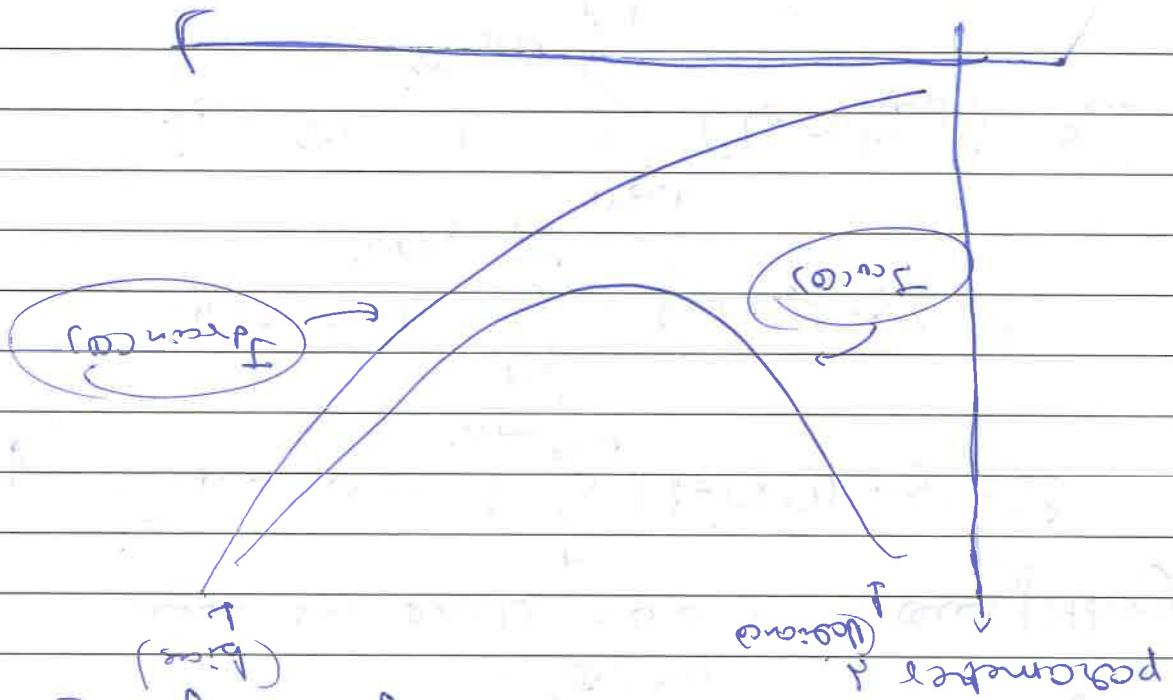
Iteration with small bias & low

Date: / /

papergrid

(P.T.O)

→ ↘



Process (Time) as a function of the performance

Let say the peak \Rightarrow $J_{avg} \rightarrow$ Total error \rightarrow J_{avg}

$$J_{avg} = \frac{1}{2} \int_{-\infty}^{\infty} J_{avg}(t) dt$$

$$J_{avg} = \frac{1}{2} \int_{-\infty}^{\infty} \frac{1}{2} \left[\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(t-\mu)^2}{2\sigma^2}} + \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(t-\mu)^2}{2\sigma^2}} \right] dt$$

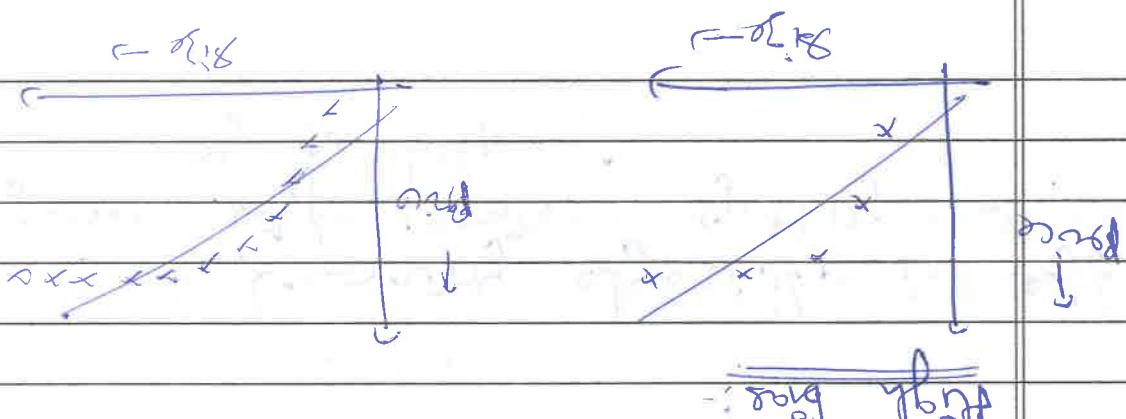
$$J_{avg} = \frac{1}{2} \left[\frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt + \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt \right]$$

Model :- $J_{avg} = 0.0 + 0.2x + 0.2x^2$

Choosing the model which is performed -

Date: / /

papergrid



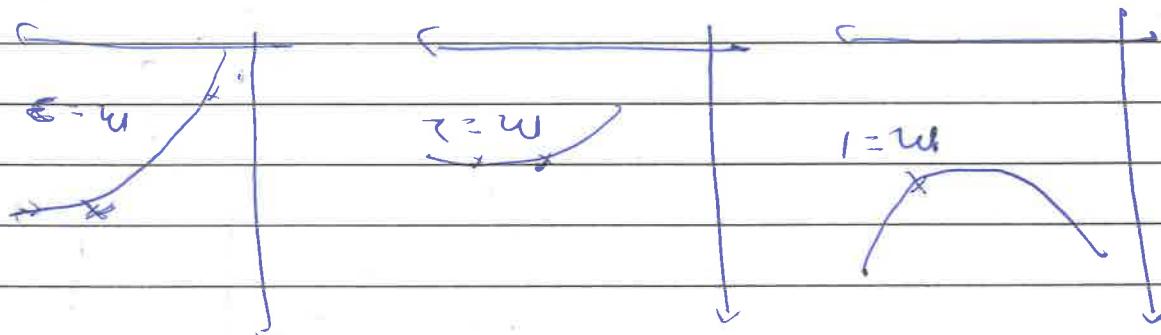
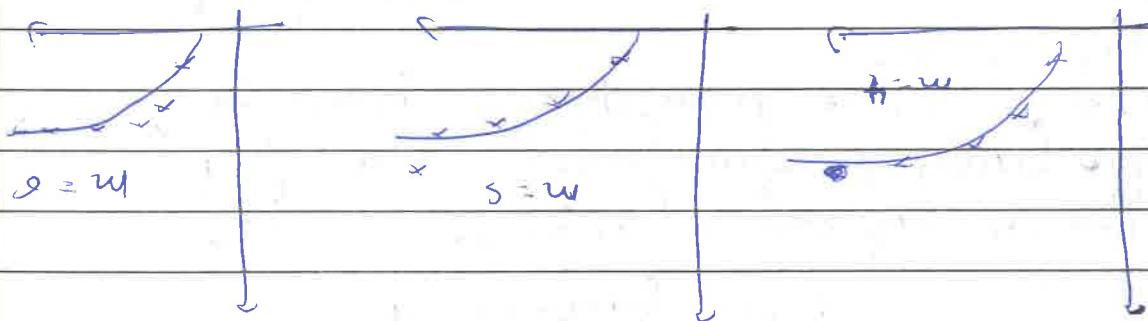
\rightarrow (deflecting end)

$I_{\text{deflection}}$

I_{cu}

I_{cu}

So we can see that as m increases the deflection curve of the beam



$$h_{02} = Q_0 + Q_1 x + Q_2 x^2$$

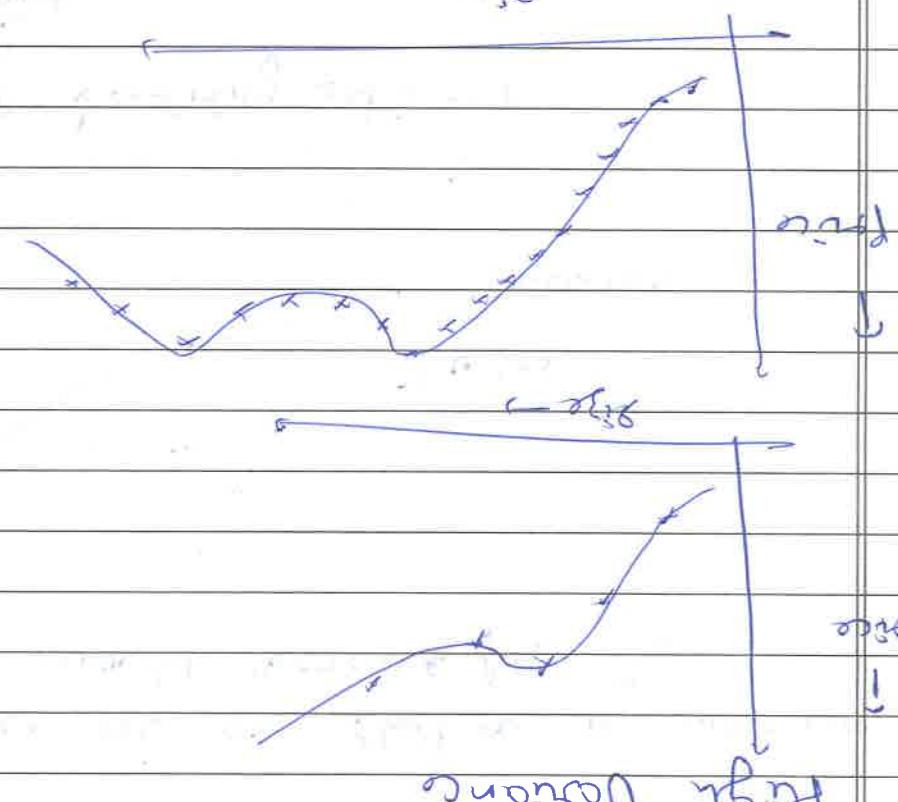
deflecting curve

Date: / /

papergrid

→ a learning algorithm is for different types of learning tasks, getting more data from the task to help the model learn.

→ $\alpha_{1,2}$



out of following

High variance

→ a learning algorithm is suffering from high bias, getting more data will not help much as we can see that the curve is following the training data.

→ $\alpha_{2,3}$ for α deep learning

out of following

Deep learning

Deep learning

get stuck to be flat



Date: / /

papergrid

Locality

Aug 18 Business Plan

A hand-drawn graph on lined paper. A blue curve starts at a point labeled 'top' and increases towards the right. The curve is labeled 'sales over time'.

m (describing)

(20)

20

20

20

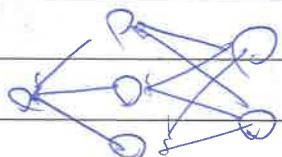
20

• x = feature of small specimens (or not) $y =$ specimen ID \rightarrow x : character 100 \rightarrow $y =$ specimen ID \rightarrow x : character 100 \rightarrow the q of specm / not specm. \rightarrow food, body, dietcond \rightarrow indicate specm \rightarrow dietcond, heat \rightarrow not specm

(containing the following words)

Chitosan polymers: more prone to some
of the same "age" and nutritional

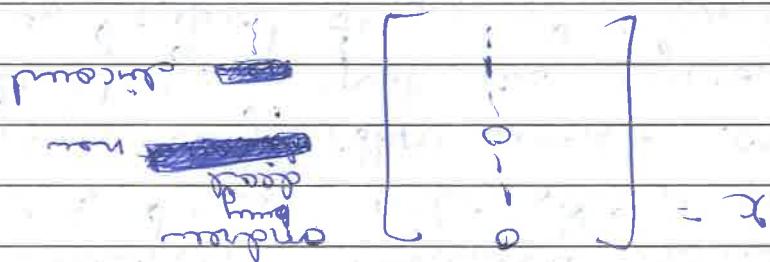
Contour of the mouth



more prone to cognitive bias.
Since neural network of some people
tend to make more mistakes and supporting

some word!
and "discuss". It's directed to the
music body. e.g. should "discuss"
↳ Develop soft/hatched footprints for
outdoor discourses
→ effect multiplications (e.g. multigenerations)
→ Develop soft/hatched always thinner than the
signature (from small to large)
footprints based on small trampling
→ Develop soft/hatched
→ could lost of dots
how to see
How to spread your things to make it!

like 100 words
in drawing set, together than mainly
with drawing in which c 10,000 to 50,000
Note: → in practice, take most figures



End of the week! See you,

See you next!

To:

(- 20m)

Fig. 9. 9 spam e-mail

Date: 11

papergrid

Chores :- 31

Household work: 23

House: 4

Chores: 12 miss drawing

Household work: 23 some of that

Chores: 12 draw them correctly

Household work: 23 house holded the

Chores: 12 (jacket) out

Household work

Household work: 23

Household work: 23 small if it is

Household work: 23 house hold on

Household work: 23 house hold on

Household work: 23 house hold on

Household work: 23

Date: / /

papergrid

just past it (y = 0) more often.
it might just mean that now  can't do this
we get a 3-1 away, it's also closest
but supporters who do some changes so that
they can't go to the away game, this also doesn't

work, as quick as the original number
that we have imagined double the number of
but it's only 0.5%. patients then it seems
that set  more supporters we get 1.1. after an

(then  can't $y = 0$ otherwise)

(curse a classification example)

ask of how will we implement the new example
but, so that (0) is no longer a good option
that would now significantly for the first
that we may just do some changes of our
new footnotes by reclassifying the first set
that (0) is because it's the default
that the first date which is now the complete
date and the community is in (0) rather
so far calculating the date using the cause of last action
The recommended approach the prediction

numerical (spomining) prediction :- 32

numerical medical smoking :- 16

(magazine, news, etc)

Directorate health

Highly recall, however, people in
the world do
not do $0.3 \leq 0.5$
so we may do
this they don't think they might do
this how more and we the
only we are absolutely sure be cause
now is we want that we want the



↳ Definition of segregation :- $Q \leq 0.5$



Two pairs difference + False negative

Two pairs difference =

out of two pairs

Percent :- Two pairs difference

Two pairs difference + False positive

Two pairs difference =

bad了好 file

Proportion :- Two pairs difference

(1) Actual Class		(2) Predicted Class	
Actual Class	False Positive	Predicted Class	True Positive
0	1	1	0
1	0	0	1

Proportion (Percent)

Date: / /

papergrid

the software
we bring together & solve

$$E_1 \cdot 8\text{core} = Q \cdot (P + R)$$

Then we come up with $E_1 \cdot 8\text{core}$ *surplus*

actually what also would be very nice
would be ~~position but a good value of~~
OR ~~we would get considerable loss~~
~~that would be very nice.~~

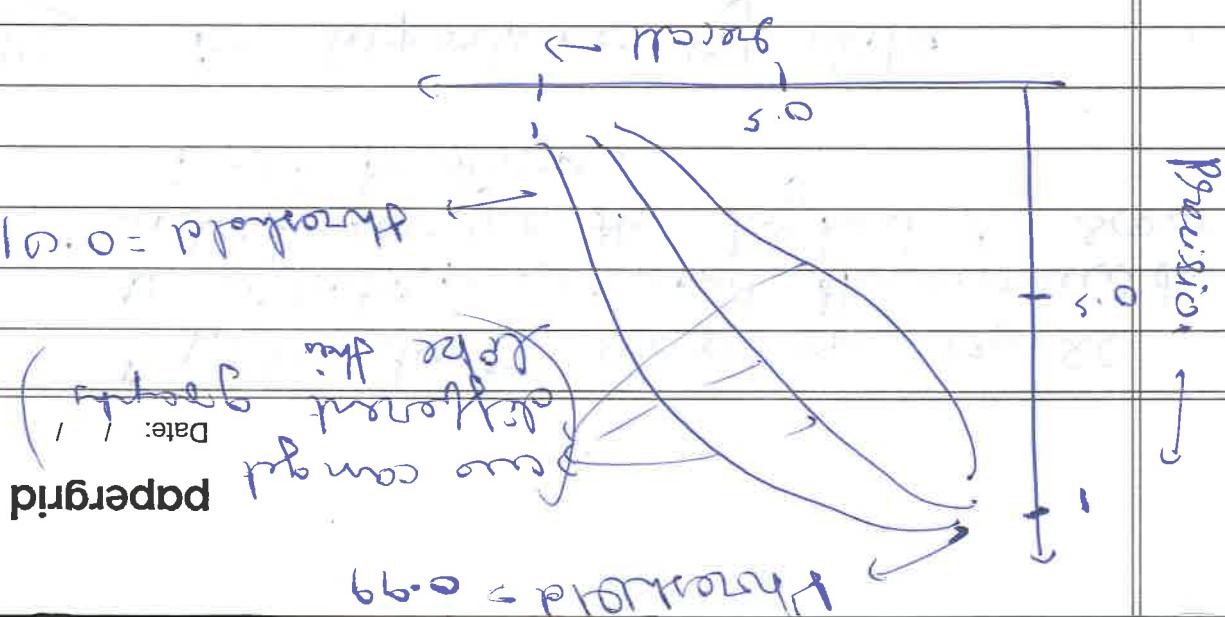
~~we get considerable loss value~~
~~but a good value of power~~
~~that would be very nice.~~

~~value~~
~~result~~

Practical (P)	Practical (R)	Algorithm 1	Algorithm 2	Algorithm 3	0.0352
1.0	0.2	0.7	0.8	0.9	0.125
1.0	0.3	0.7	0.8	0.9	0.444
0.4	0.5	0.5	0.6	0.7	0.808

How do we compute measure/actual

$$E_1 \cdot 8\text{core} (E_1 \cdot 8\text{core})$$



many body. It is world well known and
then how much more companies are
able to produce.

for food. for fuel
for all kinds of
for all kinds of
not be able to help in
then run the ~~delays~~ ~~delay~~ ~~delay~~
but it will have enough power.

88

For all of these reasons there is no
future of all of these things that are
available.

small

they just (1) would also become

so $I_{deion}(0) \approx I_{far}(0)$

because then we would get

the result.

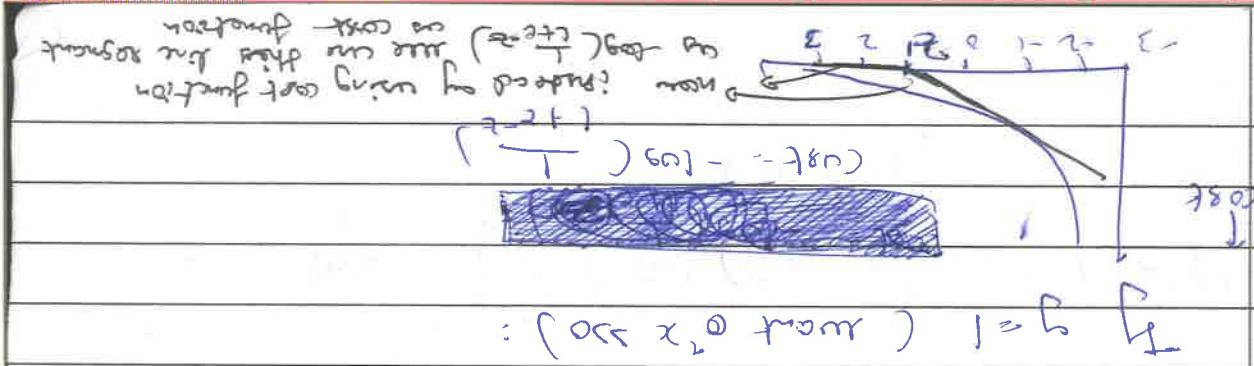
But then it will then if we can be able to
determine the $I_{deion}(0)$ \approx $I_{far}(0)$
small.

large approximation.

Then we would get
total number with many difficulties.
large as approximation with many difficulties
many possibilities (e.g. logistic, sigmoidal)
if we use a learning algorithm with

The information of data

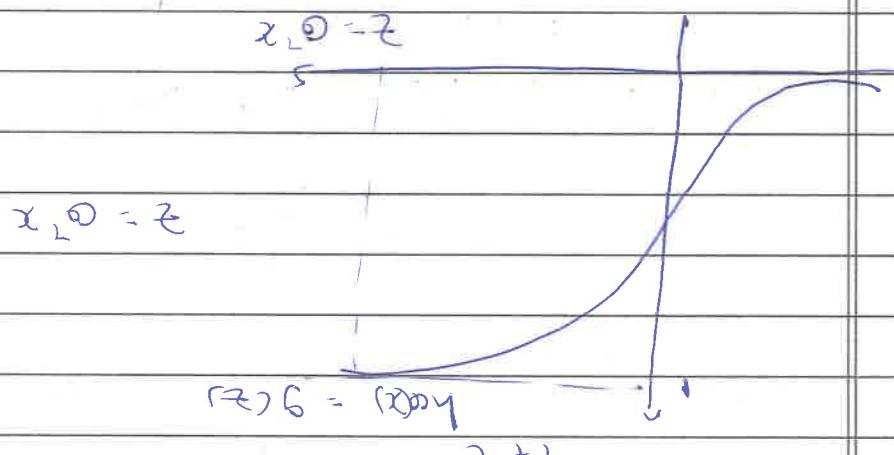
should be chosen.
which we get the highest E score
selection of and the one with
value of E score on $Cross$



$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Example: $-\log(\text{softmax}(x_i)) = -\log(\frac{e^{x_i}}{\sum_j e^{x_j}}) = \log(\sum_j e^{x_j}) - \log(e^{x_i}) = \log(\sum_j e^{x_j}) - x_i$

$y=1$, we want $\text{log}(x) \approx 0$, $\text{softmax}(x) \approx 1$, $\text{softmax}(x) \approx 0$



Difference with logit activation

~~softmax with margins~~

and why so we need more parameters
 that cannot produce the right output
 need more features, (x_1, x_2, \dots) .
 the output and is not probabilistic so
 others that would be able to predict
 should who be predicted with the person
 we could do. This that is a human

Date: / /

papergrid

$$\text{Cost} = \min \left(\left[\sum_{i=1}^n f_i (0.87, (0.7x_{i1}) + (1-y_{i1})0.87, (0.7x_{i2})) \right] \right)$$

what is the goal of
 $A + B$ for such
 but now we should do

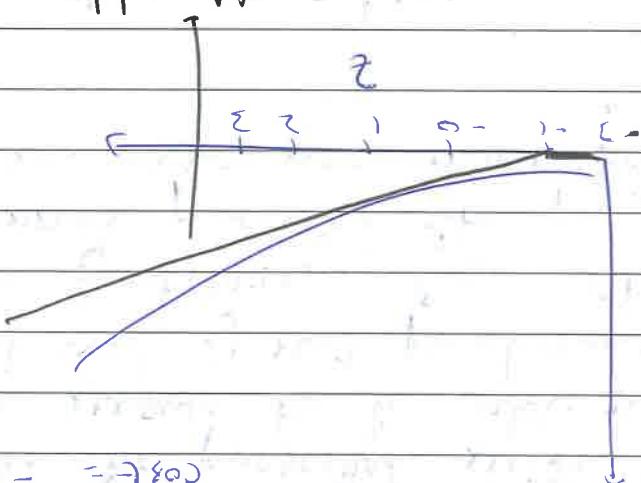
$A + B$ what is the goal of
 before for solving problem we do

working of the sum in machine

$$\text{Cost} = \sum_{i=1}^n y_{i1} (0.87, (0.7x_{i1}) + (1-y_{i1})0.87, (0.7x_{i2}))$$

for solving sum in machine

use this as cost function



$$\text{Cost} = -\log(1 - 1 + e^{-x})$$

$$y = e^{-x} (\text{cost of } x \text{ cost of } y)$$

Lazy migration (distribution) :-

Hot code { i.e. often used }
if $\Theta(x) \leq \epsilon$

Date: / /

papergrid

» [a, b] = Square and Cube (5)

» a

$$\text{a} = 25$$

» b

$$\text{b} = 125$$

Function $J = \text{cost function}(X, y, \theta)$

1. X is the "design matrix" containing our training example.

1. y is the class labels.

$m = \text{size}(X, 1)$; 1. number of training examples.
predictions = $X \cdot \theta$; 1. predictions of hypothesis
on all m

examples

$$\text{sqError} = (\text{Predictions} - y)^2$$

$$J = \frac{1}{2} (1/m) \cdot \text{sum}(\text{sqError});$$

```
>> indices = 1:10;
>> for i=indices,
    disp(i);
end;
```

>> i = 1;

```
>> while i <= 5,
    v(i) = 100;
    i = i + 1;
end;
```

>> i = 1;

```
>> while true,
    v(i) = 999;
    i = i + 1;
    if i == 6,
        break;
    end;
end;
```

>> v(1) = 2;

```
>> if v(1) == 1,
    disp('The value is one');
else if v(1) == 2,
    disp('The value is two');
else
    disp('The value is not one or two');
end;
```

The value is two

function y = squareTheNumber(x)

$$y = x^2$$

function [y1, y2] = SquareAndCube(x)

$$y1 = x^2;$$

$$y2 = x^3;$$

```
→ y2 = cos(2 + pi * 4 * t);  
⇒ [redacted] hold on;  
⇒ plot(t, y2, 'y');  
⇒ xlabel('t');  
⇒ [redacted] ylabel('y');  
⇒ title('---'); legend('sin', 'cosine');  
⇒ figure();  
⇒ plot(t, y1); [redacted] figure(1);  
⇒ plot(t, y2); figure(2);  
⇒ cd('C:\Users\...');  
⇒ print -dpng 'myplot.png'  
⇒ close y. This could cause the figure to  
go away.  
⇒ subplot(1, 2, 1); y. divides plot a 1x2  
grid, uses first element.  
⇒ plot(t, y1);  
⇒ subplot(1, 2, 2); y. uses 2nd element  
⇒ plot(t, y2);  
⇒ plot(t, y2);  
⇒ [redacted] help axis  
⇒ axis([0.5, -1, 1])  
⇒ clf;  
⇒ A = magic(5)  
⇒ imagesc(A)  
⇒ imagesc(A), colorbar, colormap gray.  
⇒ V = zeros(10, 1)  
⇒ for i = 1:10,  
    V(i) = 2^-i;  
    end  
⇒ V  
V =  
 0.5  
 0.25  
 0.125  
 0.0625  
 0.03125  
 0.015625  
 0.0078125  
 0.00390625  
 0.001953125  
 0.0009765625
```

» let $A =$

$$\begin{matrix} 8 & 1 & 6 \\ 3 & 5 & 7 \\ 4 & 9 & 2 \end{matrix}$$

↓ diag

» $\max(A, \text{col}, 1)$ ↴ returns maximum of each column

Ans =

$$\begin{matrix} 8 & 9 & 7 \end{matrix}$$

» $\max(A, \text{col}, 2)$ ↴ returns maximum of each row.

Ans =

$$\begin{matrix} 8 \\ 7 \\ 9 \end{matrix}$$

» $\max(A(:))$

» $A = \text{magic}(9)$

» $\text{sum}(A, 1)$ ↴ sums up each column of A

» $\text{sum}(A, 2)$ ↴ sums up each row of A

Now if we want to find the sum of the diagonal

» ~~$\text{sum}(\text{sum}(A \cdot \text{eye}(9)))$~~

Ans = 369

» $\text{eye}(3)$

$$\begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix}$$

» $\text{flipud}(\text{eye}(9))$ ↴ this would flip the matrix vertically

$$\begin{matrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{matrix}$$

» $t = [0:0.01:0.98];$

» $y_1 = \sin(2\pi t + 4\pi t);$

» $\text{plot}(t, y_1);$

new

new

new

new

» $a = [1 \ 15 \ 2 \ 0.5]$

» $c < 3$

» $\begin{matrix} 1 & 0 & 1 & 1 \\ 1 & 3 & 4 \end{matrix}$

↳ index of elements less than 3.

» $A = \text{magic}(3)$

$A =$

$\begin{matrix} 8 & 1 & 6 \\ 3 & 5 & 7 \\ 4 & 9 & 2 \end{matrix}$

(sum of every row, column and diagonal adds up to the same number)

» $\text{sum}(a)$

18.500

» $\text{prod}(a)$

15

» $\text{floor}(a)$

$\begin{matrix} 1 & 1 & 5 \\ 2 & 0 & 0 \end{matrix}$

» $\text{ceil}(a)$

$\begin{matrix} 1 & 1 & 5 \\ 2 & 1 & 1 \end{matrix}$

» $\text{max}(\text{rand}(3), \text{rand}(3))$

0.72763

0.78777

0.93872

:

:

:

This would return us a ~~matrix~~ a 3×3 matrix in which each element is the maximum of 2 randomly generated matrix.

$$\text{let } A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix}$$

$$B = \begin{bmatrix} 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{bmatrix}$$

$$\gg C = [A \ B]$$

$$C = \begin{bmatrix} 1 & 2 & 11 & 12 \\ 3 & 4 & 13 & 14 \\ 5 & 6 & 15 & 16 \end{bmatrix}$$

$$\gg C = [A : B]$$

$$C = \begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 11 & 12 \\ 13 & 14 \\ 15 & 16 \end{bmatrix}$$

$$\gg A \leftarrow B$$

$\gg A \leftarrow B$ i.e. element wise multiplication of 2 matrices

$\gg A \leftarrow A^2$ i.e. element wise squaring of A

$\gg 1./A$ i.e. element wise reciprocal

$\gg \log(u)$

$\gg \exp(u)$ i.e. e^u (element wise)

$\gg \sin(u)$

$\gg A'$ (would return transpose of A)

$\gg \max(v)$ i.e. return maximum value of v

$\gg [v, i] = \max(v)$ i.e. returns maximum value as well as index of maximum value

>> $U = \text{vector_name}(1:10) \text{ } \therefore \text{ now}$
we will have the first 10 elements
of vector-name into U.

>> clear \therefore would all the variables.

>> $A = [1 \ 2 \ 3 \ 4 \ 5 \ 6]$

>> $A(3, 2)$ \therefore element in 3rd row 2nd column
6

>> $A(2, :)$ \therefore every element in 2nd row.

>> $A(:, 2)$ \therefore every element in 2nd column.

>> $A([1 \ 3], :)$ \therefore every element in 1st & 3rd row.

>> $A(:, 2) = [10 \ 11 \ 12]$

so now A would become:

1 10
3 11
5 12

>> $A = [A, [100, 101, 102]]$.

now A would become:

1 10 100
3 11 101
5 12 102

>> $A(:)$ \therefore put all elements of A into
a single vector

1
3
5
100
101
102

- rents
- planned
- ind
- ..
- 2) $\text{A} = [3 2; 1 3]$
- 3) $\text{size}(\text{A})$
1 2
- 4) $\text{size}(\text{A}, 1)$ ✓ would return the number of rows.
- 5) $\text{size}(\text{A}, 2)$ ✓ would return the number of columns.
- 6) $\text{U} = [1 2 3 4]$
- 7) $\text{U} =$
1 2 3 4
- 8) $\text{length}(\text{U})$
ans = 4
- 9) $\text{length}(\text{A})$ ✓ returns the longest dimension.
ans = 3
- 10) pwd ✓ would tell us the current directory path.
- 11) ~~cd~~
- 12) $\text{cd} \dots$ ✓ would change directory.
- 13) ls ✓ list files of current directory.
- 14) $\text{load}('name of the file')$ ✓ this would load the file into octave
- 15) who ✓ this would tell us what variables we have instantiated.
- 16) $\text{size}(\dots)$
- 17) whos ✓ would give us the detailed view of all the variables.
- 18) clear variable name.

- onion
- Raw banana
- Buffalo

papergrid

Date: / /

~~keep~~ out of Rs 100000 is invested in 2 shares. first one interest 9% and second 11% per annum. Total interest is 9.3%. final amt invested in each shares?

1000

$C = \text{ones}(1, 3)$

$C = 2^{\text{nd}} \text{ ones}(1, 3)$

$C = \text{zeros}(1, 3)$

ones(1, 3)

$C = \text{rand}(3, 3)$

$C = \text{randn}(1, 3)$

(normal distribution) \rightarrow variance = 1

mean = 0

number of numbers returned

\rightarrow 3x3 matrix of random numbers from 0 to 1

$w = -6 + \text{sqrt}(10) * (\text{randn}(1, 10000))$

here mean would be -6

Variance = 1

\rightarrow this would return 10000 numbers

eye(3)

1 0 0
0 1 0
0 0 1

help eye

① Walk - melan - 2

green grapes.

buffalo meat.

- 1/2 (lotion - 2)

②

ginger.

③

methi sevai.

④

cheese.

⑤

raw - banana

⑥