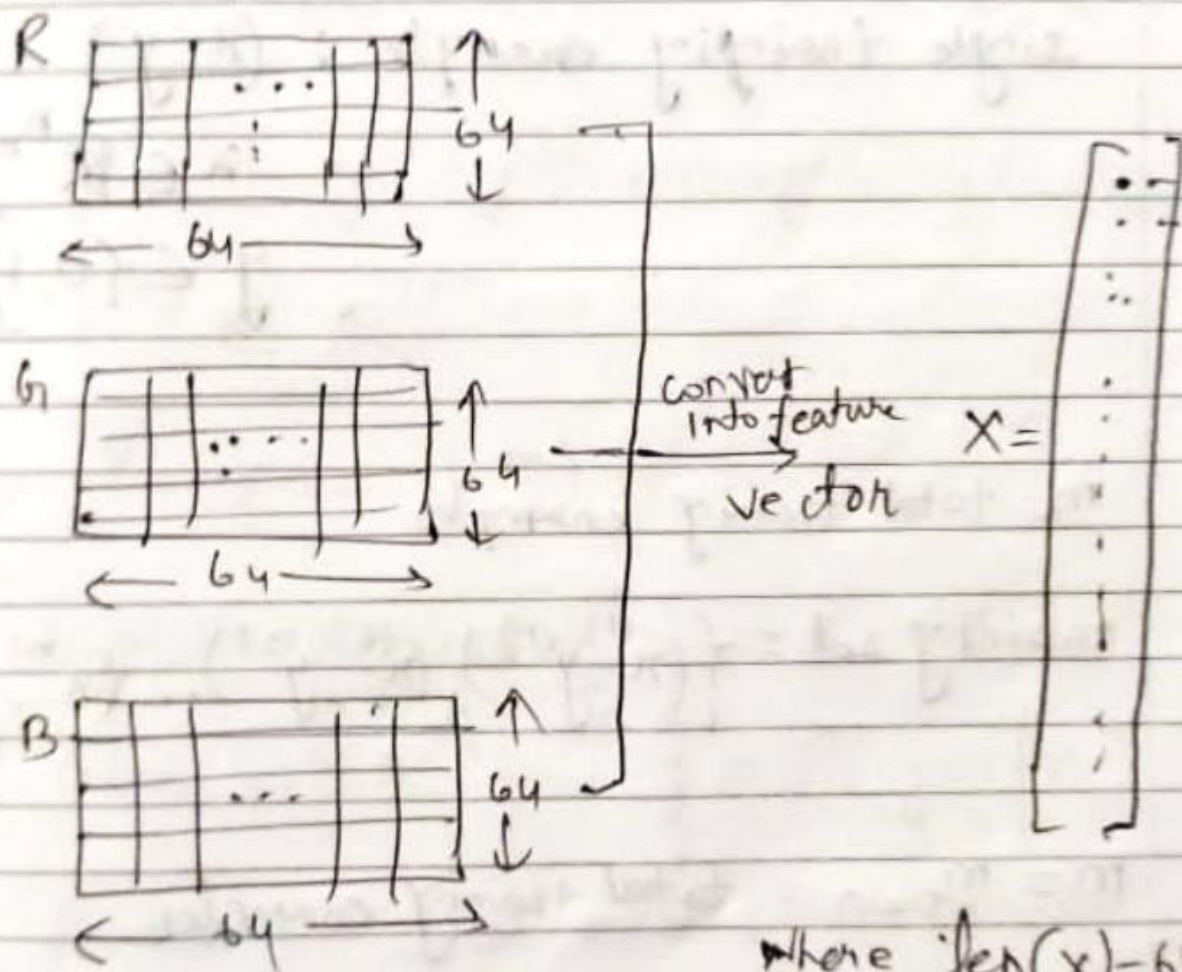
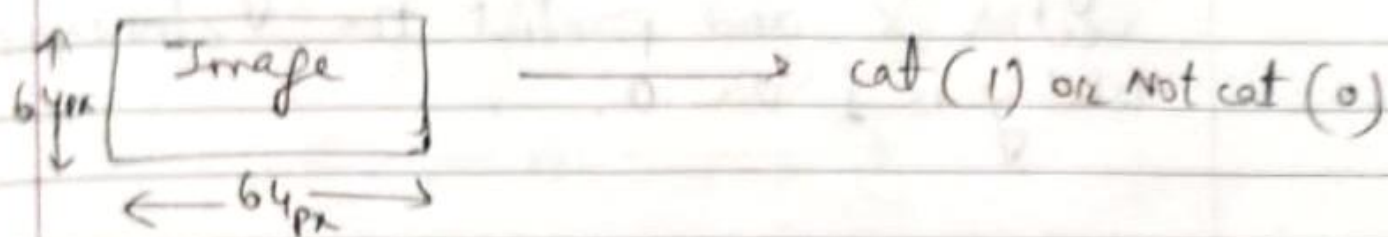


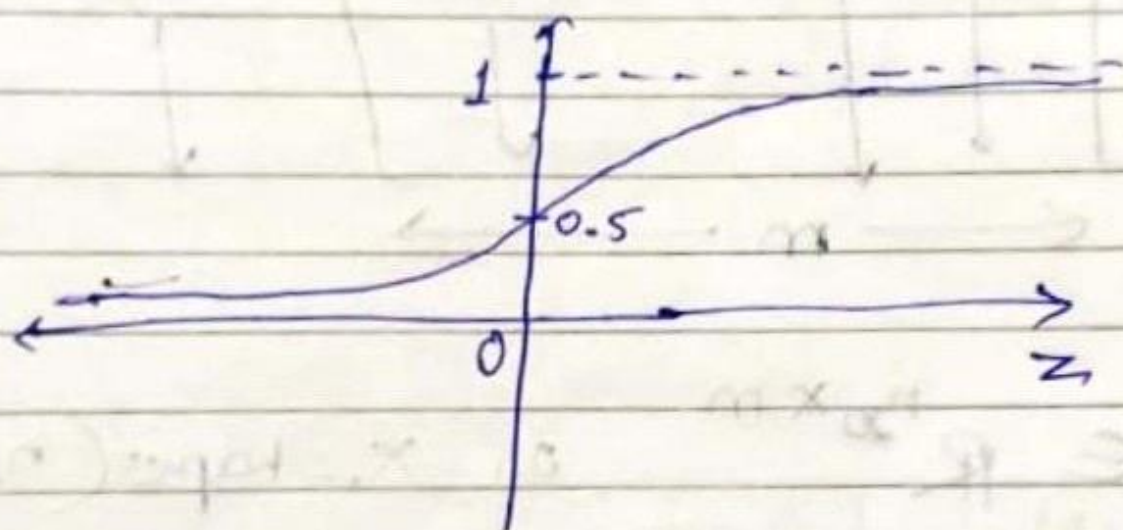
Binary classification :-



$$\text{Where } \text{len}(X) = 64 \times 64 \times 3$$
$$n_X = 12288$$

$$\text{RGB} = (64, 64, 3)$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(w^T x + b)}}$$



Sigmoid
function

Another alternative to b vector is

$x_0 = 1$ so that now $x \in \mathbb{R}^{n+1}$

$$\hat{y} = \sigma(\theta^T x)$$

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

w

will not
use here
in this
class

set of all training examples

PAPERWILL

Date :

Page no.

$$X = \begin{bmatrix} \uparrow & \uparrow & \dots & \uparrow & \uparrow \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} & n_n \\ \downarrow & \downarrow & \dots & \downarrow & \downarrow \\ \leftarrow m \rightarrow & & & & \end{bmatrix}$$

$$X \in \mathbb{R}^{n_n \times m}$$

$$\text{or } X.\text{shape} = (n_n, m)$$

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$

Logistic Regression

$$\text{Given } x \in \mathbb{R}^{n_n}, \hat{y} = P(y=1|x)$$

Parameter θ or w = parameter vector

$$w \in \mathbb{R}^{n_n}, b \in \mathbb{R}$$

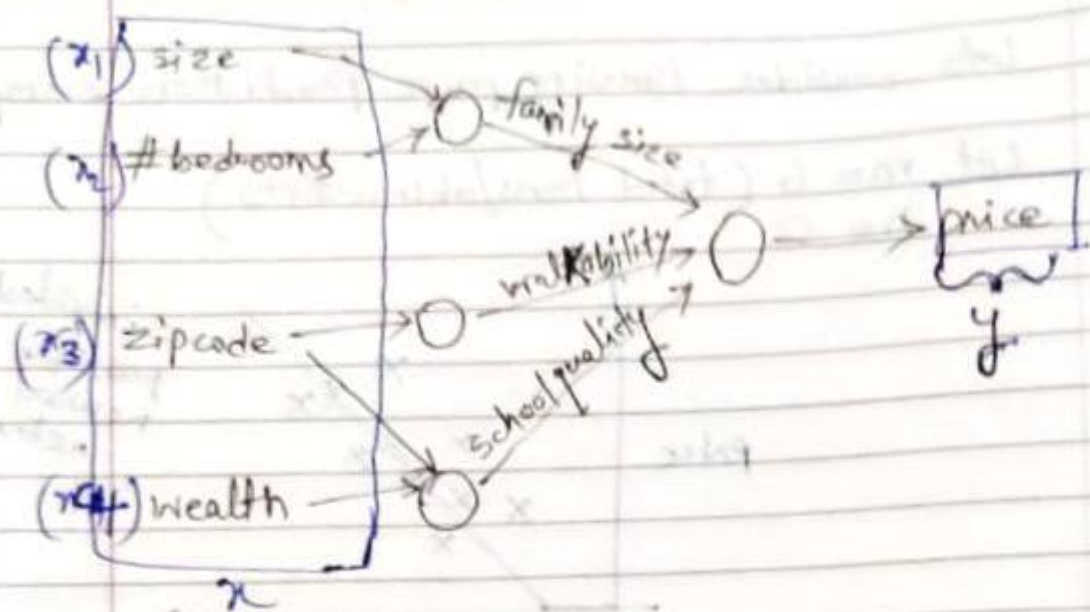
$$\text{Output } \hat{y} = w^T x + b \rightarrow \text{will not work good}$$

$$\text{as } -\infty \leq \hat{y} \leq \infty$$

but we want $\hat{y} \in [0, 1]$

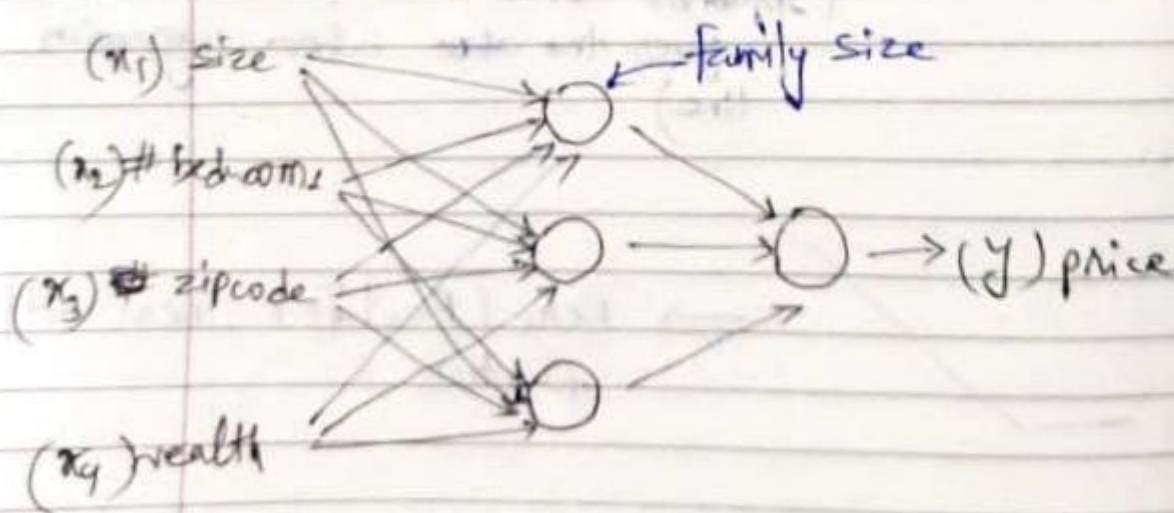
So

$$\hat{y} = \sigma(z) = \sigma(w^T x + b) = \frac{1}{1 + e^{-z}}$$



All we need to give is input x & output y in training set and all the intermediate levels / hidden layers will be taken care automatically.

↓ turns out to become

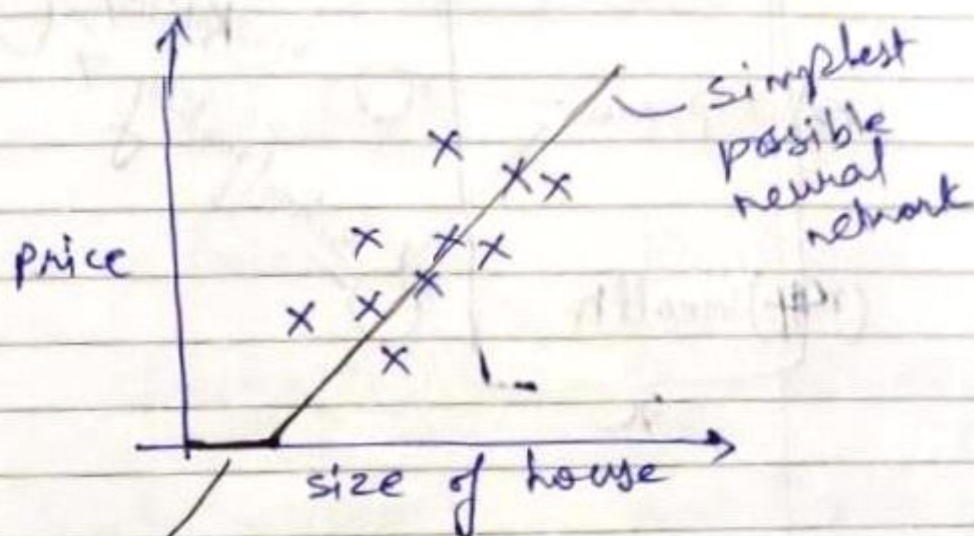


Although family size is just dependent on size of house & # bedrooms we will say well neural network

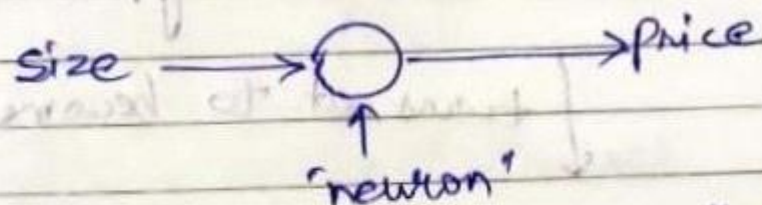
Q What is Neural Network?

Let's consider Housing price prediction example.

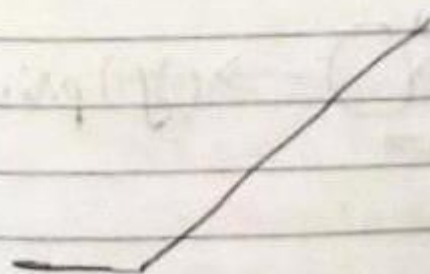
Let $m = 6$ (total rows/observations)
 $n = 2$



prices cannot be negative



(Implement above linear fn that gives the above regression line)



→ Relu (Rectified Linear Unit) function

$$J = \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

Note: $a^{(i)} = \hat{y}^{(i)} = \text{sigmoid}(h_0(x))$

$w_1, w_2, x_1, x_2 \in \mathbb{R}$

Lets Initialize

$J=0; dw_1=0; dw_2=0; db=0;$

Let total features ($n=2$) i.e. x_1, x_2

For $i=1$ to m

$$\left\{ \begin{aligned} z^{(i)} &= w^T x^{(i)} + b \end{aligned} \right.$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J = J + \left[-y^{(i)} \log a^{(i)} - (1-y^{(i)}) \log (1-a^{(i)}) \right]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 = dw_1 + x_1^{(i)} dz^{(i)}$$

$$dw_2 = dw_2 + x_2^{(i)} dz^{(i)}$$

$$db = db + dz^{(i)}$$

may need to write the for loop if
lets say $n=k$ features

for $j=1$ to k

$$\left\{ \begin{aligned} dw_j &= dw_j + x_j^{(i)} dz^{(i)} \end{aligned} \right.$$

$$J = J/m \rightarrow \text{gives } \left(\frac{1}{m} \right) \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$dw_1 = dw_1/m \rightarrow \text{gives } \left(\frac{1}{m} \right) \sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})$$

$$dw_2 = dw_2/m$$

$$db = db/m$$

}

you decide whatever this node "family size" should be

* Supervised learning with neural network :-

| Input (x) | Output (y) | Application |
|--------------------|------------------------|---------------------|
| ① Home features | Price | Real Estate |
| ② Ad, User Info | Click on ad? (0/1) | Online Advertising |
| ③ Image | Object (1, ..., 1000) | Photo tagging |
| ④ Audio | Text transcript | Speech Recognition |
| ⑤ English | Hindi | Machine translation |
| ⑥ Image/Radar Info | Position of other cars | Autonomous driving |

→ Slightly different types of neural networks are used for different applications.

For above examples

- ① Real state (Home features → price)
 - ② Online Advertising (Ad → Click or Not)
 - ③ ~~Image~~ Photo tagging (Image → lots of Images)
 - ④ Speech Recognition → RNN
 - ⑤ Machine translation → RNN
 - ⑥ ~~Image~~ Autonomous driving → Hybrid NN
- } Standard Neural Network

Gradient Descent

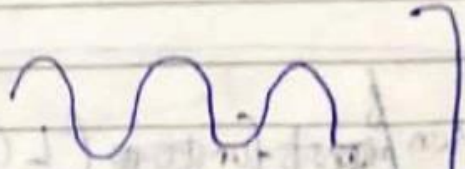
$$\hat{y} = \sigma(w^T x + b), \quad \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

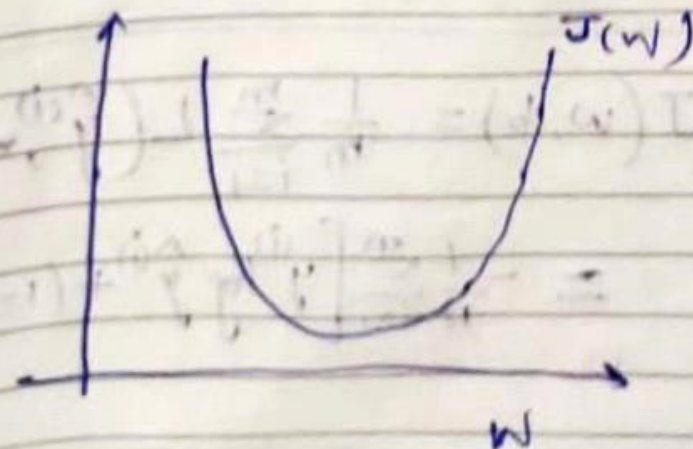
$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

Find w, b that minimizes $J(w, b)$

$J(w, b) \rightarrow$ convex function (U)

[Non-convex f : 

Let's ignore b for now, for simplicity



$$h_{\theta}(x) = \hat{y} \quad y = g(\theta^T x) = \sigma(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

Logistic Regression cost function :-

$$\hat{y} = \sigma(w^T x + b) \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

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

want $\hat{y}^{(i)}$ vs $y^{(i)}$

Prediction for
ith ~~training~~ ~~the~~
training example

label for ith
training example

$$\hat{y} \equiv h_{\theta}(x)$$

Loss function / ~~cost function~~ $(L(\hat{y}, y)) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$
for single training example

Cost function $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$
for entire training set
 $= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$

will try to find w & b that minimizes cost J

So in binary classification our goal is to learn a classifier that can input an image ~~representation~~ represented by feature vector x and predict the corresponding label y is 1 or 0.

Notation we will use :-

single training example : (x, y)

$$x \in \mathbb{R}^{n_x}$$

$$y \in \{0, 1\}$$

m = total training examples

$$\text{training set} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

$m = m_{\text{train}}$ = total training examples

m_{test} = total test examples

Now we have to come up with a cost function that is applicable to logistic regression & is convex f .

→ cost function for logistic regression:-

$$\text{cost}(h_0(x), y) = \begin{cases} -\log(h_0(x)) & \text{if } y=1 \\ -\log(1-h_0(x)) & \text{if } y=0 \end{cases}$$

$\text{cost}(h_0(x), y)$ refers to the cost/penalty the learning algo has to pay when it predicts $h_0(x)$ & actual label is y .

So FINAL vectorized Logistic Regression Implementation

$$Z = W^T X + b \quad (\text{i.e. } \text{np.dot}(W^T, X) + b)$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dW = \frac{1}{m} X dZ^T$$

$$db = \frac{1}{m} \text{np.sum}(dZ)$$

So now Gradient descent code will be

$$W := W - \alpha dW$$

$$b := b - \alpha db$$

↳ With this we just implemented a single iteration of gradient descent without using a single 'for' loop.

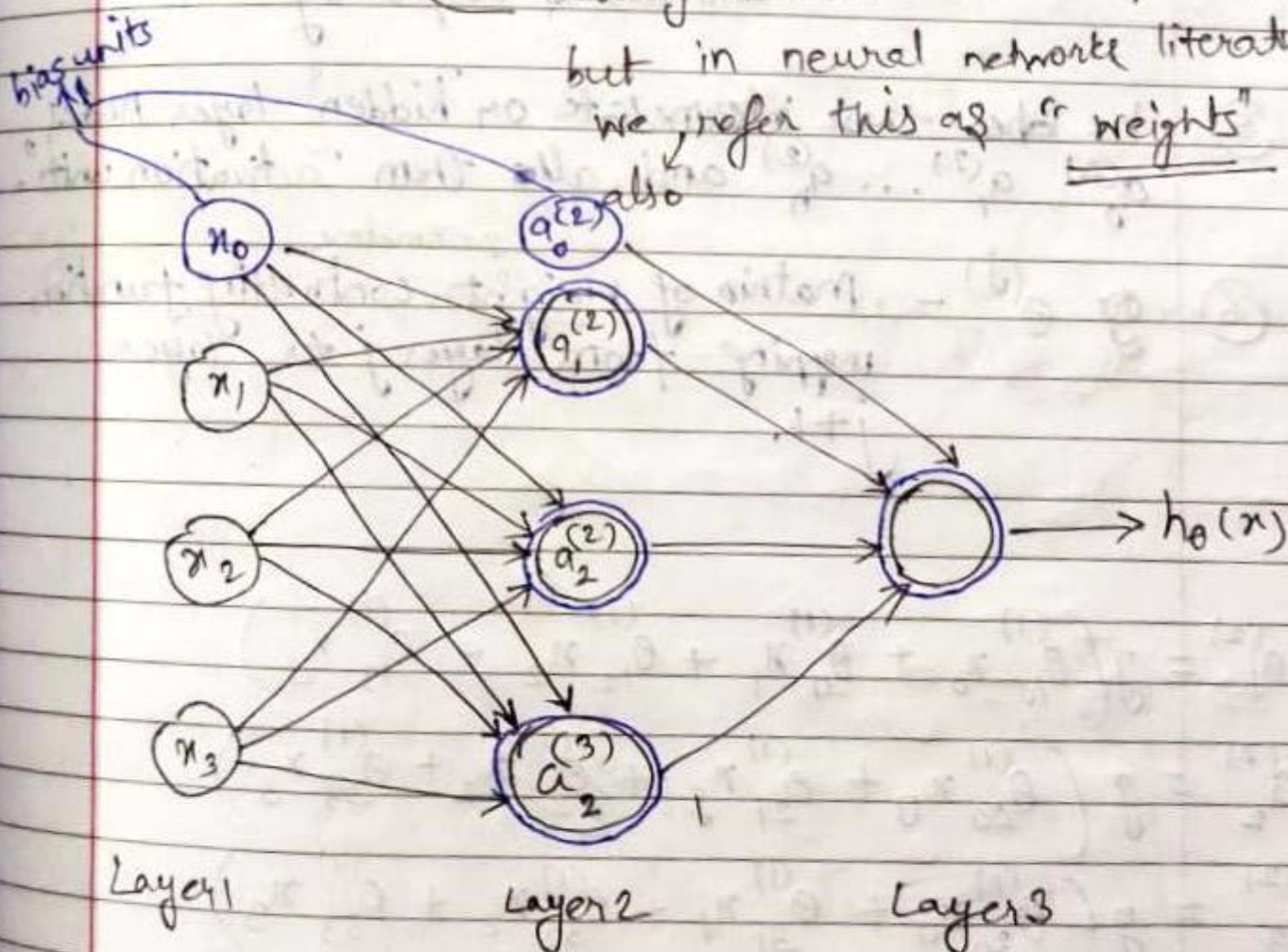
→ But to implement multiple iterations of gradient descent we still need a 'for' loop even in vectorized implementation, for the number of iterations.

bias unit that usually outputs value 1.

→ Sometimes when we talk about neural networks sometimes we will say that this is a neuron with a ~~same~~ sigmoid/logistic activation function.

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

(usually we call this as parameters but in neural network literature we refer this as "weights")



- In neural networks
- ① The first layer (Layer 1) is called input layer
 - ② The final layer (Layer 3) is called output layer because that layer has the neuron that ~~computes~~ outputs the final value computed by the hypothesis.

Neural network:- Much better way to learn complex ~~hypothesis~~ non linear hypothesis even when your input feature space is large

Neuron model:-

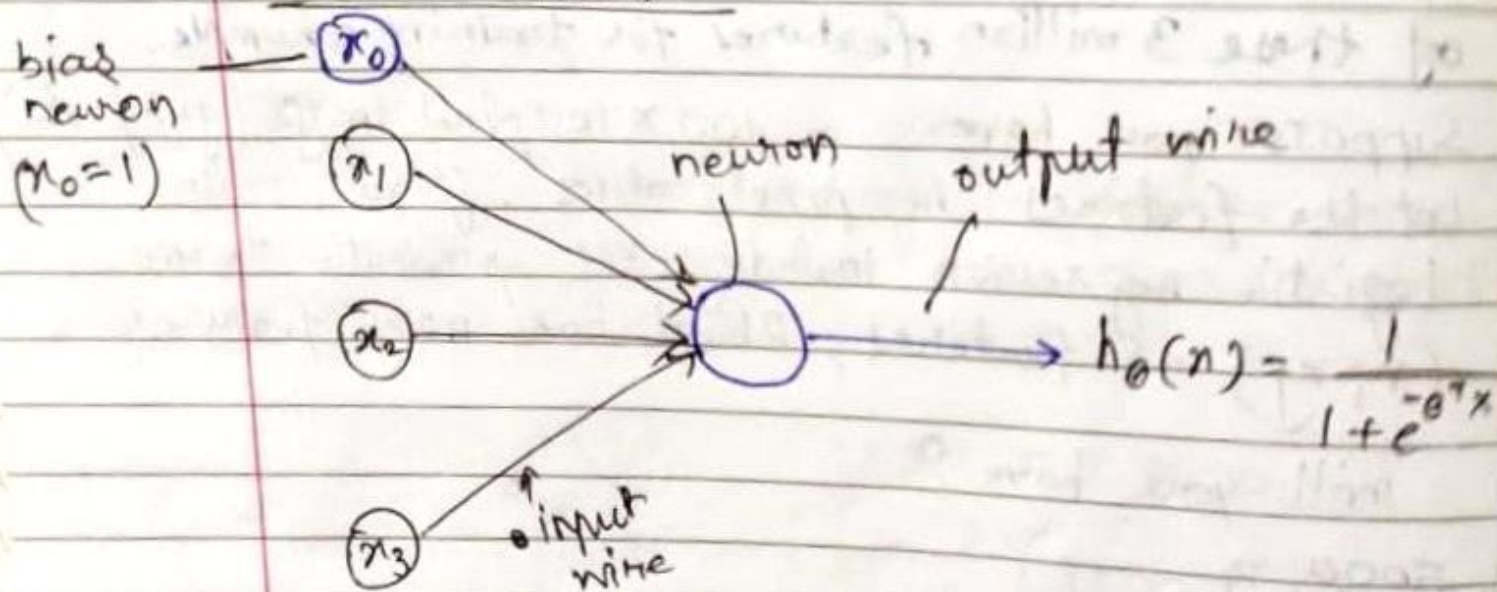
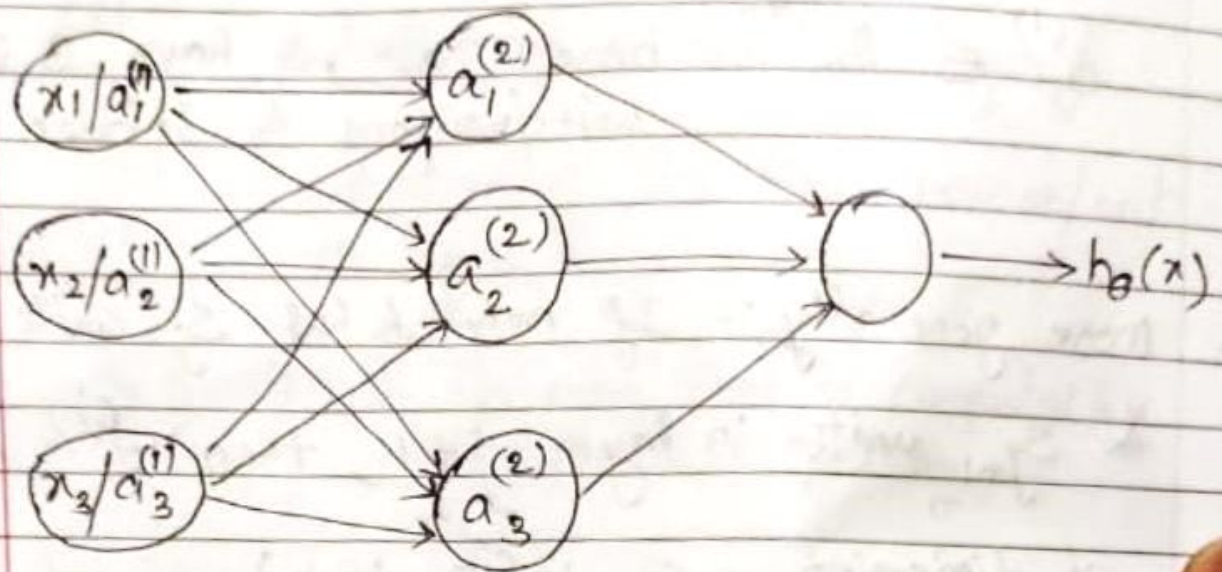


Diagram with single neuron

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

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

Forward Propagation: Vectorized Implementation



Let

$$z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3$$

$$z_2^{(2)} = \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3$$

$$z_3^{(2)} = \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3$$

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

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

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

So these $z_j^{(i)}$ values are just ~~var~~ linear combination of the input values x_0, x_1, x_2, x_3 that go into a particular j th ~~var~~ neuron of i th level

We have 3 input units/neurons
3 hidden units

Date
DELTA Pg No.

③ → Layer 2 is called the hidden layer. We can have more than 1 hidden layers also.

④ $a_i^{(j)}$ — Activation of neuron i in layer j

eg. $a_1^{(2)}$ — activation of the first neuron/unit in layer 2.

Activation → The value that is computed by and that is output by

⑤ We label these intermediate or "hidden" layer nodes $a_0^{(2)}, a_1^{(2)}, \dots, a_n^{(2)}$ and call them "activation units".

⑥ $\theta^{(j)}$ = Matrix of ^{parameters} weight controlling function mapping from layer j to layer $j+1$.

$$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_\theta(x) = a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

clearly from the curve we can observe that

since we draw for

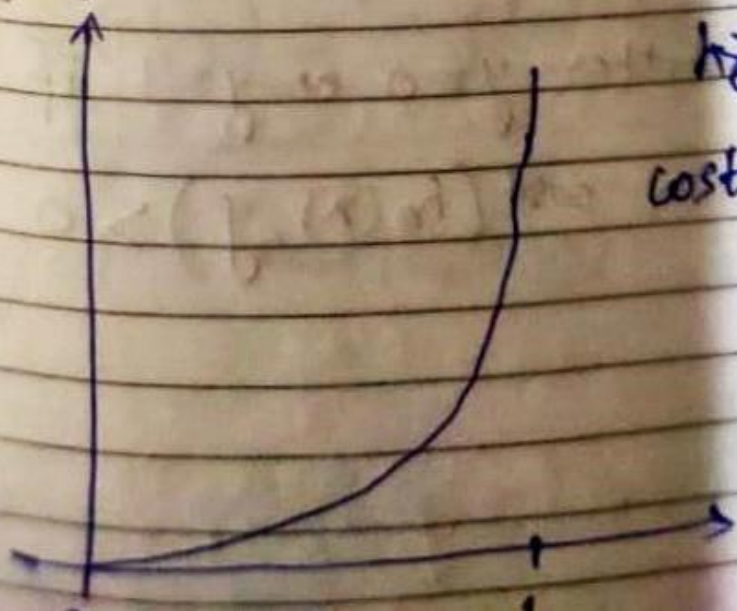
$$\text{cost}(h_0(x), y) = -\log(h_0(x)) \text{ for } y=1$$

so when $h_0(x)=1 \wedge y=1 \Rightarrow \text{cost}(h_0(x), y)=0$

$h_0(x) \rightarrow 0 \wedge y=1 \Rightarrow \text{cost}(h_0(x), y) \rightarrow \infty$

so if we go and say the patient that he has malignant tumor but it is actually benign then the learning algo has to pay ∞ cost/penalty for wrong prediction.

→ If $y=0$:-



$\text{cost}(h_0(x), y) = -\log(1-h_0(x))$
for $y=0$