

# DEEP LEARNING (deeplearning.ai / Andrew Ng)

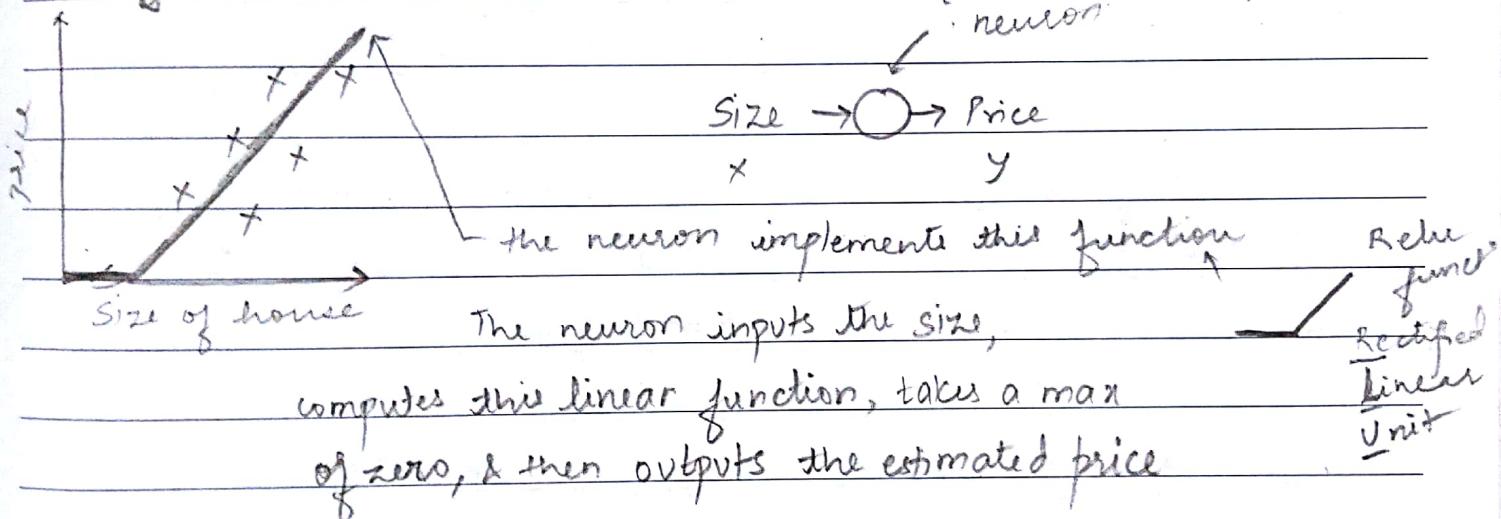
## SPECIALIZATION

### COURSE 1: Neural Networks and Deep Learning

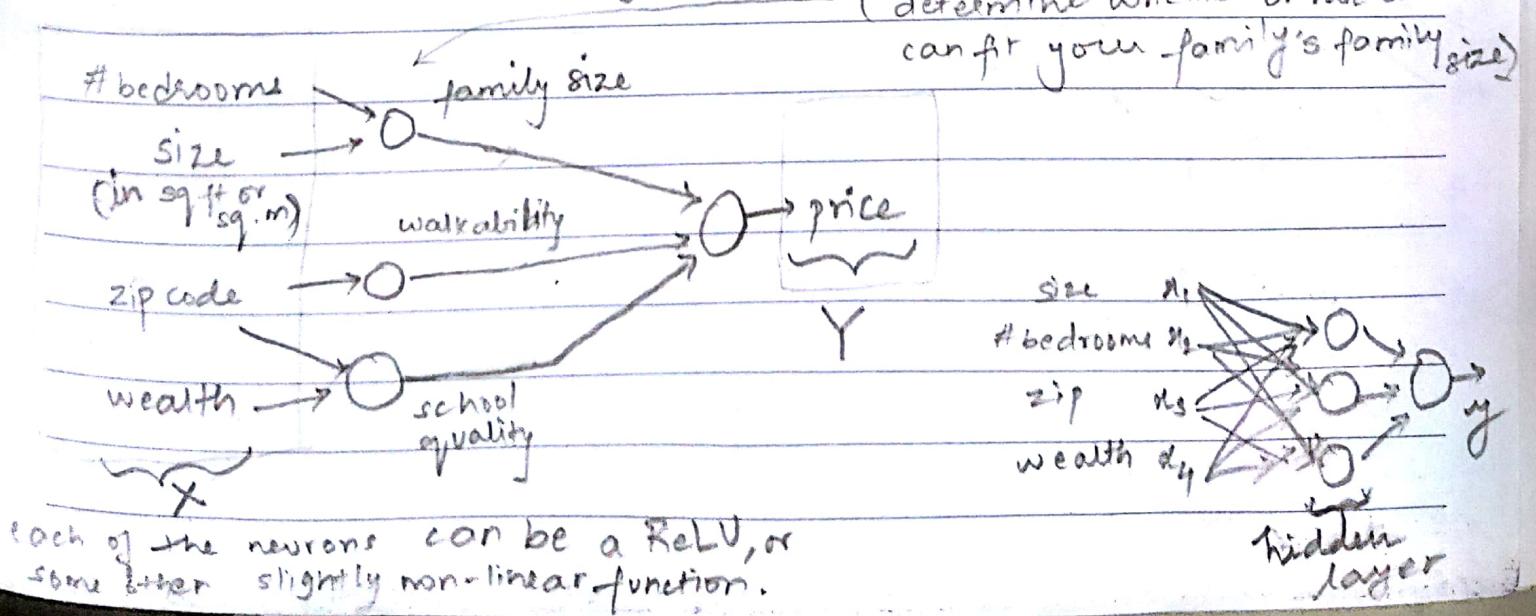
#### Week 1:

- What is neural network?

#### Housing Price Prediction



Stacking many neurons together to get a neural network. e.g:



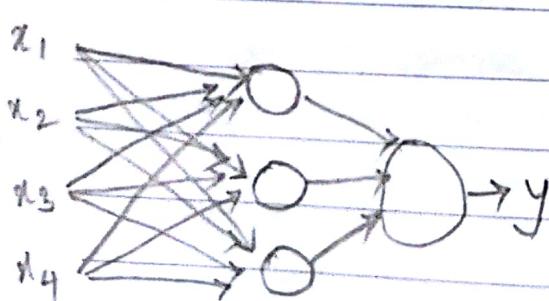
## Supervised Learning with Neural Networks

You have  $i/p(x)$  & want to predict the  $o/p(y)$ , e.g.:

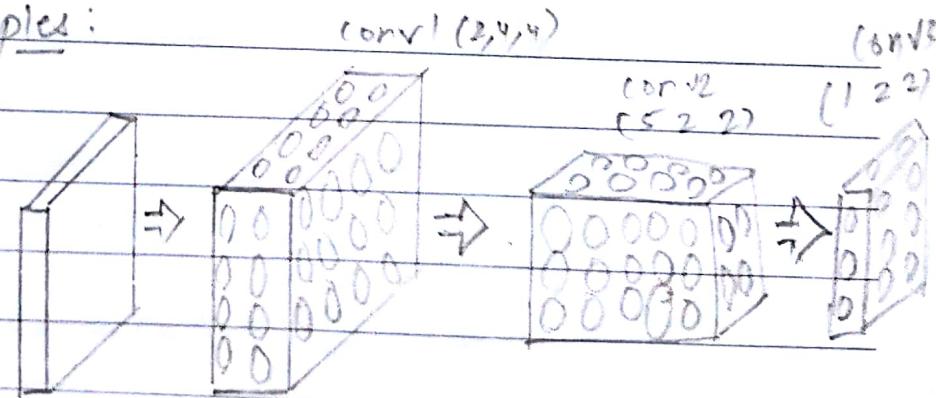
Input ( $x$ )	Output ( $y$ )	Application
Home features	Price	Real Estate
Ad, user info	click on ad? (0/1)	Online Advertising
Image	Object (1..., 1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

Real Estate	standard	Phototagging - Convolutional NN (CNN)
online Ad.	Neural N/w	Speech & Machine Transl": Recurrent NN (RNN)
Autonomous driving	- custom / Hybrid NN	

## Neural Network examples:



Standard NN



Convolutional NN

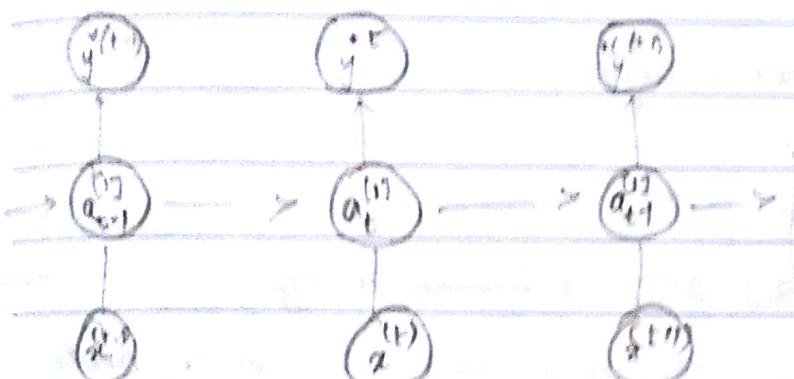


१४ / PAUSA

819 / Shaka

## जनवरी / JANUARY

बुधवार / WEDNESDAY



### Recurrent N.N



<u>Size</u>	<u># bedrooms</u>	<u>Price</u>
2104	3	400
:		
3000	4	510

## Audio, Image, Text

User Age	Ad 1d	1d	...	click
41	93242			1
80	93283			0
37	71244			1

## Week 2 : Neural Networks - Basic

### Logistic Regression as a Neural N/w :

#### Binary Classification :

In Binary Classification the goal is, given an i/p, predict to which of two classes this i/p belongs to. (e.g) Given a picture of a cat, we want to classify it between cat (1) or non-cat (0). Image is stored in 3 separate matrices corresponding to red, green, and blue color channels of this image.

For a 64 pixel  $\times$  64 image,

Blue	255 134 93 222	$\times$	255	Total dimension of this vector =
Green	255 134 202	=	231	
Red	255 231	(feature vector)	:	$64 \times 64 \times 3 = 12288$
64			255	(R,G,B)
64			134	
			:	
			255	
			134	

So, our goal is to learn a classifier that can i/p an image represented by a feature vector  $\mathbf{x}$  and predict whether the corresponding label  $y$  is 1 (cat) or 0 (non-cat).

#### Notation :-

single training example is represented as

$(\mathbf{x}, y)$  where  $\mathbf{x} \in \mathbb{R}^{n_x} \Rightarrow n_x$ -dimensional feature vector  
 $y \in \{0, 1\} \Rightarrow$  label.

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

$m_{train} = m$  = no. of training examples

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots \\ | & | & | & | \end{bmatrix}_{n_x \times m} \quad X \in \mathbb{R}^{n_x \times m}$$

$\leftarrow m \rightarrow$

In python, command to find shape of matrix  
 $X.shape = (n_x, m)$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] ; Y \in \mathbb{R}^{1 \times m}$$

In python,  $Y.shape = (1, m)$

## Logistic Regression:

Given  $X$ , we want a prediction of the probability that the class is either one or zero

Given  $X$ , want  $\hat{y} = P(y=1|X)$

$$X \in \mathbb{R}^{n_x \times n_x}$$

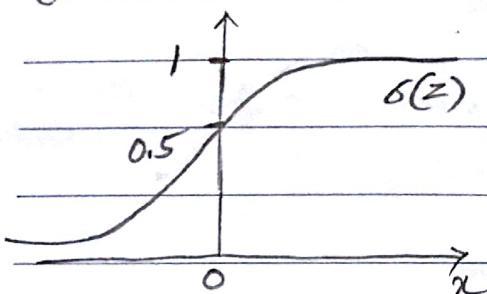
Parameters :  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$   
weights & bias

Output  $\hat{y} = w^T x + b$ . This is good for linear regression but not for binary classification since here we want  $0 \leq \hat{y} \leq 1$  (Probability). But  $w^T x + b$  can be much higher than one or it can even be negative. A probability can neither be negative or greater than one.

So, in logistic regression, our output is going to be the sigmoid function applied to that linear equation.

$$\therefore \text{output } \hat{y} = \sigma(w^T x + b) \quad | w \in \mathbb{R}^{n_x} \text{ & } b \text{ is a real number}$$

Sigmoid func<sup>n</sup>:



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

If  $z$  is large<sup>+ve</sup>,  $e^{-z} \approx 0 \therefore \sigma(z) = \frac{1}{1+0} = 1$

If  $z$  large negative,  $e^{-z} \approx \infty \therefore \sigma(z) = \frac{1}{1+\infty} = 0$

So, while implementing logistic regression, we try to learn parameters

$w$  &  $b$ , so that  $\hat{y}$  becomes a good estimate of the chance of  $Y$  being equal to one.

In programming, we use  $x_0 = 1$ , then  $X \in \mathbb{R}^{n_x+1}$

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix} \quad \left. \begin{array}{l} \{ b \} \\ \{ \} \\ \{ w \} \end{array} \right\} \begin{array}{l} \text{we use} \\ \text{them} \\ \text{separately} \end{array}$$

No need to worry about this!

## Logistic Regression Cost Function:

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

To train the parameters  $w$  and  $b$ , we need to define a cost function. Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$ .

Loss (error) function: To measure how good our output  $\hat{y}$  is when the true label is  $y$ .

Traditionally, the most used loss function would be:

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

However, this function may not be convex, making us get stuck in local minima. For this reason, our loss function is going to be instead:

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

Since we want loss function to be small,

→ If  $y=1$ , loss function will be  $-\log \hat{y}$ . We want  $\log \hat{y}$  to be large  $\Rightarrow \hat{y}$  should be large. Since  $\hat{y}$  is sigmoid function, it can never be greater than one. So we want  $\hat{y}$  to be close to one as well.

→ If  $y=0$ , then loss function will be  $-\log(1-\hat{y})$ . Since we want loss func to be small, we want  $\log(1-\hat{y})$  to be large or  $1-\hat{y}$  to be large. This implies that our  $\hat{y}$  must be very small. OR we want our  $\hat{y}$  to be as close to zero as possible.

The loss function was defined with respect to a single training example. Now, we define cost function, which measures how well you're doing on entire training set.

So, cost funct' is avg of sum of loss fun' applied to all m training

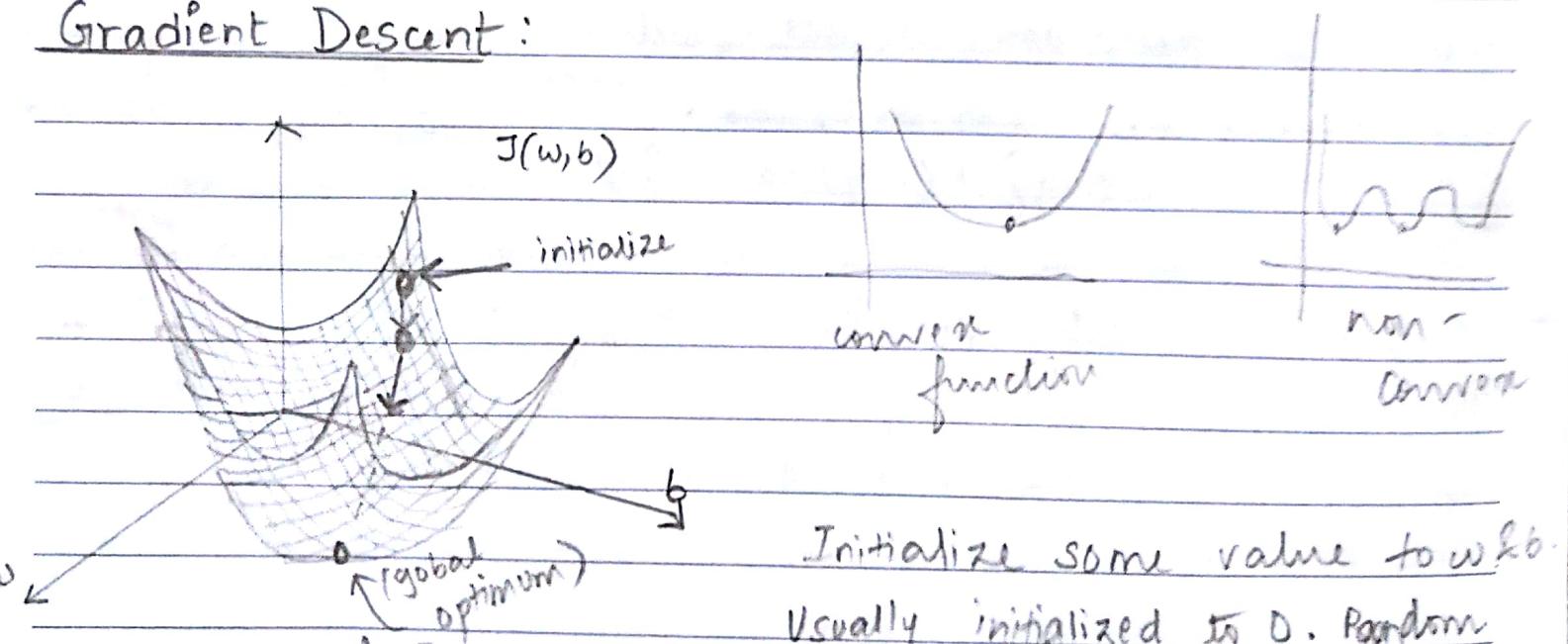
examples

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

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

This is going to be cost function, whereas the loss function refers to the single example. Our goal is then to find  $(w, b)$  that minimise  $J(w, b)$ .

### Gradient Descent:



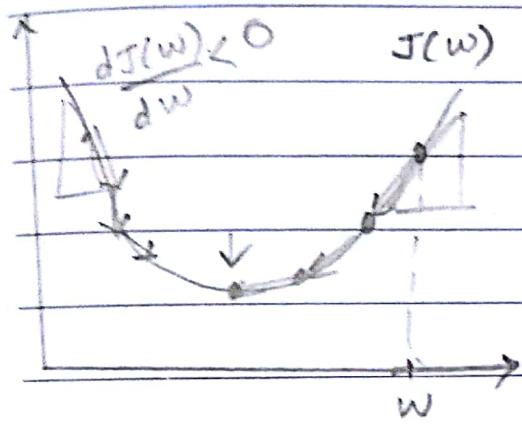
Initialize some value to  $w, b$

Usually initialized to 0. Random

initialization works but ppl. don't usually do

that for logistic regression. But since this funct' is convex, you should get to the same point or roughly the same pt. irrespective of where you initialize.

Gradient descent will start from initial value for  $J$  and makes step in the direction of the biggest gradient (steepest direction in the curve), until converging to the global optimum (or close by).



Ignoring 'b' for now, just to get 1-D plot. We are going to repeatedly update the value of  $w$ .

learning rate (how big a step we take on each iteration of gradient descent)

Repeat

$$w := w - \alpha \frac{dJ(w)}{dw}$$

} until convergence

The gradient's meaning is the slope of the curve. Then you assuming that we're starting from the right side of the curve (pt.  $w$ ), its gradient is going to be positive (slope  $> 0$ ), so the G.D step is going to move our pt. in the opposite direction of the slope, which is towards the center, where the minimum is. The same applies if we start with a  $J$  that is in the left side of the curve (slope  $< 0 \Rightarrow$  step in the opposite direction, towards the center).

On right side of slope :  $w := w - \alpha \frac{\text{slope}}{(+ve)}$  ... GD moves left  
 if  $\frac{\text{slope}}{(-ve)} : w := w - \alpha \frac{\text{slope}}{(-ve)}$  ... GD moves right

All these work if  $J$  is a funct<sup>n</sup> of  $w$  only.

If  $J$  also depends on  $b$ , i.e.  $J = J(w, b)$ , then we have two updates:

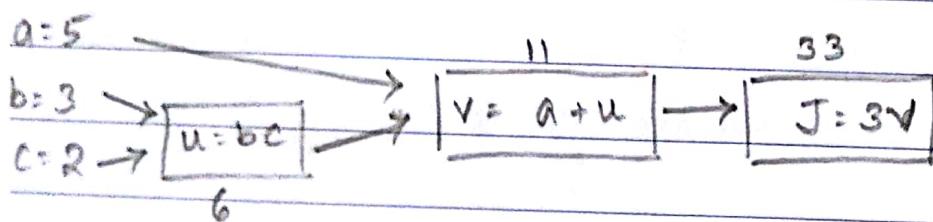
$$w := w - \alpha \frac{\delta J(w, b)}{\delta w}$$

$$b := b - \alpha \frac{\delta J(w, b)}{\delta b}$$

- $\delta$  implies partial derivative.

Derivatives with a computation Graph:

$$\text{let } J(a, b, c) = 3(a + bc)$$



Now we calculate  $\frac{dJ}{dv} = ?$  ... if we change 'v', how does the 'J' change?

$J = 3v$ .  $v$  is currently 11  $\therefore J = 33$ .

$$\begin{aligned} v &= 11 \rightarrow 11.001 \\ J &= 33 \rightarrow 33.003 \end{aligned} \quad \left\{ \begin{array}{l} J \text{ goes up 3 times.} \\ \therefore \frac{dJ}{dv} = 3 \end{array} \right.$$

$$\begin{aligned} \frac{dJ}{da} &=? & a &= 5 \rightarrow 5.001 \\ && v &= 11 \rightarrow 11.001 \\ && J &= 33 \rightarrow 33.003 \end{aligned} \quad \left\{ \begin{array}{l} J \text{ goes up} \\ 3 \text{ times} \end{array} \right. \quad \therefore \frac{dJ}{da} = 3$$

So changing  $a$  changes  $v$  which in turn changes  $J$

$$\therefore \frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da} \quad \therefore \frac{dv}{da} = 1$$

To easily implement this in python, will use the variable names as follows:  $\frac{dJ}{da} = da$ ;  $\frac{dJ}{dv} = dv$ ;  $\frac{dJ}{db} = db$ , and so on.

$$\begin{array}{c}
 \text{a=5} \quad \frac{da}{da} = 3 \\
 \text{b=3} \quad \frac{db}{db} = 6 \\
 \text{c=2} \quad \frac{dc}{dc} = 9
 \end{array}
 \rightarrow
 \begin{array}{c}
 V = a + u \quad \frac{dv}{dv} = 3 \\
 J = 3v \quad \frac{dJ}{dv} = 3
 \end{array}
 \quad
 \frac{dJ}{du} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 1$$

$$\frac{dJ}{db} = db = \frac{dJ}{du} \cdot \frac{du}{db} = \frac{3}{3} \cdot \frac{2}{2} = 6 \quad
 \begin{array}{l}
 b = 3 \rightarrow 3.001 \\
 u = b \cdot c = 6 \rightarrow 6.002 \\
 J = 33 \rightarrow 33.004
 \end{array}
 \quad
 \frac{du}{db} = \frac{0.002}{0.001} = 2$$

$$\rightarrow \frac{dJ}{db} = db = 6$$

$$\rightarrow \frac{dJ}{dc} = \frac{dJ}{du} \cdot \frac{du}{dc} = \frac{3}{3} \cdot \frac{2}{3} = 2 \quad
 \begin{array}{l}
 c = 3 \rightarrow 3.001 \\
 u = b \cdot c = 9 \rightarrow 9.003
 \end{array}
 \quad
 \frac{du}{dc} = \frac{0.003}{0.001} = 3$$

$$\therefore \frac{dJ}{dc} = dc = 9$$

## Logistic Regression Gradient Descent:

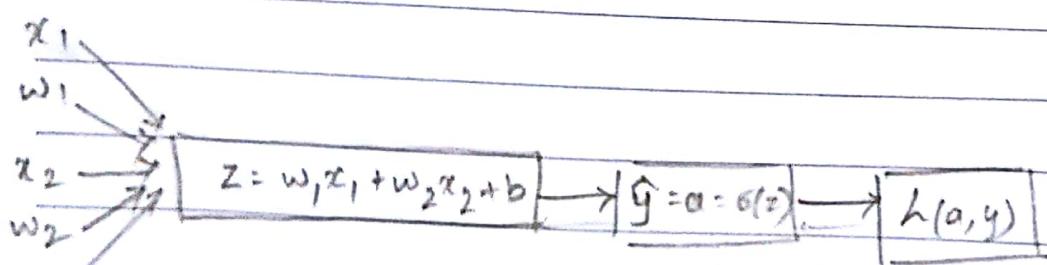
$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a, y) = -(y \log a + (1-y) \log(1-a)) \quad // \text{for one example}$$

Let's write this out as a computation graph and let's assume we have two features  $x_1$  and  $x_2$ .

So, to get  $z$ , we'll need inputs  $w_1, w_2$  &  $b$  in addition to feature values  $x_1, x_2$ .



$$dz = \frac{dL(a, y)}{dz}$$

$$da = \frac{dL(a, y)}{da}$$

$$= \frac{-y}{a} + \frac{1-y}{1-a}$$

$$= \frac{dL}{da} \cdot \frac{da}{dz}$$

$$\rightarrow dz = a - y$$

Then the final step is to go back & compute how much you need to change  $w_{k,b}$ .

$$\frac{\partial L}{\partial w_1} = dw_1 = x_1 dz ; \quad dw_2 = x_2 dz \quad \& \quad db = dz$$

then  $w_1 := w_1 - \alpha dw_1$        $b := b - \alpha db$   
 $w_2 := w_2 - \alpha dw_2$

Gradient Descent on m examples:

$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y)$       It turns out that the derivative of the loss w.r.t a certain weight  $w_j$ , for all the examples is going to be simply the avg. of the derivative of the cost for each variable w.r.t a certain weight  $w_j$ .

$$\frac{\partial J(w, b)}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})}_{d w_1^{(i)}} - (x^{(i)}, y^{(i)})$$

Let's formalise all this in a concrete algorithm.

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

For  $i = 1$  to  $m$

$$z^{(i)} = w^T x^{(i)} + b$$

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

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

} inside the  
for loop

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

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

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

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

} for n=2  
another for loop  
over all the  
features

$J / m$ ; //since we are computing avg.

$dW_1 / m$ ;  $dW_2 / m$ ;  $db / m$

$$w_1 = w_1 - \alpha * dW_1$$

$$w_2 = w_2 - \alpha * dW_2$$

$$b = b - \alpha * db$$

$dW_1, dW_2$ , do not have superscript (t) as these are accumulators + sum over the entire training set.

$$dW_1 = \frac{\partial J}{\partial W_1}$$

... till here we have implemented one step of gradient descent. We'll need to run this multiple times in order to take multiple steps of gradient descent.

We'll use 2 for loops and thus the code runs less efficiently. So, we use a technique called vectorization that allows you to get rid of these for loops.

## Python and Vectorization:

$$Z = w^T x + b$$

non-vectorized:

$$Z = 0$$

for i in range(0 to n):

$$Z += w[i] * z[i]$$

$$Z += b$$

vectorized: (in python or numpy)

$$import numpy as np$$

$$Z = np.dot(w, z) + b$$

$$w^T x$$

more examples:

whenever possible, avoid explicit for-loops.

i)  $U = A \cdot V$  then matrix multiplication is done as  $U_i = \sum_j A_{ij} V_j$

Non-vectorized: ( $N \cdot V$ )

$U = np.zeros((n, 1))$

for  $i \dots$

for  $j \dots$

$U[i] = A[i][j] * V[j]$

Vectorized: (python)

( $V$ )

$U = np.dot(A, V)$

ii) say you need to apply the exponential operation on every element of a matrix / vector

$$V = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}; \text{ then we want } U = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

Non-vec:

$U = np.zeros((n, 1))$

for  $i \in \text{range}(n):$

$U[i] = \text{math.exp}(V[i])$

import numpy as np

$U = np.exp(V)$

more examples:

$\text{np.log}(V)$ ; // element-wise log

$\text{np.abs}(V)$ ; // element-wise abs

$\text{np.maximum}(V, 0)$  // element wise max w.r.t 0

$V^{\star 2}$ ; // square of each element

$1/V$

→ Let's use vectorization to compute derivatives of our logistic regression

non vectorized form:

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for  $i = 1$  to  $m$ :

$$z^{(i)} = w^T x^{(i)} + b$$

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

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

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

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

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

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

If  $n_x > 2$ , then we require another for loop  
for  $j = 1$  to  $n_x$   
 $dw_j += ...$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m, db = db/m$$

vectorizing our second for-loop:

→ our  $dw$  will be a vector now -

$$dw = np.zeros((n_x, 1)) \quad // \text{initialisation}$$

then,

$$dw += x^{(i)} dz^{(i)}$$

$$\text{and } dw /= m$$

Now, let's try to get rid of the main for-loop as well.

Vectorizing Logistic Regression: (forward propagation, ie computing  $z$ ,  $a$ )

for  $m$  training examples, we need to compute

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad \dots$$

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} \quad X \text{ is a } (n \times m) \text{ matrix}$$

Now, let's compute  $z^{(1)}, z^{(2)}, z^{(3)} \dots$  all in one line of code

$$Z = [Z^{(1)} \ Z^{(2)} \ \dots \ Z^{(m)}] = w^T X + [b \ b \ \dots \ b]$$

↓      ↓      1xm matrix  
~~Z~~  $\begin{bmatrix} \longrightarrow \\ \text{row vector} \end{bmatrix} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(m)} \end{bmatrix} + [b \ b]$   
 $= \underbrace{[w^T x^{(1)} + b]}_{1xm} \ \underbrace{[w^T x^{(2)} + b]}_{1xm} \ \dots \ \underbrace{[w^T x^{(m)} + b]}_{1xm}$

In python,

$$Z = \text{np.dot}(w.T, x) + b$$

Next, we'll compute  $a^{(1)}, a^{(2)}, \dots, a^{(m)}$   
all at the same time.

python automatically adds 'b' to each element of the matrix  $w^T x$ . This is called 'broadcasting'

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \text{sigmoid}(z)$$

## Vectorizing Logistic Regression's Gradient Output

Previously, we computed the derivatives for each training example's prediction as  $dZ^{(1)} = a^{(1)} - y^{(1)}$ ,  $dZ^{(2)} = a^{(2)} - y^{(2)}$  and so on.

$$dZ^{(1)} = a^{(1)} - y^{(1)}, dZ^{(2)} = a^{(2)} - y^{(2)} \dots$$

$$\begin{aligned} dZ &= \begin{bmatrix} dZ^{(1)} & dZ^{(2)} & \dots & dZ^{(m)} \end{bmatrix}^T = A - Y \\ &= \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots \end{bmatrix}^T \end{aligned}$$

in our prev. page,  
 $A = [a^{(1)} \dots a^{(m)}]^T$   
 $& Y = [y^{(1)} \dots y^{(m)}]^T$

We had another for loop where we calculated  $d$ :

1939

1

माघ / MAGHA

शक / Shaka

**21**

जनवरी / JANUARY

रविवार / SUNDAY

$$dw = 0$$

$$dw += X^{(1)} dZ^{(1)}$$

$$dw += X^{(2)} dZ^{(2)}$$

$$dw /= m$$

$$db = 0$$

$$db += dZ^{(1)}$$

$$db += dZ^{(2)}$$

$$db /= m$$

vectorising the second for-loop:

$$\nabla \frac{1}{m} X^T dZ^T$$

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)}$$

$$dw = \frac{1}{m} X^T dZ^T \quad *$$

$$= \frac{1}{m} np.sum(dZ) \quad *$$

$$\begin{aligned} &= \frac{1}{m} \begin{bmatrix} X^{(1)} & \dots & X^{(m)} \end{bmatrix} \begin{bmatrix} dZ^{(1)} \\ \vdots \\ dZ^{(m)} \end{bmatrix} \\ &= \frac{1}{m} \begin{bmatrix} X^{(1)} dZ^{(1)} + \dots + X^{(m)} dZ^{(m)} \end{bmatrix} \end{aligned}$$

$n \times 1$

putting it all together:

$$\# A = 1/(1 + np \cdot \exp(-z))$$

Non-vectorised

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for  $i = 1$  to  $m$ :

$$z^{(i)} = w^T x^{(i)} + b$$

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

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

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

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

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

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

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$$

$$db = db/m$$

Vectorised:

$$Z = w^T X + b$$

$$= np \cdot \text{dot}(w.T, X) + b$$

$$A = \text{sigmoid}(Z)$$

$$dZ = A - Y$$

$$dw = \frac{1}{m} \times dZ^T$$

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

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

But we will need to call this multiple times for multiple iterations of gradient descent. So we will need a 'for' loop. There may not be a way to get rid of this for loop.

with this we have just implemented a single elevation of gradient descent to logarithmic degrees

## Broadcasting in Python:

↳ is a technique used to speed up computation.

(eg)	Apples	Beef	Eggs	Potatoes
carbs	56.0	0.0	4.4	68.0
Proteins	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

our goal is to calculate the percent. of calories from C, P, & F for each of the 4 foods

A: 59 cal

so. carbs =  $\frac{56}{59} = 94.9\%$  . . . and so on for each of the 4 foods

Can we do this without a for-loop? Yes!

→ let the above matrix be A

- ① With 1 line of code, we'll calculate sum of each columns.  $(3 \times 4)$
- ② divide each of the columns with their value of corresp. sum

python:

```
import numpy as np
```

```
A = np.array([[56.0, 0.0, 4.4, 68.0], [1.2, 104.0, 52.0, 8.0], [1.8, 135.0, 99.0, 0.9]])
```

③ `cal = A.sum(axis=0)` // calculate sum of all columns  
`print(cal)` // output: [59. 239. 155.4 76.9]

④ `percentage = 100 * A / cal.reshape(1, 4)` // divide by  $1 \times 4$  mat  
`print(percentage)` // o/p: [[94.9152 0. 2. 88. ]]  
[. . . . ]

$\Rightarrow \text{cal} = A \cdot \text{sum}(\text{axis}=0)$   
 $\text{axis} = 0 \dots \rightarrow \text{sum vertically}; \text{axis} = 1 \dots \text{sum horizontally}$

$\Rightarrow \text{percentage} = 100 * \text{A} / \text{cal} \cdot \text{reshape}(1, 4)$ .  
since our 'cal' is already a  $1 \times 4$  matrix, the 'reshape' part is a bit redundant. We have divided a  $(3 \times 4)$  matrix by a  $(1 \times 4)$  matrix.

### More examples of broadcasting:

$\Rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 = ?$ . Python will auto expand '100' into a  $(4 \times 1)$  vector as well:  $\begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix}$

The result is  $\begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$ .

$\Rightarrow$  Similarly, if we add:  
 $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}$ . Python will expand matrix B as  
 $\begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$   
then result is  $\begin{bmatrix} 101 & 102 & 103 \\ 104 & 105 & 106 \end{bmatrix}$

$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 101 & 102 \\ 200 & 201 & 202 \end{bmatrix}$   
broadcasting

$\Rightarrow$  General Principle: python converts it into  
 $(m, n) \pm (1, n) \rightsquigarrow (m, n)$   
matrix  $\times (m, 1) \rightsquigarrow (m, n)$

$\Rightarrow (m, 1) \pm (1, m) \rightsquigarrow (1 \times 1)_{\text{matrix}} \rightsquigarrow (m, 1)$

## A note on python/numpy vectors:

```
import numpy as np
```

```
a = np.random.randn(5)
```

```
print(a) // o/p: [0.5... -0.2... 0.95... -0.82... -1.42...]
```

```
print(a.shape) // o/p: (5,) ... so this is neither a row vector
```

transpose → or a column vector.

```
print(a.T) // the o/p of this ends up looking the same as  
print(a). ie [0.5... -0.2... 0.95... -0.82... -1.42...]
```

```
print(np.dot(a, a.T)) // o/p: 4.0657... now we may think that this  
should result in a matrix but we get a number instead.
```

So, to solve this we should do:

```
a = np.random.randn(5, 1) → // o/p will be a 5x1 matrix
```

```
print(a)
```

```
print(a.T) // o/p will be a 1x5 matrix.
```

```
print(np.dot(a, a.T)) // o/p will be a outer product of 2 matrices... a  
5x5 matrix
```

In detail,

→ { a = np.random.randn(5) ... this is basically a 'rank 1 array'.  
so never → a.shape = (5, ) ←

use this

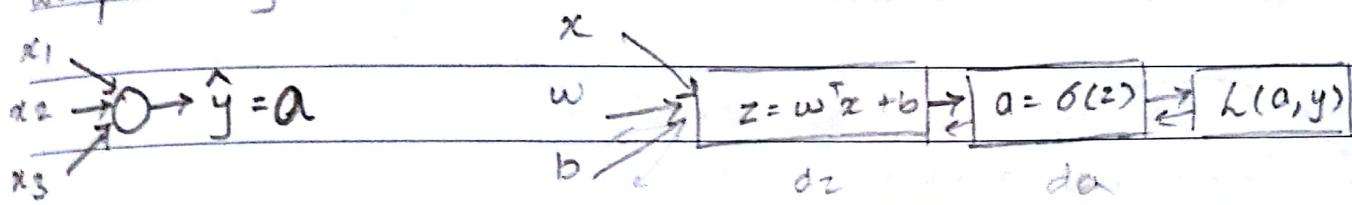
```
a = np.random.randn(5, 1) → a.shape = (5, 1) ... column vector
```

```
assert (a.shape == (5, 1)) ← check what this does
```

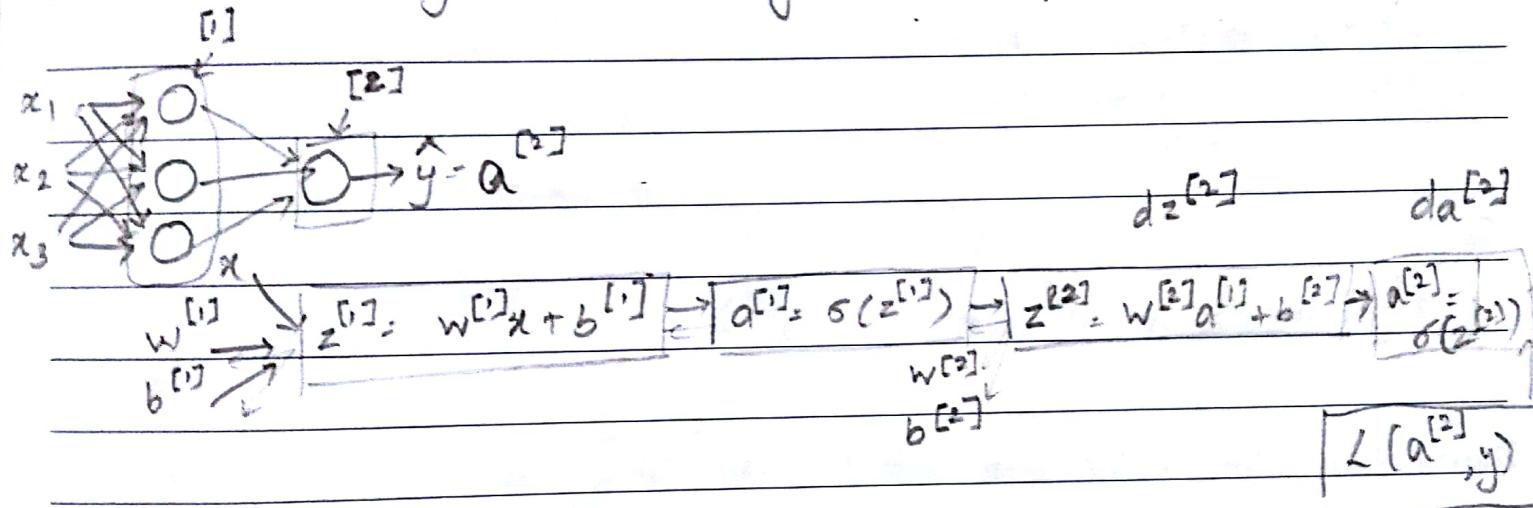
```
a = a.reshape((5, 1)) // we can use reshape to convert a  
rank matrix to a (5x1) or  
(1x5) matrix
```

## Week 3: Shallow Neural Networks:

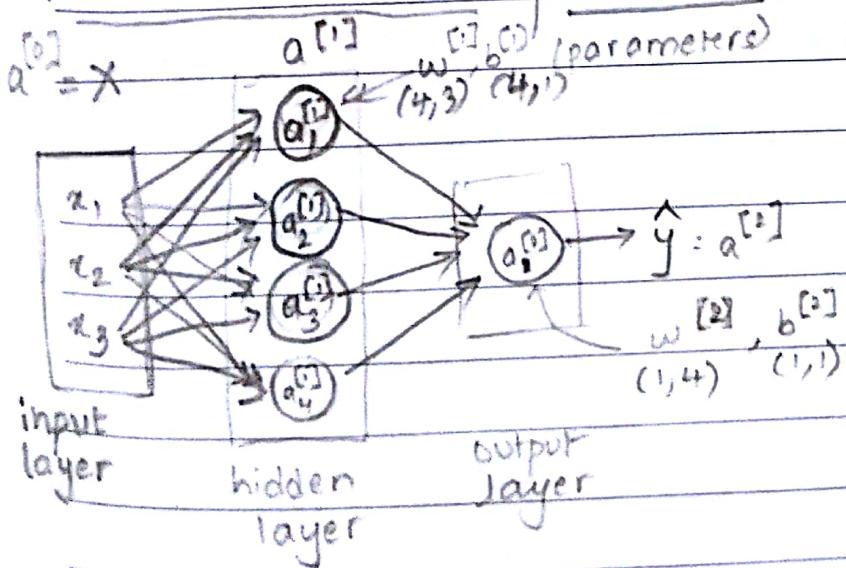
we previously had:



with 1 hidden layer: (or 2 layer neural n/w)



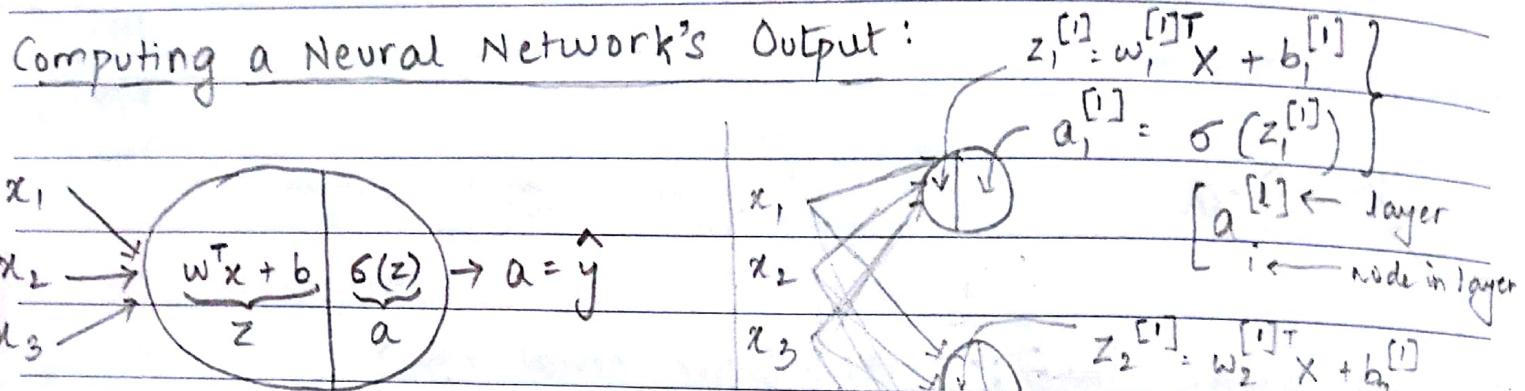
### Neural Network Representation:



$a^{[i]}$  = activations of the i/p layer

$$a^{[i]} = \begin{bmatrix} a_1^{[i]} \\ a_2^{[i]} \\ \vdots \\ a_n^{[i]} \end{bmatrix}$$

2 layer N.N (1-hidden layer)  
2 - o/p layer



$z = w^T x + b$  } a single sigmoid  
 $a = \sigma(z)$  } neuron does these  
 two things

So  $z_1^{[1]} = w_1^{[1] T} x + b_1^{[1]}$ ,  $a_1^{[1]} = \sigma(z_1^{[1]})$   
 $z_2^{[1]} = w_2^{[1] T} x + b_2^{[1]}$ ,  $a_2^{[1]} = \sigma(z_2^{[1]})$   
 :

Doing this with a for loop will be very inefficient. So, we vectorize.

$$\begin{bmatrix} -w_1^{[1] T} \\ -w_2^{[1] T} \\ -w_3^{[1] T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_1^{[1] T} x + b_1^{[1]} \\ \vdots \\ \vdots \\ w_4^{[1] T} x + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$\underbrace{\phantom{w^{[1]}}}_{W^{[1]}}$        $\underbrace{\phantom{b^{[1]}}}_{b^{[1]}}$

then we calculate 'a' as

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

So, given input  $x$ :

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

second layer

$$\begin{cases} z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{cases}$$

(output layer)      (1,1)      (1,1)

vectorizing across multiple examples:

for a 2 layer N.N :

$x \rightarrow a^{[2]} = \hat{y}$ , This is for one training example.

If there are  $m$  training examples, then

$x^{(1)} \rightarrow a^{[2](1)} = \hat{y}^{(1)}$

$x^{(2)} \rightarrow a^{[2](2)} = \hat{y}^{(2)}$

$x^{(m)} \rightarrow a^{[2](m)} = \hat{y}^{(m)}$

$a^{[2](i)}$  example  $i$   
layer 2

So, for all training examples,

for  $i = 1$  to  $m$ ,

$$z^{[1](i)} = w^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = w^{[2]}a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

Non-vectorized form

1939

15

शक / Shaka  
B-E-A-V

4

फरवरी / FEBRUARY

रविवार / SUNDAY

In vectorized form, the computation will be.

$$x = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}^T$$

↑ diff. features

$$Z^{[1]} = \begin{bmatrix} 1 & z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}^T$$

$$(n_{in}, m) \quad A^{[0]}$$

$$A^{[1]} = \begin{bmatrix} 1 & a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \end{bmatrix}^T$$

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

diff. training examples

hidden layers

Explanation for vectorized implementation:

justification for vectorized implementation -

$$z^{[0]}(1) = w^{[0]}x^{(1)} + b^{[0]} \rightarrow \text{for 1st training example}$$

$$z^{[0]}(2) = w^{[0]}x^{(2)} + b^{[0]} \rightarrow \text{for 2nd training example}$$

$$w^{[0]} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}; w^{[0]}x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad \text{so are } w^{[0]}x^{(2)} \text{ & } w^{[0]}x^{(3)}$$

column vector

our  $x = \begin{bmatrix} 1 & x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 & 1 \end{bmatrix}$ , then  $w^{[0]}x = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$

$w^{[0]}x^{(3)}$

$$= \begin{bmatrix} 1 \\ z^{[0]}(1) \\ z^{[0]}(2) \\ z^{[0]}(3) \end{bmatrix} = \begin{bmatrix} 1 \\ z^{[0]}(1) + b^{[0]} \\ z^{[0]}(2) + b^{[0]} \\ z^{[0]}(3) + b^{[0]} \end{bmatrix} = \begin{bmatrix} 1 \\ z^{[0]}(1) \\ z^{[0]}(2) \\ z^{[0]}(3) \end{bmatrix}$$

so,  $w^{[0]}x^{(i)} = z^{[0]}(i)$

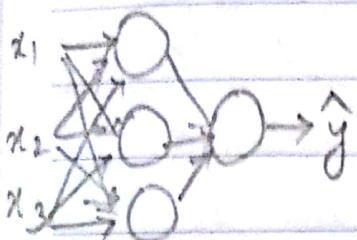
## Activation functions:

$$\text{given } x: z^{[1]} = w^{[1]T}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

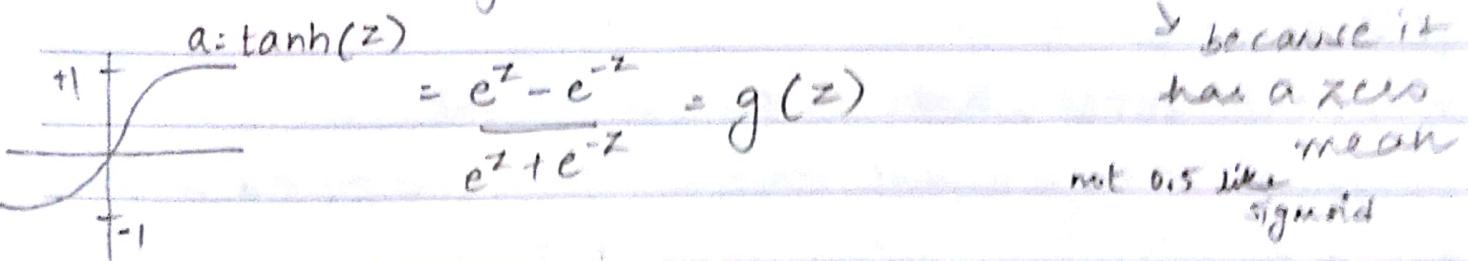
$$z^{[2]} = w^{[2]T}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$



Sigmoid function is one activation function. We can have more:

1) tanh (almost always works better than sigmoid):



Sigmoid can be used in binary classification at the output layer as in B.C., where  $y \in \{0, 1\}$ . So it makes sense to have sigmoid where  $0 \leq \hat{y} \leq 1$  instead of tanh, where  $-1 \leq \hat{y} \leq 1$ .

Disadvantage of sigmoid & tanh is that if  $z$  is very large, or very small, then the slope is close to zero, so this can slow down gradient descent.

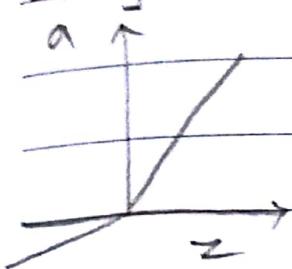
## 2) ReLU

$a = \max(0, z)$  derivative is 1 as long as  $z$  is positive  
and " " is 0  $\rightarrow z$  is negative  
when  $z=0$ , slope is not defined.  
But in computer, we get derivative as 0.0000..

Some thumb rules:-

- When output is 0 or 1, sigmoid function is natural choice for output layer and for other layers ReLU is default choice.
  - When not sure what to use, use ReLU
  - Advantage of ReLU is that for a lot of space of  $z$ , the derivative or slope of activation function is very different from zero. So in practice using ReLU, your NN will often learn much faster than tanh or sigmoid as slope of the function doesn't much go to zero which slows down learning. Though half of the range of  $z$ , the slope of ReLU is zero but in practice, enough of your hidden units will have  $z$  greater than 0 so learning can still be fast
  - You can try all the activation functions & choose the best.
- Pros & cons:
- ▷ Sigmoid: Never use except in o/p layer of binary classification

### 3) Leaky ReLU:



$$a = \max(0.01z, z)$$

works better than ReLU

→ why do you need non-linear activation functions?

Given  $x$ :

$$z^{[1]} = W^{[1]}x + b^{[1]} \quad a^{[1]} = z^{[1]} \text{ or}$$

$a^{[1]} = g^{[1]}(z^{[1]})$  ... so instead of this let's say we have  $g(z) = z$   
then  $g(z)$  is a linear activation function. or

identity activation function i.e it outputs whatever is the input

so  $\hat{y}$  is a linear function of your input features  $x$ .

$$\begin{aligned} a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \quad \& \quad a^{[2]} = z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ &= W^{[2]}(W^{[1]}x + b^{[1]}) + b^{[2]} \\ &= (W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]}) \\ &= Wx + b' \end{aligned}$$

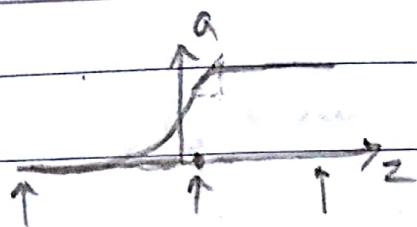
so o/p is just a linear funct<sup>n</sup> of the i/p.

whatever hidden layers you add, the activation will be always linear like logistic regression (so it is useless in a lot of complex problems).

You might use linear activation function in one place - in the o/p layer if the output is real nos. (like price of the house \$0 to ... \$10000... ; basically regression problem). But even in this case if the o/p value is non-negative you could use ReLU instead.

## Derivatives of activation function:

### Sigmoid function:



$$g(z) = \frac{1}{1+e^{-z}} \quad \text{then slope: } \frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left( 1 - \frac{1}{1+e^{-z}} \right)$$

$$\text{or } g'(z) = g(z)(1-g(z))$$

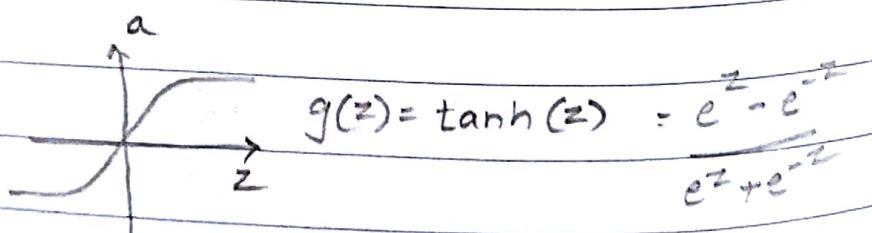
Sanity check:

If  $z=10$ ,  $g(z) \approx 1$ :  $\frac{d}{dz} g(z) \approx 1(1-1)=0$  ... correct, as if  $z$  is large,  
slope is close to zero

$z=-10$ ,  $g(z) \approx 0$ .  $\frac{d}{dz} g(z) \approx 0(1-0)=0$  ... correct, as if  $z$  is very small,  
slope is close to zero

$z=0$ , then  $g(z)=\frac{1}{2}$   $\frac{d}{dz} = 1.5(1-0.5)=0.25$  ... correct slope.

Since  $a: g(z) = \frac{1}{1+e^{-z}}$ , then  $g'(z) = a(1-a)$



$$g(z) = \tanh(z) \quad \text{then } \frac{d}{dz} g(z) = 1 - (\tanh(z))^2$$

Sanity check:  $z=10$ ,  $\tanh(z) \approx 1$ ;  $g'(z)=0$

$z=-10$ ,  $\tanh(z) \approx -1$ ;  $g'(z)=0$

$z=0$ ,  $\tanh(z) \approx 0$ ;  $g'(z)=1$

$$a: g(z) \text{ then } g'(z) = 1 - a^2$$

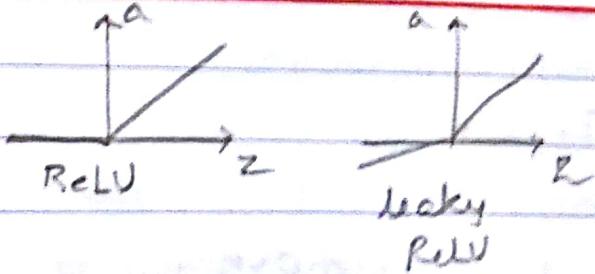
## ReLU and Leaky ReLU:

$$g(z) = \max(0, z) \dots \text{ReLU}$$

$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

undefined if  $z=0$  ... but when implemented on computer you can set it to 1, no problem. so

$$g'(z) = 1 \text{ if } z \geq 0$$



$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

## Gradient descent for Neural Networks:

1939

22

माघ / MAGHA

शक / Shaka

**11**

फरवरी / FEBRUARY

रविवार / SUNDAY

Parameters :  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

Input layer,  $n_2 = n^{[0]}$

Hidden layer,  $n^{[1]}$

Output layer,  $n^{[2]} = 1$

$$\text{So, } w^{[1]} = (n^{[1]}, n^{[0]}) ; b^{[1]} = (n^{[1]}, 1) ; w^{[2]} = (n^{[2]}, n^{[1]}) ; b^{[2]} = (n^{[2]})$$

Cost function :

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}, y)$$

gradient descent :

Repeat {

compute predictions ( $\hat{y}^{(i)}$  for  $i=1\dots m$ )

compute derivatives ( $dW^{(1)}, db^{(1)}, \dots, dW^{(2)}, db^{(2)}$ )

update  $w^{(1)} := w^{(1)} - \alpha \cdot dW^{(1)}$  } also for  $w^{(2)} \& b^{(2)}$   
 $b^{(1)} := b^{(1)} - \alpha \cdot db^{(1)}$

This is one iteration of gradient descent. So repeat this until convergence

Formulas for computing derivatives:

Forward propagation steps:

$$Z^{(1)} = W^{(1)}X + b^{(1)}$$

$$A^{(1)} = g^{(1)}(Z^{(1)})$$

$$Z^{(2)} = W^{(2)}X + b^{(2)}$$

$$A^{(2)} = g^{(2)}(Z^{(2)})$$

Backpropagation steps:

$$dZ^{(1)} = A^{(1)} - Y \quad ; \quad Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dW^{(2)} = \frac{1}{m} dZ^{(2)} A^{(2)T}$$

$$db^{(2)} = \frac{1}{m} \text{np.sum}(dZ^{(2)}, \text{axis}=1, \text{keepdims}=\text{True})$$

$$dZ^{(1)} = W^{(1)T} dZ^{(2)} \times g^{(1)'}(Z^{(1)}) \quad (n^{(1)}, m) \text{ matrix}$$

element wise product

( $n^{(1)}, m$ ) matrix

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

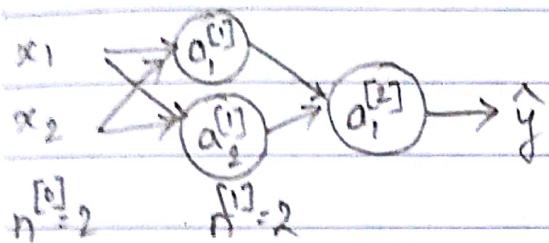
$$db^{(1)} = \frac{1}{m} \text{np.sum}(dZ^{(1)}, \text{axis}=1, \text{keepdim}=\text{True})$$

( $n^{(1)}, 1$ ) vector

## Random Initialization:

For logistic regression, it was OK to initialize the weights to zero.

For N.N., to initialize the weight to zero and then applying gradient descent, it won't work. Let's see why.



$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{let's say we initialize with all zeros} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- Turns out, initialising  $b$  to zero is OK, but initializing  $w$  to all 0s is a problem.

With this formalization, for any example you give it, we have  $a_1^{[1]} = a_2^{[1]}$  (as both are computing the exact same function.)

During backpropagation,  $dZ_1^{[1]} = dZ_2^{[1]}$

And by proof of induction, it turns out that after every single iteration of training your two hidden units are still computing exactly the same function.

our  $dw$  looks like this:

$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$  ... every row takes on the same value.

so when we perform:  $w^{[1]} = w^{[1]} - \alpha dw$  first row is equal to the second row.

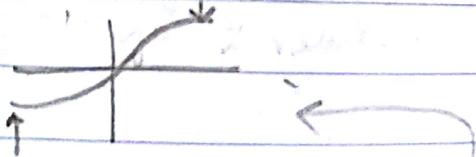
All hidden units become symmetric.

so, there's really no point in having more than one hidden unit as they are all computing the same thing.

Solution : initialize your parameters randomly.

You can set  $w^{[1]} = np.random.randn(2, 2) * 0.01$   $\uparrow$  because we want very small values  
 $b^{[1]} = np.zeros((2, 1))$

Large values of weights imply that  $z$  will be very large or very small.



In this case, we are more likely to end up at these points, where the slope of the gradient is very small. Meaning that the gradient descent will be very slow!

Also, at 0<sup>th</sup> layer we can have  $w^{[0]} = np.random.randn(1, 2) * 0.01$

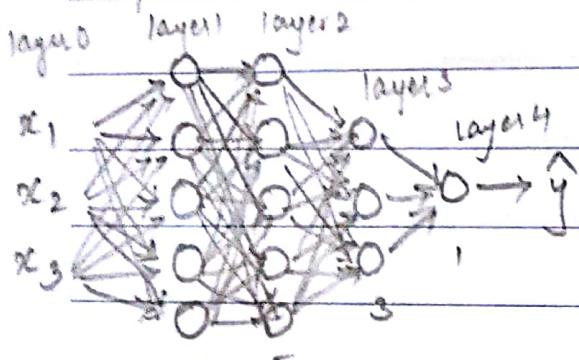
When working with very deep N.N, we might want to pick a different constant than 0.01.

## Week 4 : Deep Neural Networks

Deep L-layer neural network :

Shallow N.N.  $\rightarrow$  one or 2 layers

Deep N.N.  $\rightarrow$  three or more layers



'4' layer N.N.

Notations:

$L$  = no. of layers in the network,  
here  $L = 4$ .

$n^{[l]}$  = no. of units in  $l$

$n^{[0]} = 3$ ;  $n^{[1]} = 5$ ;  $n^{[2]} = 5$ ;  $n^{[3]} = 3$ ;  $n^{[4]} = 1$

$n^{[0]} = n_x = 3$

$n^{[L]}$

$a^{[l]}$  = activations in layer  $l$

$a^{[l]} = g^{[l]}(z^{[l]})$ ,  $w^{[l]}$  = weights for  $z^{[l]}$

$x = a^{[0]}$        $a^{[L]} = \hat{y}$

Forward Propagation in Deep Neural Network :-



$$x: z^{[1]} = w^{[1]} X + b^{[1]}$$

$a^{[0]}$

layer 1:  $a^{[1]} = g(z^{[1]})$

layer 2:  $z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$

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

layer 4:  $z^{[4]} = w^{[4]} a^{[3]} + b^{[4]}$

$a^{[4]} = g(z^{[4]})$

-  $y$

generally,

$$\begin{aligned} z^{[l]} &= w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(z^{[l]}) \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{all for a single training eg.}$$

vectorised way for the whole training set :-

layer 1:  $z^{[1]} = w^{[1]} X + b^{[1]}$

$A^{[1]} = g^{[1]}(z^{[1]})$

layer 2:  $z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$

$A^{[2]} = g^{[2]}(z^{[2]})$

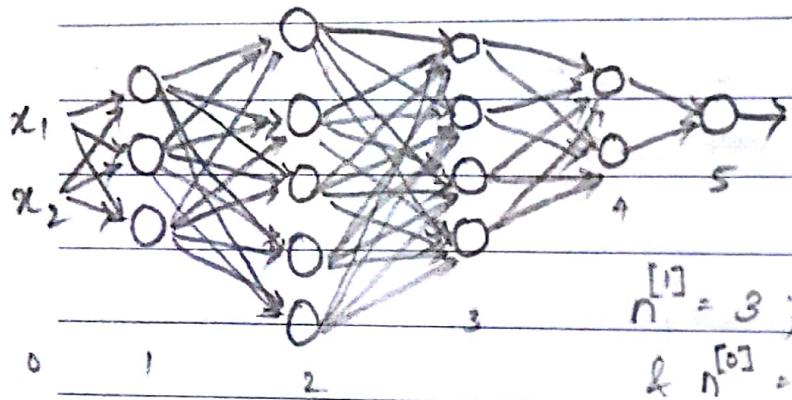
$$z = \begin{bmatrix} z^{[1](1)} & z^{[2](1)} & \dots & z^{[m](m)} \end{bmatrix} \quad \text{similarly } A = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}$$

then,

$$\hat{y} = g(z^{[4]}) = A^{[4]}$$

while performing forward propagation, it is perfectly OK to have a for-loop to calculate activation function for layers 1, 2, 3 & 4.

Getting your matrix dimensions right: (best way is to use paper & pencil)



$$L=5$$

$$z^{[l]} = w^{[l]} \cdot x + b^{[l]}$$

$$n^{[1]} = 3; n^{[2]} = 5; n^{[3]} = 4; n^{[4]} = 2; n^{[5]} = 1$$

$$\& n^{[0]} = n_x = 2$$

$$\text{our } z^{[l]} = w^{[l]} \cdot x + b^{[l]} \quad (3,1) \quad (2,1) \quad (2,1) \\ (n^{[l]}, 1) \quad (n^{[l]}, 1)$$

So, we need a matrix  $w^{[l]}$  to be something that when we multiply an  $(n^{[0]}, 1)$  matrix, we get an  $(n^{[l]}, 1)$  vector.  $z^{[1]}$  is  $(3, 1)$  &  $x$  is  $(2, 1)$ . So  $w^{[1]}$  will be  $(3, 2)$ . Or generally,  $w^{[l]}$  is  $(n^{[l]}, n^{[0]})$  matrix.

$$w^{[1]}: (n^{[1]}, n^{[0]}) \dots \text{generally } w^{[l]}: (n^{[l]}, n^{[l-1]})$$

$$\text{so, } w^{[2]}: (5, 3) \text{ or } (n^{[2]}, n^{[1]})$$

$$z^{[2]} = w^{[2]} \cdot a^{[1]} + b^{[2]} \quad (5,1) \quad (5,3) \quad \uparrow \quad (5,1) \quad (3,1)$$

$$w^{[3]}: (4, 5)$$

$$w^{[4]}: (2, 4)$$

$$w^{[5]}: (1, 2)$$

generally,

$$b^{[l]}: (n^{[l]}, 1)$$

$$dw: w : (n^{(1)}, n^{(l-1)})$$

$$db: b: (n^{[l]}, 1)$$

$$z^{[l]} = g^{[l]}(a^{[l]})$$

so  $z \& a$  should have the same dimension

When we have a vectorized implementation, that looks at multiple examples at a time.

The dimensions of  $w_b$ ,  $d_w$ , and  $d_b$  will stay the same. But the dimensions of  $z$ ,  $a$ , as well as  $x$  will change a bit.

Previously we had,

$$z^{[l]} = w^{[l]} \cdot x + b^{[l]} \quad (n^{[l]}, 1)$$

$(n^{[l]}, 1) \quad (n^{[l]}, n^{[l]}) \rightarrow (n^{[l]}, 1)$

In vectorised implementation:

$$z^{[l]} = w^{[l]} \cdot x + b^{[l]} \quad \dots \quad z^{[l]} = \begin{bmatrix} z^{[l](1)} & z^{[l](2)} & \dots & z^{[l](m)} \end{bmatrix}$$

$(n^{[l]}, m) \quad (n^{[l]}, n^{[l]}) \quad (n^{[l]}, m) \quad (n^{[l]}, 1)$

$\downarrow$

$(n^{[l]}, m)$  matrix by python broadcasting

$z^{[l]}, a^{[l]} : (n^{[l]}, 1)$  ... previously,  
now in vectorized,

$$z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

\* when  $l=0$ ,  $A^{[0]} = x : (n^{[0]}, m)$

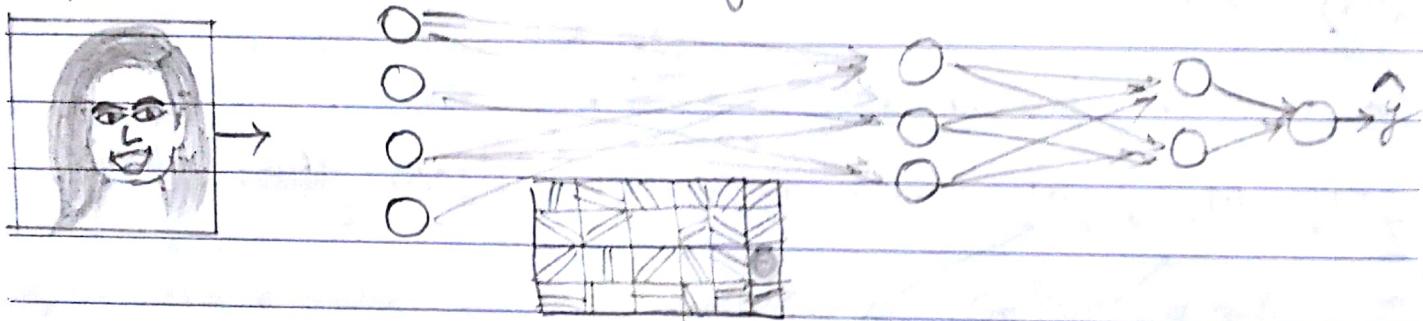
$$dz^{[l]}, dA^{[l]} : (n^{[l]}, m) \dots \text{same as } Z \text{ & } A.$$

## Why deep representations?

Why deep networks work well?

First, what is deep network computing?

- If you're building a system for face recognition, here's what a deep neural network could be doing.



The first layer of the N.N, you can think of as maybe being a feature detector or an edge detector. If I have 20 hidden units, then the hidden units are trying to figure out where the edges of that orientation in the image, some other hidden unit maybe trying to figure out where are the horizontal lines in the image, etc. We can think of the first layer of the neural nw as looking the picture and try to figure out where are the edges in the picture

Now let's think about where are the edges in this picture. by grouping together pixels to form edges. It can then detect the edges and group the edges together to form parts of faces. For eg, you might have a neuron trying to see if it's finding an eye or a diff. neuron trying to find that part of the nose. And so by putting together lots of edges, it can start to detect diff. parts of faces.

And then finally, by putting together different parts of faces, like an eye or a nose or a chin, it can then try to recognise or detect different types of faces. So intuitively, we can think of the earlier layers of the N.N as detecting simple functions like edges. And then composing them together in the later layers of a N.N. so that it can learn more & more complex functions. The edge detectors are looking in relatively small areas of an image, and facial detectors look at much larger areas of image. Main takeaway: find simpler things like edges; composing them together to find complex things & composing those together to find even more complex things.

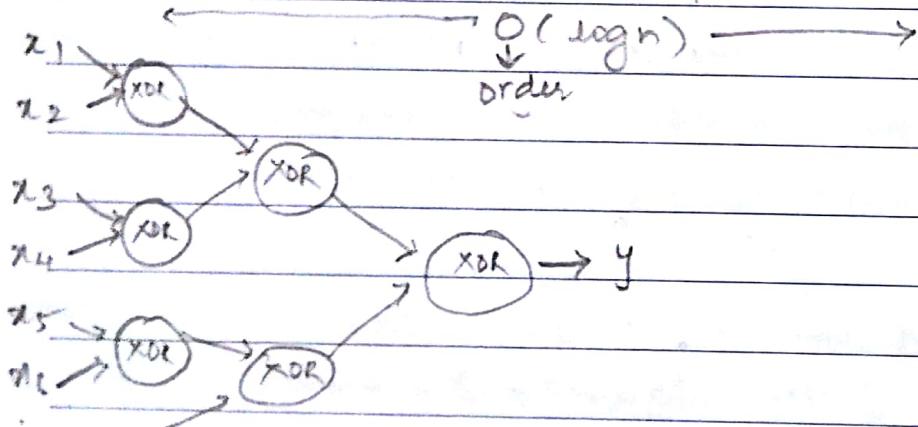
In Speech recognition, if you input an audio clip, the first level of a N.N might learn to detect low level audio wave features,

1939 6 फाल्गुन / FALGUN such as is this tone going up? or down? is it  
शक / Shaka white noise or slithering sound & what is the pitch.  
**25** फरवरी / FEBRUARY And then composing them, you'll learn to  
रविवार / SUNDAY detect basic units of sound.(phonemes). Then composing these together maybe learn to recognise words in the audio. Then compose those together in order to recognise entire phrases or sentences.

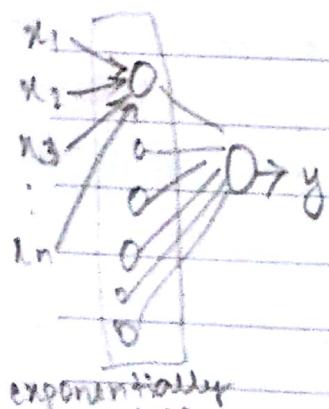
## Circuit theory and deep learning:

Informally: There are functions you can compute with a 'small' L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

(eg):  $x_1 \oplus x_2 \oplus x_3 \oplus x_4 \dots \oplus x_n = y$



Without a N.N. (only 1 hidden layer):

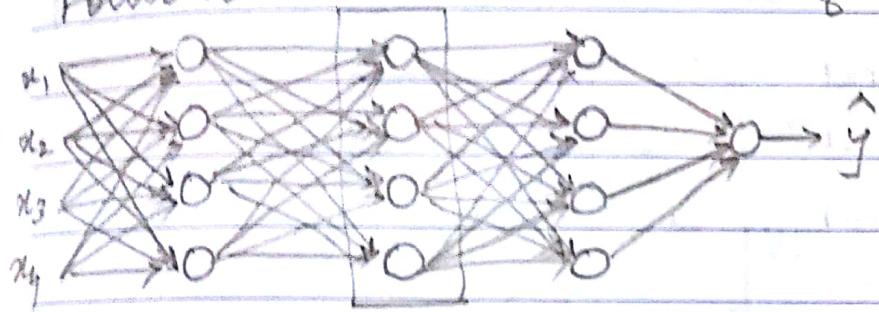


then in order to compute this XOR function, this hidden layer will need to be exponentially large because essentially, you need to exhaustively enumerate over  $2^n$  possible configurations of the input bits that result in XOR being either 1 or 0.

<sup>large</sup>  
 $O(2^t)$  So, there are mathematical functions, that are much easier to compute with deep N.N than with shallow networks.

## Building blocks of deep neural networks

### Forward and backward functions:



Let's pick one layer & look at the computations focussing on just that layer for now.

Layer  $l$ :  $w^{[l]}, b^{[l]}$

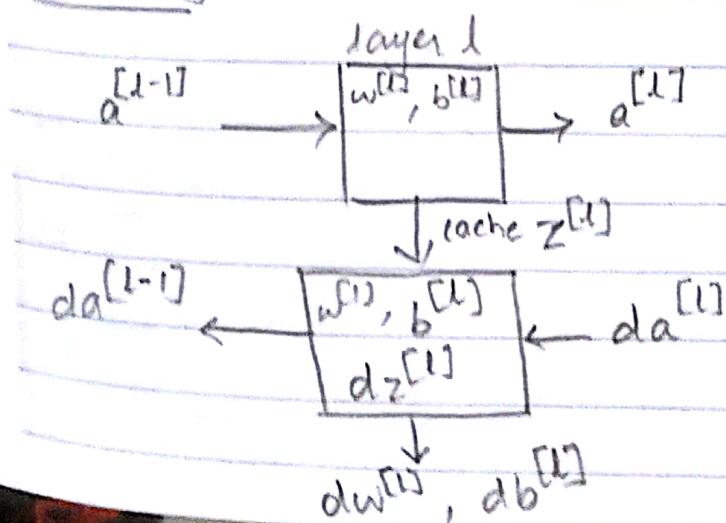
forward: Input:  $a^{[l-1]}$ ; output  $a^{[l]}$   
Back prop

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \quad \text{and} \quad a^{[l]} = g^{[l]}(z^{[l]})$$

We also store the value  $z^{[l]}$  as it will be useful for back prop step later.  
(or cache)

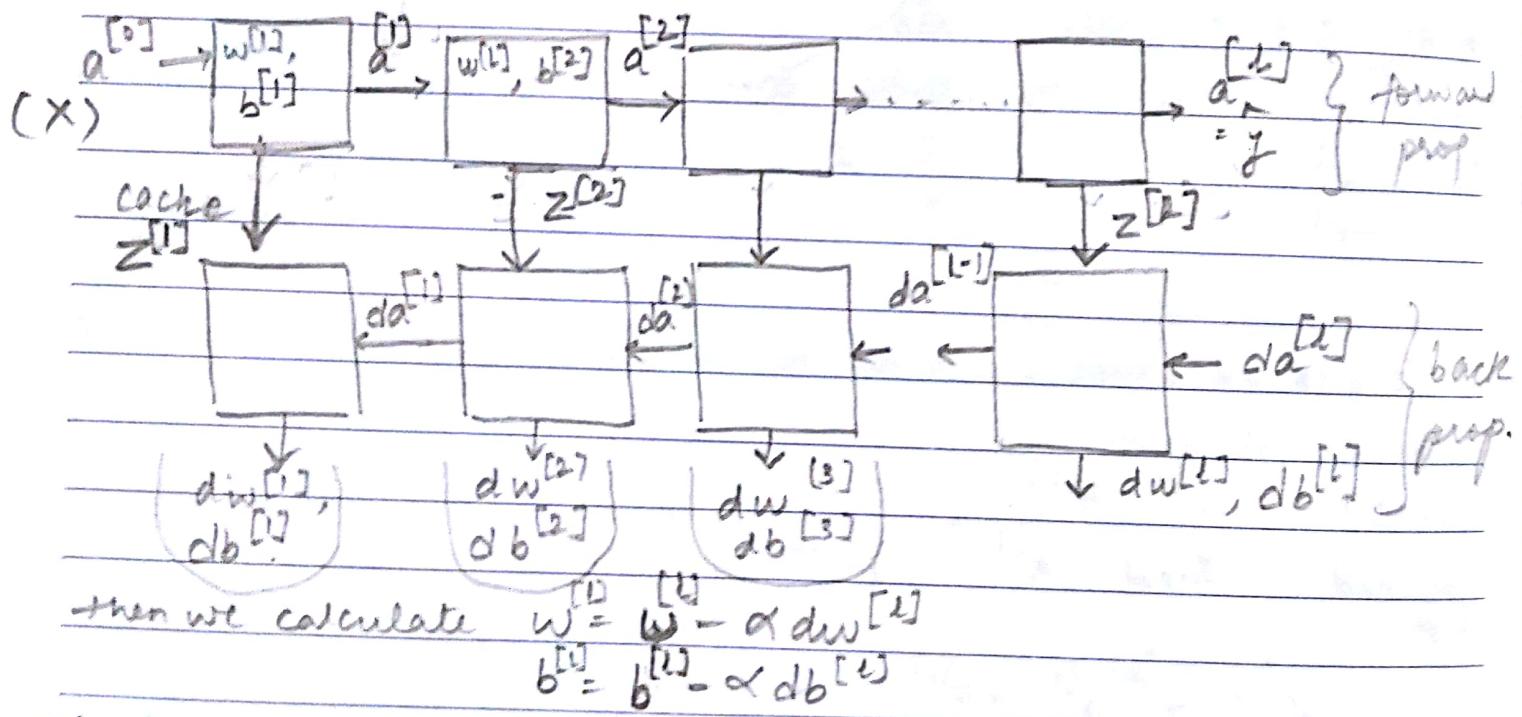
Back prop: Input:  $da^{[l]}$ ; output  $da^{[l-1]}$   
cache:  $z^{[l]}$        $dw^{[l]}$   
                           $db^{[l]}$

Summarising:



2018

If you can implement these two functions (forward prop), then the basic computation of the NN will be as follows:



This is just one iteration of gradient descent for your N.N.  
Though it is written cache  $z^{[1]}, z^{[2]}, \dots$  but actually while  
programming, we cache  $z^{[0]}, w^{[1]}, b^{[1]}, z^{[2]}, w^{[2]}, b^{[2]}, \dots$ , etc..

## Forward and Backward Propagation:

We previously saw the basic blocks of implementing a deep N.N - a fwd prop for each layer, and a corresponding back prop step. Let's see how you actually implement these steps. Starting with fwd prop

Forward propagation for layer l :

Input  $a^{[l-1]}$

Output  $a^{[l]}$ , cache ( $z^{[l]}$ )

$$\text{eq} \quad z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

$$a^{[0]} = x$$

vectorized

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

$$A^{[0]} = X$$

Backward propagation for layer l

Input  $da^{[l]}$

Output  $da^{[l-1]}, dW^{[l]}, db^{[l]}$

element-wise product

vectorized,

$$dz^{[l]} = da^{[l]} \times g^{[l]}(z^{[l]})^T$$

$$dW^{[l]} = 1/m \cdot dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = 1/m \cdot \text{np.sum}(dz^{[l]}, \text{axis}=1, \text{keepdim=True})$$

$$dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

$$\text{eq} \quad da^{[l-1]} = W^{[l]T} \cdot dz^{[l]} \quad \text{②}$$

putting ② in ①, we get

$$dz^{[l]} = W^{[l+1]T} \cdot dz^{[l+1]} \times g^{[l]}(z^{[l]})$$

11

ફાળગુણ / FALGUN

શક્તિ / Shakti

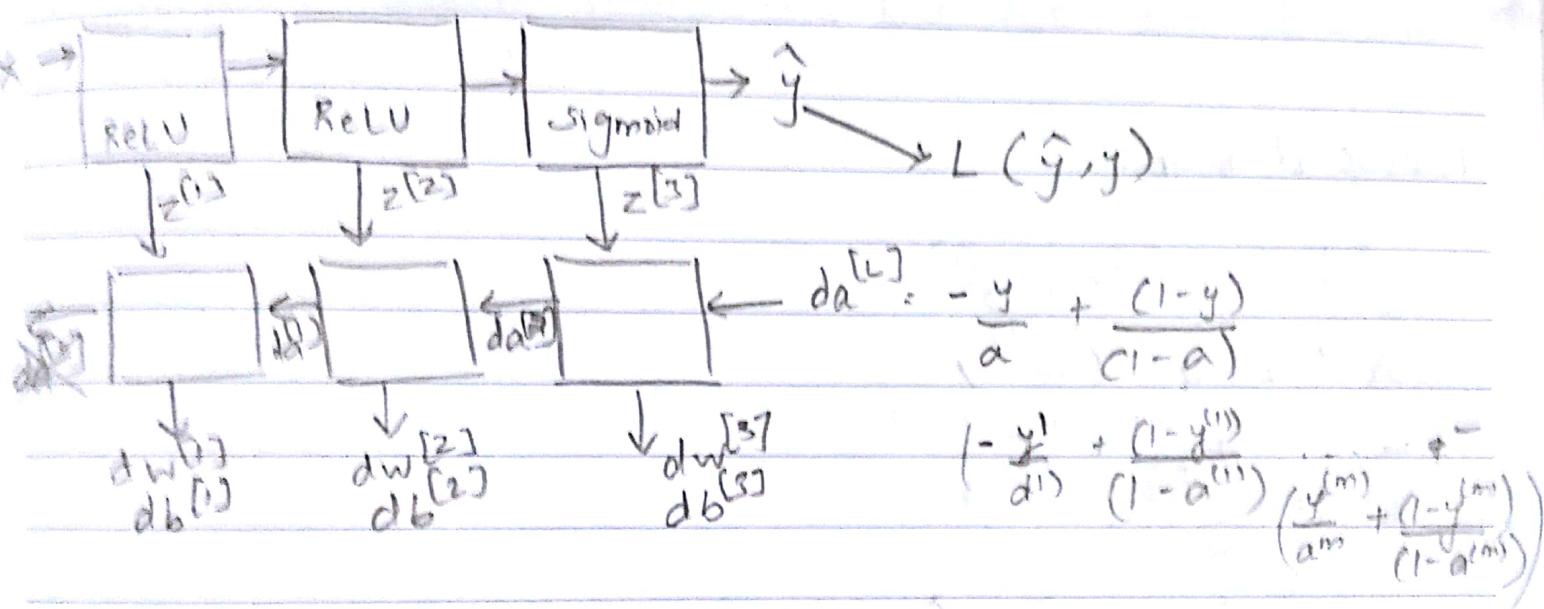
2

માર્ચ / MARCH

શુક્રવાર / FRIDAY



Summary :



### Parameters vs Hyperparameters:

Parameters:  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots$

Hyperparameters: Learning rate,  $\alpha$ ,  
 # iterations of gradient descent,  $L$ ,  
 # hidden layers,  $n^{[1]}, n^{[2]}, \dots$ ,  
 choice of activation function } all these  
 control our parameters  $w$  &  $b$

There are more hyperparameters that we'll learn later.

X

10

→ Forward and Backward propagation equations :

Forward prop:

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

:

$$A^{[L]} = g^{[L]}(Z^{[L]}) = \hat{Y}$$

1939

13

फाल्गुन / FALGUN

शक / Shaka

4

मार्च / MARCH

रविवार / SUNDAY

Backward prop:

$$dZ^{[L]} = A^{[L]} - Y$$

$$dW = (1/m) dZ^{[L]} A^{[L]}^T$$

$$db^{[L]} = (1/m) \text{np.sum}(dZ^{[L]}, \text{axis}=1, \text{keepdims=True})$$

$$dZ^{[L-1]} = dW^{[L]}{}^T dZ^{[L]} g'^{[L]}(Z^{[L-1]})$$

$$dZ^{[1]} = dW^{[1]}{}^T dZ^{[2]} g'^{[1]}(Z^{[1]})$$

$$dW^{[1]} = (1/m) dZ^{[1]} A^{[1]}{}^T$$

$$db^{[1]} = (1/m) \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$$