

# Courese 1 : Neural Networks and Deep Learning..

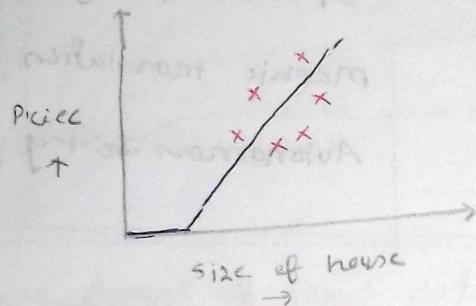
start Date  
30-11-2016  
3:20 PM.

## WEEK 1 :- Introduction to Deep Learning.

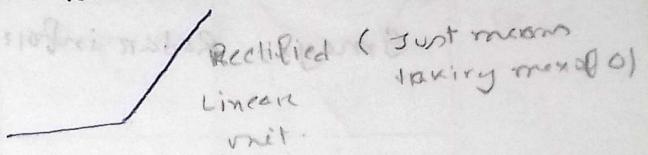
### L-1 : what is neural networks. 11.05 AM (1-12-18)

Neural network is a powerful learning algorithm inspired by how the brain works. Deep learning refers to train neural network sometimes very large neural network.

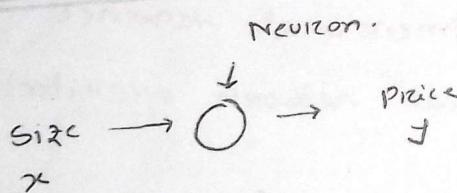
Housing price prediction.



An price can never be negative so we are creating a function called rectified Linear unit (ReLU).



### Simple Network.

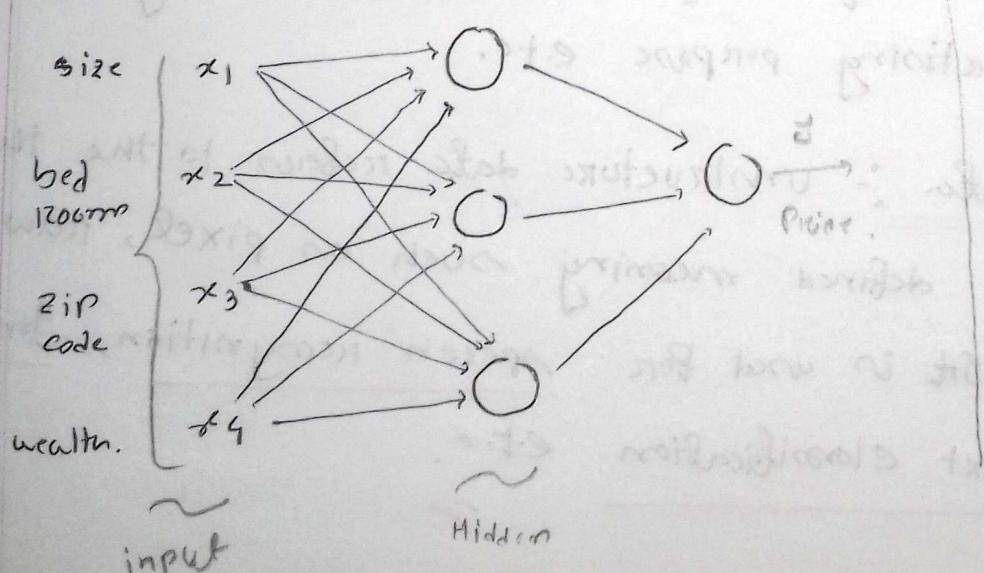


$x$  (input)  $\rightarrow$  size of house.

$y$  (output)  $\rightarrow$  Price.

The neuron implements the function ReLU.

### Large network.



magic of neural network is that it automatically generate the hidden units we only need to give input  $x$  and output  $y$ .

## L-2: Supervised Learning with Neural Networks (11:35 AM 01-12-18)

In supervised learning, we are given a dataset and already know what our correct output should look like. Some Examples.

Input ( $x$ )	Output ( $y$ )	Application.
Home features.	Price	Real Estate.
Ad, user info.	Click on Ad? (0/1)	online Advertising.
Image	object detection.	photo tagging.
Audio	Text transcript	speech recognition
English.	Chinese	Machine translation
Image, Radar info.	position of other cars.	Autonomous driving.

\*\* standard NN used for standard neural network application.

\*\* CNN used for image application.

\*\* RNN used for one dimensional sequence data and temporal element. { Hybrid neural network architecture is also used }

Structure data :- Structure data refers to the things that has a defined meaning such as price, age, etc. Structure data in database of data. It is used for housing price prediction, advertising purpose etc.

Unstructure data :- Unstructure data refers to the things that has not defined meaning such as pixel, raw audio, text. It is used for speech recognition, image recognition, Text classification etc.

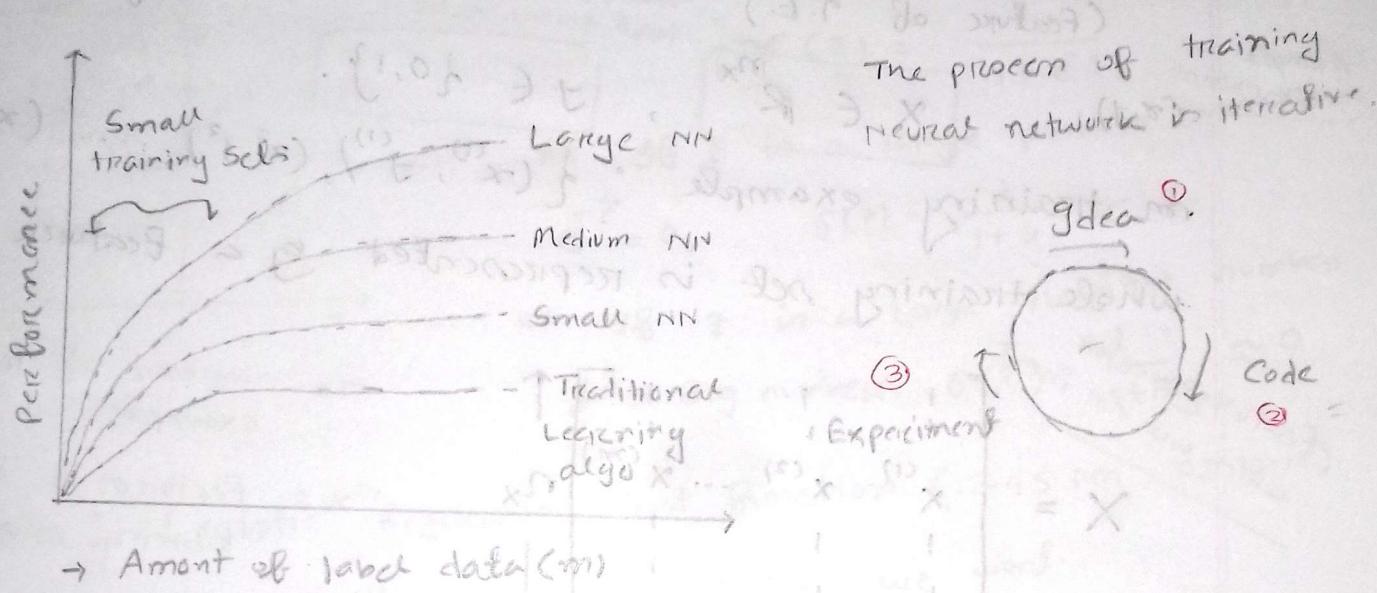
L 3 : why Deep Learning taking off? (12.10 pm 1-12-18)

Deep learning is taking off due to,

- ① Large amount of data
- ② Faster computation.

Scale drove deep learning Program.

- ③ Innovations in NN algorithm



L 4 : About this course. - (12.25 pm 1-12-18)

Week 2:

## P1 : Logistic Regression as a Neural Network

E1 : Binary classification. (12:05 pm 02-12-18)

In binary classification result is discrete value output.

Here,  $x \rightarrow y$ . (corresponding label of training example)

(Feature of T.E.) - dimension of input feature vector

A single training example.

$(x, y)$

$$x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

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

whole training set is represented by a Feature matrix

feature vector which will be input.

$$X = \begin{bmatrix} & & & \\ & & & \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ & & & \\ & & & \end{bmatrix} \quad \begin{array}{c} \uparrow \\ n_x \\ \downarrow \\ \text{number of features} \end{array}$$

$m = \text{number of training examples}$   
 $n_x = \text{features of each training examples}$

$\leftarrow m \rightarrow$

label of the training example is denoted by a row vector.

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

$$X.\text{Shape} = (n_x, m).$$

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

$X \in \mathbb{R}^{n_x \times m}$  is the input matrix.

$x^{(i)} \in \mathbb{R}^{n_x}$  is the  $i^{\text{th}}$  example represented as a column vector.

$Y \in \mathbb{R}^{n_y \times m}$  is the label matrix.

$y^{(i)} \in \mathbb{R}^{n_y}$  is the output label for the  $i^{\text{th}}$  example

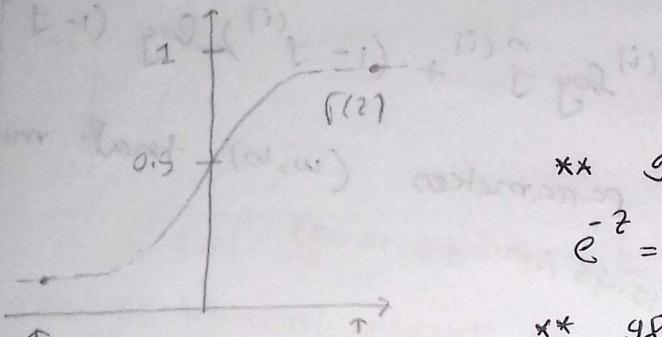
— (can be viewed as a very very small neural network)

## L2: Logistic Regression. (12.30 pm 02-12-18).

Given  $X$ , we want to predict  $\hat{y} = p(y=1|x)$ .  $[0 \leq \hat{y} \leq 1]$ .

$X \in \mathbb{R}^{n_x}$ , Parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ .

output  $\hat{y} = \sigma(w^T x + b)$ . Now:  $w^T x + b = z$



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

\* if  $z$  is a large positive number.

$$e^{-z} = 0, \quad \sigma(z) = \frac{1}{1+0} \approx 1.$$

\* if  $z$  is a large negative number.

$$e^{-z} = \text{Big number}, \quad \sigma(z) = \frac{1}{1+\text{Big}} \approx 0$$

## L3: Logistic Regression Cost Function. (12.45 pm 02-12-18)

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

$\hat{y}^{(i)} \approx y^{(i)}$  [Predicted function  $\hat{y}^{(i)}$  to be as close as actual output  $y^{(i)}$ ]

Loss (error) Function:  $\ell(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$  [squared error function]

we can not use this squared error function. Because if we use it we do not get global optima because there are multiple local optima.

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

$$\ell(\hat{y}, y) = -\log \hat{y} \leftarrow \text{want } \log \hat{y} \text{ large, want } \hat{y} \text{ large.} \\ [\text{close to 1}]$$

if  $y=1$ :  $\ell(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{want } \log(1-\hat{y}) \text{ large...}$

if  $y=0$ :  $\ell(\hat{y}, y) = -\log \hat{y} \leftarrow \text{want } \log \hat{y} \text{ small.} \\ [\text{close to 0}]$

Cost Function calculation on a single training set.

Cost Function measures how well our algorithm is doing in entire training set.

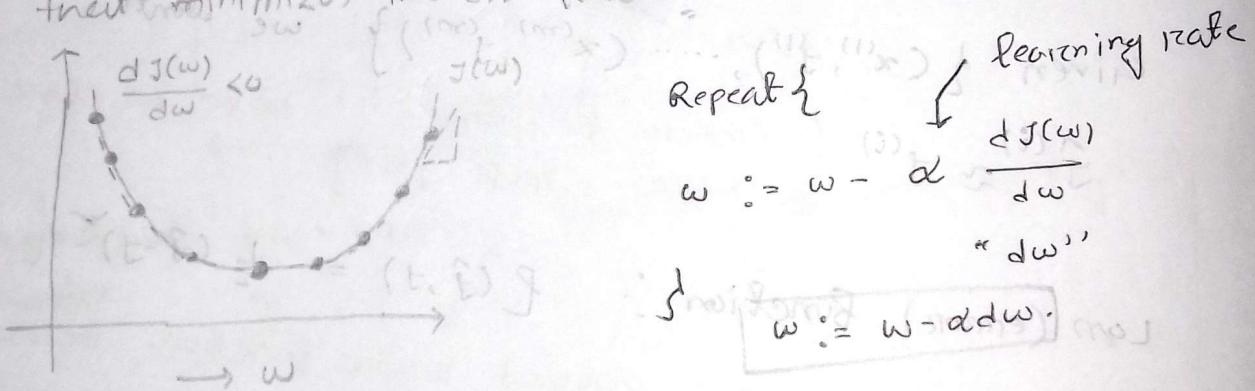
$$\text{Cost function: } J(w, b) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)})$$

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

We are going to find parameters  $(w, b)$  that minimizes Cost Function.

L9: Gradient Descent. (video w/ previous slide)

\* Gradient descent is an algorithm that helps to find  $(w \text{ and } b)$  that minimizes the cost function.



The amount by which you update  $w$

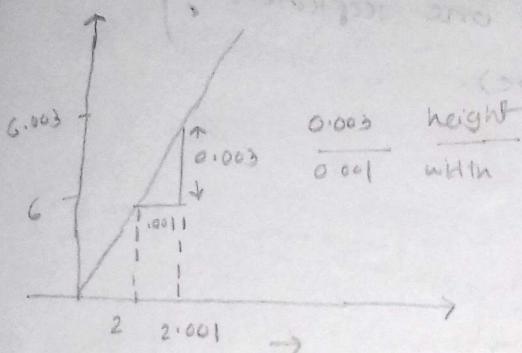
$$(t+1)^{\text{th}} (w_1) = w_1 - \alpha \frac{\partial J(w_1)}{\partial w} \rightarrow dw$$

$$J(w_1, b) \rightarrow b - \alpha \frac{\partial J(w_1, b)}{\partial b} \rightarrow db$$

The amount by which you update  $b$

L 4 : Derivatives. (11.10 Am 03-12-18).

In mathematics, the derivative is a way to show rate of change. The amount by which a function is changing at one given point.



$$f(a) = 3a, \quad a=2, \quad f(a)=6$$

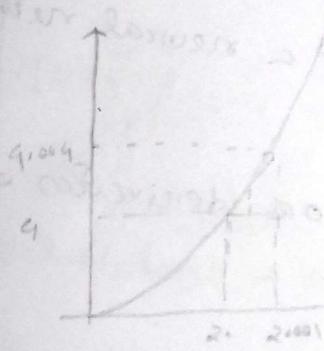
$$(a+0.001) = (2.001) \quad a=2.001 \quad f(a)=6.003$$

slope (derivative) of  $f(a)$  at  $a=2$  is 3.

$\left[ \frac{df(a)}{da} \right]$  The amount of change in  $(a)$  cause change in  $f(a)$ .

L 5 : More derivative examples. (11.40 Am 03-12-18)

Derivatives can be different at different points of a line.



$$f(a) = 4$$

$$f(a) = 4.004$$

slope (derivative) of  $f(a)$  at  $a=2$  is 4.

slope of  $f(a)$  at  $a=5$  is 10.

Ans 3 for Derivative example, 92cm after problem no.

$$f(a) = a^2 \quad \frac{d}{da} f(a) = 2a$$

$$f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2$$

$$f(a) = \ln(a) \quad \frac{d}{da} f(a) = \frac{1}{a}$$

L6 : Computation Graph. (12.15 pm 03-12-18)  
 Help to find more complex derivatives.

In neural network computation can be represented by a data flow graph called computation graph. Computation graph consists of a set of nodes each one representing an operation.

Example :-  $J(a, b, c) = 3(a+b+c)$ .

$$U = b + c$$

$$V = a + U$$

$$J = 3V.$$

— — —

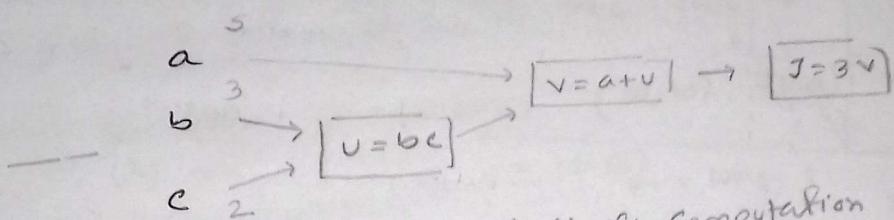


Fig:- A simple computation graph.

\*\* Computation graph is useful when we want to optimize a special variable in this case that is  $J$ .

A neural network is organized in terms of

- Forward Pass. (Computes output of a neural network) [it is a left to right pass]
- Backward Pass (computes gradients or derivatives of a neural network) [it is a right to left pass]

L7 : Derivatives with a computation graph. (3.50 pm 16-12-18)

$\frac{dJ}{dv} = ?$  [it means that if we change the value of  $v$  how does it effect the value of  $J$ ].

[Chain rule goes from concept  $a \rightarrow v \rightarrow J$ ] [Successive derivatives].

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

The change in  $J$  due to change in  $a$  can be calculated by.

First computing how much  $v$  changes due to changing  $a$ . multiplied by how much  $J$  changes due to change in  $v$ .

\* blue lines compute the output of neural network.

\* Red lines compute the derivative of neural network.

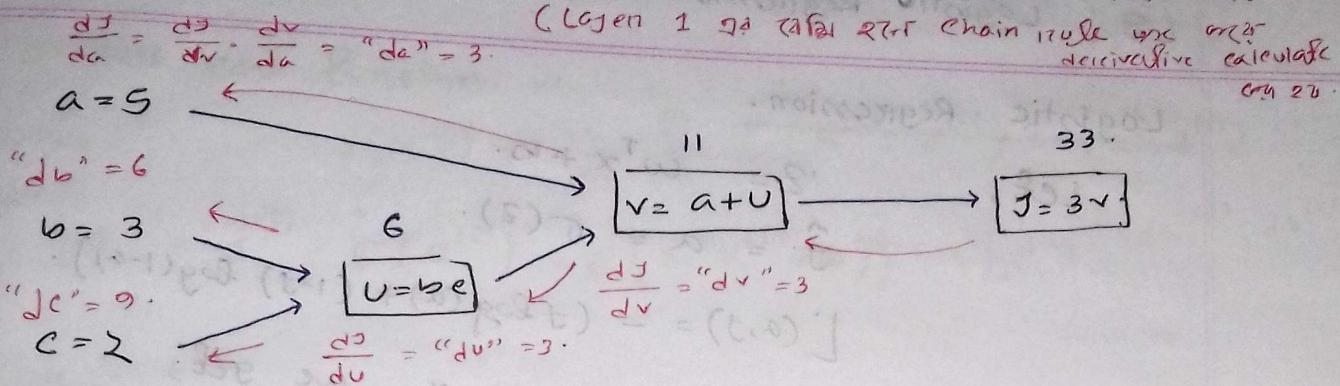


Fig:- Computation graph.

notation:- we have to compute the derivative of final output with respect to some intermediate variable. It can be represented by,

$$\frac{d \text{final output var}}{d \text{var}} = ("d \text{var}")$$

$$\frac{dJ}{dv} = 3.$$

- we will write

[chain rule.  
unit layer a value change  
outgoing layer g for change  $\Delta J$ ]

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da} = 3 \cdot 1 = 3.$$

$$\frac{dJ}{db} = \frac{dJ}{dv} \cdot \frac{dv}{db} = 3 \cdot 1 = 3.$$

$$\frac{dJ}{dc} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{dc} = 3 \cdot 1 \cdot 2 = 6$$

$$\frac{dJ}{dc} = \frac{dJ}{dv} \cdot \frac{dv}{du} \cdot \frac{du}{dc} = 3 \cdot 3 \cdot 3 = 9.$$

Example  
value put  
var? mark

" $db$ " or " $dc$ " mean derivatives of final output with respect to intermediate variables " $b$ " and " $c$ ".

L8 : Logistic Regression Gradient Descent. [ 4.30 PM  
16.12.18 ]

Logistic Regression.

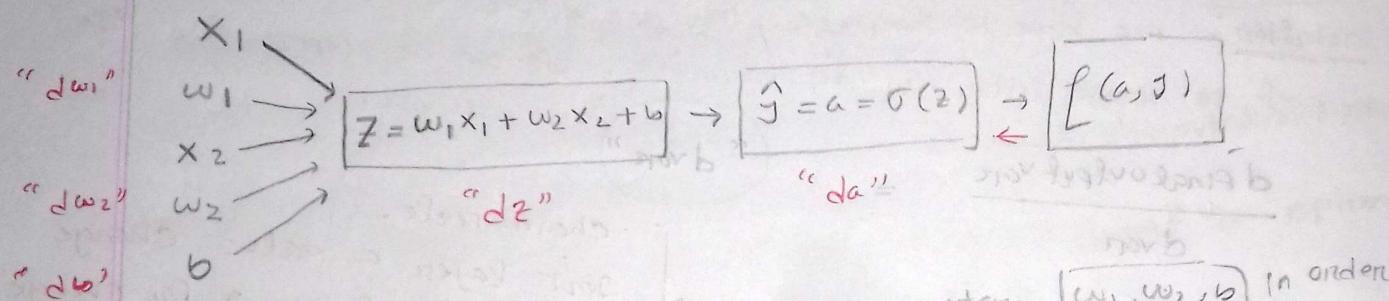
Let,

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

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

$$L(a, j) = - (j \log(a) + (1-j) \log(1-a))$$

Now considering only two variables  $x_1, x_2$  we get,



\*\* Our objective is to modify the parameters

to reduce  $L(a, j)$

Now,  $dL(a, j) = \frac{dL(a, j)}{da} = -\frac{j}{a} + \frac{1-j}{1-a}$

$$dL(a, j) = \frac{dL(a, j)}{dz} = \frac{dL(a, j)}{da} \cdot \frac{da}{dz} = \left(-\frac{j}{a} + \frac{1-j}{1-a}\right) \cdot a \cdot (1-a)$$

$$= a - j$$

$$dL(a, j) = \frac{dL}{dw_1} = x_1 \cdot dL$$

$$dL(a, j) = \frac{dL}{dw_2} = x_2 \cdot dL$$

$$dL(a, j) = dL$$

one step of gradient descent is,

$$w_1 := w_1 - \alpha \cdot dL$$

$$w_2 := w_2 - \alpha \cdot dL$$

$$b := b - \alpha \cdot dL$$

L9: Gradient Descent of m Examples.

[ 5:00 pm  
16.12.18 ]

Overall cost function is.

Prediction for

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m \ell(a^{(i)}, j^{(i)}). (x^{(i)}, j^{(i)})$$

In example

$$[a^{(i)} = j^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^T x^{(i)} + b)] \quad [d\omega_1^{(i)}, d\omega_2^{(i)}, db^{(i)}]$$

Code :- Let

$$j=0, d\omega_1=0, d\omega_2=0, db=0$$

For  $i=1$  to  $m$

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

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

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

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

Features  $x_1, x_2$

Loop for  $i$

Up to  $m$ .

$$\left. \begin{array}{l} d\omega_1+ = x_1^{(i)} dz^{(i)} \\ d\omega_2+ = x_2^{(i)} dz^{(i)} \end{array} \right\} \begin{array}{l} n=2 \\ \text{for } i \\ \text{feature} \end{array}$$

vectorization via  
Easily implement  
using Numpy.

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

$$j/m$$

$$d\omega_1/m, d\omega_2/m, db/m.$$

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

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

[This is simple code for  
one iteration]

## P2 : python and vectorization

L1 : vectorization. (10.30 Am, 17-12-18)

\* Faster implementation of our vectorized version (10.21.18) for 225.  
vectorization just matrix multiplication w.r.t. for 225.

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

$$w \in \mathbb{R}^{n_x} \\ w = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \quad x \in \mathbb{R}^{n_x} \\ x = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

Non-vectorized.

$$z = 0$$

For i in range (n\_x)

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

$$z += b.$$

time required = 475 ms

vectorized.

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

time required = 1.5 ms.

\* Built-in function use core function  
output 1.34 ms or built-in function  
use core Try opt.

## L2 : more vectorization Example (10.45 Am, 17-12-18)

\*\* for neural network programming whenever possible avoid explicit for loop. (Built-in function)

Let,

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \text{ we have to find } v = \begin{bmatrix} e^{v_1} \\ \vdots \\ e^{v_n} \end{bmatrix}.$$

non-vectorized.

$$v = np.zeros(n)$$

for i in range (n):

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

vectorized.

import numpy as np

$$v = np.exp(v).$$

\*\* Elementwise operation

(10.21.18) or 225.

$$\rightarrow np.log(v)$$

→ Elementwise log.

$$\rightarrow np.abs(v)$$

→ Elementwise absolute value

$$\rightarrow np.maximum(v, 0)$$

→ Maximum of v.

$$\rightarrow v ** 2$$

→ Elementwise square.

### L3: Vectorizing Logistic Regression. (11:00 AM 17-12-18)

During Logistic regression predictions are,

$$z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(3)} = w^T x^{(3)} + b$$

$$a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)})$$

This implementation can be vectorized,   
 ( $n_x$  = number of features  
 $m$  = number of Training examples)

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

$$\begin{bmatrix} z^1 & z^2 & \dots & z^m \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & \dots & w_{n_x} \end{bmatrix} \times \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix} + \begin{bmatrix} b & b & \dots & b \end{bmatrix}$$

$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & \dots & w^T x^{(3)} + b \end{bmatrix}$$

vectorized

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

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(Z)$$

computing

Forward propagation

### L4: Vectorizing Logistic Regression Gradient output (11:15 AM 17-12-18)

$$A = [a^{(1)} \ \dots \ a^{(m)}]$$

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

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

$$dZ = A - Y$$

$$(1) \quad dw = 0$$

$$(2) \quad dw = x^{(1)} dz^{(1)}$$

$$(3) \quad dw = x^{(2)} dz^{(2)}$$

vectorized

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

$$= \frac{1}{m} \left[ x^{(1)} \ \dots \ x^{(m)} \right] \begin{bmatrix} dZ^{(1)} \\ \vdots \\ dZ^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \left[ x^{(1)} dZ^{(1)} + \dots + x^{(m)} dZ^{(m)} \right]$$

$$\begin{aligned}db &= 0 \\db + &= dz^{(1)} \\db + &= dz^{(2)} \\&\vdots \\db + &= dz^{(m)} \\db/m &= \end{aligned}$$

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

$$db = \frac{1}{m} np \cdot \text{sum}(dz)$$

Gradient or vectorize code

Vectorized implementation of logistic regression.

$$j=0, dw_1=0, dw_2=0, db=0$$

For i = 1 to m:

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

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

vectorization

for m training

Example

for iter in range(1000):

1000 iteration of gradient descent.

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

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

$$A = \sigma(z)$$

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

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

$$dz = A - Y$$

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

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

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

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

$$db = \frac{1}{m} np \cdot \text{sum}(dz)$$

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

for m iterations  
for n features

$$\left. \begin{array}{l} w := w - \alpha dw \\ b := b - \alpha db \end{array} \right\} \text{update}$$

This the code of iteration

## L9: Broadcasting in Python. [10:00 pm , 17-12-18]

Notebook  
JAN 2018

\* Fair loop PER 42213 OR 0720 225 3<sup>rd</sup> year Comp technique ENGG 2018.

$$cal = A \cdot \text{sum}(axis=0)$$

Axis 0 → Sum vertically ↑  
Axis 1 → Sum horizontally →

$$\text{percentage} = 100 \times A / \text{cal.reshape}(1, 4)$$

• reshape matrix size from 4x14 into 4x213 OR 225.

General principle of broadcasting.

1. given (m,n) matrix as well as (1,n) matrix of operation  
or a given (1,m) matrix broadcast 225, (m,n) matrix 225

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 \rightarrow \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix} \rightarrow \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}. \quad \begin{array}{l} \text{Broadcasting normally} \\ +, -, *, /, \text{etc.} \\ \text{No error.} \end{array}$$

2. given (m,n) matrix as well as (m,1) matrix of operation  
or a given (m,1) matrix broadcast 225 (m,n) matrix 225.

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + 100 \rightarrow \begin{bmatrix} 100 & 100 & 100 \end{bmatrix} \rightarrow \begin{bmatrix} 101 & 102 & 103 \end{bmatrix}.$$

L5: A note on python/numpy vector. [10:25 17-12-18]

\* Do not make a rank 1 array. → This will lead to error,

$$\left\{ \begin{array}{l} a = np.random.rand(5) \\ a.shape = (5, 1) \end{array} \right\} \text{ v.v. } \quad a = np.random.rand(5). \quad \begin{array}{l} \text{Column vector} \end{array}$$

$$a = np.random.rand(1, 5) \quad \begin{array}{l} \text{Row vector} \end{array}$$

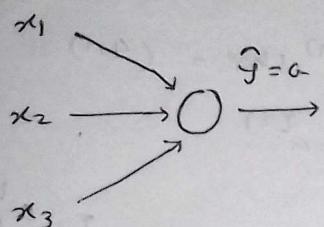
correct ( $a.shape == (5, 1)$ )

④ 225 rank 1 Array 225 0720.5 reshape into  
row or column vector 225 OR 0720 225.

week 3 : Shallow neural networks. [other mathematical function can be used]

# L1: Neural Networks overview.

Logistic regression is a simple neural network, it is.



$$x \quad \backslash$$

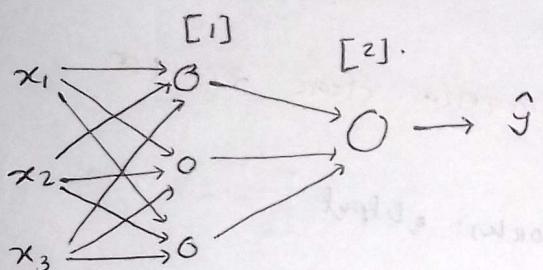
$$\omega \quad \quad \quad \boxed{z = \omega^T x + b} \rightarrow \boxed{a = g(z)} \rightarrow \boxed{L(a, y)}$$

$$b \quad /$$

$\downarrow z \quad \quad \quad \downarrow a$

Fig :- steps of logistic regression

Typically a neural network has multiple hidden layers, output of one layer is taken as input of another layer.

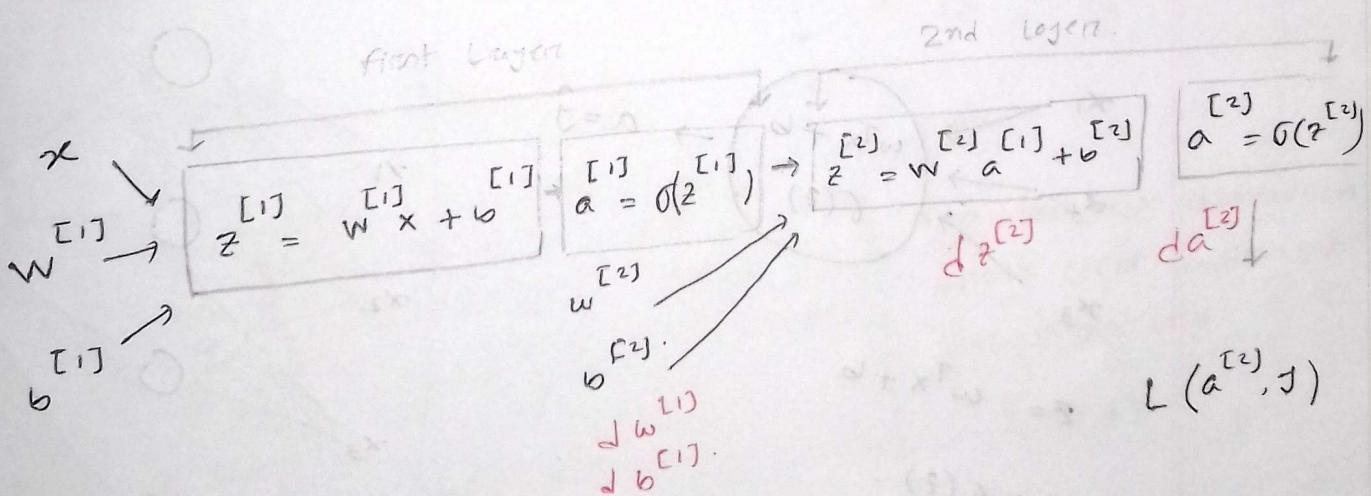


$x^{(i)}$  → This means  $i^{\text{th}}$  training example

$\times [i]$   $\rightarrow$  This means it's hidden layer.

square bracket [i] represents  
 with hidden layer of  $NN =$

Simple calculation of this crucial network,



$w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ . ~~for~~ weight and bias. Neural network backpropagation go ~~and~~ go weight ~~as~~  $474125cm$   $(25^{\circ}C)$ .

## L2: Neural Network Representation [12:06 Am 19-12-18]

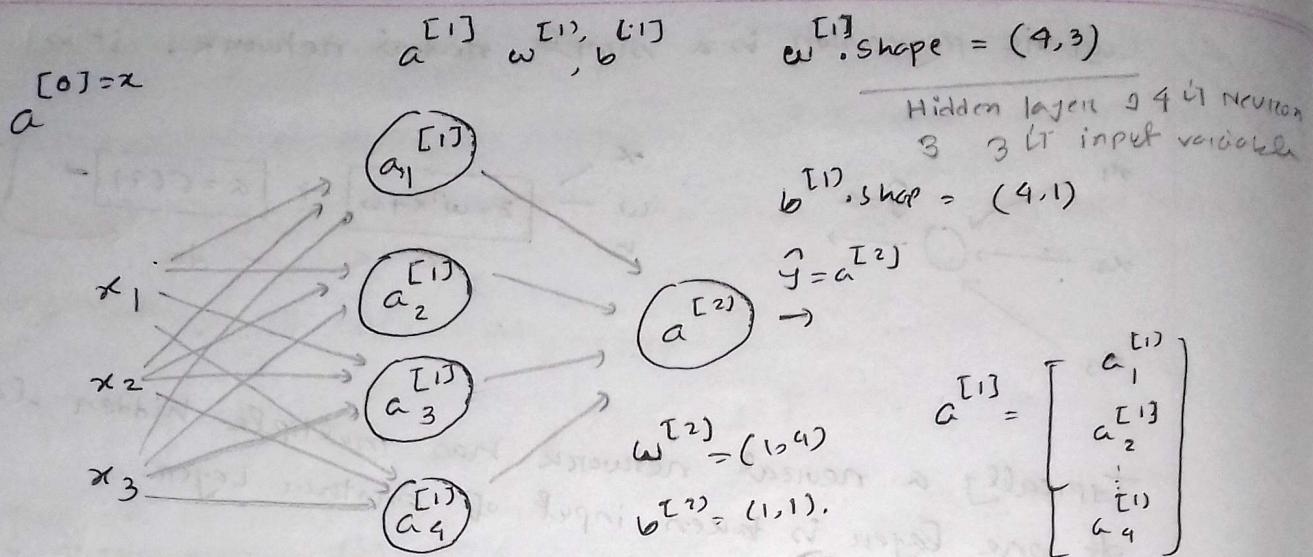


Fig :- 2 layer NN.

\* Layer 0: dimension & not matrix clear 270 275

## L3: Computing a Neural Networks' Output [12:20 Am 19-12-18]

Computation on a single neuron,

Neuron mit 3x Dimensionen

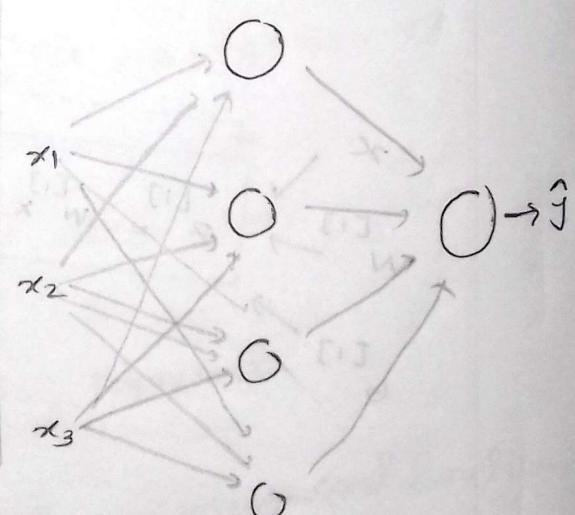
$x_1 \rightarrow$

$x_2 \rightarrow$

$x_3 \rightarrow$

$z = w^T x + b$

$a = f(z)$



→ simply a neural network do this computation for each neuron in each layer.

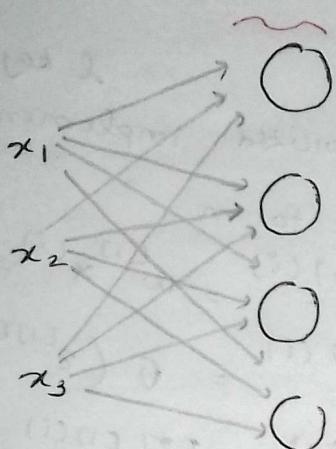
(Easy concept)

fig :- 2 layer NN

Given  $w_i^{[1]} \rightarrow$  dimension 2x1  
After taking input & bias, sum  
 $w_1^{[1]} = (3, 1)$

$[l] \leftarrow$  Layer No.  
 $a_i \leftarrow$  node in the layer

Now Computation of each part  
Layer 1.



$$z_1^{[1]} = w_1^{[1] \cdot T} + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1] \cdot T} + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1] \cdot T} + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1] \cdot T} + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})$$

vectorizing

$$\begin{matrix} w^{[1]} \\ \hline -w_1^{[1] \cdot T} \\ -w_2^{[1] \cdot T} \\ -w_3^{[1] \cdot T} \\ -w_4^{[1] \cdot T} \end{matrix} \quad \begin{matrix} x = a^{[0]} \\ \hline x_1 \\ x_2 \\ x_3 \end{matrix} \quad \begin{matrix} b^{[1]} \\ \hline b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{matrix} \quad \begin{matrix} z^{[1]} \\ \hline z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{matrix} \quad \begin{matrix} a^{[1]} \\ \hline a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{matrix}$$

$$\begin{matrix} x = a^{[0]} \\ \hline x_1 \\ x_2 \\ x_3 \end{matrix} + \begin{matrix} b^{[1]} \\ \hline b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{matrix} = \begin{matrix} z^{[1]} \\ \hline z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{matrix}$$

$$\sigma(z^{[1]}) = \begin{matrix} a^{[1]} \\ \hline a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{matrix}$$

Equations of 2-layer NN :-

$$z^{[1]} = w^{[1]} \cdot a^{[0]} + b^{[1]}$$

$$\rightarrow (4, 1) \quad (4, 3) \quad (3, 1) \quad (4, 1)$$

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

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

$$(1, 1) \quad (1, 4) \quad (4, 1) \quad (1, 1)$$

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

$$(1, 1) \quad (1, 1)$$

\* matrix of dimension  
and how to implement  
 $(4, 3) \times (3, 1) = (4, 1)$

(4, 1)

Easy

L 4: Vectorizing Across multiple Examples. [ 6.15 PM 19.12.18 ]

For a training example what we get from neural network,

$$x \longrightarrow a^{[2]} = \hat{y} \quad a^{[l](i)} \rightarrow i^{\text{th}} \text{ training example}$$

For m training examples,

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

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

— — — — —

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

unvectorized implementation.

for  $i=1$  to  $m$ ,

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

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

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

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

$\rightarrow$  Just write equation for each training example separately.

\*\* vectorized implementation.

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

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

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

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

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

( $n \times m$ )  
number of input features

$$z = \begin{bmatrix} \rightarrow & & & & \\ | & | & | & & | \\ z^{(1)(1)} & z^{(1)(2)} & \dots & z^{(1)(m)} \\ | & | & & | \end{bmatrix}$$

$$A = \begin{bmatrix} \downarrow & & & & \\ | & | & | & & | \\ a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(2)} \\ | & | & & | \end{bmatrix}$$

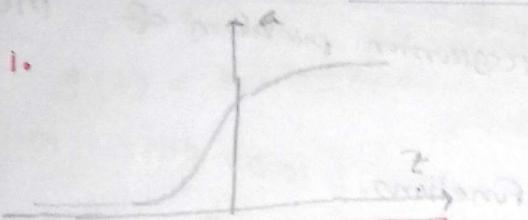
row # Number of training examples  
col # Number of units in a hidden layer

## L5:- Explanation for vectorized implementation. [6.35 pm 19.12.18]

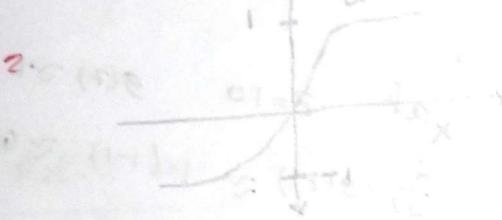
[ Vectorize our implementation or we can implement our code in a better way and avoid many loops = [Easy] ]

## L6: Activation Functions. [6.45 pm 19.12.18]

→ The activation function of a node defines the output of that node given an input or net of inputs. Several types of activation functions are available.



sigmoid :  $a = \frac{1}{1+e^{-z}}$



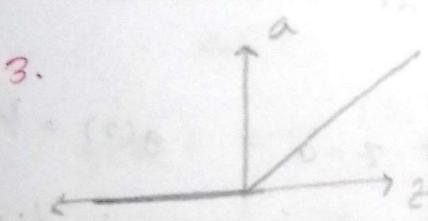
\* Almost never use it because it is slow.

\*\* If the problem is a binary classification problem it can be used in final layer.

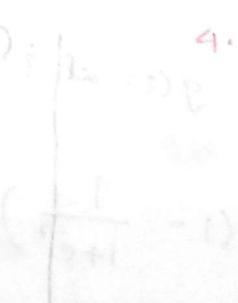
$$\tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

\*\* It is superior than sigmoid.

→ The main drawback of gradient descent of sigmoid and tanh is that if  $z$  is very small or very large gradient descent will be very slow.



ReLU =  $\max(0, z)$



Leaky ReLU =  $\max(0.01z, z)$ .

\*\* In both of these algorithms our network will learn much faster. (It is better to use ReLU or Leaky ReLU).

\*\* We can use different activation function at different layers of neural network.

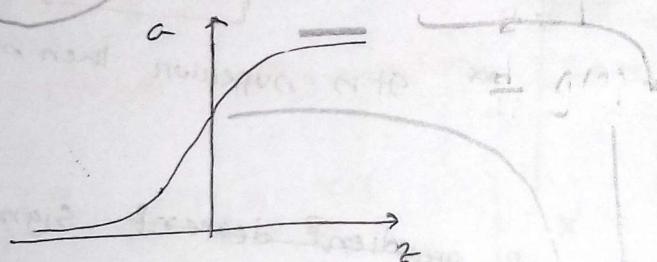
(Ques & Ans for)

L7: why do you need non-linear activation functions [7.20 pm 19.12.18]

- If we use linear activation function then hidden layers in the network will be useless because composition of two linear function is a linear function.
- Unless we throw some non-linearity in there we do not computing more interesting functions even if we go deeper in network.
- Linear function is used in regression problem of ML

L8: Derivatives of Activation Functions. [10.00 pm 19.12.18]

Sigmoid activation Function.



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

$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$

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

$$\Rightarrow g(z) (1 - g(z))$$

$$g'(z) = a(1-a)$$

$$[g'(z) = a(1-a), a = \frac{1}{1+e^{-z}}]$$

$$\begin{array}{l} \text{at } z=10 \quad g(z) \approx 1 \\ \text{at } z=-10 \quad g(z) \approx 0 \end{array}$$

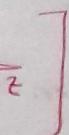
$$\frac{d}{dz} (g(z)) \approx 1 (1-1) \approx 0$$

$$\text{at } z=0 \quad g(z) = \frac{1}{2}$$

$$\frac{d}{dz} (g(z)) \approx 0. (1-0) \approx 0$$

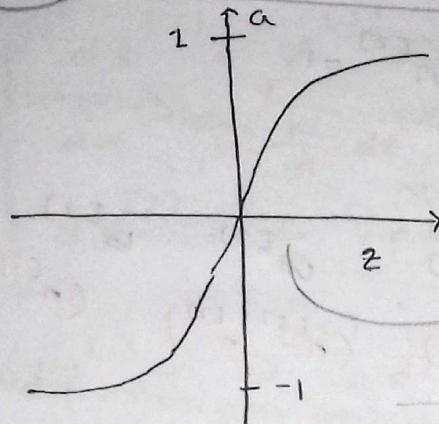
$$\text{at } z=0 \quad g(z) = \frac{1}{2}$$

$$\frac{d}{dz} (g(z)) = \frac{1}{2} (1-\frac{1}{2}) = \frac{1}{4}$$



tanh

activation function.



at  $z=10$ ,  $\tanh(8) \approx 1$   
 $g'(z) \approx 0$

at  $z=0$   $\tanh(0) = 0$

$g'(0) = 1$

at  $z=-10$   $\tanh(-8) \approx 0$

$g'(-8) = 1$

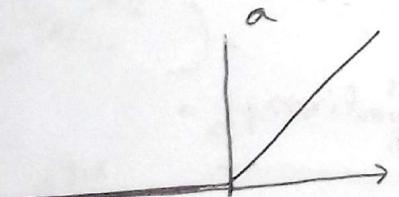
$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned} g'(z) &= \frac{d}{dz}(g(z)) = \text{slope of } g(z) \text{ at } z \\ &= 1 - (\tanh(z))^2 \\ &= 1 - a^2. \end{aligned}$$

$$g'(z) = 1 - a^2$$

$$\begin{aligned} a &= \tanh(z) \\ &= \frac{e^z - e^{-z}}{e^z + e^{-z}} \end{aligned}$$

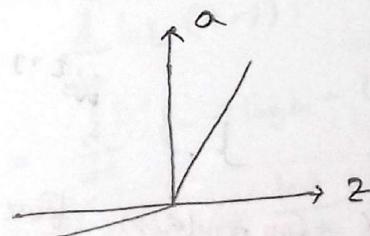
ReLU and Leaky ReLU.



ReLU

$$g(z) = \max(0, z)$$

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

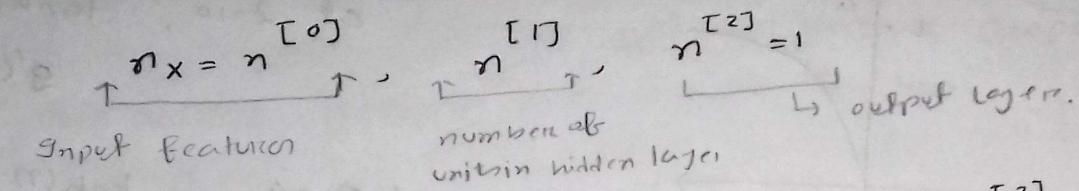


$$g(z) = \max(0.01z, z)$$

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

L7: Gradient descent for Neural Networks. [10.20 19.12.18]

Consider a neural network with one hidden layer,



Parameters :  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$ .  
Dimensions  $\rightarrow (n^{[1]}, n^{[0]}) \quad (n^{[2]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[0]}, 1)$ .

\* In neural network we have to initialize parameters randomly.

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

Gradient descent,

Repeat {

Compute prediction  $(\hat{y}^{(1)}, \hat{y}^{(2)}, \dots, \hat{y}^{(m)})$

$$dw^{[1]} = \frac{dJ}{dw^{[1]}}, \quad db^{[1]} = \frac{dJ}{db^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha \cdot dw^{[1]}.$$

$$b^{[1]} = b^{[1]} - \alpha \cdot db^{[1]}.$$

$$w^{[2]} = \dots \quad b^{[2]} = \dots$$

Forward propagation  
to calculate cost.

dt week 24  
Hazari GMP (featur)

\*\* formulas for computing derivatives.

forward propagation :-

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

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

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

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

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

$g(z)$

↓  
Activation function.

binary classification 2-class  
sigmoid function  
probability output

and sigmoid function

probability output

probability output

probability output

~~Backpropagation :-~~

Similar  
as logistic  
regression:

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} \approx -\frac{1}{m} dZ^{[2]} \cdot A^{[1] T}$$

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

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

prevent python to  
produce rank 1 array

$$dZ^{[1]} = W^{[2], T} \cdot dZ^{[2]} * g^{[1]} \cdot (Z^{[1]})$$

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

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

L10 : Backpropagation Intuition.

Logistic Regression.

$$\begin{array}{c}
 \begin{array}{ccc}
 x & \rightarrow & \boxed{z = w^T x + b} \\
 w & \rightarrow & \boxed{a = g(z)} \\
 b & \rightarrow & \boxed{L(a, y)}
 \end{array}
 \\ 
 \begin{array}{l}
 "dZ" = \frac{\partial L}{\partial a} \cdot \frac{da}{dz} \\
 dw = dZ \cdot x \\
 db = dZ
 \end{array}
 \quad
 \begin{array}{l}
 da = \frac{d}{da} L(a, y) \\
 = \frac{d}{da} (-y \log a - (1-y) \log(1-a)) \\
 = -\frac{y}{a} + \frac{1-y}{1-a}
 \end{array}
 \end{array}$$

*(off) gradient based error metric from 21.*

*activation function.*

$"dZ" = da \cdot g'(z)$

$$x_1 \rightarrow \textcircled{1} \rightarrow \hat{y} = 0.6$$

$$x_2 \rightarrow \textcircled{2} \rightarrow \hat{y} = 0.3125$$

$$x_3 \rightarrow \textcircled{3} \rightarrow \hat{y} = 0.3125$$

④ Neural Network

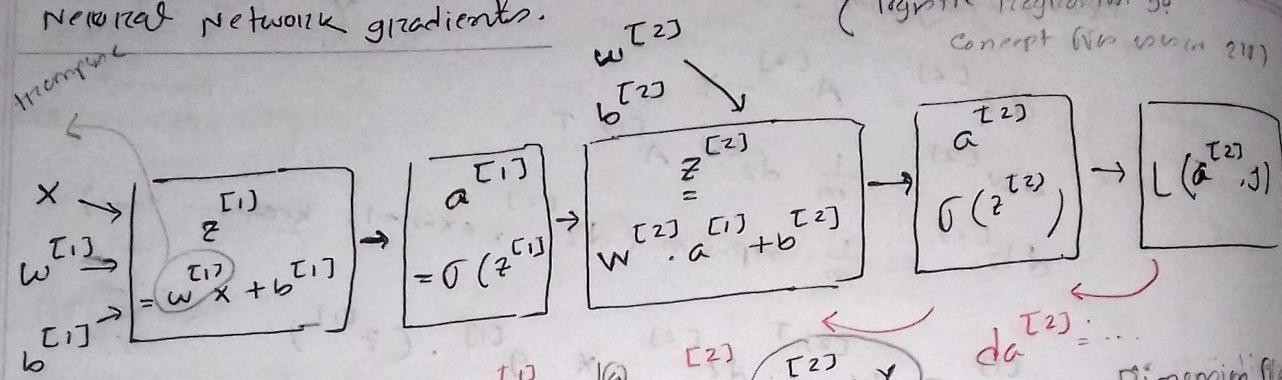
→  $\hat{y}$ -calculation

\* Most important,  
matrix dimension 2x2x2 4x4x4 2x2.

$$\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \begin{array}{c} = \\ = \\ = \end{array} \begin{array}{c} o-j \\ o-j \\ o-j \end{array}$$

$$n^{[0]} n^{[1]} n^{[2]} = 1$$

Neural network gradients.



$$\textcircled{1} \quad dz^{[1]} = w^{[2] \cdot T} dz^{[2]} \quad da^{[1]} = \dots$$

Dimension matching.

$$w^{[2]} \rightarrow (n^{[2]}, n^{[1]})$$

$$z^{[1]}, dz^{[1]} \rightarrow (n^{[1]}, 1)$$

$$z^2, dz^{[2]} \rightarrow (n^{[2]}, 1)$$

$$\therefore dz^{[1]} = [n^{[1]}, n^{[2]}] \cdot (n^{[2]}, 1) * [n^{[1]}, 1]$$

$$= (n^{[1]}, 1)$$

$$\textcircled{2} \quad dw^{[1]} = dz^{[1]} \cdot x^T$$

$$\textcircled{3} \quad db^{[1]} = dz^{[1]}$$

Summary.

$$dz^{[2]} = a - y$$

$$dw^{[2]} = dz^{[2]} \cdot a^{[1]} \cdot x^T$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = w^{[2] \cdot T} dz^{[2]} * g'(z^{[1]})$$

$$dw^{[1]} = dz^{[1]} \cdot x^T$$

$$db^{[1]} = dz^{[1]}$$

(logistic regression by  
concept (WV) from 211)

$$\begin{aligned} \textcircled{4} \quad dz^{[2]} &= a - y \\ \textcircled{5} \quad dw^{[2]} &= dz^{[2]} \cdot a^{[1]} \cdot x^T \\ \textcircled{6} \quad db^{[2]} &= dz^{[2]} \end{aligned}$$

Dimension filter  
4 12x3, 3x1  $\rightarrow$   
4 12x3 2x3.

2x3  $\rightarrow$   
2x3

logistic regression error  
-0.6f

$$dz^{[2]} = A^{[2]} - Y$$

$$dw^{[2]} = Y^T dz^{[2]} \cdot A^{[1] \cdot T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dz^{[2]}, axis=1, keepdim=True)$$

$$dz^{[1]} = w^{[2] \cdot T} dz^{[2]} * g'(z^{[1]})$$

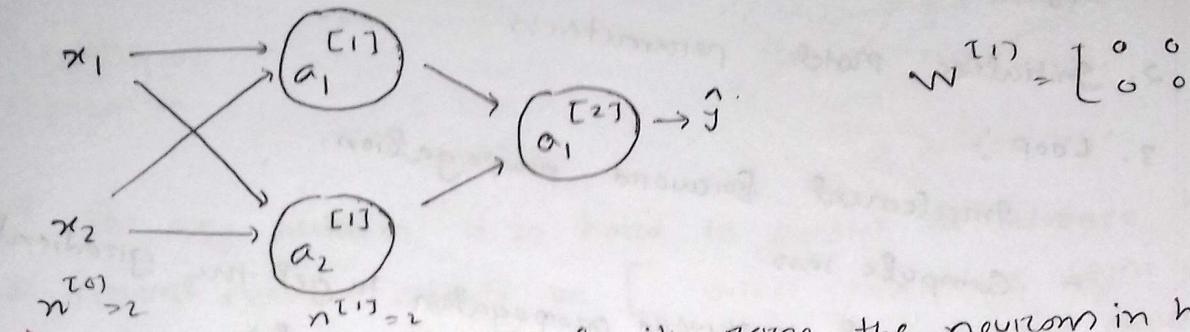
$$dw^{[1]} = \frac{1}{m} dz^{[1]} \cdot x^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \dots)$$

✕ Forward propagation (4 equations)  
 ✕ Backward propagation (6 equations) > Programming implement up to part.

L11:- Random initialization. [11:50 pm 19.12.18]

\*\* Neural Network will not work if we initialize weights with  $\underline{\underline{0}}$



$n^{[0]} = 2$        $n^{[1]} = 2$   
 → If we initialize weight with zeros, the neurons in hidden unit will compute symmetric because they calculate same function.

$$\text{So, } a_1^{[1]} = a_2^{[1]} \dots = a_n^{[1]}$$

And the derivatives will also be same.  $d^2_1 = d^2_2 \dots d^2_n =$

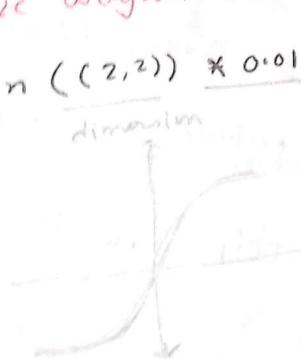
\*\* The solution is to initialize weights with random values.

$$w^{[1]} = \text{np.random.randn}((2, 2)) * 0.01$$

$$b^{[1]} = \text{np.zeros}((2, 1))$$

$$w^{[2]} = \dots$$

$$b^{[2]} = \dots$$



→ we initialize weights with small values because if we use large values learning rate of gradient descent will be very slow.

→ Programming Assignment: [will be able to].

1. Develop an intuition of back-propagation.

2. Recognize that more hidden layers can compute more complex structures.

3. Built all helper functions to implement a full model with one hidden layer.

Programming Assignment [week 3] [11.30 AM 26-12-18]

→ General methodology to build a neural network.

1. Define Neural Network structure (# no of inputs, # of hidden units, etc)
2. Initialize model's parameters.
3. Loop :-
  - Implement forward propagation.
  - Compute loss.
  - Implement backward propagation to get the gradients.
  - update parameters. (gradient descent).

\*\* don't helper function

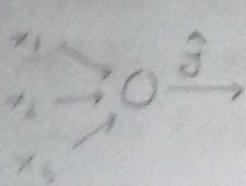
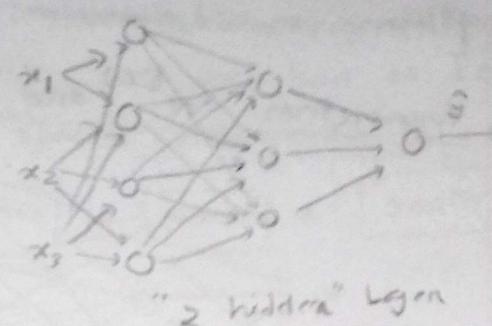
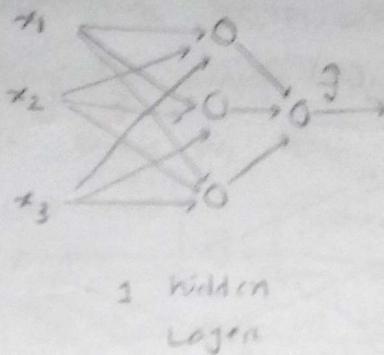
don't helper function build

- \* Refer Hidden layer gradients Accuracy (good)
- \* (Future good)
- \* don't Helper Function (good)

(good)

L-1

## Deep L-layers neural network. [9.10 AM 26-12-18]

logistic  
Regression

→ for any problem, it is hard to predict in advance how deep a neural network should be. [number of layers  $\geq 3$  to implement  $\text{exp}(x)$ ].  
Layer & number of neuron change

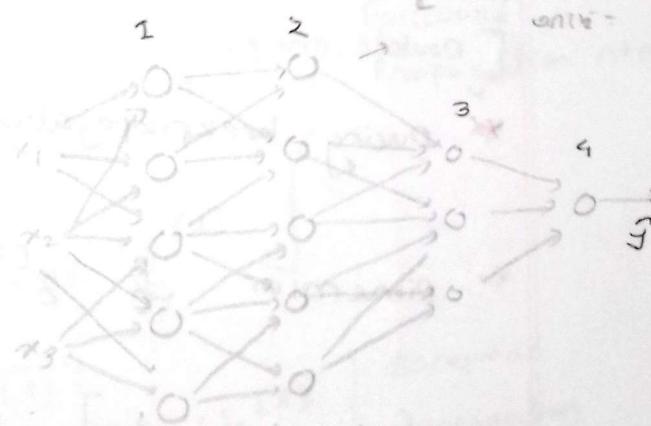
Notation:

 $L$  = number of hidden layers. $n^{[l]}$  = number of units in layer  $l$ . $a^{[l]}$  = activation in layer  $l$ . $w^{[l]}$  = weights in layer  $l$ . $b^{[l]}$  = bias in layer  $l$ . $x = a^{[0]}$  = input features. $\hat{o} = a^{[L]}$  = outputs.

## L-2 Forward Propagation in a Deep Network

[9.25 PM 26-12-18]

[softmax unit]



$$\begin{cases} z^{[i]} = w^{[i]} A^{[i-1]} + b^{[i]} \\ A^{[i]} = g^{[i]}(z^{[i]}) \end{cases}$$

$L=4$

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

for  $i=1, 2, \dots, L$ .

→ make explicit for loop for

go hidden layer by layer with

implement on to LSTM

L-3 Getting your matrix dimensions right [ 9:40 PM, 26-12-18 ]

[ note : matrix dimension after giving Debugging easy 22 ]

Parameters  $W^{[l]}$  and  $b^{[l]}$ . [  $\oplus L=5$  .

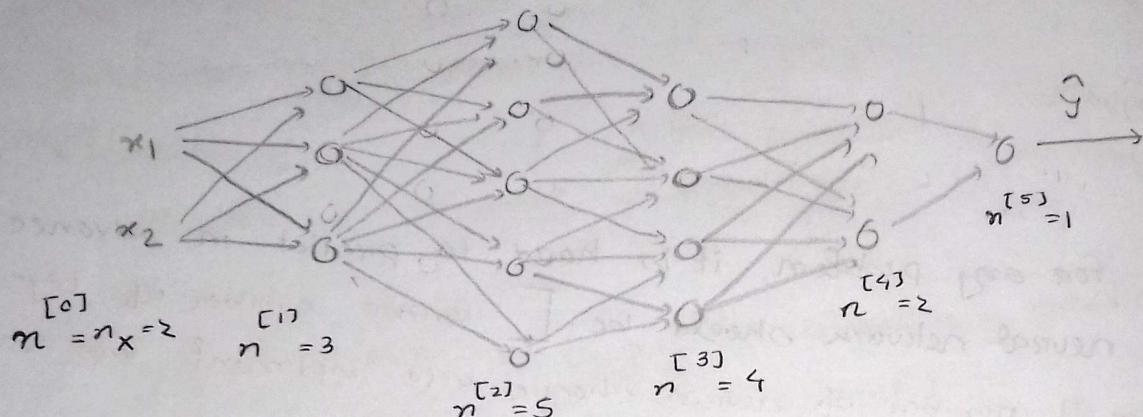


fig:- Neural Network.

\*\* Dimension of  $W^{[l]}$  is :  $(n^{[l]}, n^{[l-1]})$

$W^{[1]} = (3, 2)$ ,  $W^{[2]} = (5, 3)$ ,  $W^{[3]} = (4, 3)$ ,  $W^{[4]} = (2, 4)$ ,  $W^{[5]} = (1, 2)$ .

\*\* Dimension of  $b^{[l]}$  is :  $(n^{[l]}, 1)$

$b^{[1]} = (3, 1)$ ,  $b^{[2]} = (5, 1)$ ,  $b^{[3]} = (4, 1)$ ,  $b^{[4]} = (2, 1)$ ,  $b^{[5]} = (1, 1)$ .

[ Double check of here matrix dimension check orgy implement v610211 ]

\*\* During backpropagation dimension

$$W^{[l]} = dW^{[l]}, b^{[l]} = db^{[l]}$$

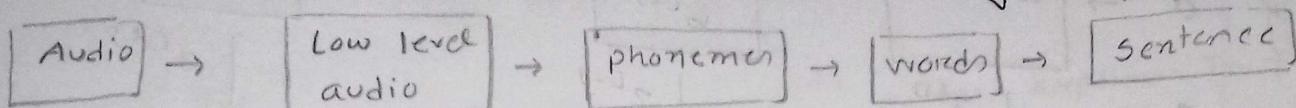
\*\* Dimension of  $Z^{[l]}$  for m training example is,  $(n^{[l]}, m)$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \vdots & x^{(m)} \\ | & | & | & | \end{bmatrix} \quad \left| \quad Y = \begin{bmatrix} y^{(1)} & \cdots & y^{(m)} \end{bmatrix} \quad \right| \quad Z = \begin{bmatrix} | & | & | \\ z^{(1)} & z^{(2)} & \vdots & z^{(m)} \\ | & | & | & | \end{bmatrix}$$

④ number of hidden unit 3 hidden layer  
orgy dimension change 22

## L-4 why deep representation? [10:00 pm 26-12-18]

Neural Network is a complex hierarchical representation where at first it computes simple function and use this to learn complex features. Such as for "speech recognition"



\* profit take off next another layer

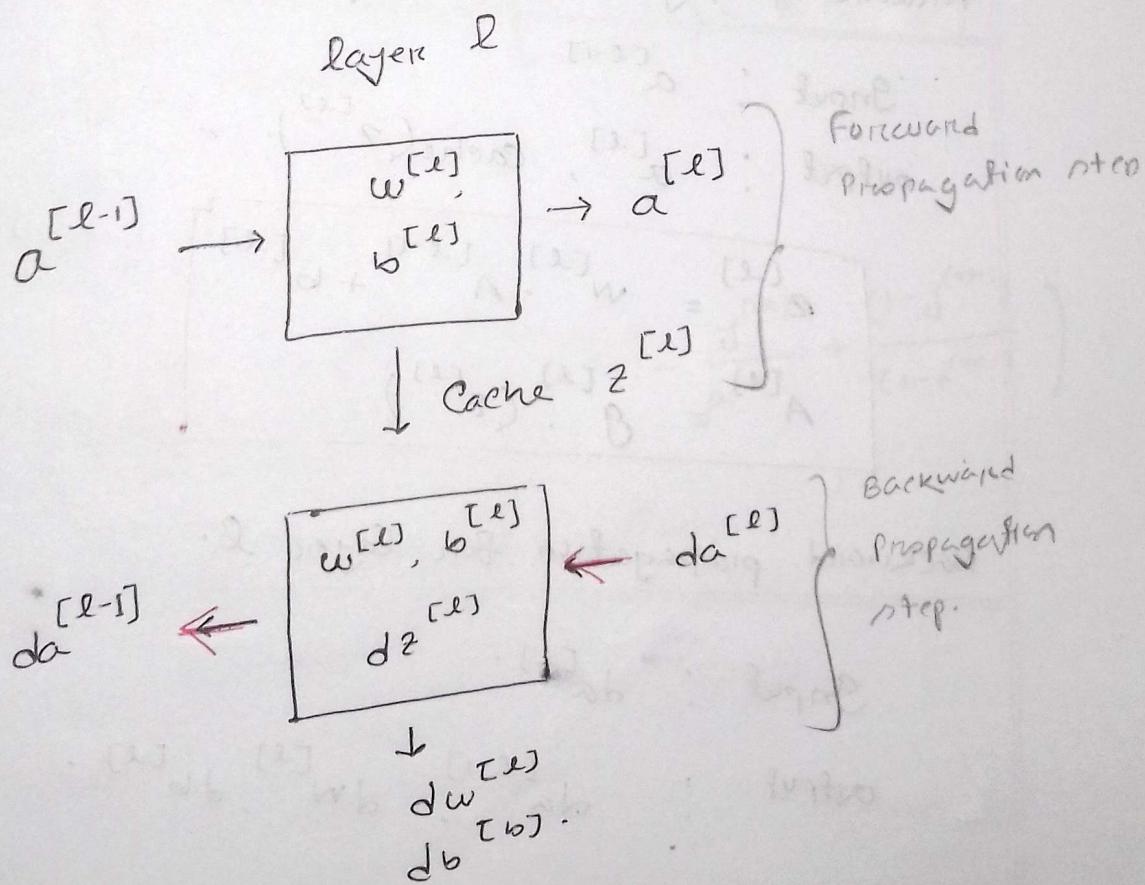
Informally: There are functions you can compute with a "small" L-layer deep neural network that shallow networks require exponentially more hidden units to compute.

[Lecture video to Example from [Google](#) ([video](#))

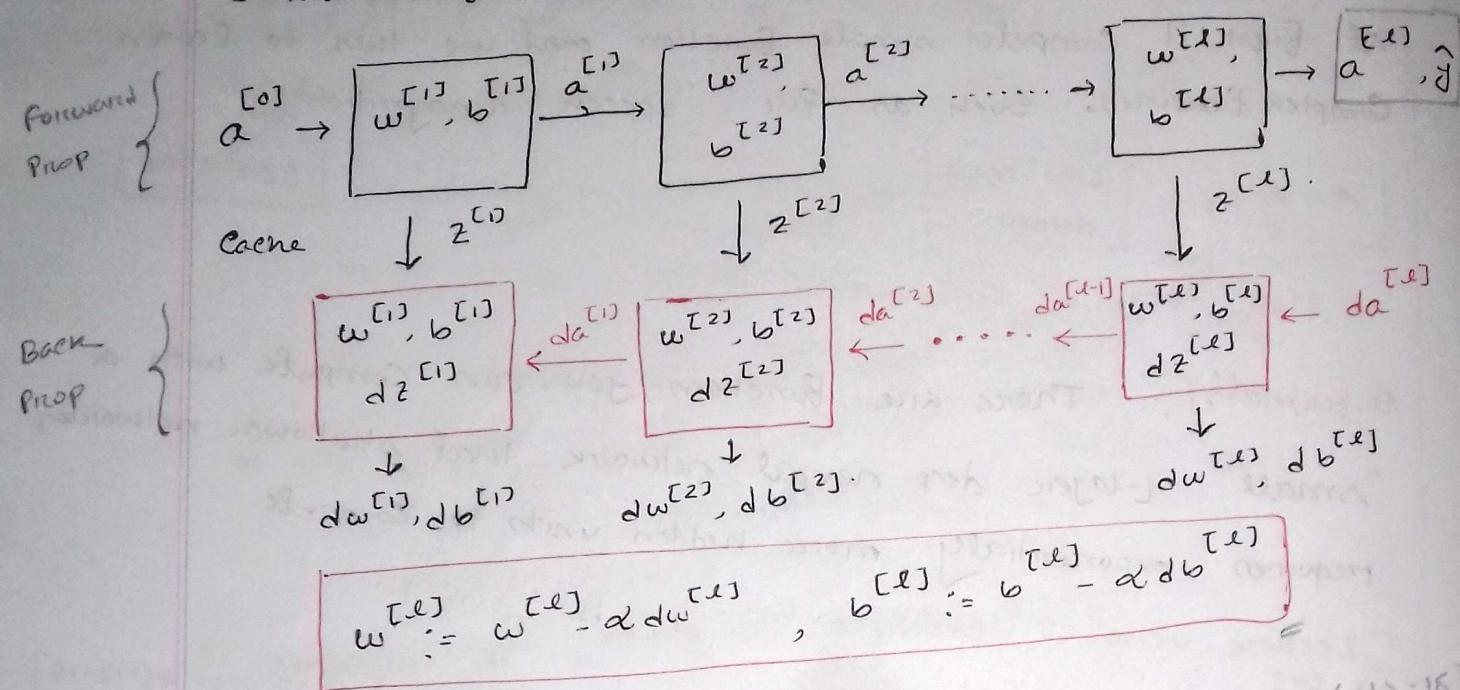
[11:30 pm 26-12-18]

## L-5 Building blocks of Deep neural networks.

Forward and Backward Functions:



Iteration of a neural network :-



## L-6 Forward and backward propagation [11.50 pm 26-12-18]

Forward propagation for layer  $\ell$ .

Input :  $a^{[\ell-1]}$

Output :  $a^{[\ell]}, \text{cache}(z^{[\ell]}), w^{[\ell]}, b^{[\ell]}$

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

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

[equation for implementation]

Backward propagation for layer  $\ell$ .

Input :  $da^{[\ell]}$

Output :  $da^{[\ell-1]}, d\omega^{[\ell]}, db^{[\ell]}$

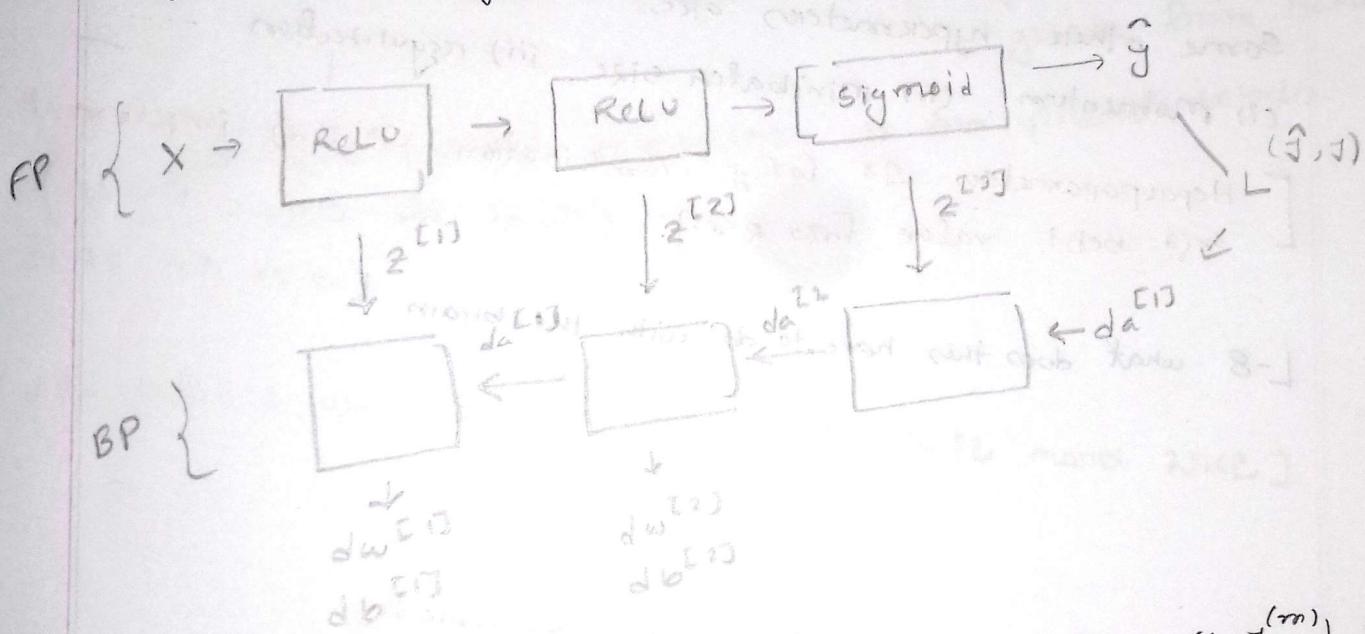
$$dz^{[l]} = dA^{[l]} * g'(z^{[l]})$$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1].T}$$

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

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

By summarizing we get.



$$da^{[l]} = \left( -\frac{j^{(l)}}{a^{(l)}} + \frac{(1-j^{(l)})}{(1-a^{(l)})} + \dots + \frac{j^{(m)}}{a^{(m)}} + \frac{(1-j^{(m)})}{(1-a^{(m)})} \right)$$

initialize  
[ implement softmax function for y ]

## L-7 Parameters vs hyperparameters. [12.10 AM 27-12-18]

Parameters :  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, w^{[3]}, b^{[3]} \dots$

Hyperparameters : → learning rate  $\alpha$ .

→ number of iterations.

→ hidden layer  $L$ .

→ number of hidden units

Choice of activation function.

GMP

[ $\alpha$  hyperparameter change  $\Rightarrow$  computation]

change paper length

27/12

[GMP, FOR PAPER WRITING].

\* hyperparameter controls parameters.

Some other hyperparameters are,

(i) momentum, (ii) minibatch size, (iii) regularization

[Hyperparameters  $\Rightarrow$  safest combination  $\Rightarrow$  best implement  
best value into  $\alpha$ ]. (Trial and error process)

[12.20 AM 27-12-18]

L-8 what does this have to do with the brain

[The brain has to learn more about your environment (ML)]

→ Programming Assignment. [will be able to].

1. Use non-linear units like ReLU to improve your model.

2. Build a deeper neural network.

3. Implement easy to use neural network class.

[Exercise is easier Theory that uses GMP and GMP 27/12].

V.V.Y.

## \*\* General Methodology for Deep Learning.

1. Initialize parameters / Define hyperparameters.
2. Loop for num\_iterations :
  - a) Forward propagation.
  - b) Compute cost function.
  - c) Backward propagation.
  - d) update parameters (using parameters and grads from Backprop)
3. Use trained parameters to predict labels.