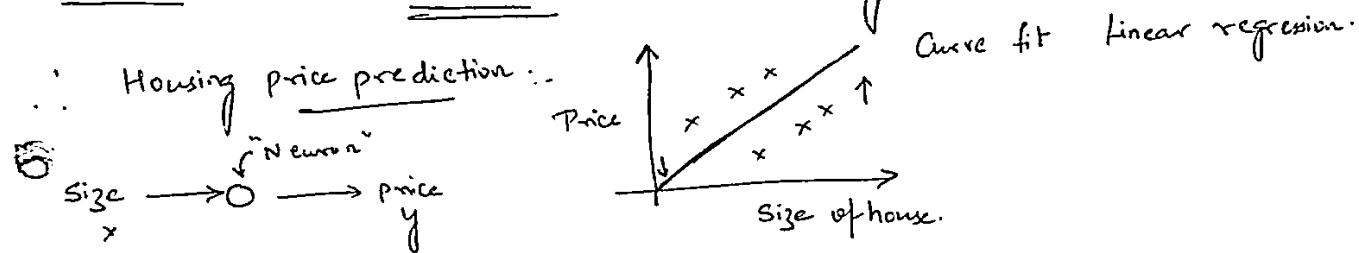
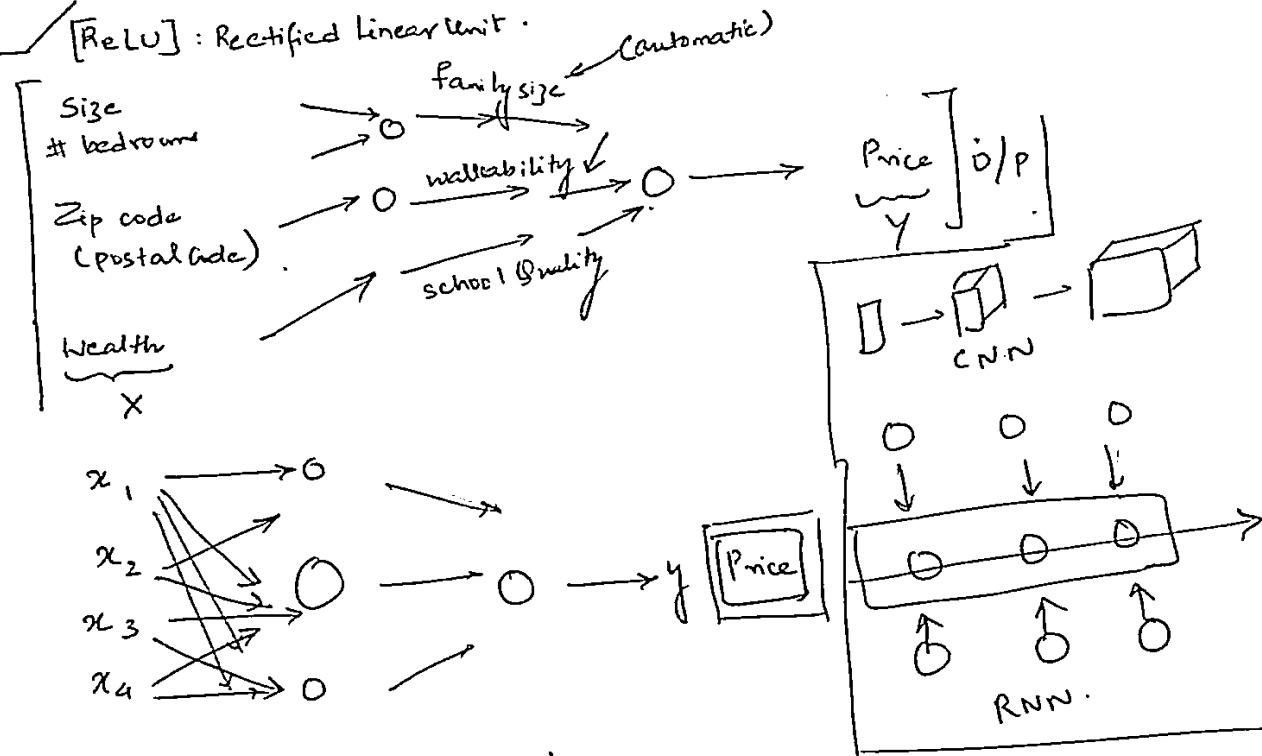


- \* Courses:-
1. Neural Networks and Deep Learning
  2. Im
  - 3.
  - 4.
  - 5.

\* Course 1: Neural Networks & Deep learning.



[ReLU]: Rectified Linear Unit.



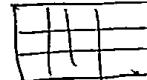
II Supervised learning :  $I/p(x)$   $O/p(y)$  Application

Homefeat..	Price	Real Estate
Ad, Userinfo	Click	Online Adver
Image	Obj (...)	Photo tagging
Audio	Text transcript	Speech Recognition
—	—	Audi RNN (Recurrent)
Image/A-	—	PTC-translations
—	—	Autonomous Driving
—	—	Custom/Hybrid

## Supervised learning :

- 2 -

Structured data ::



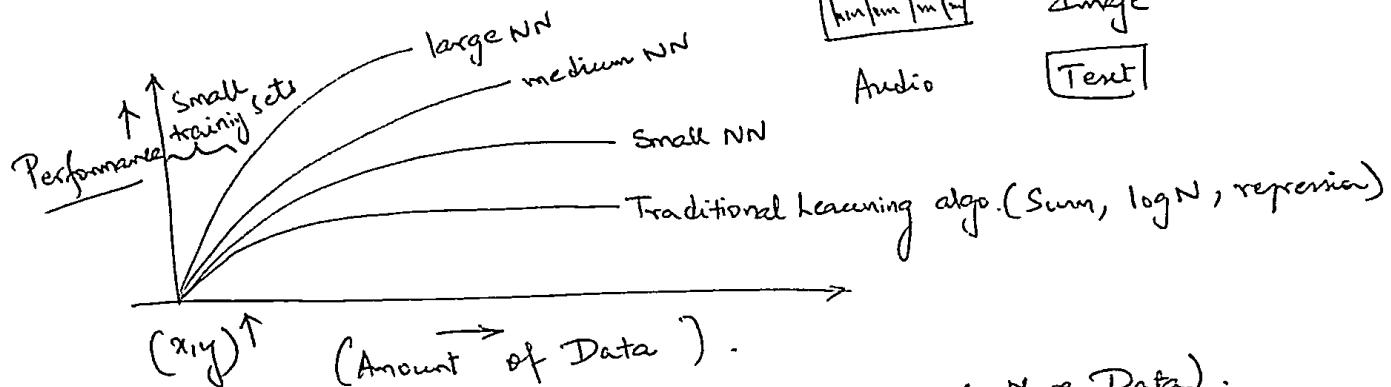
## Unstructured Data

format [img]

Image

Audio

Text



Digitization of Data (Collecting More & More Data).

Performance increases with large NN.  $\rightarrow$  medium  $\rightarrow$  small NN.

Data sigmoid  
Using it in MC learning, learning becomes very slow.

Computation  $\rightarrow$  medium  
ReLU  $\rightarrow$  small NN.  
func  
gradient is imp.  
idea  
Code  
Experiment

## INTRODUCTION TO NEURAL NETWORKS.

### 1. Neural Networks & Deep Learning

Week 1 : Intro.

Week 2 : Basics of NN programming

Week 3 : One hidden layer NN

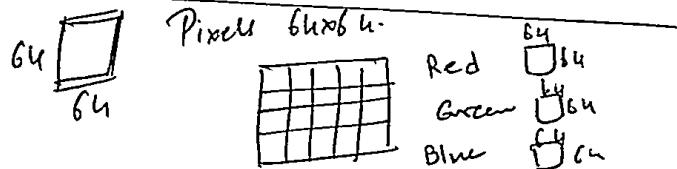
Week 4 : Deep Neural Networks.

### Quiz :

- AI is the new electricity :- Similar to electricity starting about 100 years ago, AI is transforming multiple industries.
- Deep learning requires multiple iterations ; train models, experience is 2<sup>nd</sup> order
- Increasing the size of NN, training set does not hurt algorithm performance

### \* Binary Classification :-

FPP / BPP : Forward/ Backward propagation



$$x = [:] = 64 \times 64 \times 3 = 12288 \quad \gamma = \gamma_x = 12288$$

$(x, y)$   $x \in \mathbb{R}^n$ ,  $y \in \{0, 1\}$  m training examples :  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$x = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(n)}] \quad [ \uparrow \downarrow \dots \uparrow \downarrow ] \quad [= x^{(1)} \dots x^{(n)}]$$

$x \in \mathbb{R}^{n \times m}$

$$x.\text{shape} = (n, m) ; y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}].$$

## Logistic Regression

$y \in \mathbb{R}^{l \times m}$

$$y.\text{shape} = (l, m).$$

Given  $x$ , want  $\hat{y}$

$$\hat{y} = P(y=1|x)$$

$x \in \mathbb{R}^{n \times 1}$

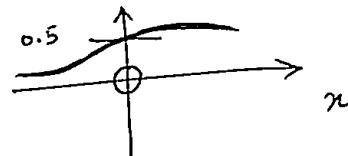
( $x$  dim vector).

Par.  $w \in \mathbb{R}^{n \times 1}$   
 $b \in \mathbb{R}$ .

Output

$$\hat{y} = w^T x + b$$

$$0 \leq \hat{y} \leq 1$$



Sigmoid func.:  $\hat{y} = \sigma(w^T x + b)$

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

if  $z$  large  $\sigma(z) \approx \frac{1}{1+0} = 1$ .

if  $z$  large negative num,  $\sigma(z) = \frac{1}{1+e^{-z}} \approx$

$$\approx \frac{1}{1+\text{Big num}} \approx 0.$$

alternative notation:-

$$x_0 = 1, x \in \mathbb{R}^{n \times 1}, \hat{y} = \sigma(\theta^T x).$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_R \end{bmatrix} \begin{matrix} \left. \right\} b \leftarrow \\ \left. \right\} w \leftarrow \end{matrix}$$

Summary Parameters :  $w \in \mathbb{R}^{n \times 1}$ ,  $b \in \mathbb{R}$ ,  $\hat{y} = w^T x + b$ ,  $\hat{y} = \sigma(w^T x + b)$

$\hat{y} = \sigma(\theta^T x); \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_R \end{bmatrix} \begin{matrix} \left. \right\} b \leftarrow \\ \left. \right\} w \leftarrow \end{matrix}$

## Logistic Regression Cost function

$$\hat{y} = \sigma(w^T x + b) \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}. \quad z^{(i)} = w^T x^{(i)} + b \quad \begin{matrix} x^{(i)} \\ y^{(i)} \\ z^{(i)} \end{matrix}$$

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$ .

\* loss (error) func :  $L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

(efficiency of Alg.)

how good  $\hat{y}$  is ?? given by

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

\*  $L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$  if  $y=0: L(\hat{y}, y) = -\log(1-\hat{y})$

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

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

Summary:

w, b minimize cost func;  
loss func for efficiency of algorithm

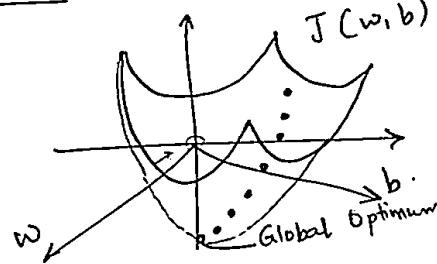
Logistic  
Regression

Note: The loss func computes the error for a single training example, the cost func is the average of the loss funcs of the entire training set.

∴  $J(w, b) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)}))$  O

Gradient Descent: - (train / learn w, b on training set).

until now:



Cost func:  $J(w, b)$

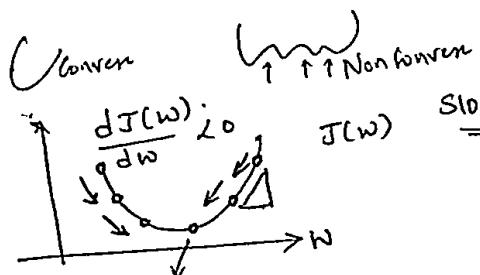
Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$

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

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

Want to find w, b that minimize  $J(w, b)$  ??? \*\*\*

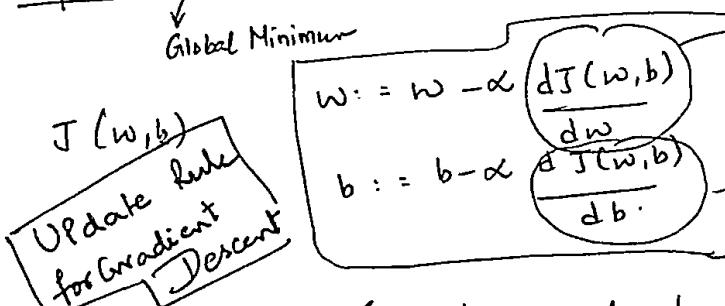
O



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

learning rate  
update "dw" on w.

$w := w - \alpha \partial w$  Gradient Descent Update.



$$\frac{\partial J(w, b)}{\partial w} \quad \frac{\partial J(w, b)}{\partial b}$$

$\partial$  = partial derivative

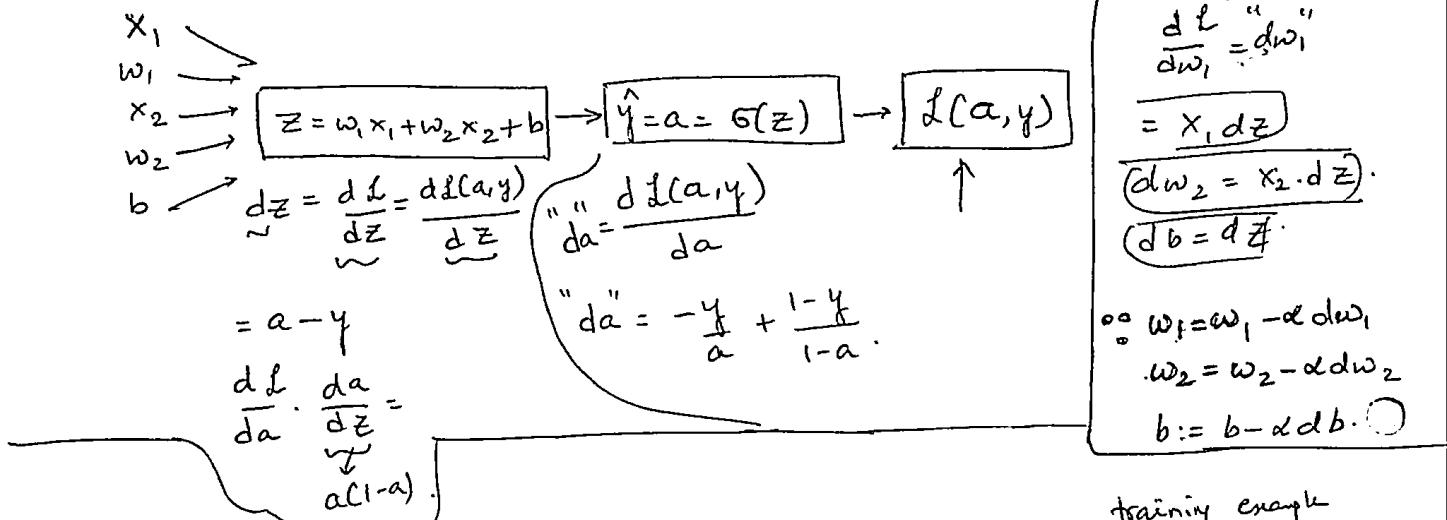
∴ A convex func has single local optima



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

$$\leftarrow \hat{y} = a = \sigma(Z)$$

$$\text{Loss func: } L(a, y) = -(y \log(a) + (1-y) \log(1-a)).$$



\* Gradient Descent on m Examples:

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

$$a^{(i)} = \hat{y}^{(i)} = \sigma(Z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})}_{dw_1^{(i)} - (x^{(i)}, y^{(i)})}$$

training example  
 $(x^{(i)}, y^{(i)})$ .  
 $dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$ .

lets say  
 $J=0, dw_1=0, dw_2=0,$   
 $db=0$ .  
for  $i=1 \text{ to } m$

$$Z^{(i)} = w^T x^{(i)} + b ; a^{(i)} = \sigma(Z^{(i)}) ; J += -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dZ^{(i)} = a^{(i)} - y^{(i)} ; dw_1 += x_1^{(i)} dZ^{(i)} \uparrow n=2 ; J /= m ; dw_1 /= m ; dw_2 /= m$$

$$dw_2 += x_2^{(i)} dZ^{(i)} \uparrow i=1 \text{ to } m$$

$$db += dZ^{(i)} \uparrow db /= m$$

$$\boxed{\begin{aligned} dw_1 &= \frac{\partial J}{\partial w_1} \\ w_1 &:= w_1 - \alpha dw_1 \\ w_2 &:= w_2 - \alpha dw_2 \\ b &:= b - \alpha db \end{aligned}} ; \quad \text{2 for loops} \rightarrow dw_1, dw_2 \dots \text{features.}$$

= Vectorization . . . to avoid for loops.

## \* Derivation of $\frac{dL}{dZ}$

8, 6, 10, 7, 9, 11, 6, 3, 7,

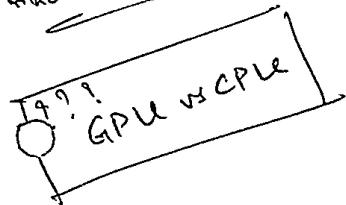
Time 12:00

$$\frac{dL}{dZ} = a - y$$

## \* Vectorization:-

$$z = w^T \times b.$$

if vectorize your code  
Deep learning runs  
much faster



GPU

CPU

} SIMD - single instruction multiple data.

## AVOID Explicit "FOR" Loops.

$$u = Av,$$

$$u_i = \sum_j A_{ij} v_j$$

$$\text{Normal: } u = np.zeros(n, 1).$$

for i ...  
for j ...

$$u[i] += A[i][j] * v[j]$$

Vectorize

that

$$u = np.dot(A, v)$$

## \* Vectors & Matrix valued functions

Apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$\therefore u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = math.exp(v[i])$$

import numpy as np.

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

np.log(v); np.abs(v); np.maximum(v, 0)

$$v^{**2}, 1/v.$$

$$\therefore dw = np.zeros((n, 1)) \quad (\text{dimensional vector})$$

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

$$dw = dw/m$$

$$|dw| = m$$

## Deep learning. AI.

Normal

Vectorize

300 times faster than normal

```
a = np.random.rand(1000000)
b = "
tic = time.time()
for i in range(100000):
    c += a[i]*b[i]
toc = time.time()

for c = time.time()
```

481.31102039 ms | 1.501 ms.

GPU (Graphical P.U. has many cores & used for parallel computing (Simult. many threads))

CPU (few cores & less parallelization)

Note : if you use In-Built func it enables python to efficient parallelism]

# for training examples Vectorizing Logistic Regression

$$\rightarrow z^{(1)} = w^T x^{(1)} + b \quad z^{(2)} = w^T x^{(2)} + b \quad z^{(3)} = w^T x^{(3)} + b \\ a^{(1)} = \sigma(z^{(1)}) \quad a^{(2)} = \sigma(z^{(2)}) \quad a^{(3)} = \sigma(z^{(3)}) \\ \therefore x = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_m^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_m^{(m)} \end{bmatrix} \in \mathbb{R}^{n_x \times m}$$

$(z = np.dot(w.T, x) + b)$  ("Broadcasting" in Python)

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

VECTORIZING LOGISTIC REGRESSIONS GRADIENT DESCENT

$$\frac{\partial z^{(1)}}{\partial w} = a^{(1)} - y^{(1)} \quad \frac{\partial z^{(2)}}{\partial w} = a^{(2)} - y^{(2)} \dots$$

$$\frac{\partial z}{\partial w} = [d z^{(1)} \ d z^{(2)} \ \dots \ d z^{(m)}].$$

$$\frac{\partial z}{\partial b} = A - Y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots].$$

$$\begin{aligned} \frac{\partial b}{\partial b} &= \frac{1}{m} \sum_{i=1}^m d z^{(i)} \\ &= \frac{1}{m} \text{np.sum}(d z). \end{aligned}$$

$$\frac{\partial w}{\partial w} = \frac{1}{m} X \cdot d z^T.$$

$$= \frac{1}{m} \left[ \begin{array}{c|c} x_1^{(1)} & \dots & x_1^{(m)} \\ \hline x_2^{(1)} & \dots & x_2^{(m)} \\ \vdots & \ddots & \vdots \\ x_n^{(1)} & \dots & x_n^{(m)} \end{array} \right] \left[ \begin{array}{c|c} d z^{(1)} \\ \vdots \\ d z^{(m)} \end{array} \right]$$

$$= \frac{1}{m} \left[ \begin{array}{c|c} x_1^{(1)} d z^{(1)} & \dots & x_1^{(m)} d z^{(m)} \\ \vdots & \ddots & \vdots \\ x_n^{(1)} d z^{(1)} & \dots & x_n^{(m)} d z^{(m)} \end{array} \right].$$

$$= \begin{bmatrix} \cdots & \cdots & \cdots & \cdots \\ | & | & | & | \\ \cdots & \cdots & \cdots & \cdots \end{bmatrix}$$

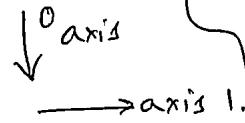
Broadcasting in Python (faster).

$$A = np.array([[[\dots, \dots, \dots], [\dots, \dots, \dots], [\dots, \dots, \dots]]])$$

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

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

Matlab/Octave: bsxfun



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

Ex ①

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}.$$

GENERAL PRINCIPLE of BoC

$$(m, n) \xrightarrow{*} (1, n) \rightsquigarrow (m, n)$$

$$(m, 1) \xrightarrow{*} (R \rightsquigarrow (m, 1))$$

②  $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 1,000 & 2,000 & 3,000 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 204 & 205 & 206 \end{bmatrix}. (m, 1)$

buit Python makes it  $(m, n)$

$$(1, m) \xrightarrow{*} (R \rightsquigarrow (1, m))$$

A Note on python/numpy vectors.

never use data structures

$a = np.random.rand(5)$   $\times$  rank 1 array

Use in NN  $a = np.random.rand(5, 1)$  -  $\begin{bmatrix} [ ] \\ [ ] \\ [ ] \end{bmatrix}$

$\xrightarrow{\text{print}(a.T)} \text{a transpose.}$

$\xrightarrow{\text{np.dot}(a, a.T)}$  value outerproduct of a vector

Note Never use rank 1 array.  
Always use vector

$a = np.random.rand(5, 1) \rightarrow a.shape = (5, 1)$  column vector

$a = np.random.rand(1, 5) \rightarrow a.shape = (1, 5)$  row vector

reshape an array  $a = a.reshape((5, 1))$

assert ( $a.shape = (5, 1)$ ) inexpensive assert

---

## Tour of Jupyter

### Explanation of logistic regression Cost func

\*  $y | x$  y given x. Note: Minimizing the loss corresponds with maximizing  $\log P(y|x)$ .

### Exercises: Week 2

Q1: A neuron computes a linear func ( $Z = Wx + b$ ) followed by an activation func.  
We generally say that the output of a Neuron is a  $g(Wx + b)$  where  $g$  is the activation func  $\hat{y} = \sigma(W^T x + b)$ ;  $\sigma(z) = \frac{1}{1+e^{-z}}$ .  
~~(Sigmoid, tanh, ReLU, ...).~~

Q2. Logistic Loss :-  $L^{(i)}(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)}))$ .

Q3:  $\text{img} = (32, 32, 3)$  array  $\approx x = \text{img.reshape}((32*32*3, 1))$ , column vector

Q4  $a = np.random.rand(2, 3); c = a+b; c.shape = (2, 3)$  Broadcasting of "b". to match "a".

Q5  $a = np.random.rand(4, 3); b = " " (3, 2); c = a * b;$  "Error" The "\*" operator in Numpy, indicates element-wise multiplication.

If it is different from "np.dot()". If you would try "c = np.dot(a, b)" you would get  $c.shape = (4, 2)$

Q6. Suppose you have  $n_x$  input features per example.  $X = [x^{(1)} \ x^{(2)} \dots x^{(m)}]$ .

What is the dimension of  $X$ ?  $(n_x, m)$

???

Q7  $a = np.random.rand(12288, 150)$   $c = np.dot(a, b)$ ;  $c.shape = (12288, 45)$ .  
 $b = \dots$   $(150, 45)$ .

Q8.  $a.shape = (3, 4)$ . for  $i$  in range (3):

$b.shape = (4, 1)$  for  $j$  in range (4)

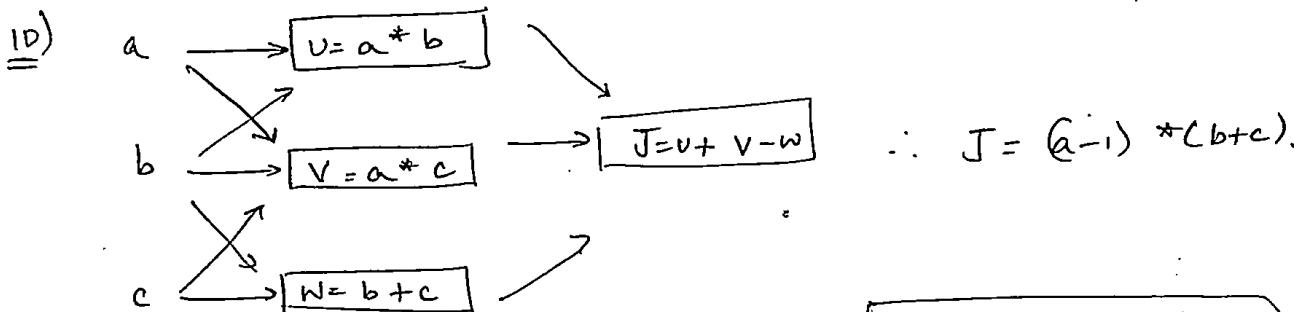
$$c[i][j] = a[i][j] + b[j].$$

Vectorize  $C = a + b^T$ .

Q9.  $a = np.random.rand(3, 3)$ .  $c = a * b$

$b = np.random.rand(3, 1)$ ,

This will invoke broadcasting, so  $b$  is copied three times to become  $(3, 3)$  and  $*$  is an element-wise product so  $c.shape$  will be  $(3, 3)$ .



\* Programming exercise. Q. Build a func returns

$$\text{Sigmoid}(x) = \frac{1}{1+e^{-x}}$$

Logistic func

`math.exp(x)`

Note: We rarely use "math" library in deep learning because the inputs of the func are real numbers. In deep learning we mostly use matrices and vectors.

This is why "numpy" is more useful.

Furthermore, if  $x$  is a vector, then Py operation such as  $S = x + 3$  or  $S = \frac{1}{x}$  will output  $S$  as a vector of the same size as  $x$ .

Note: ML/DL Normalize our data because gradient Descent converges faster after Normalization.

\* gradient of Sigmoid func = Slope of Sigmoid func.  $\frac{x}{\|x\|} \equiv np.linalg.norm(x, axis=1, keepdims=True)[S]$

Broadcasting Changes matrix dimensions.  $\begin{cases} np.dot() \rightarrow \text{matrix-matrix multiplication} \\ np.multiply() \rightarrow \text{element wise multiply} \end{cases}$

# Programming Exercise Part 2

— 10 —

Centre & Standardize your Data Set in Machine Learning

## Logistic Regression to identify cats.

255 (maximum value of a pixel channel).

\* m-train images labelled as Cat ( $y=1$ ) & non-cat ( $y=0$ )

\*  $(\text{num\_px}), (\text{num\_px}, 3)$  3 → 3 channels (RGB).

\* Logistic regression is actually a very simple neural network

$$w^T x^{(i)} + b \mid \sigma$$

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

\* The main steps for building a Neural Networks are:

1. Define the model structure (such as number of input features).

2. Initialize the model's parameters.

3. loop: - Calculate current loss (forward propagation).

- calculate current gradient (backward propagation).

- update parameters (gradient descent).

$$W = np.zeros((dim, 1))$$

You often build 1-3 separately and integrate them into one func we call `model()`.

Propagate():

$W, b, X, Y$ .

$$A = \sigma(W^T X + b) = (a^{(1)}, a^{(2)}, \dots, a^{(m)})$$

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

$$\frac{\partial J}{\partial W} = \frac{1}{m} X(A-Y)^T$$

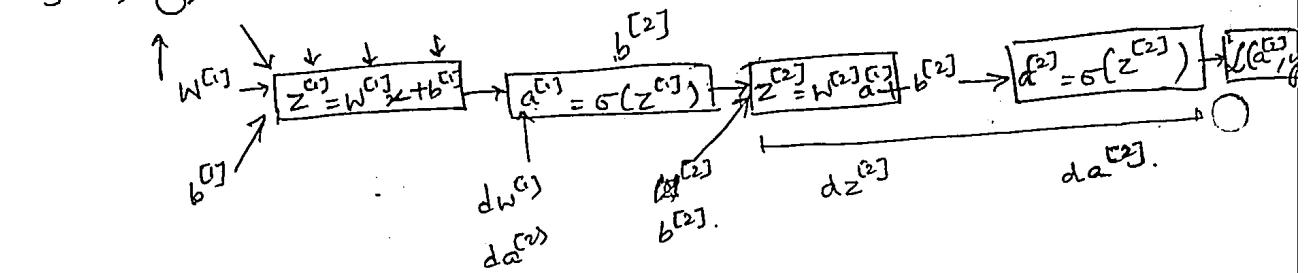
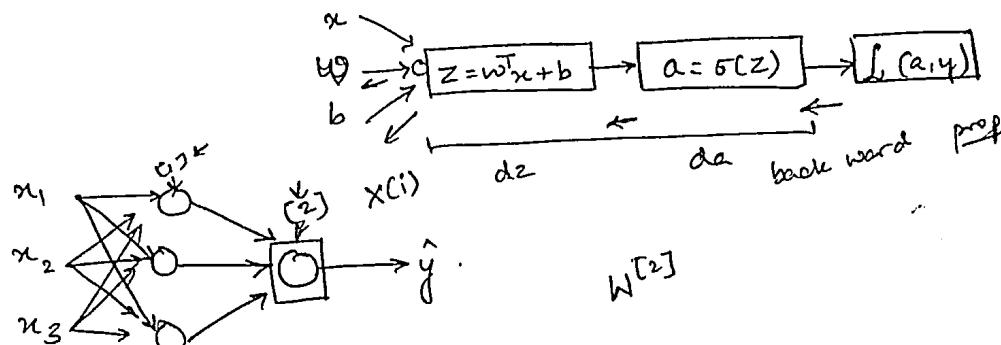
$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

\*  $\Theta = \Theta - \alpha \Delta \Theta$

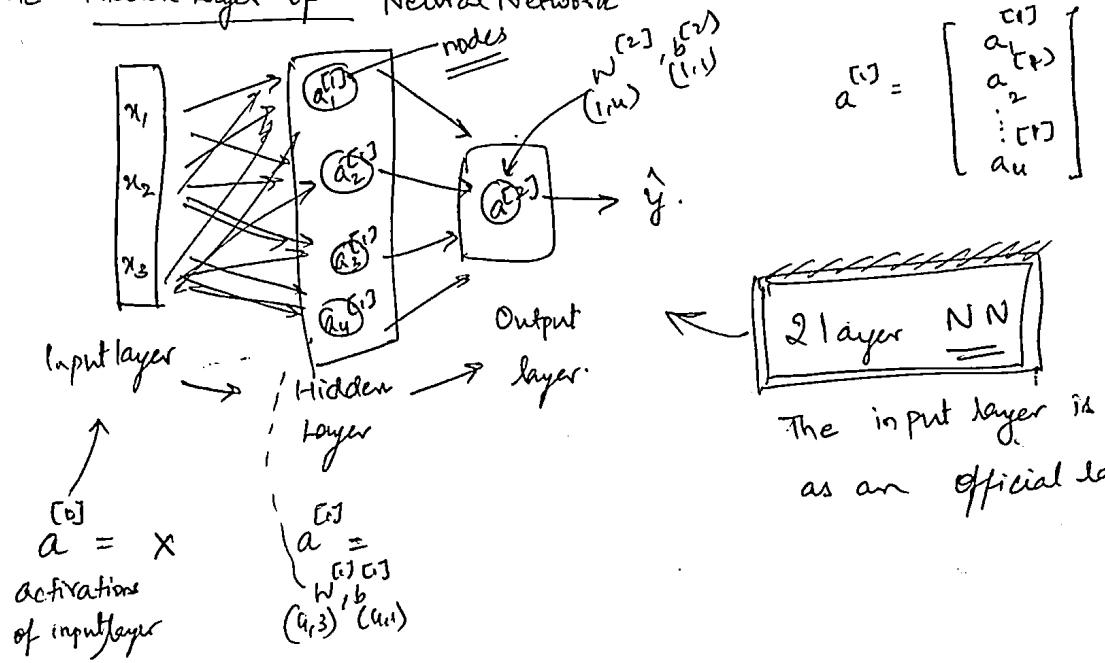
$\alpha$  = learning rate

Glimpse

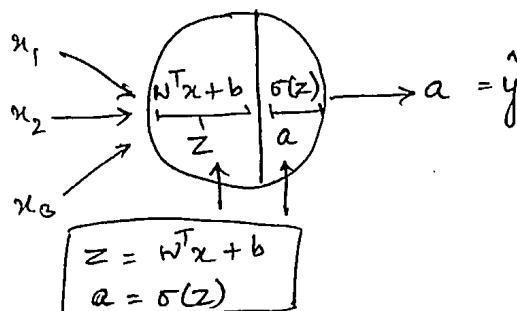
$$x_1 \rightarrow O \rightarrow \hat{y} = \varphi.$$



### One Hidden Layer of Neural Network

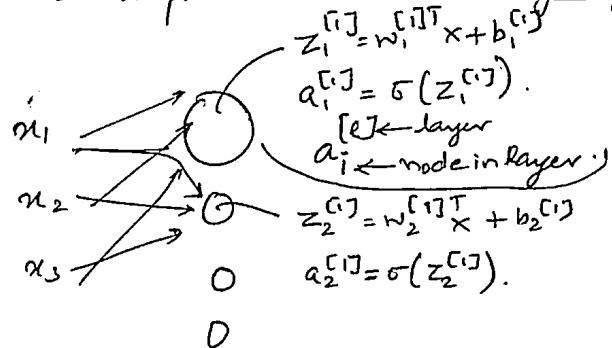


### Computing a NN output : - 2 steps



1. Compute Z

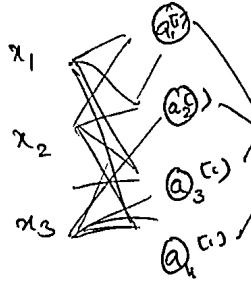
2. Compute Activation(a) from sigmoid of Z



$$x_1 \rightarrow z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$w_2 \rightarrow a_2^{[1]} = \sigma(z_2^{[1]}).$$

### \* Neural Network Representation



$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]}). \end{aligned}$$

O

$$z^{[1]} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}}_{(4,1)} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{(3,1)} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{(4,1)} = \underbrace{\begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}}_{(4,1)} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

O

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

$$\rightarrow z^{[1]} = w^{[1]T} x + b^{[1]}.$$

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

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

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

$$\boxed{z = w^T x + b.}$$

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

7, 10, 11, 12, 13, 14

\* Vectorize across multiple examples :-

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

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

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

$\vdots$   $\vdots$

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

$$X = \begin{bmatrix} 1 & x^{[1]} & x^{[2]} & \dots & x^{[m]} \end{bmatrix}$$

↑      ↑      |      |      |  
(n\_x, m)

training examples  
hidden units.

$$z^{[i]} = \begin{bmatrix} z^{[1](i)} \\ z^{[2](i)} \\ \vdots \\ z^{[m](i)} \end{bmatrix}$$

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

$$\left\{ \begin{array}{l} z^{[i]} = w^{[i]} x + b^{[i]} \\ a^{[i]} = \sigma(z^{[i]}) \\ z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}). \end{array} \right.$$

for  $i = 1 \text{ to } m,$

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

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

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

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

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

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

\* Explanation for Vectorized implementation

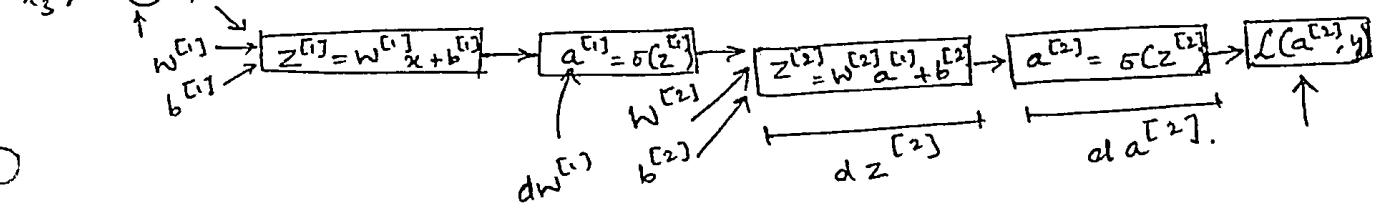
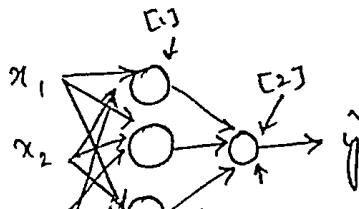
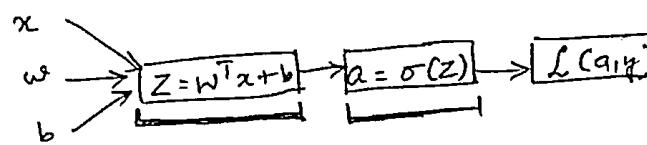
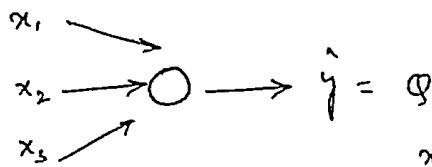
∴ first training example

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

$$z^{[1](2)} = w^{[1]} X^{(2)} + b^{[1]},$$

$$w^{[1]} = \begin{bmatrix} \dots \\ \dots \end{bmatrix}.$$

○



### \* Activation (Network) functions:-

Better  
f.actn → Sigmoid f.nc

$$z^{[1]} = w^{[1]}x + b^{[1]} \quad \text{until now Activation func}$$

$$a^{[1]} = \sigma(z^{[1]}) \rightarrow g(z^{[1]}) \quad \text{is not sigmoid func}$$

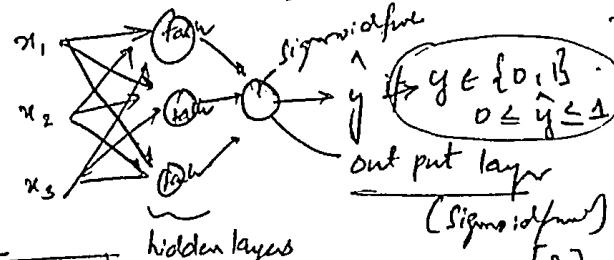
$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}.$$

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

Nonlinear func (Not sigmoid).

Now on Activation func is tanh func

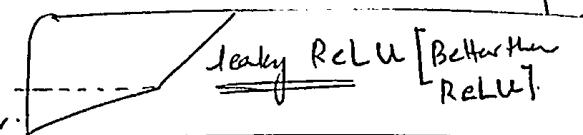
$$g(z^{[1]}) = \tanh(z^{[1]}).$$



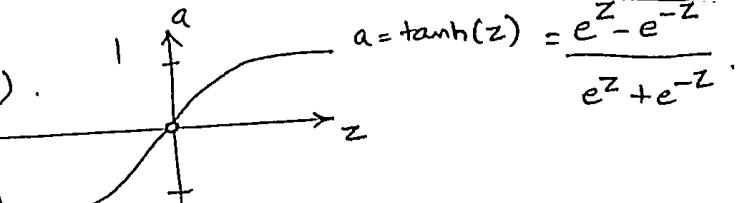
Note

hidden layers have data closer to zero mean so, tanh better for Output layer (Sigmoid) better.

[1].



Leaky ReLU [Better than RELU]



a = max(0, z).

ReLU

Rectified  
Linear Unit

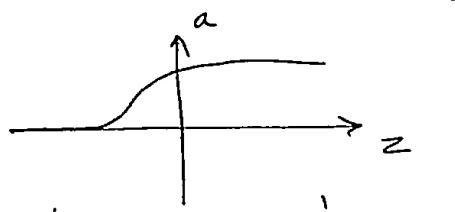
(Default  
choice of Activation

func for hidden  
layers if not  
sure then ...  
tanh Activ func  
for y)

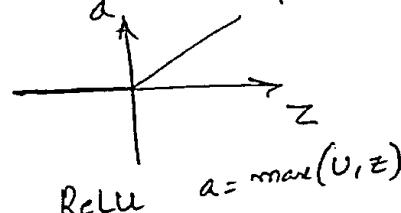
- NN is much faster when using ReLU Activation func than tanh or sigmoid Activation func

## \* Pros & Cons of Activation func.

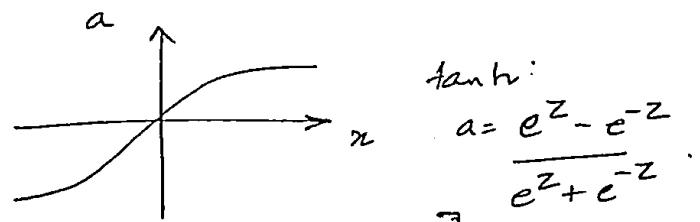
- 2 - 23, 24, 25, 26, 27, 28, 29, 30, 1, 2, 3, 4, 5, 6, 7.  
15 days.



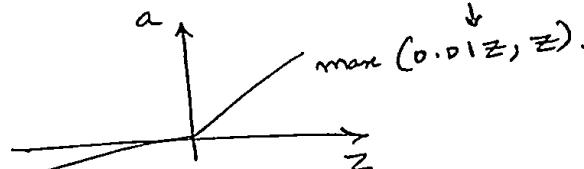
Sigmoid:  $a = \frac{1}{1 + e^{-z}}$   
 [Never for hidden layers & always for Output layers].



ReLU  $a = \max(0, z)$   
 (Default) for hidden layers.



[best for hidden layers]



leaky ReLU  $a = \max(0.01z, z)$   
 [sometimes better than ReLU].

## \* Why do we need a N.L Activation func. - :

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

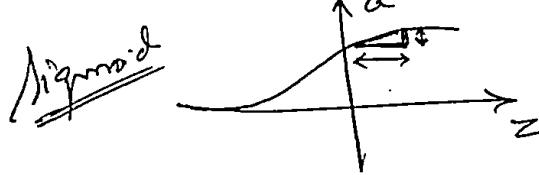
$$a^{[i]} = g^{[i]}(z^{[i]})$$

" $g(z) = z$   
 linear orientation function"

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

∴ A linear Hidden layer is more or less useless. N.L is better

## \* Derivatives of Activation Functions



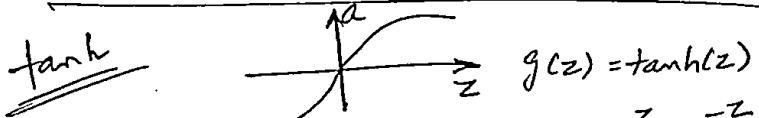
$$g(z) = \frac{1}{1 + e^{-z}}$$

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

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

$$= g(z)(1 - g(z)).$$

$$= a(1 - a).$$



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

$$g'(z) = \frac{d}{dz}(g(z))$$

$$= 1 - (\tanh(z))^2$$

$$z = 10 \quad \tanh(z) \approx 1$$

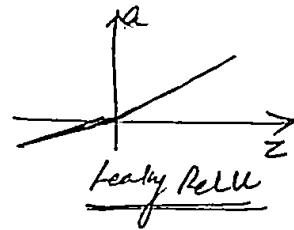
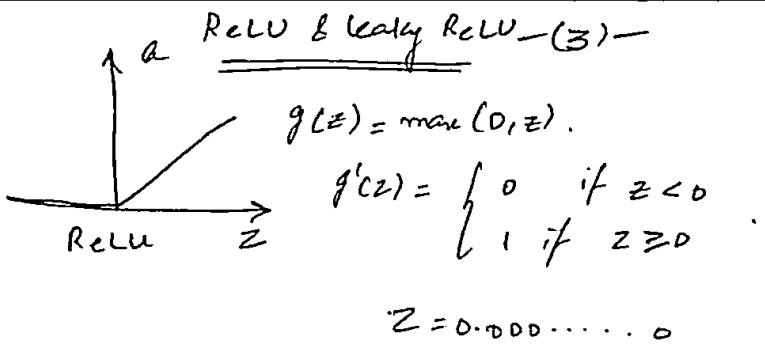
$$g'(z) \approx 0.$$

$$z = -10 \quad \tanh(z) \approx -1$$

$$g'(z) \approx 0.$$

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

$$g'(z) = 1.$$



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

## Gradient Descent for Neural Networks:-

Parameters :-  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$ .  
 $(n^{[1]}, n^{[0]}) (n^{[1]}, 1) (n^{[2]}, n^{[1]}) \underset{n}{\underbrace{(n^{[2]}, 1)}}$ .

Cost func :  $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}_i, y_i)$ .

### ○ Gradient Descent :

Repeat { Compute predicts  $(\hat{y}^{(i)}, [-1, \dots, m])$ .

$$dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}, db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$$

Update rule  $\boxed{W^{[1]} := W^{[1]} - \alpha dW^{[1]}, b^{[1]} := b^{[1]} - \alpha db^{[1]}}$  Update Rule for Gradient Descent

### ○ Formulas for computing derivatives:

#### Forward Propagation.

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

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

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

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

Correct vectorised Implementation  
of forward propagation for layer  $i$ ;  $1 \leq i \leq L$ :

$$Z^{[i]} = W^{[i]}A^{[i-1]} + b^{[i]}.$$

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

Backward Propagation

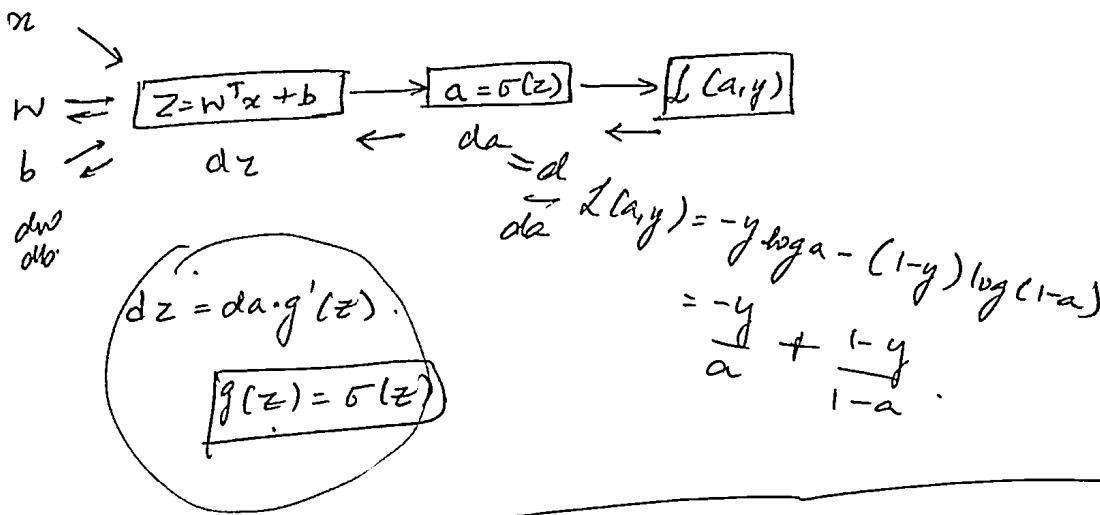
$$\begin{aligned} dz^{[2]} &= A^{[2]} - y \\ dW^{[2]} &= \frac{1}{m} dz^{[2]} A^{[1]T}, \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis}=1, \text{keepdims=True}). \end{aligned}$$

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

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

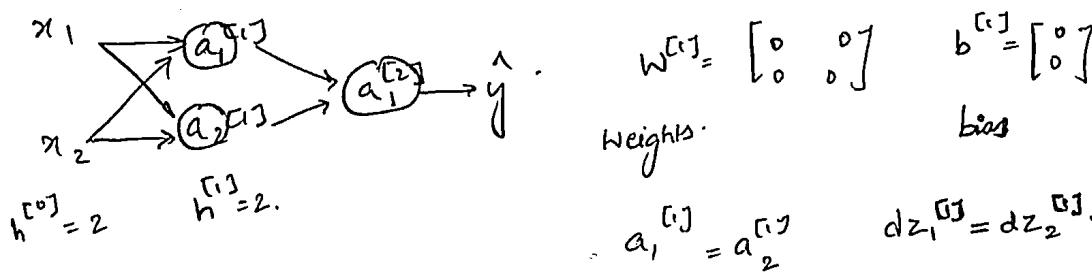
$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$

## Logistic regression



## Random Initialization :-

\* What happens if you initialize weights to zero?



Symmetry:-

$$dw = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad w^{[1]} = w^{[2]} - \alpha dw. \quad w^{[1]} = \begin{bmatrix} \dots \end{bmatrix}.$$

all the values of  $w$ :

∴ if you initialize  $w$  with zero,  $w$  after every iteration, will have the first row equal to the second row. ∵ They all are computing the same thing

Initialize Randomly :-

$$\therefore w^{[1]} = np.random.rand((2, 2)) * (0, 0.1)$$

if it's too large then ?!

Initialize to very small random values.

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

∴  $w$  is random, so no longer there is a problem of symmetry breaking.

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

$$b^{[2]} = 0.$$

- ① Initialize parameters
- ② Predict using forward prop.
- ③ Gradient Descent using Backward prop

if it is very large,  $w$  will be large, correspondingly  $Z^{[1]} = w^{[1]} x + b^{[1]}$

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

$Z$  will be very large or small → the slope of the gradient will be very small so Gradient Descent will be very small & slow

LabPlanar Data - Classification with one hidden layer v6C.

np.random.seed(1).

2-class dataset

Mathematically,

For one example,

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

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

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

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

$$y_{\text{prediction}} = \begin{cases} 1 & \text{if } a^{[2](i)} > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$J = -\frac{1}{m} \sum_{i=0}^m \left( y^{(i)} \log(a^{[2](i)}) + (1-y^{(i)}) \log(1-a^{[2](i)}) \right)$$

General Net<sup>4</sup> to build NN

- ① Define NN structure  
(# of inp units, # of hidden units, etc..)

- ② Initialize the Model parameters

- ③ Loop:
  - Implement forward propagation
  - Compute loss
  - Implement backward prop. to get the gradients
  - Update parameters.  
(gradient descent)

- Define NN structure

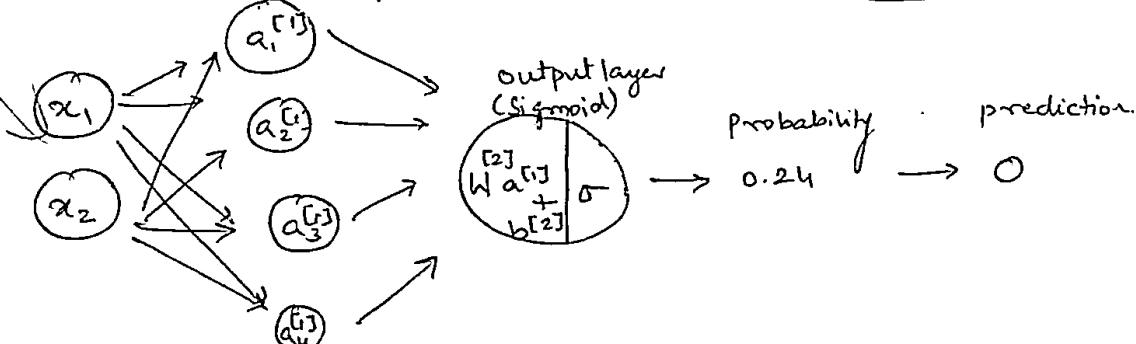
Define 3 variables:

↳  $n_x$ : the size of the input layer↳  $n_h$ : the size of the hidden layer (set this to 4)↳  $n_y$ : the size of the output layer

- Use shapes of X & Y to find  $n_x$  and  $n_y$

Also, hard code the hidden layer size to be 4.

hidden layer size 4



Sigmoid()  
np.tanh()

$$J = -\frac{1}{m} \sum_{i=1}^m \left( y^{(i)} \log(a^{(2)}[i]) + (1-y^{(i)}) \log(1-a^{(2)}[i]) \right)$$

Implement `compute_cost()` to compute the value of the cost:

$$-\sum_{i=0}^m y^{(i)} \log(a^{(2)}[i]):$$

$$\text{log_probs} = \text{np.multiply}(\text{np.log}(A2), Y)$$

$$\text{cost} = -\text{np.sum(log_probs)}.$$

Hardest Part in Deep Learning is Backward Propagation

17

Summary of Gradient Descent:-

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

$$dW^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

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

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

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

①

Functions Used

X, Y

\* `layer_sizes(X, Y)`

\* `initialize_parameters(n_x, n_h, n_y)`

\* `forward_propagation(X, parameters)`

\* `compute_cost(A2, Y, parameters)`

\* `backward_propagation(parameters, cache, X, Y)`

\* `update_parameters(parameters, grads, learning_rate = 1.2)`

# Summary

## Forward propagation

$$z^{[l]} = W^{[l]}x + b^{[l]}$$

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

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

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

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

○

Implementation info

## Cost function:

Compute Cost func to check if your model is actually learning:-

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

⇒ Backward Propagation:-

$$\text{○ } dW^{[l]} = \frac{\partial J}{\partial W^{[l]}} = \frac{1}{m} \sum_{i=1}^m dz^{[l]} A^{[l-1]T}$$

$$db^{[l]} = \frac{\partial J}{\partial b^{[l]}} = \frac{1}{m} \sum_{i=1}^m dz^{[l]}(c_i)$$

$$dA^{[l-1]} = \frac{\partial L}{\partial A^{[l-1]}} = W^{[l]T} dz^{[l]}$$

## Backward Propagation

$$dz^{[l]} = A^{[l]} - y$$

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

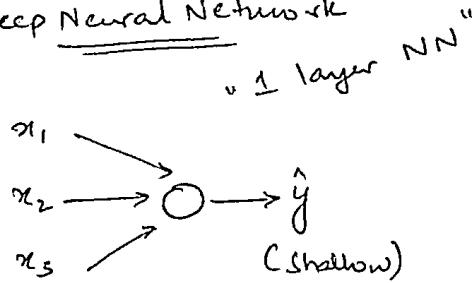
$$db^{[l]} = \frac{1}{m} \text{np. sum}(dz^{[l]}, \text{axis}=1, \text{keepdims=True})$$

$$dz^{[l-1]} = W^{[l]T} \cdot dz^{[l]} * g'(z^{[l]})$$

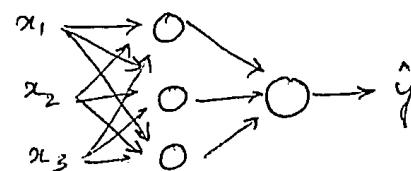
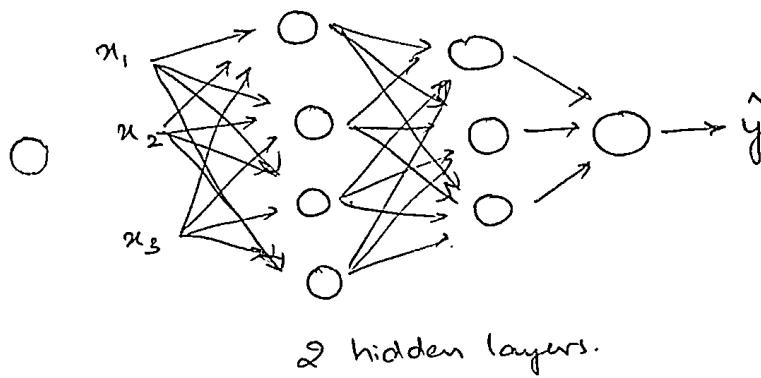
$$dz^{[1]} = aW^{[2]} dz^{[2]} * g'(z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dz^{[1]} A^{[0]T}$$

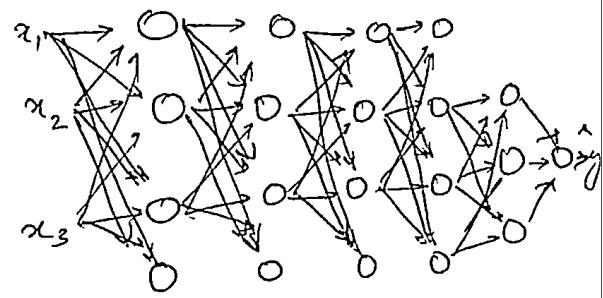
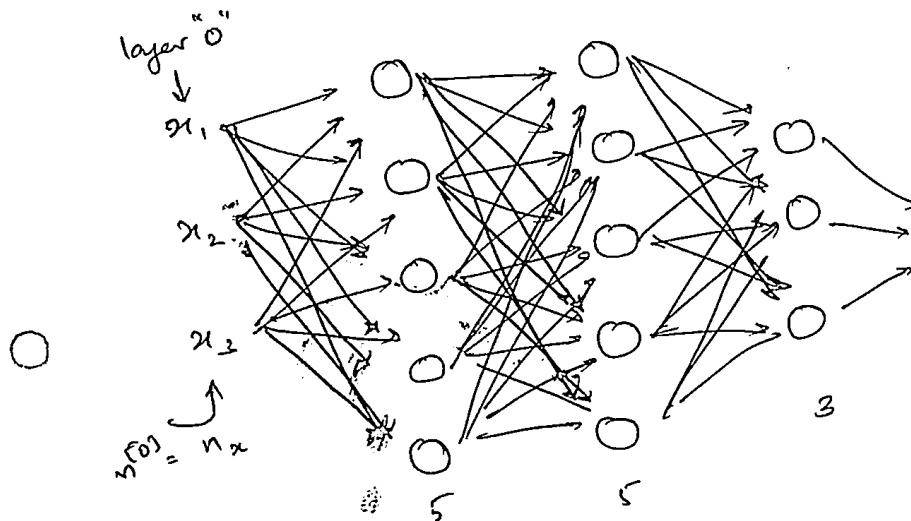
Note that  $A^{[0]T}$  is another way to see previous page.

Deep Neural Networks - Deep L-layer Neural Network# Deep Neural Network

logistic regression "Shallow"

1 hidden layer.  
 (2 Layer) NN

2 hidden layers.

"deep".  
 5 - hidden layers.

4 layer NN

$$O \rightarrow \hat{y} = a^{[L]}$$

$L = 4$  (# layers);  $n^{[l]} = \# \text{units in layer } l$ .  $n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3,$   
 $n^{[4]} = 3, n^{[4]} = n^{[L]} = 1$

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

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

$$a^{[l]} = g(z^{[l]}), \quad w^{[l]} = \text{weights for } z^{[l]}. \quad b^{[l]} = \text{biases }$$

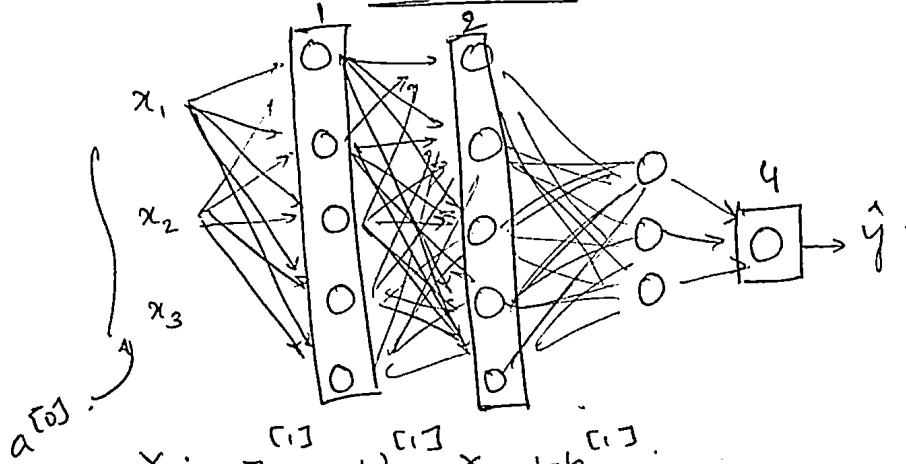
Activation index

finally,

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

Input layer &amp; also activation at index 0.

## \* Forward propagation in a deep network



General Rule

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

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

General forward propagation Eqn.

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

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}.$$

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

activation func

$$z^{[l]} = \underbrace{W^{[l]} a^{[l-1]} + b^{[l]}}_{\substack{\text{Parameters} \\ \text{from } [l-1] \text{ layer}}} \quad \underbrace{\text{activation from} \\ \text{previous layer}}_{\substack{\text{bias from } [l] \text{ layer}}} \quad \text{vector}$$

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

Estimated Output  $\hat{y}$  is computed.

• Vectorized -

$$z^{[1]} = W^{[1]} \times A^{[0]} + b^{[1]} \quad X = A^{[0]}.$$

Seems like  
for  
 $l=1 \dots 4$

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

But

there is  
no way

to compute  
for layers

1, 2, 3, ...

Perfectly ok

$$Y = g(z^{[4]}) = A^{[4]}.$$

Stacking columns  
for m training  
examples.

$$\begin{aligned} z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]}, \\ A^{[2]} &= g^{[2]}(z^{[2]}). \end{aligned} \quad \left. \begin{array}{c} z^{[2](1)} \\ | \\ z^{[2](2)} \\ | \\ \dots \\ | \\ z^{[2](m)} \end{array} \right\} = \underline{\underline{z^{[2]}}}.$$

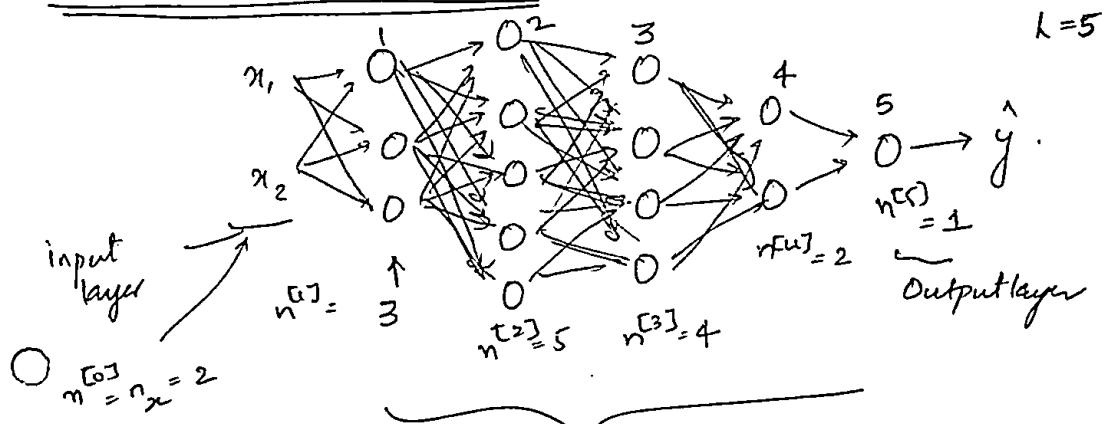
Correction for Video

$$a^{[e]} = g^{[e]}(z^{[e]}). \quad a, z \text{ have dimensions } (n^{[e]}, 1).$$

Week 3

Getting your matrix dimensions right

Parameters  $W^{[e]}$  and  $b^{[e]}$ .



$$z^{[1]} = W^{[1]} \cdot x + b^{[1]}$$

$$(3,4) \leftarrow (3,2)(2,1) + (3,1)$$

$$(n^{[1]}, 1) (n^{[0]}, n^{[1]}) + (n^{[1]}, 1)$$

$$\begin{bmatrix} : \\ : \\ : \end{bmatrix}_{3 \times 1} = \begin{bmatrix} : \\ : \\ : \end{bmatrix}_{3 \times 2} \begin{bmatrix} : \\ : \\ : \end{bmatrix}_{2 \times 1} + \begin{bmatrix} : \\ : \\ : \end{bmatrix}_{3 \times 1}$$

from Matrix Multiplication

$$W^{[1]} : (n^{[1]}, n^{[0]})$$

$$W^{[2]} : (5, 3) : (n^{[2]}, n^{[1]})$$

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

$$(5,1) \leftarrow (5,3)(3,1) + (5,1)$$

$$W^{[3]} : (4, 5) \quad (n^{[3]}, n^{[2]})$$

$$W^{[4]} : (2, 4), \quad W^{[5]} : (1, 2)$$

Dimensions of your matrices.

$$W^{[e]} : (n^{[e]}, n^{[e-1]}), \quad b^{[e]} : (n^{[e]}, 1)$$

for Back propagation:

$$dW^{[e]} : (n^{[e]}, n^{[e-1]})$$

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

same as same as

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

## Vectorized Implementation :-

Normally (Nonvector)

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

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

for Vectorized (Stacking them).

$$\left[ \begin{array}{cccc} Z^{[l](1)} & Z^{[l](2)} & \dots & Z^{[l](m)} \end{array} \right]$$

$(n^{[l]}, m)$      $(n^{[l]}, n^{[l]})$      $(n^{[l]}, m)$      $(n^{[l]}, 1)$      $\downarrow$   
 through Python Broadcasting  
 $(n^{[l]}, m)$

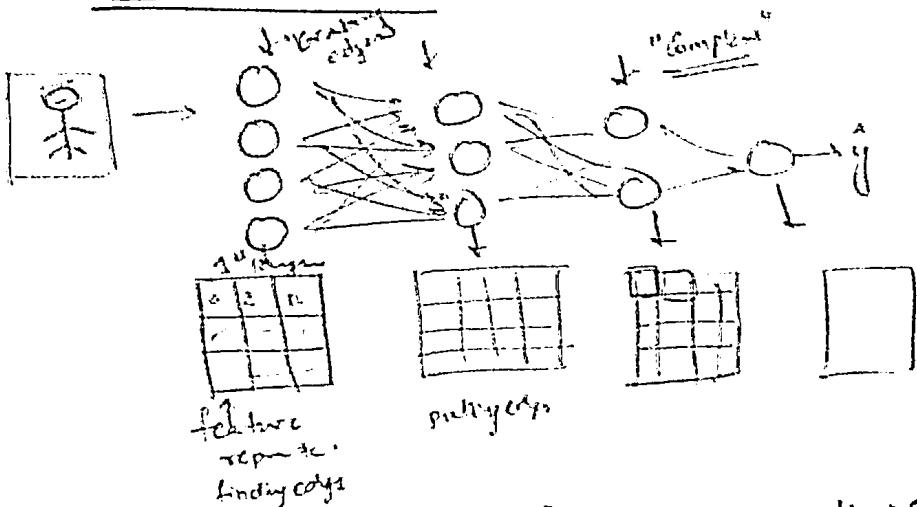
$$Z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

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

Capital ~~vector~~  
 if  $l=0 \rightarrow A^{[0]} = X = (n^{[0]}, m)$

$$dZ^{[l]}, dA^{[l]} : (n^{[l]}, m).$$

## Intuition about deep representation



Audio  $\rightarrow$  low level  $\rightarrow$  Phonemes  $\rightarrow$  words  $\rightarrow$  sentences/phrases  
 Audio waveform  $\rightarrow$  CAT

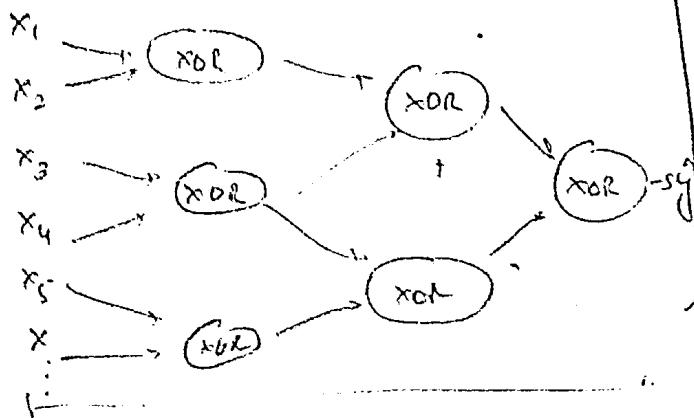
## Circuit Theory and Deep Learning

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

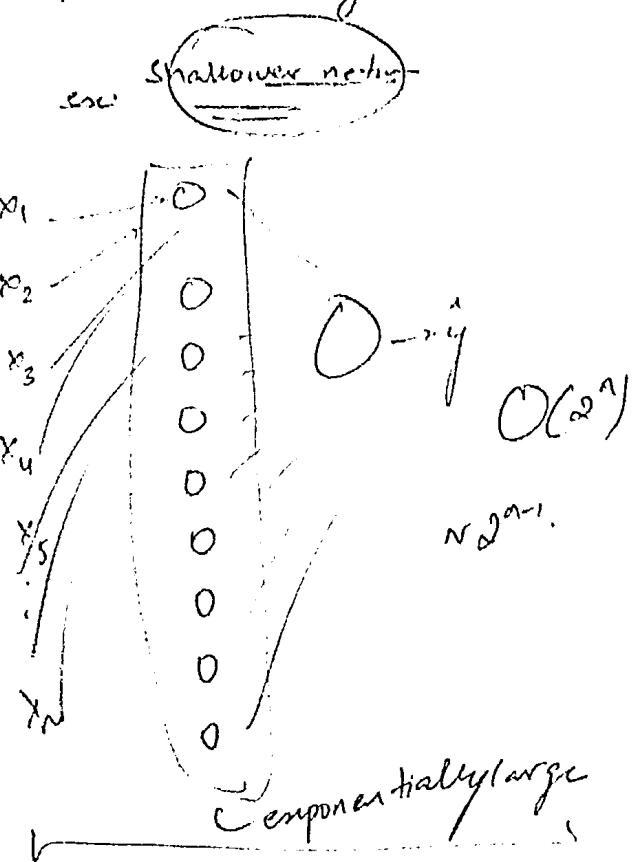
### Deep Netw.

$$y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n.$$

$O(\log n)$



More efficient



## \* Building blocks of Deep Neural Networks:

$$\begin{matrix} x_1 & 0 & 0 & 0 \\ x_2 & 0 & 0 & 0 \\ x_3 & 0 & 0 & 0 \\ x_4 & 0 & 0 & 0 \end{matrix}$$

$$0 \rightarrow \hat{y}$$

layer  $l: W^{[e]}, b^{[e]}$

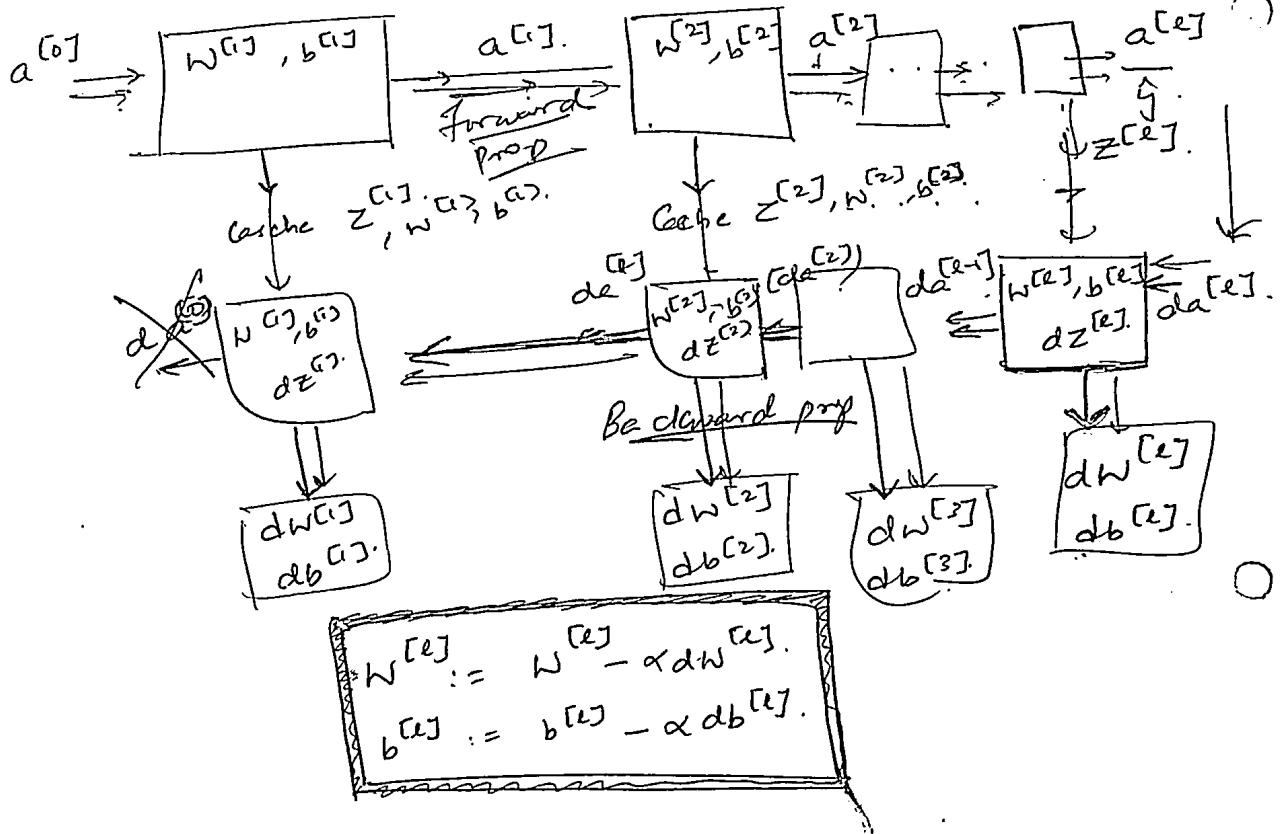
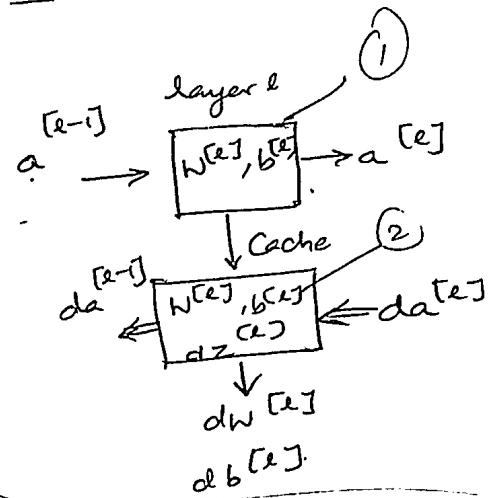
forward  $\circ$  input  $a^{[e-1]}$ , output  $a^{[e]}$ .

$$\underline{z}^{[e]} = W^{[e]} a^{[e-1]} + b^{[e]}. \quad \text{Cache } z^{[e]}.$$

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

Backward  $\circ$  input  $d a^{[e]}$ , output  $d a^{[e-1]}$ .  
Cache  $(z^{[e]})$

$$\begin{aligned} d a^{[e-1]} \\ d W^{[e]} \\ d b^{[e]} \end{aligned}$$



Correct Video

$$dW^{[e]} = dZ^{[e]} * a^{[e-1]T}$$

## # Forward propagation for layer $i$

→ Input  $a^{[e-i]}$ .

→ Output  $a^{[e]}$ , Cache ( $z^{[e]}$ ).

$$z^{[e]} = w^{[e]} \cdot a^{[e-i]} + b^{[e]}$$

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

$w^{[e]}, b^{[e]}$ .

Vectorized

$$z^{[e]} = w^{[e]} \cdot A^{[e-i]} + b^{[e]}$$

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

$$X = A^{[0]} \rightarrow \square \rightarrow \square \rightarrow \square \rightarrow \dots$$

## Backward Propagation for layer $i$



→ Input  $da^{[e]}$ .

→ Output  $da^{[e-i]}, dw^{[e]}, db^{[e]}$ .

$$\frac{dz^{[e]}}{da^{[e]}} = g^{[e]}'(z^{[e]})$$

$$dw^{[e]} = dz^{[e]} * a^{[e-i]T}$$

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

$$da^{[e-i]} = w^{[e]T} \cdot dz^{[e]}$$

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

Vectorized

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

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

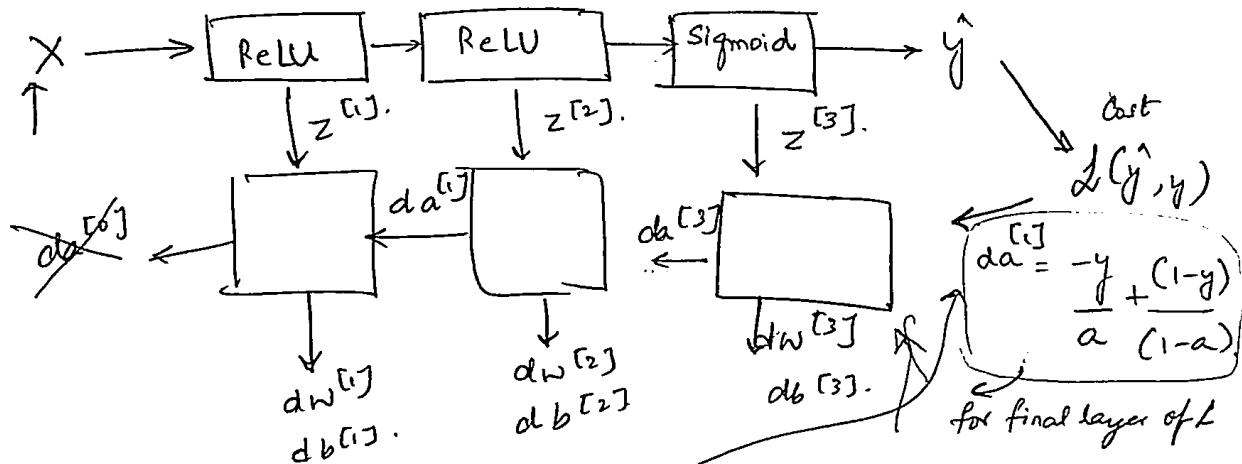
$$db^{[e]} = \frac{1}{m} \text{np.sum}(dz^{[e]}, \text{axis}=1)$$

$$dA^{[e-i]} = w^{[e]T} \cdot dz^{[e]}$$

Backpropagation



## SUMMARY



This is how you  
Initialize the vectorized version of the  
Back propagation -

$$dA^{[e]} = \left( -\frac{y^{(1)}}{a^{(1)}} + \frac{(1-y^{(1)})}{(1-a^{(1)})} \dots \right) \dots \left( -\frac{y^{(m)}}{a^{(m)}} + \frac{(1-y^{(m)})}{(1-a^{(m)})} \right)$$

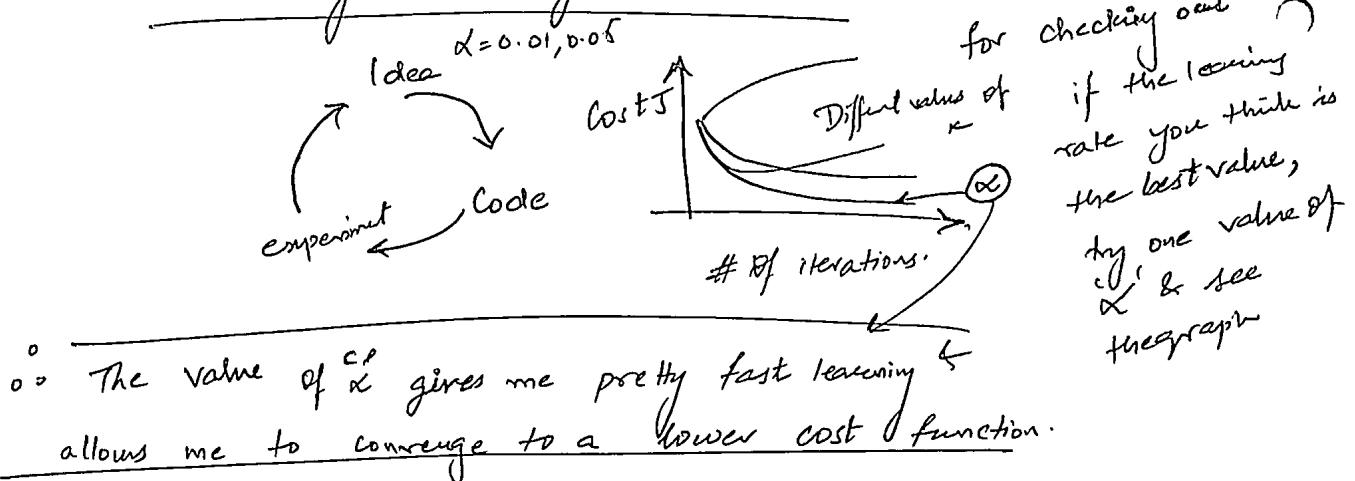
\* Parameters vs Hyperparameters:

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

Hyperparameters:  
learning rate  $\alpha$   
 $\#$  iterations  
 $\#$  hidden layers  $L$   
 $\#$  hidden units  $n^{[1]}, n^{[2]}, \dots$   
Choice of activation function.

Later: Momentum, mini batch size, regularizations...

Note: Applied deep learning is a very empirical process.



// ① Try out many different values & check out which is best.

②

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

$$dW^{[e]} = \frac{1}{m} dZ^{[e]} A^{[e-1]T}$$

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

$$dZ^{[e-1]} = W^{[e]T} dZ^{[e]} * g'(z^{[e-1]}).$$

: element wise multiplication

$$dZ^{[1]} = W^{[2]} dZ^{[2]} * g'(z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} A^{[0]T}$$

Note that  $A^{[0]T}$  is another way to denote the input features, which is also written as  $X^T$ .  $db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$ .