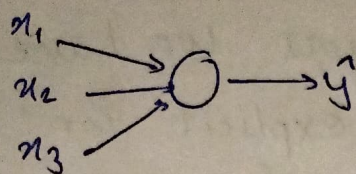
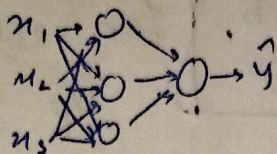


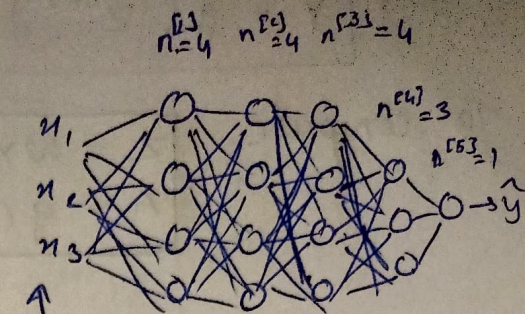
Deep Neural Networks



Logistic Regression.
(Shallow NN)



"2 layer NN"
1 hidden layer



"deep" NN
4 hidden layers.

(5-layer NN)
11
L

$L \rightarrow$ no. of layers

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

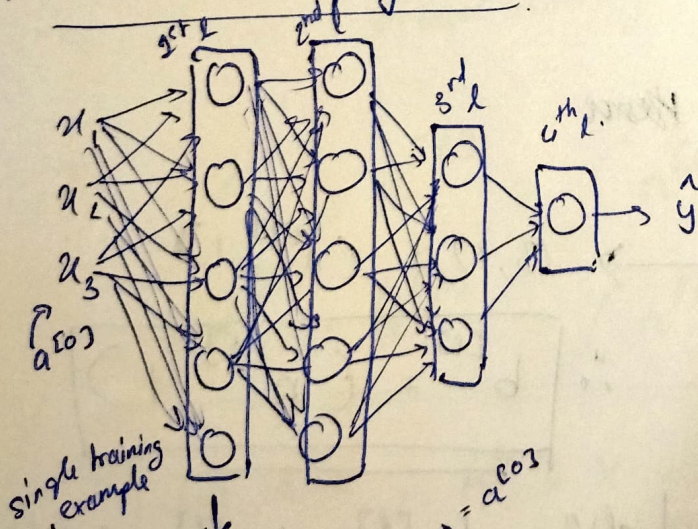
$a^{[l]} =$ activations in layer l .

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

$a^{[0]} = x$ (Input layer = 0)

$a^{[L]} = \hat{y}$ (Output layer = L)

Forward Propagation in a deep network.



vectorized:

$$z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

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

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

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

for training examples

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

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

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

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

General:

$$\begin{cases} z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

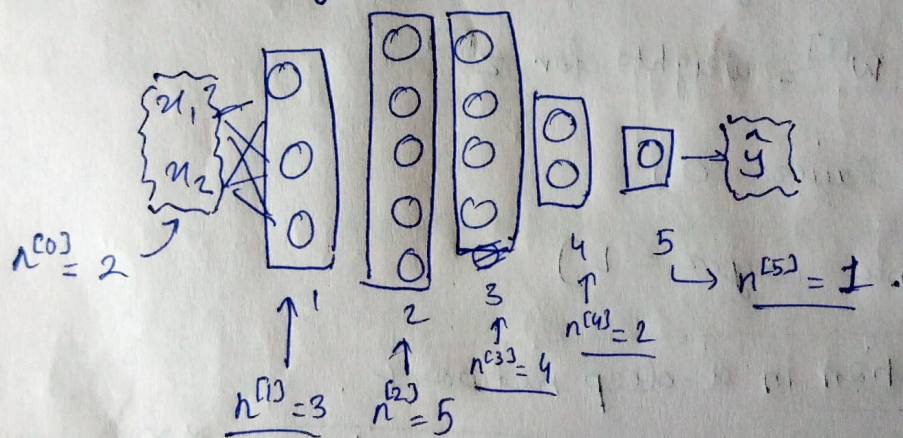
Here, $\hat{y} = g(z^{[L]}) = A^{[L]}$

To compute $z = wx + b$, $A = g(z)$ for all layers we have to use explicit 'for' loop. No other way to do it.

* One of the bugs in your neural network could be occur because of matrix dimensions.

• Parameters $w^{[L]}$ and $b^{[L]}$

We have a 5-layer NN.



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

$$(3,1) \quad (3,2) \quad (2,1)$$

$$(n^{[L]}, 1) \quad (n^{[L]}, n^{[L-1]}) \quad (n^{[L-1]}, 1)$$

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}_{3 \times 2} \begin{bmatrix} 0 \\ 0 \end{bmatrix}_{2 \times 1}$$

Hence

$$(3,1) = (n^{[L]}, 1)$$

$$\therefore b^{[L]} = (n^{[L]}, 1)$$

$$\therefore w^{[L]} = (n^{[L]}, n^{[L-1]})$$

$$w^{[2]} = (5, 3) \quad w^{[5]} = (1, 2)$$

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

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

also, $dw^{[L]} = w^{[L]} = (n^{[L]}, n^{[L-1]})$

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

• $Z^{[L]}$, X , $A^{[L]}$

for 1 training example:- $Z^{[L]} = W^{[L]} \cdot X + b^{[L]}$
 $(n^{[L]}, 1) \quad (n^{[L]}, n^{[0]}) \quad (n^{[0]}, 1) \quad (n^{[L]}, 1)$

for m training examples:-

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

m training examples

$$Z^{[L]} = W^{[L]} \cdot X + b$$

$(n^{[L]}, m) \quad (n^{[L]}, n^{[0]}) \quad (n^{[0]}, m) \quad (n^{[L]}, 1)$

remains same

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

no. of examples

no. of features
(x_1, x_2)

$\underbrace{W \cdot X}_{(n^{[L]}, m)} + \underbrace{b}_{(n^{[L]}, 1)}$
 \downarrow
 by python broadcasting $\rightarrow (n^{[L]}, m)$

$\therefore Z^{[L]}, a^{[L]} = (n^{[L]}, 1) \rightarrow$ for 1 training example

$Z^{[L]}, A^{[L]} = (n^{[L]}, m) \rightarrow$ for m training examples.

{ for $l=0$, $A^{[0]} = X = (n^{[0]}, m)$ }

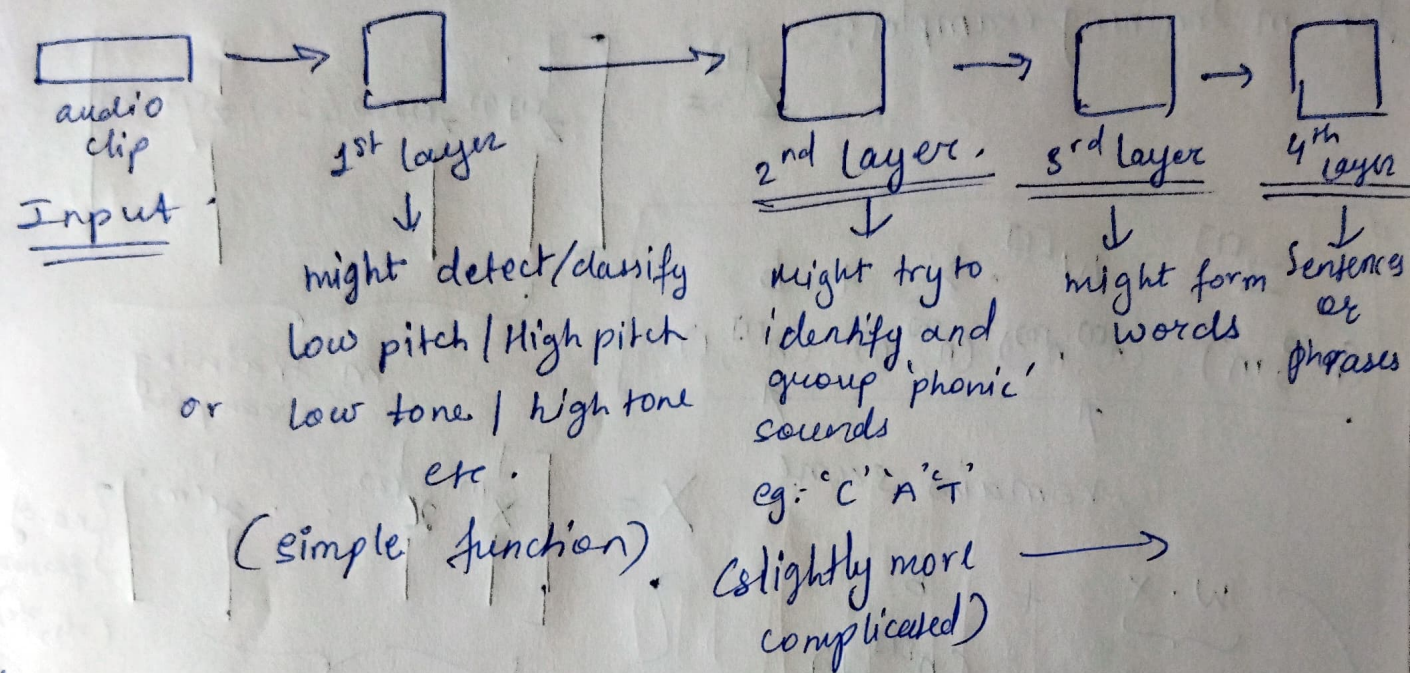
in back propagation

$dz^{[L]}, dA^{[L]} = (n^{[L]}, m)$

WHY Deep Representation?

What does a deep representation do?

eg: In case of audio recognition,



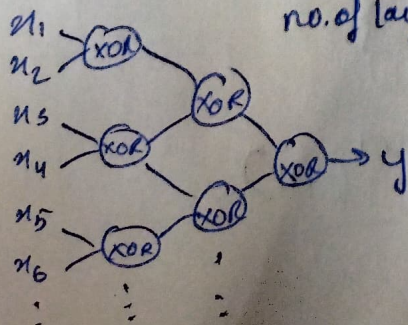
Circuit theory

There are some functions which we can compute with a "small" L-layer deep neural network, that shallower and "bigger" networks require exponentially more hidden units to compute.

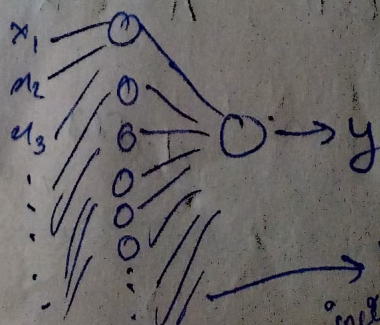
eg:- $y = x_1 \text{ XOR } x_2 \text{ XOR } x_3 \dots \text{ XOR } x_n$

small but deep NN

no. of layers = $O(\log n)$



big but shallow NN



no. of units $\approx 2^n$
increases exponentially

* forward and backward functions.

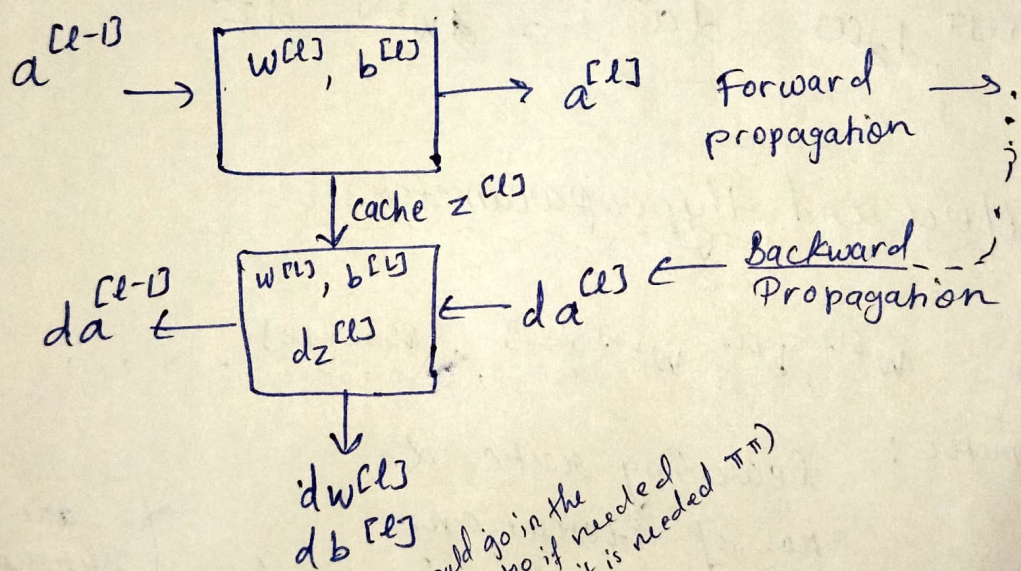
13

for a layer l : $w^{[l]}, b^{[l]}$

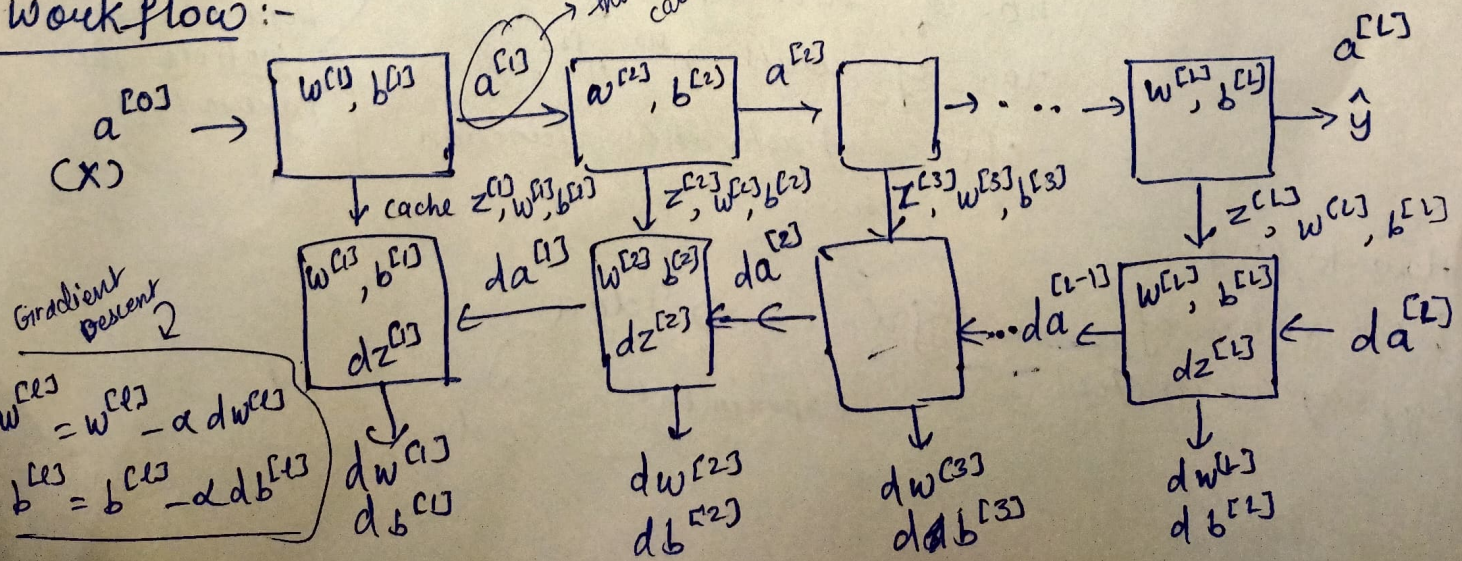
→ forward: Input $a^{[l-1]}$, output $a^{[l]}$
 $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$ cache $z^{[l]}$
 $a^{[l]} = g^{[l]}(z^{[l]})$

→ Backward: input $da^{[l]}$, output $da^{[l-1]}$
cache $(z^{[l]})$ $\frac{dw^{[l]}}{dz^{[l]}}$
 $\frac{db^{[l]}}{dz^{[l]}}$

layer l .



Workflow:-



Forward Propagation:- Input $a^{[L-1]}$ Output $a^{[L]}$, cache ($z^{[L]}$)

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

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

Backward Propagation:- Input $da^{[L]}$ Output $da^{[L-1]}$, $dw^{[L]}$, $db^{[L]}$

vectorized

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

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

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

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

$$dz^{[L]} = da^{[L]} * g^{[L]'}(z^{[L]})$$

$$dw^{[L]} = dz^{[L]} \cdot a^{[L-1]} \rightarrow \text{obtained from cache.}$$

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

$$da^{[L-1]} = dw^{[L]T} \cdot dz^{[L]}$$

* Parameters and Hyperparameters

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

Hyperparameters :- learning rate α .

no. of iterations

no. of hidden layers L

no. of units $n^{[1]}, n^{[2]}, \dots$

choice of activation function etc.

are called hyperparameters because they control the parameters

How to find best value of hyperparameter?

eg] α .

Idea
experiment
code