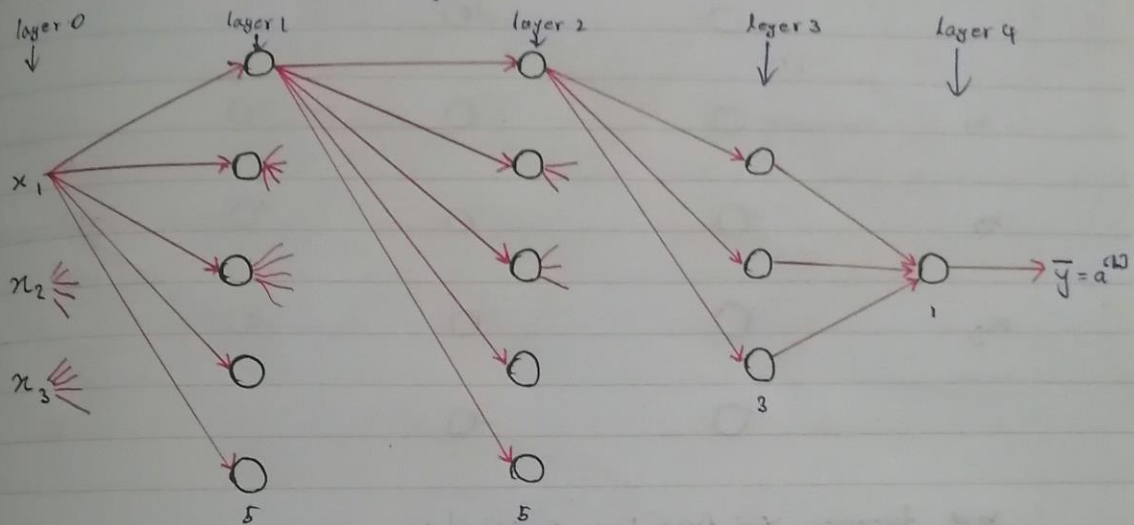


Week 4

Deep Neural Network

Deep L-layer Neural Network



$L = 4$ (Number of layers)

$n^{[l]}$ = units in layer l

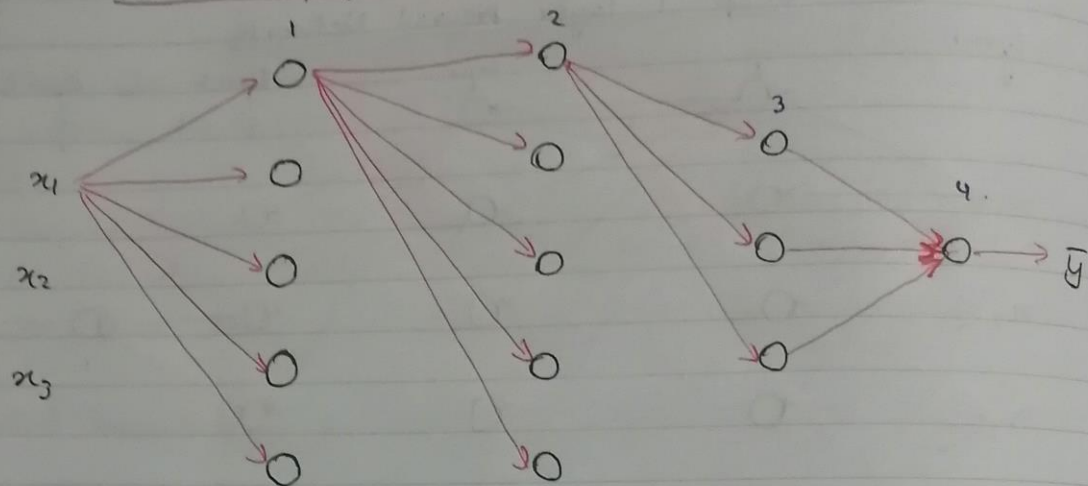
$a^{[l]}$ = activations in layer l

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

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

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

Forward Propagation in a Deep Network



~~x~~ ~~training~~ x : training example

layer 1

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

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

layer 3

$$z^{[3]} = w^{[3]}a^{[2]} + b^{[3]}$$

$$a^{[3]} = g(z^{[3]})$$

layer 2

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

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

layer 4

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

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

formula

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

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

Now we look into vectorized method for above calculations.

layer 1

$$\mathbf{Z}^{[1]} = \mathbf{W}^{[1]} \mathbf{A}^{[0]} + \mathbf{b}^{[1]}$$

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

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

layer 2

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

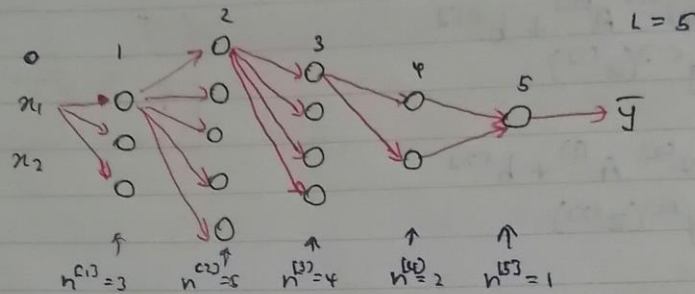
$$\mathbf{A}^{[2]} = g^{[2]}(\mathbf{Z}^{[2]})$$

$$\mathbf{Y} = g^{[3]}(\mathbf{Z}^{[2]}) = \mathbf{A}^{[2]}$$

vectorized version formula

$$\begin{aligned} \mathbf{Z}^{[1]} &= \mathbf{W}^{[1]} \mathbf{A}^{[0]} + \mathbf{b}^{[1]} \\ \mathbf{A}^{[1]} &= g^{[1]}(\mathbf{Z}^{[1]}) \end{aligned}$$

Getting your Matrix Dimension right



layer 1

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

Dimensions: $(3,1)$ $(3,2)$ $(2,1)$

Dimensions: $(n^{(1)}, 1)$ $(n^{(1)}, n^{(0)})$ $(n^{(1)}, 1)$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

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

$$W^{(2)}: (5, 3) \quad (n^{(2)}, n^{(1)})$$

layer 2

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

Dimensions: $(5,1)$ $(5,3)$ $(3,1)$

likewise \rightarrow

$$W^{(3)}: (4, 5)$$

$$W^{(4)}: (2, 4)$$

$$W^{(5)}: (1, 2)$$

Vectorized implementation above model

$$\begin{matrix} \mathbf{z}^{[1]} \\ (n^{[1]}, m) \end{matrix} = \begin{matrix} \mathbf{W}^{[1]} \\ (n^{[1]}, n^{[0]}) \end{matrix} \cdot \begin{matrix} \mathbf{X} \\ (n^{[0]}, m) \end{matrix} + \begin{matrix} \mathbf{b}^{[1]} \\ (n^{[1]}, m) \end{matrix}$$

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

note

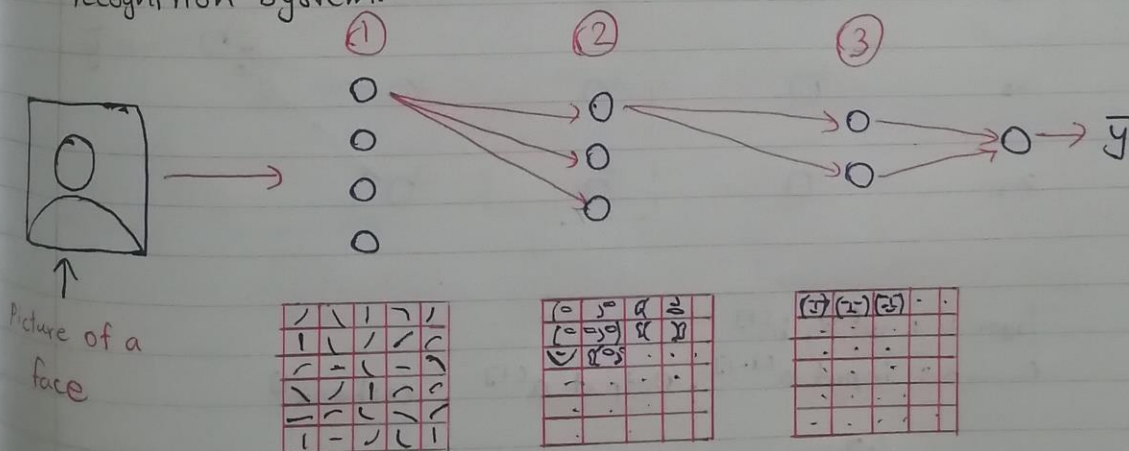
$$\mathbf{z}^{[1]}, \mathbf{a}^{[1]} : (n^{[1]}, 1)$$

$$\mathbf{Z}^{[1]}, \mathbf{A}^{[1]} : (n^{[1]}, m)$$

$$d\mathbf{Z}^{[1]}, d\mathbf{A}^{[1]} : (n^{[1]}, m)$$

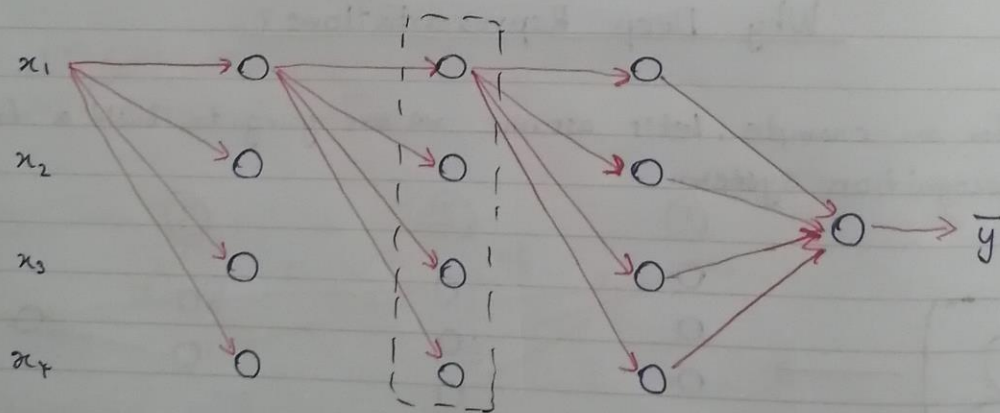
Why Deep Representations?

* For an example, let's assume we are going to build a face recognition system.



- ① → You enter a picture of a face and then first layer of the neural network, we can think may be use for feature detector of an edge detector. given all square bones represent edges of the given picture.
- ② → In this layer, edges and group edges detected together and it represent parts of the faces (nose, eyes, ears, etc.)
- ③ → In this layer, it putting together different parts of faces and try to recognize or detect different typer of faces.

Building Blocks of Deep Neural Networks



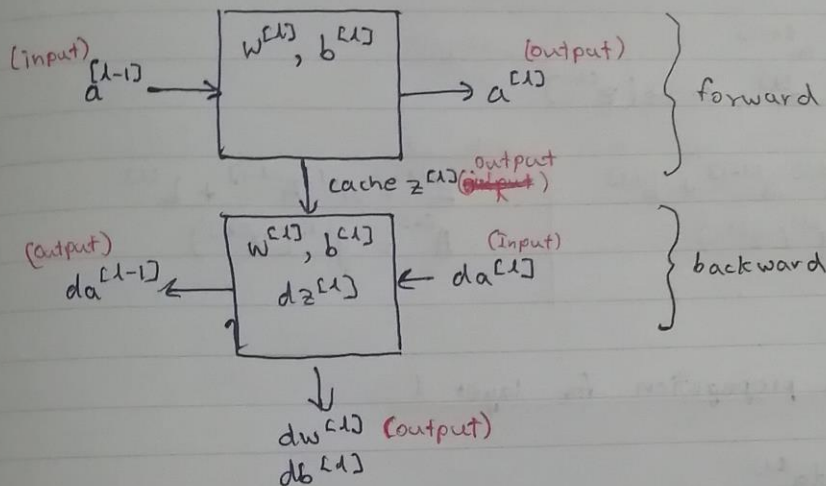
layer l : $w^{[l]}$, $b^{[l]}$
 forward : input $a^{[l-1]}$, Output $a^{[l]}$, cache $z^{[l]}$

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

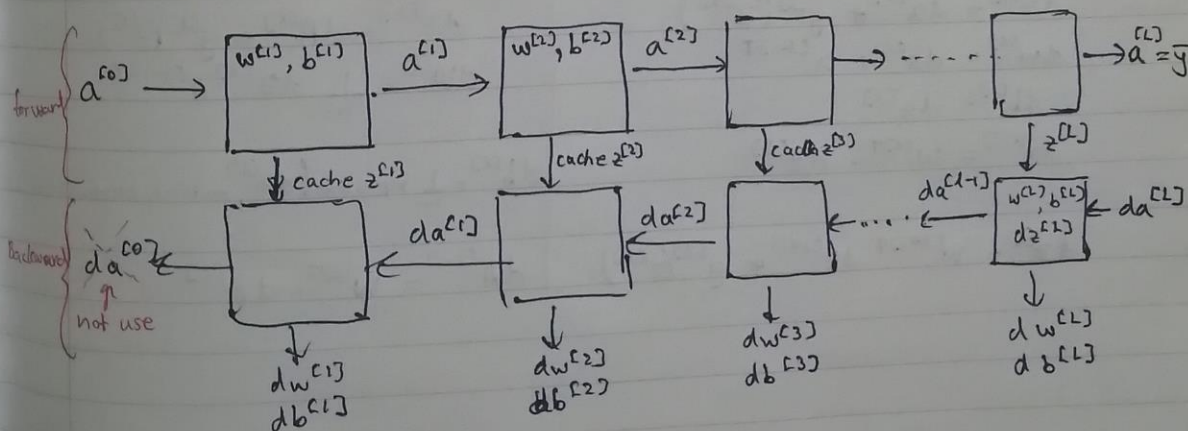
$$a^{[l]} = g^{[l]}(z^{[l]})$$

Backward : input $da^{[L]}$, output $da^{[L-1]}$, cache $z^{[L]}$

layer 1



* After applying above method into neural network



Forward and Backward Propagation

Forward propagation for layer l

Input $a^{[l-1]}$
Output $a^{[l]}$, cache($z^{[l]}$) $\leftarrow w^{[l]}, b^{[l]}$

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

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

vectorized

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

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

Backward propagation for layer l

Input $da^{[l]}$

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

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

$$dW^{[l]} = dz^{[l]} * a^{[l-1]T}$$

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

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

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

vectorized

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

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

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

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

Parameters vs Hyperparameters.

What are hyperparameters?

Parameters $\therefore w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, w^{(3)}, b^{(3)}, \dots$

Hyperparameters \therefore learning rate α (ex: number of iterations of gradient descent you carry out $\frac{1}{2}$)

\therefore number of hidden layers (L)

\therefore number of hidden units ($n^{(1)}, n^{(2)}, \dots$)

\therefore choice of activation function

- * hyperparameters need to ~~test~~ control parameters above mentioned ($w^{(1)}, b^{(1)}, \dots$)
- * hyperparameters determine final values of w and b .
- * In later we will learn these hyperparameters \therefore momentum term, mini batch size, various forms of regularization parameters, \dots