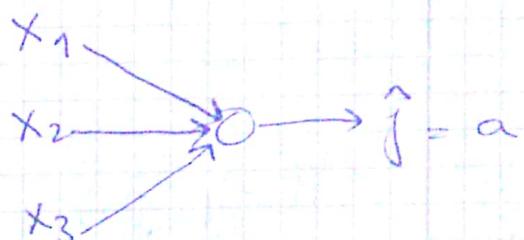


Week # 3:-

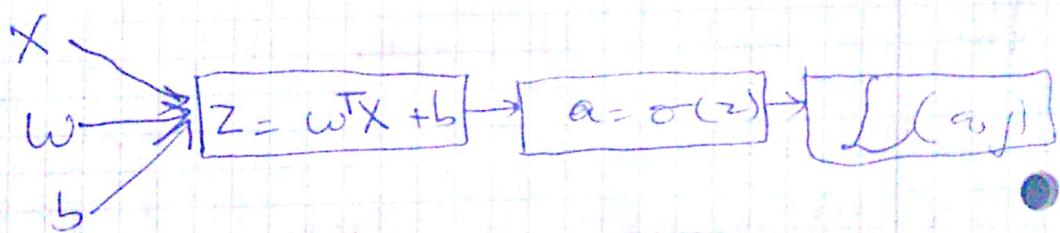
Shallow Neural Network:-

Neural network overview:-

What is neural network



Neural network can be formed by stacking together a lot of sigmoid units

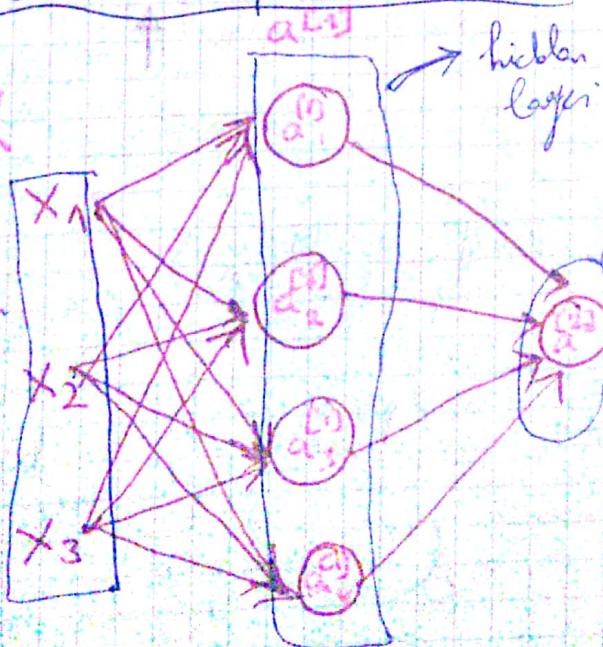


Neural Network Representation:-

Ip. feature (x)

Input layer

to neural network



2-layer NN

output layer \Rightarrow predicted value

$$\hat{j} = a_{21}$$

$$a_{21}$$

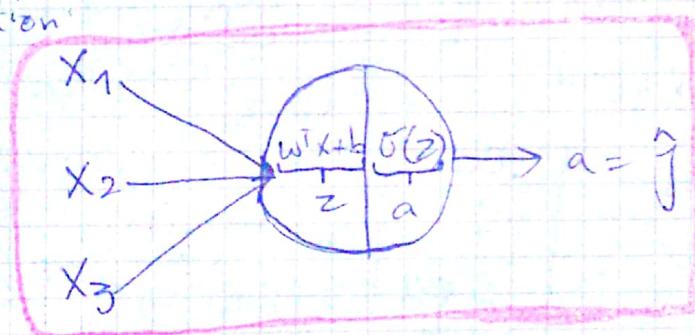
$w^{[1]}$ here will be 4×3 matrix. 4 comes from the fact that we have 4 units inside the hidden layer and 3 comes from the fact that we have 3 features.

$$b^{[1]} \in 4 \times 1$$

$$w^{[2]} \in R^{1 \times 4} \quad b^{[2]} \in R^{1 \times 1}$$

Computing a Neural Network output:-

- Now let's recall we said logistic regression consists of 2 steps of computation first you compute z then you compute sigmoid as activation function



- In the last 2 layered network, let's look at the first node first layer. Same as above figure it will compute:

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]} \text{ then } a_1^{[1]} = \sigma(z_1^{[1]}) \Rightarrow \text{activation:}$$

So convention here is

$a^{[l]}$ \leftarrow layer

$a_i^{[l]}$ \leftarrow node in layer

Second node in first layer:

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

3rd node 1st layer

$$z_3^{(1)} = w_3^{(1)T} x + b_3^{(1)}; a_3^{(1)} = \sigma(z_3^{(1)})$$

4th node 1st layer

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}; a_4^{(1)} = \sigma(z_4^{(1)})$$

So now as we have finished our first layer
Computation now to implement it we will not
Use for loops we will use vectorization

→ Stack w_i 's in a matrix.

$$\rightarrow z = \begin{bmatrix} w_1^{(1)T} x + b_1^{(1)} \\ w_2^{(1)T} x + b_2^{(1)} \\ w_3^{(1)T} x + b_3^{(1)} \\ w_4^{(1)T} x + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(1)T} \\ w_2^{(1)T} \\ w_3^{(1)T} \\ w_4^{(1)T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_n^{(1)} \end{bmatrix} = \sigma(z^{(1)}).$$

→ Now for the second layer

$$\rightarrow z^{(2)} = (w^{(2)} a^{(1)} + b^{(2)})^{(1 \times 4)}$$

$$\rightarrow a^{(2)} = \sigma(z^{(2)})^{(1 \times 1)}$$

Vectorizing across multiple examples:

- Let's see how to vectorize across multiple training examples:

$X^{[2]}$ → $a^{[2]} = \hat{y}$
 Example:
 $a^{[2](i)}$ ← Layer 2

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

$$X^{[1]} \rightarrow a^{[1]} = \hat{y}^{[1]}$$

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

$$\vdots$$

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

* for $i=1$ to m :

$$z^{(i)} = w^{(2)} x^{(i)} + b^{(2)}$$

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

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

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

We want to vectorize the whole computation, in order to get rid of the main loop *

→ So Recall, our training examples stored in a matrix like this

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

Note the vectorized version of $A^{[l]}$ $x = a^{[l]}, X = a^{[l] \times [l+1]}$

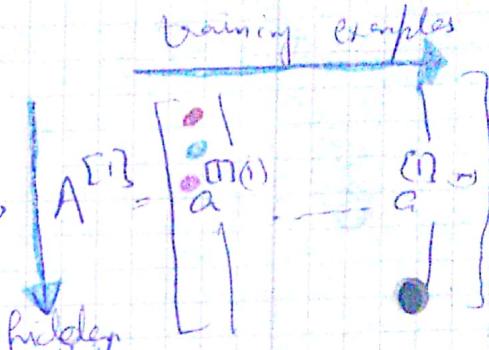
$$\rightarrow Z^{[1]} = w^{[1]} X + b^{[1]} \quad \leftarrow w^{[1]} A^{[0]} + b^{[1]}$$

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

$$\rightarrow Z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

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

$$Z^{[1]} = \begin{bmatrix} | & | & | & | \\ z^{(1)(1)} & z^{(1)(2)} & \dots & z^{(1)(m)} \end{bmatrix}; \quad A^{[1]} = \begin{bmatrix} | & | & | \\ a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(n)} \end{bmatrix}$$



- Corresponds to the activation value of 1st hidden unit on first training example in 1st layer
- Corresponds to the activation value of 2nd hidden unit on first training example in 1st layer
- activation value \Rightarrow 3rd hidden unit, 1st layer (hidden)
- 1st training example

Explanation for Vectorized implementation:-

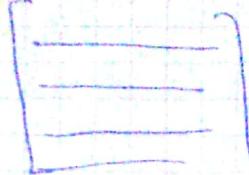
Let's say we have 3 training examples

then Computation would be

$$Z^{1} = w^{1} X + b^{[1]}, \quad Z^{[1](2)} = w^{[1](2)} X + b^{[1]}, \quad Z^{[1](3)} = w^{[1](3)} X + b^{[1]}$$

For the sake of simplicity let's say that all b 's are zero

So $w^{[1]}$



$$w^{[1]} x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}; w^{[2]} x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}; w^{[3]} x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$w^{[1]}$

$$\begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots \end{bmatrix}$$

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

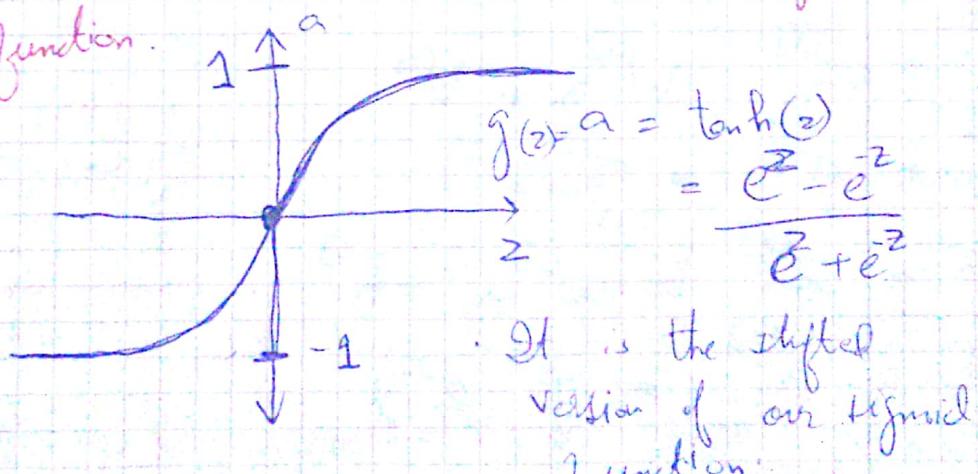
$$= \begin{bmatrix} z^{(1)(1)} & | & z^{(1)(2)} & | & z^{(1)(3)} \\ 1 & | & 1 & | & 1 \end{bmatrix} = z^{(1)}$$

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

Activation functions:-

So far, we have been using sigmoid activation function, now we look other options too.

→ Other function that we can use is tangent hyperbolic function.

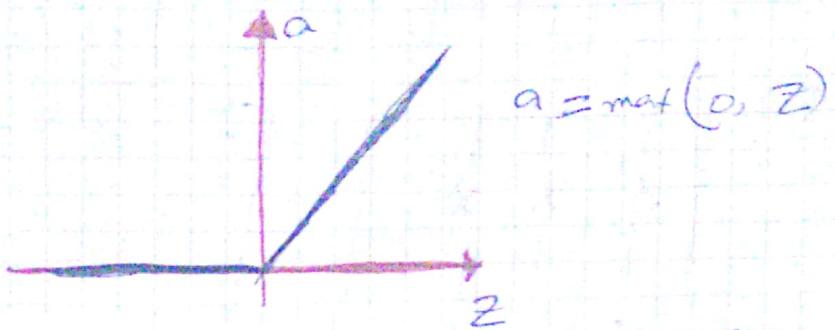


It is the shifted version of our sigmoid function.

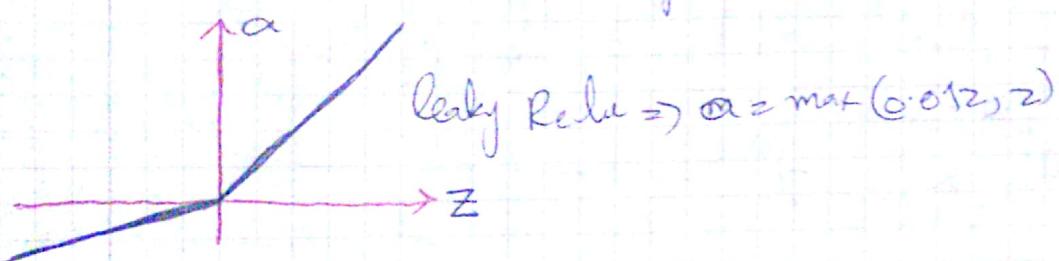
Q) "tanh" works better (always) with the exception of the output layer where you want the probability to be between 0 and 1.

S. in layered NN $g^{[1]}(z^{[1]}) = \tanh(z^{[1]})$
 $g^{[2]}(z^{[2]}) = \delta(z^{[2]})$.

One other function is ReLU (Rectified linear unit function)



One disadvantage of ReLU is that it has derivative zero with z being less than zero.



→ Why do you need a non-linear activation function?

From the algorithm that we established earlier

Given x :

1. $\bullet z^{[1]} = W^{[0]}x + b^{[0]}$
2. $\bullet a^{[1]} = g^{[1]}(z^{[0]})$ $g(z) = z$
3. $\bullet z^{[2]} = W^{[1]}a^{[0]} + b^{[1]}$
4. $\bullet a^{[2]} = g^{[2]}(z^{[1]})$

Let's think for a moment that the function $g(z)$ is the Identity function (that is whatever goes in comes out)

then:

$$\begin{aligned}
 a^{[1]} &= z^{[1]} = w^{[1]}x + b^{[1]} \\
 a^{[2]} &= g(z) = w^{[2]}_{a^1} a^{[1]} + b^{[2]} \\
 a^{[3]} &= w^{[3]}(w^{[1]}x + b^{[1]}) + b^{[3]} \\
 &= \underbrace{(w^{[2]} w^{[1]})x}_{w^*} + \underbrace{(w^{[3]} b^{[1]} + b^{[3]})}_{b^*} \\
 a^{[2]} &= w^*x + b^*
 \end{aligned}$$

If we were to use Identity activation function then our neural network will be using linear function and outputting it i.e. linear output

→ So later when we will use network with many hidden layers, then with identity activation function we will be outputting linear function's output and it's like not using a network with many layers.

So linear hidden layer is more or less useless.

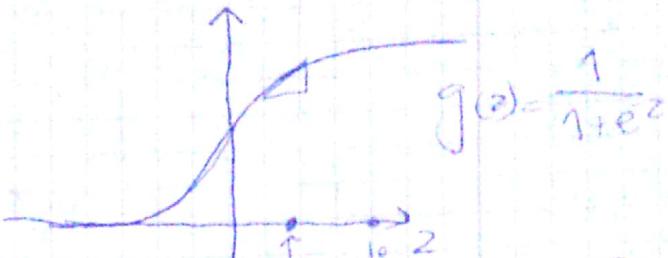
→ There is one place, where you might use linear activation function $g(z) = z$ and that is if you are doing machine learning on a regression problem. I.e. if $y \in \mathbb{R}$. not only from 0 to 1

In this case we use Identity or linear activation function only in the output layer and hidden layers still use ReLU or tanh.

Derivatives of activation functions:-

When you implement back propagation you compute derivatives of activation functions, so let's take a look at our choices of activation functions and how you can compute the slopes of these functions.

Let's see sigmoid function:



$$\frac{d}{dz} (g(z)) = \text{slope of } g(z) \text{ at } z \quad \text{let's take a look at this point}$$

So let's take derivative

$$g'(z) = \frac{d}{dz} \left[\frac{1}{1+e^z} \right] = \frac{-1(-e^z)}{(1+e^z)^2} = \frac{1}{1+e^z} \left(1 - \frac{1}{1+e^z} \right)$$

→ look at the power rule of calculus

$$\text{so } \frac{d}{dz} g(z) = g(z)(1-g(z)) \Rightarrow a(1-a)$$

Let's take a look at $z=10 \Rightarrow g(z) \approx 1$

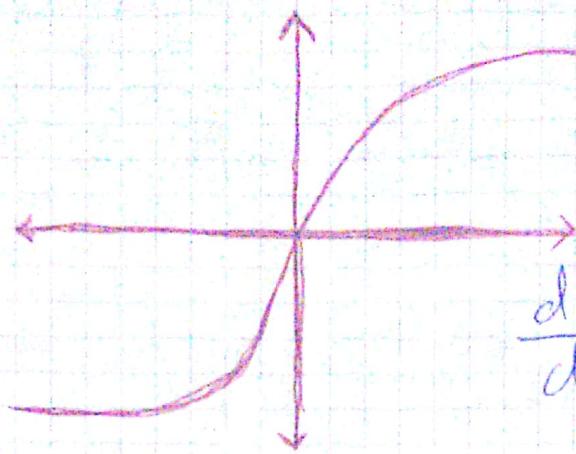
$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$$

and it's clear from the figure that as z increases rate of change decreases

→ Conversely if $z=-10 \Rightarrow g(z) \approx 0$

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

Let's now look at the tanh activation function:



$$g(z) = \tanh(z)$$

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

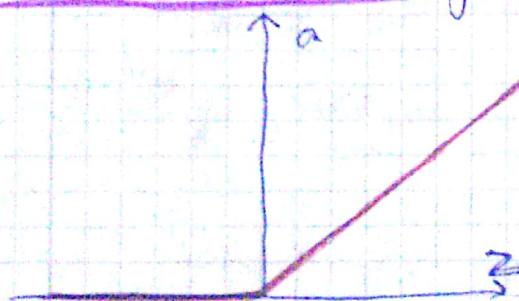
$$z=10 \quad \tanh(10) \approx 1$$

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

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

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

- ReLU and Leaky ReLU



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

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

but when undefined if $z=0$
implementing in software

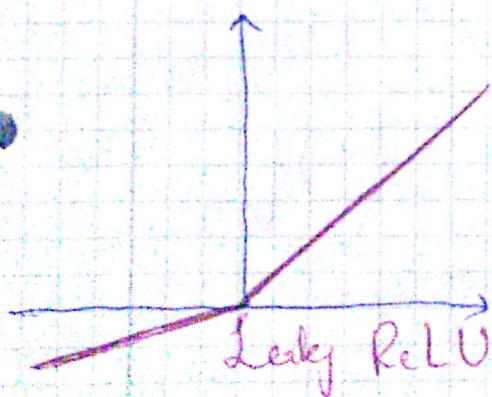
at zero we take $g'(0)=1$ or

① doesn't matter

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

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

undefined if $z=0$



* probability of $g(z)$ outputting 0.0000000000000001
Very low.

Gradient Descent for Neural network :-

Let's dive into, how to implement gradient descent on neural network with one hidden layer.

→ Neural network with single hidden layer has parameter $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

→ If you have n_x or alternatively $n^{[0]}$ input feature and $n^{[1]}$ hidden units and $n^{[2]}$ output units

then $w^{[1]} \in \mathbb{R}^{(n^{[0]} \times n^{[1]})}$ $w^{[2]} \in \mathbb{R}^{(n^{[1]} \times n^{[2]})}$
 $b^{[1]} \in \mathbb{R}^{(n^{[1]} \times 1)}$ $b^{[2]} \in \mathbb{R}^{(n^{[2]} \times 1)}$

Assuming we are doing linear classification.

Cost function is : $J(w, b^{[1]}, w^{[2]}, b^{[2]})$

$$= \frac{1}{m} \sum_{i=1}^m L(g_i, y_i)$$

$\underbrace{}_a$

Now to train the parameters, we need to implement gradient descent.

Repet

Compute predictions ($\hat{y}^i, i=1, \dots, m$).

$$\frac{dJ}{dw^{[1]}} = \frac{\partial J}{\partial w^{[1]}}, \quad \frac{dJ}{db^{[1]}} = \frac{\partial J}{\partial b^{[1]}}$$

$$w^{[1]} = w^{[1]} - \alpha \frac{dJ}{dw^{[1]}}$$

$$b^{[1]} = b^{[1]} - \alpha \frac{dJ}{db^{[1]}}$$

α : learning rate

Forward propagation:

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

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

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

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

Backward propagation:

$$d z^{[2]} = A^{[2]} - Y$$

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

$$d w^{[2]} = \frac{1}{m} d z^{[2]} A^{[1]T}$$

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

$$d z^{[1]} = \underbrace{w^{[2]T} d z^{[2]} * g'^{[2]}(z^{[1]})}_{\in (m^{[1]} \times n)} \rightarrow \in (n^{[1]} \times m)$$

element wise product

$$d w^{[1]} = \frac{1}{m} d f^{[1]} X^T$$

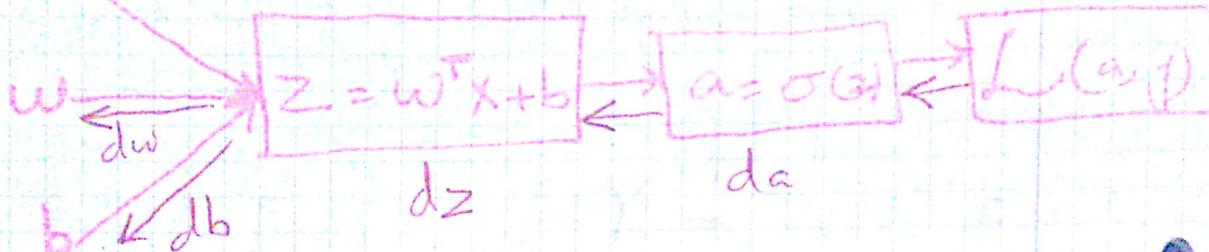
$$d b^{[1]} = \frac{1}{m} \text{np.sum}(d z^{[1]}, \text{axis}=1, \text{keepdim=True}).$$

Back propagation Intuition:-

Let's see how we derived the equation on the last page we wrote.

Logistic Regression:

X



$$\frac{d}{da} L(a, j) = -\frac{j}{a} + \frac{1-j}{1-a}$$

$$\begin{aligned} dZ &= a - j \\ dz &= da \cdot g'(z) \\ g'(z) &= \sigma'(z) \end{aligned}$$

$$dw = dz \cdot x$$

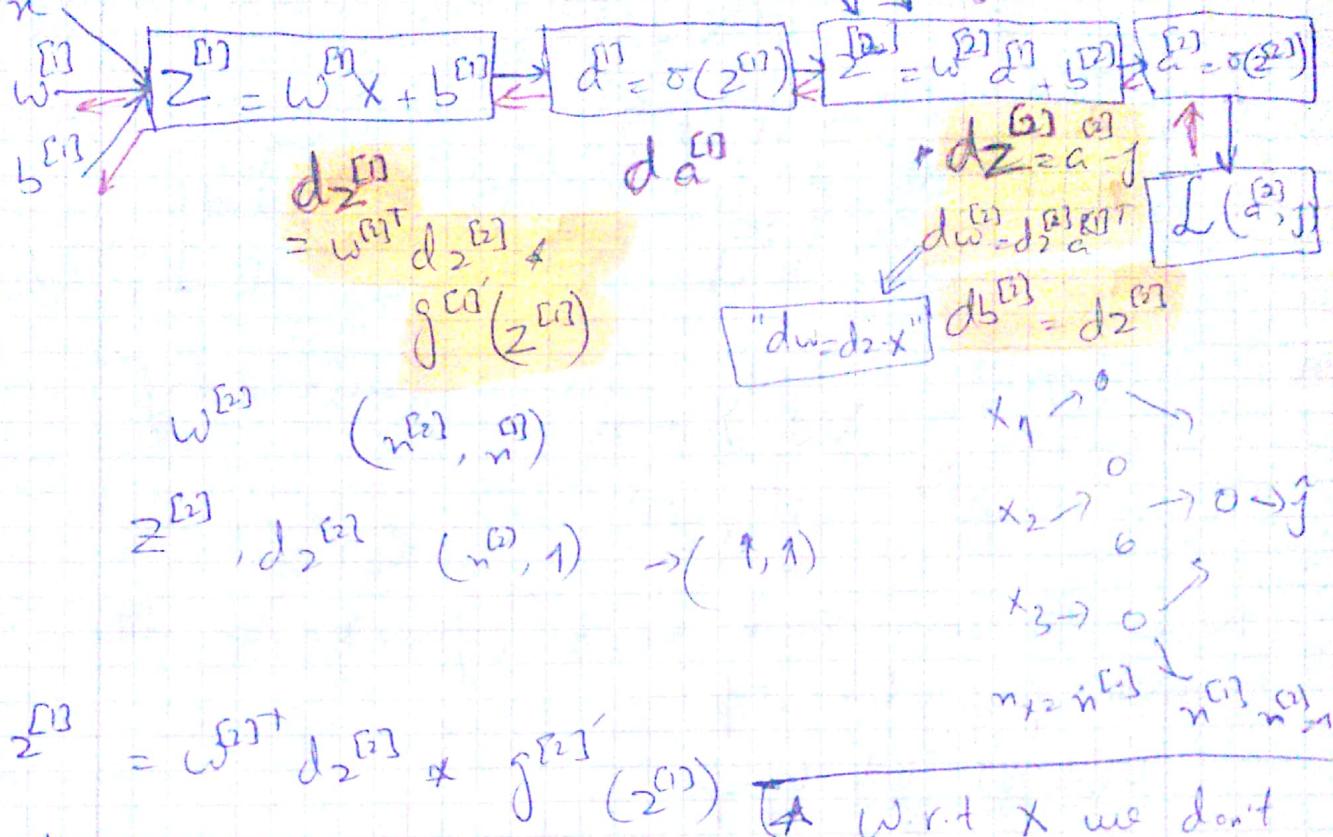
$$db = dz$$

$$\begin{aligned} \text{by chain rule} \quad \frac{dL}{dz} &= \frac{dL}{da} \cdot \frac{da}{dz} \cdot g'(z) \\ &= \left(\frac{-j}{a} + \frac{1-j}{1-a} \right) a(1-a) \\ &= -j(1-a) + (1-j)a \\ &= -j + ej + a - ej \\ \boxed{\frac{dL}{dz} = a - j} &= dz \end{aligned}$$

This is the computation for

Just one step now we are going to do it with the neural network we have defined earlier with one hidden layer that means we are going to do these computation twice.

Neural Network gradients :-



$$dZ^{[3]}$$

$$= (g^{[3]T} dZ^{[3]} * g^{[3]'}(Z^{[3]}))$$

$$dW^{[3]} = dZ^{[3]} \cdot X^T$$

$$db^{[3]} = dZ^{[3]} \cdot a^{[3]T}$$

So Summary of gradient descent:

$$dZ^{[3]} = a^{[3]} - y$$

$$dW^{[3]} = dZ^{[3]} \cdot X^T$$

$$db^{[3]} = dZ^{[3]}$$

$$dZ^{[2]} = w^{[2]T} dZ^{[3]} * g^{[2]'}(Z^{[2]})$$

$$dW^{[2]} = dZ^{[2]} \cdot X^T$$

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

w.r.t x we don't need to take derivatives because right now we are tackling a problem of Supervised learning and in these kind of problems input features remain same for particular example.

$$w^{[2]} \downarrow$$

$$b^{[2]} \downarrow$$

$$dW^{[2]} \downarrow$$

$$db^{[2]} \downarrow$$

$$da^{[2]} \downarrow$$

$$dW^{[1]} \downarrow$$

$$db^{[1]} \downarrow$$

$$da^{[1]} \downarrow$$

$$+ dZ^{[2]} = a^{[2]} - y$$

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

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

$$dW^{[0]} = dZ^{[0]} \cdot X^T$$

$$db^{[0]} = dZ^{[0]}$$

Vectorized implementations

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

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

$$z^{[1]} = \begin{bmatrix} z_1^{1} & z_1^{[1](2)} & \dots & z_1^{[1](m)} \end{bmatrix}$$

$$z^{[1]} = w^{[1]}x + b^{[1]} \quad \text{forward step.}$$

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

Now the same trick can be implemented on the back propagation step.

Vectorized B-backprop

$$dz^{[1]} = A^{[1]} - Y$$

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

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

$$dz^{[1]} = W^{[1]T} dz^{[2]} * \underbrace{g^{[1]'}(z^{[1]})}_{\text{elementwise product}}$$

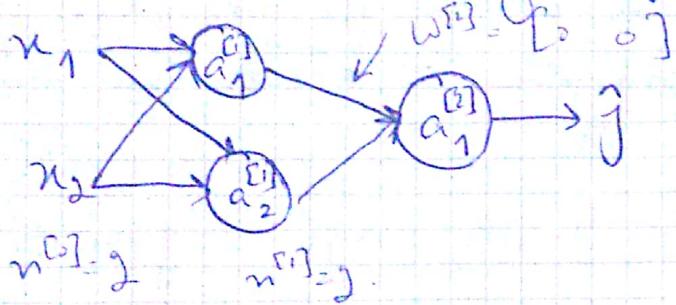
$$dW^{[1]} = \frac{1}{m} dz^{[1]} X^T$$

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

Initialization of parameters :-

- It is important to initialize your parameters randomly.
- for neural network if we initialize all parameters to zero then and then apply gradient descent, it won't work. let's find out why.

Let's consider the following:



$$w^L = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$b^{[L]} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

If we initialize weights to zero, each neuron in the first hidden layer will perform the same computation. So, even after multiple iteration each neuron in the layer will compute the same thing.

turns out initializing bias term to zero is OK! but initializing w to zero creates problems.

because:

$$a_1^{[L]} = a_2^{[L]} \text{ so in backpropagation } dZ_1^{[L]} = dZ_2^{[L]}$$

$$\text{So by symmetry } dw = \begin{bmatrix} v & v \\ v & v \end{bmatrix}$$

means $w = \begin{bmatrix} \text{---} & \text{---} \end{bmatrix} \rightarrow$ both will be same.

$$w^{[L]} = w^{[L]} - \alpha dw$$

as both are computing the same thing so there is no point in having multiple layers.

So solution to this is to initialize your parameters randomly.

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

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

$w^{[1]}$ = initialize randomly. Then input to activation function will be large so derivative will be small and optimization algorithm will be slow.

Backward Propagation Explained:

Using Chain rule to derive the derivatives of i^{th} layer in backward propagation. * Courtesy: Chirashree R

The Cost function for 1 data example is given as:

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

Total cost for all training examples :-

$$J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[l]}, b^{[l]}) = \frac{1}{m} \times \sum [L(a^{(i)}, j^{(i)})]$$

Please note that the cost J is a function of all parameters in all layers (not just the output layer) - i.e. $w^1, b^1, \dots, w^l, b^l$

$\forall i=1 \dots m$. Goal of the gradient descent is to find the values of each of these parameters using derivatives, apply to weights until it converges.

As the output has Sigmoid activation, the cost function for that is :

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

Using Chain rule of calculus:-

$$\frac{dL}{dz^{(l)}} = \frac{dL}{da^{(l)}} \times \frac{da^{(l)}}{dz^{(l)}}$$

$$\frac{dL}{da^{(l)}} = -\frac{y}{a^{(l)}} + \frac{(1-y)}{(1-a^{(l)})}$$

Simplifying this yields :-

$$\frac{dL}{da^{(l)}} = (a^{(l)} - y) / (a^{(l)}(1-a^{(l)})$$

for sigmoid activation $a^{(l)} = \frac{1}{1+e^{-z^{(l)}}}$

$$\frac{da^{(l)}}{dz^{(l)}} = \frac{\cancel{-1} \cdot e^{-z^{(l)}}}{\cancel{(1+e^{-z^{(l)}})^2}} \cdot -(1+e^{-z^{(l)}})^2 \cdot (-1) \cdot (e^{-z^{(l)}})$$

Simplifying, yields: $\boxed{\frac{da^{(l)}}{dz^{(l)}} = a^{(l)}(1-a^{(l)})}$

so Eq (1)

$$Dz^{(l)} = \frac{dL}{dz^{(l)}} = \frac{dL}{da^{(l)}} \times \frac{da^{(l)}}{dz^{(l)}} = a^{(l)} - y$$

$$\frac{dL}{dw^{(l)}} = \frac{dL}{dz^{(l)}} \times \frac{dz^{(l)}}{dw^{(l)}}$$

$$z^{(l)} = w^{(l)} \cdot a^{(l-1)} + b^{(l)}$$

$$\frac{dz^{(l)}}{dw^{(l)}} = a^{(l-1)}$$

$$\frac{dL}{dw^{(l)}} = a^{(l-1)} \cdot Dz^{(l)} = a^{(l-1)} \cdot (a^{(l)} - y)$$

$$\boxed{\frac{dL}{dw^{(l)}} = Dw^{(l)} = a^{(l-1)} Dz^{(l)}}$$

$$\frac{dL}{db^{(l)}} = \frac{dL}{dz^{(l)}} * \frac{dz^{(l)}}{db^{(l)}} \quad \because \frac{dz^{(l)}}{db^{(l)}} = 1$$

So:

$$D_L = \frac{dL}{db^{(l)}} = dz^{(l)}$$

This is for the first layer. Let's consider the 1st layer.

The key part is to apply chain rule to get the derivative of 1st parameter in terms of (1+V)th

So we apply chain rule for parameters in lth layer going through (1+V)th layer.

Applying chain rule:-

$$\frac{dL}{dw^{(l)}} = \frac{dL}{da^{(l)}} * \frac{da^{(l)}}{dz^{(l)}} * \frac{dz^{(l)}}{dz^{(l-1)}} * \dots * \frac{dz^{(l+1)}}{dz^{(l)}} * \frac{dz^{(l)}}{dw}$$

$$\text{Let } Dz^{[l+1]} = \left(\frac{dL}{da^{(l)}} \right) * \left(\frac{da^{(l)}}{dz^{(l)}} \right) * \left(\frac{dz^{(l)}}{dz^{(l-1)}} \right) * \dots * \left(\frac{dz^{(l+1)}}{dz^{(l)}} \right)$$

(as per notation)

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

Expanding for a^[l]

$$z^{[l+1]} = w^{[l+1]} * g^{[l]}(G^{[l]}) + b^{[l+1]}$$

By notation

$$Dz^{[R]} = \frac{dL}{da^{[1]}} * \frac{da^{[1]}}{dz^{[1]}} * \dots * \frac{da^{[l+1]}}{dz^{[l+1]}} * \frac{dz^{[l+1]}}{dz^{[R]}}$$

$$= Dz^{[l+1]} * \frac{dz^{[l+1]}}{dz^{[l]}}$$

$$= Dz^{[l+1]} * w^{[l+1]} * g^{[l+1]}(z^{[l]})$$

$$Dz^{[l]} = Dz^{[l+1]} * w^{[l+1]} * g^{[l+1]}(z^{[l]})$$

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

$$\frac{dz^{[l]}}{dw^{[l]}} = a^{[l-1]}$$

Therefore:

$$Dw^{[l]} = \frac{dL}{dw^{[l]}} = Dz^{[l]} * \frac{dz^{[l]}}{dw^{[l]}}$$

$$Dw^{[l]} = Dz^{[l]} * a^{[l-1]}$$

$$\frac{dL}{db^{[0]}} = \frac{dL}{dw^{[l]}} = \frac{dL}{da^{[1]}} * \frac{da^{[1]}}{dz^{[1]}} * \dots * \frac{da^{[l+1]}}{dz^{[l+1]}} * \frac{dz^{[l+1]}}{dz^{[R]}} * \frac{dz^{[R]}}{db^{[0]}}$$

$$\frac{dz^{[R]}}{db^{[0]}} = 1$$

$$\text{Therefore } Db^{[0]} = \frac{dL}{db^{[0]}} = Dz^{[R]}$$

So, for the I^{th} layer we have the following
equations -

$$Dz^{[l]} = \frac{dL}{dz^{[l]}} = Dz^{[l+1]} * w^{[l+1]}(*) g^{[l]}(z^{[l]})$$

$$Dw^{[l]} = \frac{dL}{dw^{[l]}} = Dz^{[l]} * \frac{dz^{[l]}}{dw^{[l]}}$$

$$Db^{[l]} = \frac{dL}{db^{[l]}} = Dz^{[l]}$$