Week 3 : Course 1

# One hidden layer Neural Network :



$$\begin{array}{cccc} n_1 & & & \\ n_2 & \to & \bigcirc \to \hat{y} = a & \text{(no hidden layers)} \\ n_3 & & \end{array}$$

$a^{[1]}$ . . . . $\omega^{[1]}, b^{[1]}$
$(4,3)$ , $(4,1)$

$a^{[0]} = x$

$a^{[0]} = x$

$\omega^{[2]}, b^{[2]}$
$(1,4)$ , $(1,1)$

$n_1$

$a^{[2]}_4$

$n_2 \to \bigcirc \to \hat{y} = a^{[2]}$

$n_3$

$(\hat{y} = a \text{ in Log.Regressn})$

$a^{[1]} = \begin{bmatrix} a^{[1]}_1 \\ a^{[1]}_2 \\ a^{[1]}_3 \\ a^{[1]}_4 \end{bmatrix}$

input    hidden    output (used to      "2 layer NN"
layer      layer      layer   generate $\hat{y}$)
(0)        (1)        (2)

In supervised learning, whatever calculations are
done in hidden layer are not seen during training
phase, hence the name hidden.

# Computing NN's output :



$$z^{[1]}_1 = \omega^{[1]T}_1 x + b^{[1]}_1$$

$$a^{[1]}_1 = \sigma(z^{[1]}_1)$$

$$z^{[1]}_2 = \omega^{[1]T}_2 x + b^{[1]}_2$$

$$a^{[1]}_2 = \sigma(z^{[1]}_2)$$

$[l] \leftarrow$ layer
$a_i \leftarrow$ node

$$z_1^{[1]} = \omega_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$
$$z_2^{[1]} = \omega_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$
$$z_3^{[1]} = \omega_3^{[1]} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$
$$z_4^{[1]} = \omega_4^{[1]} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

$$
Z^{[1]} =
\begin{bmatrix}
- \omega_1^{[1]T} - \\
- \omega_2^{[1]T} - \\
- \omega_3^{[1]T} - \\
- \omega_4^{[1]T} -
\end{bmatrix}
\begin{bmatrix}
n_1 \\ n_2 \\ n_3
\end{bmatrix}
+
\begin{bmatrix}
b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]}
\end{bmatrix}
=
\begin{bmatrix}
\omega_1^{[1]T} x + b_1^{[1]} \\
\omega_2^{[1]T} x + b_2^{[1]} \\
\omega_3^{[1]T} x + b_3^{[1]} \\
\omega_4^{[1]T} x + b_4^{[1]}
\end{bmatrix}
=
\begin{bmatrix}
z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]}
\end{bmatrix}
$$

$$(4 \times 3) \qquad (3 \times 1)$$

$$
a^{[1]} =
\begin{bmatrix}
a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]}
\end{bmatrix}
= \sigma(Z^{[1]})
$$

Given input $n$:

$$\text{ } z^{(1)} = w^{[1]} \overset{a^{[0]}}{\nearrow} x + b^{[1]}, \quad a^{[1]} = \sigma(z^{[1]})$$
$$(3,1) \quad (4,3)(3,1) \quad (4,1) \qquad (4,1) \qquad (4,1)$$

$$\overset{l_y}{\text{ } } \, ,$$

$$z^{[2]} = w^{[2]} \overset{a^{[1]}}{\nearrow} n + b^{[2]}, \quad a^{[2]} = \sigma(z^{[2]})$$
$$(1,1) \quad (1,4) \, (4,1) \, (1,1) \qquad (1,1) \qquad (1,1)$$

4 - 5 lined codes is all you need for a NN
with one hidden layer.

## Vectorizing across multiple examples:

For single training examples:

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

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

for 'm' training examples:

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

# $a^{[2]i}$

example $i$

layer 2

for $i = 1$ to $m$,

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}, \quad a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = W^{[2]}x^{(i)} + b^{[2]}, \quad a^{[2](i)} = \sigma(z^{[2](i)})$$

## Vectorizing across multiple examples:

$$X = \begin{bmatrix} | & | & | & & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \cdots & x^{(m)} \\ | & | & | & & | \end{bmatrix} \quad (n_x, m)$$

$$Z^{[1]} = W^{[1]} x + b^{[1]}$$
$$A^{[1]} = \sigma(Z^{[1]})$$
$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
$$A^{[2]} = \sigma(Z^{[2]})$$

} Vectorized form

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

$$A^{[1]} = \begin{bmatrix} \updownarrow & | & & \updownarrow \\ a^{[1](1)} & a^{[1](2)} & \cdots & a^{[1](m)} \\ \updownarrow & | & & | \end{bmatrix}$$

1st hidden unit on 1st training example of 1st hidden layer

# hidden units

As we move down, different activation values are stored in any particular column because there are multiple nodes for in a hidden layer

$$\longleftarrow \longrightarrow \text{ training examples}$$

$$\updownarrow \text{ hidden units}$$

Justification for vectorized implementation:

$$z^{[1](1)} = W^{[1]} x^{(1)} + b^{[1]}, \quad z^{[1](2)} = W^{[1]} x^{(2)} + b^{[1]}, \quad z^{[1](3)} = W^{[1]} x^{(3)} + b^{[1]}$$

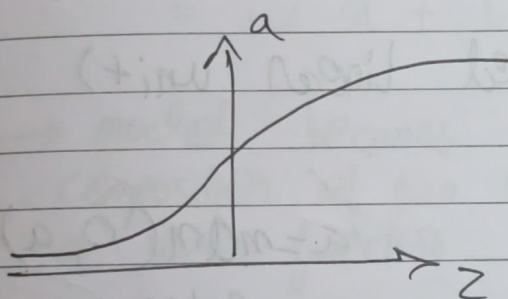$$W^{[1]} = \begin{bmatrix} \underline{\quad\quad\quad} \end{bmatrix}$$

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

$\longrightarrow$ column vectors

$$w^{[1]} \cdot \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} | & | & | \\ z^{[1](1)} & z^{[1](2)} & z^{[1](3)} \\ | & | & | \end{bmatrix}$$
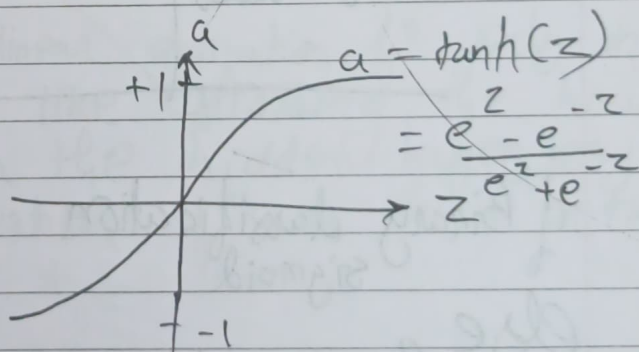
$$\underbrace{\phantom{xxxx}}_{X}$$

$$= z^{[1]}$$

## Activation function:

sigmoid is an activation $f^n$



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

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

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



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

optional

Using this function as activation $f^n$, the mean
shifts to z closer to zero hence when the
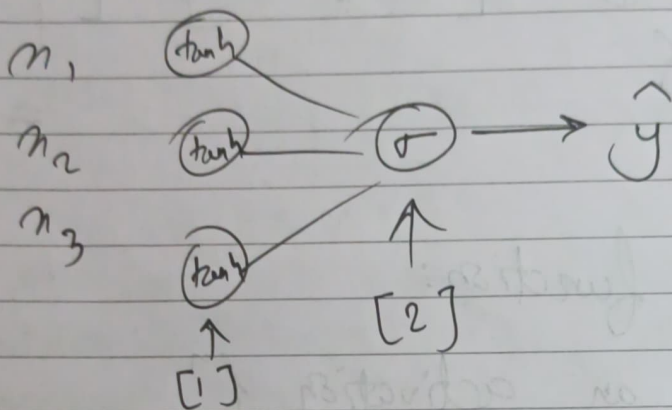algorithm requires to center the data, we would have
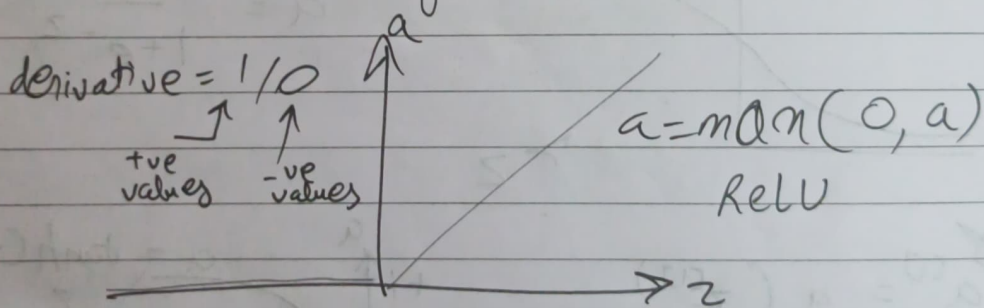zero mean.

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

When we talk about binary classification, we would want $\hat{y}$ to range between 0 & 1, for which sigmoid would be preferred as compared to ~~it~~ tanh.

Also, activation f^ns can be different for different layers like $g^{[1]}(z^{[1]}) = tanh$
$$g^{[2]}(z^{[2]}) = \sigma$$

$n_1$    (tanh)

$n_2$    (tanh) ——— ($\sigma$) ——→ $\hat{y}$

$n_3$

      (tanh)     ↑

     ↑       [2]

     [1]

✗ ReLU (Rectified linear unit)

derivative = 1/0    ↑ a

   ↗   ↑     / $a = max(0, a)$

  +ve   -ve    / ReLU
values values

            ——————→ z

if Binary classification
     sigmoid
else
   ReLU / tanh
      ↑
    more preferred because fast learning
as compared to others   as slope = 1 in most cases

Why do you need non-linear activation $f^ns$?

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

$$g(z) = z$$
"linear activation $f^n$"
/ identity activation $f^n$

$$= z^{[i]}$$

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

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

$$a^{[2]} = w^{[2]}(\underbrace{w^{[1]}x + b^{[1]}}_{a^{[1]}}) + b^{[2]}$$

$$= \underbrace{(w^{[2]} w^{[1]})}_{w'}x + \underbrace{(w^{[2]} b^{[1]} + b^{[2]})}_{b'}$$

$$= w'x + b'$$

→ model becomes useless because the composition of two linear functions is itself a linear function.

We may use linear activation $f^n$ only in output layer like estimating the house prices in which the hidden layers may contain the activation $f^ns$ as ReLU or tanh.

Derivatives of activation functions:

# sigmoid f^n

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

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

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

$$= g(z)\left(1 - g(z)\right) = a(1-a)$$

$z = 10$, $g(z) \approx 1$

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

$z = -10$, $g(z) \approx 0$

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

$z = 0$, $g(z) = \frac{1}{2}$

$$\frac{d}{dz} g(z) = \frac{1}{2}\left(1 - \frac{1}{2}\right) = \frac{1}{3}$$

# Tanh

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

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

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

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

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

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

$z = 0$, $\tanh z \approx 0$

$$g'(z) = 1$$

## # ReLU and Leaky ReLU

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

## Gradient descent for NN (1 hidden layer)

parameters: $w^{[1]}$, $b^{[1]}$, $w^{[2]}$, $b^{[2]}$

$\underset{(n^{[1]}, n^{[0]})}{} \quad \underset{(n^{[1]}, 1)}{} \quad \underset{(n^{[2]}, n^{[1]})}{} \quad \underset{(n^{[1]}, 1)}{}$

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

$n_x = n^{[0]}$ (input features)

$n^{[1]}$ (hidden units)

$n^{[2]}$ (output units)

$= 1$ (in our example)

$$= \frac{1}{m} \sum_{i=1}^{n} L(\hat{y}, y)$$

$\underset{a^{[2]}}{}$

Gradient descent :

Repeat {

$\quad\quad$ compute prediction $(\hat{y}^{(i)}, i=1,\dots, m)$

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

$$w^{[1]} := w^{[1]} - \alpha \, dw^{[1]}$$
$$b^{[1]} := b^{[1]} - \alpha \, db^{[1]}$$
$$w^{[2]} := $$
$$b^{[2]} := $$

Formulas for computing derivatives

Formulas for forward prop :

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$
$$A^{[1]} = g^{[1]} (Z^{[1]})$$
$$Z^{[2]} = W^{[2]} X + b^{[2]}$$
$$A^{[2]} = g^{[2]} (Z^{[2]})$$

$\quad\quad$ Back propogation :

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

Similar to
logistic prop: $dZ^{[2]} = A^{[2]} - Y$

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

avoids $(n^{[i]})$
vics $(n^{[i]}, 1)$

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

$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$
$(n^{[1]}, m)$     ↑ element-wise product    $(n^{[1]}, m)$

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

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

## Random Initialization



$w^{[2]}_{22} = \begin{bmatrix} 0 & 0 \end{bmatrix}$

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

$$w^{[1]} = (2, 2) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$
$(2,2)$

$$\therefore a_1^{[1]} = a_2^{[1]} \qquad dz_1^{[1]} = dz_2^{[1]}$$

i.e. the hidden units of layer 1 are exactly
identical (they are symmetric)

$$dw = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

after which, $\quad w^{[1]} = w^{[1]} - \alpha\, dw$

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

In summary, when initialized with zero,

it is actually pointless to have more than one hidden units in the hidden layer as all of them would be computing the same thing.

Solution to this is initialize the parameter randomly.

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

to make it smaller because it has to be further given to the activation functions...

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

#b can be initialized to zero as long as w is randomly initialized because it defeats the symmetry problem encountered.

$$w^{[2]} = np.random.randn((1,2)) * 0.01$$
$$b^{[2]} = 0$$

... like tanh and sigmoid so we may end 1 up with values outputs like 1 where the slope / gradient is very small ie. gradient descent will be slow. hence, learning will be slow.