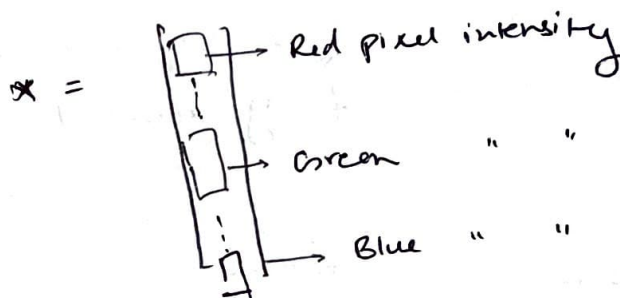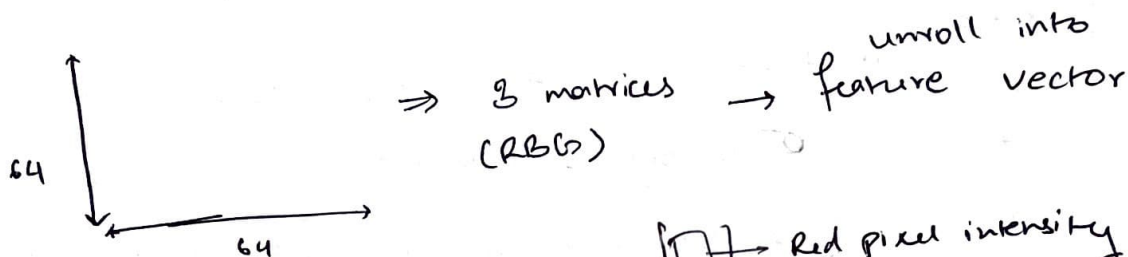THROUGHOUT THESE NOTES, I'M USING SUBSCRIPT WHERE ANDREW IS USING SUPERSCRIPT COZ I'M RACIST, ALSO, COPYRIGHT.

## CNN Notes

Binary Classification

Eg: Image → 1 (cat) vs 0 (non cat)



64

64

⟹ 3 matrices (RBG) → unroll into feature vector

$x =$

→ Red pixel intensity
→ Green " "
→ Blue " "

Dimensions of $x$:

$(64 \times 64 \times 3) \times 1 \Rightarrow 12288 \times 1$

$n = n_x = 12288$

$(x, y) \Rightarrow x \in \mathbb{R}^{n_x}, \quad y \in \{0, 1\}$

$m$ training examples: $(x^1, y^1), (x^2, y^2), \ldots\ldots, (x^m, y^m)$

$m_{train}, \quad m_{test}$

$m \to$ no. of training examples

Now define $X = \begin{bmatrix} | & | & | & & | \\ x_1 & x_2 & x_3 \cdots & x_m \\ | & | & | & & | \end{bmatrix}$

matrix $X$ is just all training examples stacked up in rows

$n_x$

$Y = [\, y_1, y_2 \cdots\cdots y_m\,]$

$X = \mathbb{R}^{n_x \times m}$

$Y = \mathbb{R}^{1 \times m}$
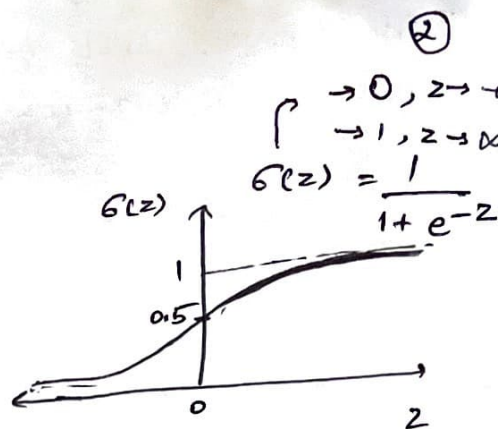
# Logistic Regression

Given : $x$, you want $\hat{y} = P(y=1|x)$

$\quad x \in \mathbb{R}^{n_x}$ parameters: $w \in \mathbb{R}^{n_x}$, $b \in \mathbb{R}$

Output $\hat{y} = w^T x + b$, but this not have range $(0,1)$.

$\Rightarrow \hat{y} = \sigma(w^T x + b)$ $\quad \sigma \Rightarrow$ sigmoid function

Our task is to have good $w$ and $b$ so that $\hat{y}$ is a very good estimation of $y$ being 1.

$$\sigma(z) \to 0, z \to -\infty$$
$$\sigma(z) \to 1, z \to \infty$$
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

---

Logistic Regression Cost Function.

Given: $\{(x_1, y_1) \cdots (x_m, y_m)\}$ we want $\hat{y}^{(i)} \approx y^{(i)}$

Loss (error) function:

$\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$ is reasonable but gradient descent goes bonkers.

$\mathcal{L}(\hat{y}, y) = -(y \log \hat{y} + (1-y)\log(1-\hat{y}))$

If $y=1$ : $\mathcal{L}(\hat{y}, y) = -\log \hat{y}$. You want small loss function so large $\hat{y}$ $\therefore \hat{y} \to 1$
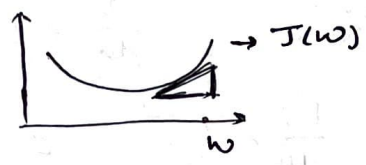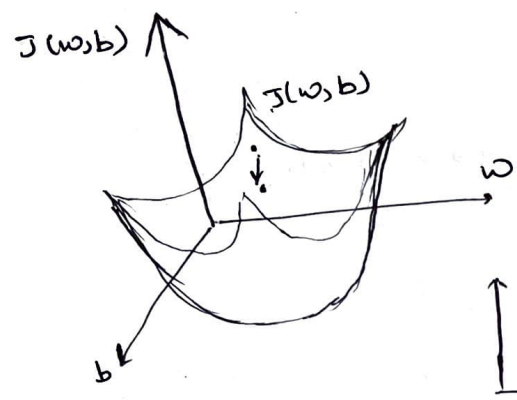
If $y=0$ : $\cdots$ $= -\log(1-y) \to \hat{y}$ small, ie $\hat{y} \to 0$

Cost function: $J(w,b) = \frac{1}{m}\sum_{i=1}^{m} \mathcal{L}(\hat{y}_i, y_i) = \frac{1}{m}\sum_{i=1}^{m}(y_i \log \hat{y}_i + (1-y_i)\log(1-\hat{y}_i))$

Now we need to find suitable $w, b$ which minimizes value of cost function $J(w,b)$.

## Gradient Descent:



$J(w,b)$

$J(w,b)$

Gradient descent takes a random $w, b$ and in iterations, moves in the steepest downward slope available.



$\rightarrow J(w)$

Repeat {

$$w := w - \alpha \frac{d\,J(w)}{dw}$$ }

Basically the intuition is just to keep approaching minima.

—
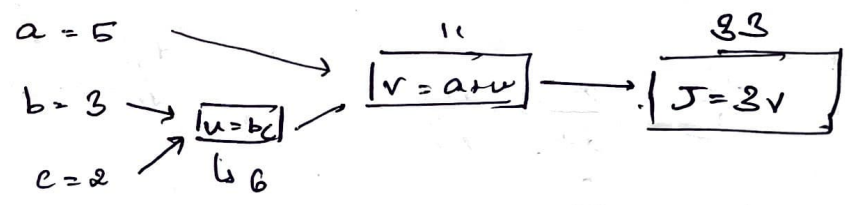
Then he explains derivatives for 17 minutes because American students are dumb
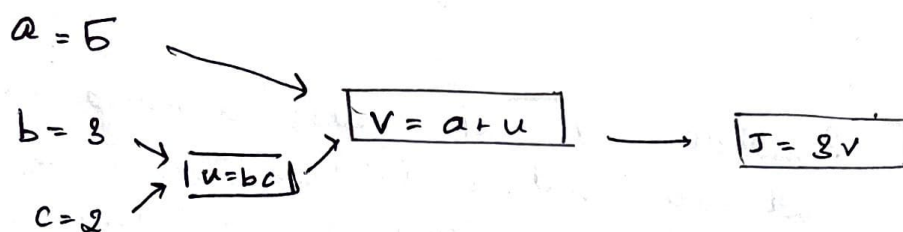
—

## Computation Graph

$J(a, b, c) = 3(a + bc)$

$u = bc \quad v = a + u \quad J = 3v$



$a = 5$

$b = 3$

$c = 2$

$u = bc$

$v = a + u$

$J = 3v$

COMPUTATION GRAPH : L → R gets output

You CAN GO BACK WARD TO GET DERIVATIVE.

$a = 5$

$b = 3$

$c = 2$

$u = bc$

$v = a + u$

$J = 3v$

Task To find $\dfrac{dJ}{da}$, $\dfrac{dJ}{db}$, $\dfrac{dJ}{dc}$ ... use chain rule,

So reverse and find $\left(\dfrac{dJ}{dv}\right)^{3} \times \left(\dfrac{dv}{da}\right)^{1}$ $\therefore \dfrac{dJ}{da} = 3$

Similarly $\dfrac{dJ}{du} = 3$, $\dfrac{dJ}{db} = \dfrac{dJ}{du} \times \dfrac{du}{db} = 3c$

## Logistic Regression Gradient Descent.

$\hat{y} = a$

$\mathcal{L}(a, y) = -(y \log a + (1 - y) \log(1 - a))$

$x_1$

$w_1$

$x_2$

$w_2$

$b$

$z = w_1 x_1 + w_2 x_2 + b$

$a = \sigma(z)$

$\mathcal{L}(a, y)$

$\dfrac{dL}{dz} = \dfrac{dL}{da} \times \dfrac{da}{dz}$

$\dfrac{d\mathcal{L}}{da} = -\dfrac{y}{a} + \dfrac{1 - y}{1 - a}$

$\phi = a - y$

$\dfrac{dL}{dw_1} = x_1 \dfrac{dL}{dz}$ ; $\dfrac{dL}{dw_2} = x_2 \dfrac{dL}{dz}$ ; $\dfrac{dL}{db} = \dfrac{dL}{dz}$

"$dw_1$"                "$dw_2$"                "$db$"

$w_1 = w_1 - \alpha \, dw_1$

$w_2 = w_2 - \alpha \, dw_2$

$b = b - \alpha \, db$

Gradient descent on m examples:

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{n} \ell(a_i, y)$$

$$a_i = \hat{y}_i = \sigma(z_i) = \sigma(w^T x_i + b)$$

$$\frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial w_1} \underbrace{\ell(a_i, y_i)}_{}$$

$$dw_i \longrightarrow (x_i, y_i)$$

for $i = 1$ to $m$

$$z_i = w^T x_i + b$$

$$a_i = \sigma(z_i)$$

$$J \mathrel{+}= -\left[ y_i \log a_i + (1 - y_i) \log(1 - a_i) \right]$$

$$\frac{d\sigma}{dz_i} = a_i - y_i$$

$$\overbrace{\phantom{xxxxxxxxxxxxxxxxx}}^{\text{for } n=2}$$

$$\frac{dJ}{dw_1} \mathrel{+}= x_{1i} dz_i \qquad \frac{dJ}{dw_2} \mathrel{+}= x_{2i} dz_i \qquad db \mathrel{+}= dz_i$$

$$\cdots \quad dw_n$$

$$w_1 = w_1 - \alpha \frac{dJ}{dw_1} \quad ; \quad w_2 = w_2 - \alpha \frac{dJ}{dw_2} \quad ; \quad w_3 = w_3 - \frac{\alpha dJ}{dw_3}$$

Vectorization

$$z = (w^T x + b$$

Vectorizing Logistic Regression.

Training examples

$$z = w^T x' + b \qquad z_2 = w^T x_2 + b \qquad \& z_3 = w^T x_3 + b$$

$$a_1 = \sigma(z_1) \qquad a_2 = \sigma(z_2) \qquad a_3 = \sigma(z_3)$$

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & --- & x_m \\ | & | & & \end{bmatrix} \quad (nx, m)$$

$$[z_1 z_2 --- z_m] = w^T X + [b \; b \; b \; -- \; b]$$
$$1 \times m.$$

$$z \overset{\rightarrow}{} [z_1 z_2 \cdots z_m] = [w_r x_1 + b \quad w_r x_2 + b \cdots w_r x_m + b]$$

$$Z = np.dot(w.T, X) + \boxed{b} \rightarrow \text{makes an appropriate}$$
matrice with all
elements b.

$$A = [a_1 \; a_2 \; --- \; a_m] = \sigma(z).$$

$$dz_1 = a_1 - y_1 \quad dz_2 = a_2 - y_2 \quad - - - -$$

$$dZ = [dz_1 \quad dz_2 \quad \cdots \quad dz_m].$$

$$1 \times m.$$

$$A = [a_1 \cdots a_m] \quad Y = [y' \cdots y^m].$$

$$dZ = A - Y. = [a_1 - y_1 \quad a_2 - y_2 \cdots \quad a_m - y_m]$$

(brace) Removed first for loop.

$$dw = 0$$
$$dw += x_1 dz_1$$
$$dw += x_2 dz_2$$
$$\vdots$$
$$dw | = m$$

$$db = 0$$
$$db += dz_1$$
$$db += dz_2$$
$$\vdots$$
$$db | = m.$$

$$db = \frac{1}{m} \sum_{i=1}^{m} dz_i \quad = \quad \frac{1}{m} \, np.\, sum(dZ)$$

$$dw = \frac{1}{m} X \, dZ^T$$

$$= \frac{1}{m} [ x_1 dz_1 + \cdots \quad x_m dz_m]$$
$$n \times m$$

## Hyper Parameters :

Parameters : $W[1]$, $b[1]$, $W[2]$, $b[2]$, $W[3]$, $b[3]$ ~~~.

Hyper parameters : Learning rate $\alpha$

       # iterations, # hidden layers $L$, # hidden units $n[1]$, $n[2]$ ...

       choice of activation function.

Later : Momentum, minibatch size, etc.

Gradient descent for neural networks:

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

$n_x = n^{[0]}, n^{[1]},$
$n^{[2]} = 1$

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

Cost function $= J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}, y)$

$\quad a^{[2]}$

Gradient descent:

For $i = 1$ to $m$:

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

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

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

- - - - - - - -

Random Initialization

w(weights) needs to be initialized randomly coz lets say

$w^{[1]} = $ zeroes $\quad a_1^{[1]} = a^{[1]}_2 \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

In this case $a^{[1]} = a_2^{[1]} \quad dz_1^{[1]} = dz_2^{[1]}$

And no matter how much you update, you'll keep computing the exact same function.

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

you keep weights small

so sigmoid/tanh doesn't tend to 1 or 0 or so

Week 4

# DEEP LAYER NEURAL NETWORK



$a_1$
$a_2$  $\rightarrow O \rightarrow \hat{y}$
$a_3$    Shallow

$x_1$
$x_2$       I hidden layer
$x_3$



$x_1$
$x_2$                    $\hat{y}$        " deep "
$x_3$

1 & 3 hidden layer

→ In

5   5   3

4 Layer NN

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

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

## Forward Propagation

$x:$ $z^{[1]} = w^{[1]}x + b^{[1]}$ ; $a^{[1]} = g^{[1]}(z^{[1]})$

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

$\vdots$

$z^{[4]} = w^{[4]}a^{[3]} + b^{[4]}$ ; $a^{[4]} = g^{[4]}(z^{[4]}) = \hat{y}$

Generic: $z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]}$
   $a^{[l]} = g^{[l]}(z^{[l]})$

Vectorized:  $Z^{[1]} = w^{[1]}A^{[0]} + b^{[1]}$
   $A^{[1]} = g^{[1]}(Z^{[1]})$
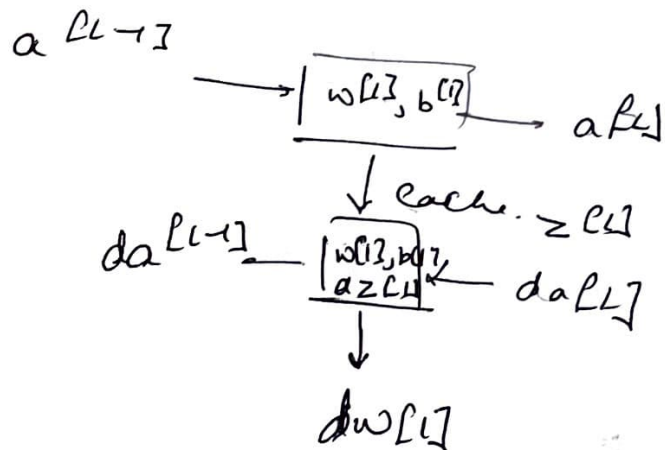   $Z^{[2]} = w^{[2]}A^{[1]} + b^{[2]}$
   $\hat{y} = g(Z^{[4]}) = A^{[4]}$

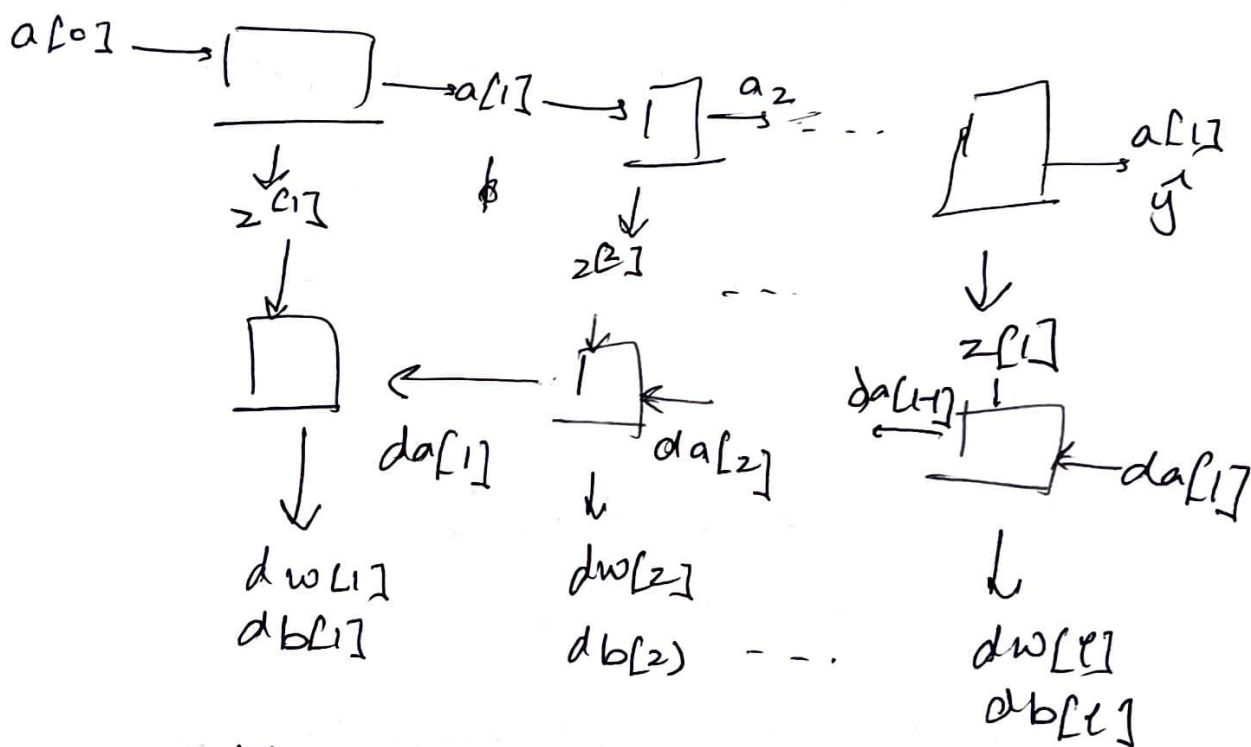for $l = 1 \cdots 4$
There is no way
to remove
this for loop

(ii)

a   Matrice dimension → get them right.

<u>backward functions</u> :

$a^{[L-1]}$

$\longrightarrow$ $\boxed{w^{[l]}, b^{[l]}}$ $\longrightarrow$ $a^{[l]}$

$da^{[L-1]}$ $\longleftarrow$ $\boxed{\begin{array}{c} w^{[l]}, b^{[l]} \\ dz^{[l]} \end{array}}$ $\longleftarrow$ $da^{[l]}$

$\downarrow$ cache $z^{[l]}$

$\downarrow$

$dw^{[l]}$

element wise multiplication $\uparrow$

$$dz^{[l]} = da^{[l]} \circledast g^{[l]'}(z^{[l]})$$
$$dw^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$
$$db^{[l]} = dz^{[l]}$$
$$da^{[l]} = w^{[l]T} \cdot dz^{[l]}$$

$a^{[0]}$ $\longrightarrow$ $\boxed{\phantom{xx}}$ $\longrightarrow$ $a^{[1]}$ $\longrightarrow$ $\boxed{\phantom{x}}$ $\xrightarrow{a_2}$ $\cdots$ $\boxed{\phantom{x}}$ $\longrightarrow$ $a^{[l]}$ $\hat{y}$

$\downarrow z^{[1]}$ $\qquad$ $b$ $\qquad$ $\downarrow$ $z^{[2]}$ $\qquad$ $\downarrow$ $z^{[l]}$

$\boxed{\phantom{x}}$ $\longleftarrow$ $\boxed{\phantom{x}}$ $\qquad$ $da^{[l-1]} \boxed{\phantom{x}}$

$\downarrow$ $da^{[1]}$ $\qquad$ $da^{[2]}$ $\qquad$ $\longleftarrow da^{[l]}$

$\downarrow$ $\qquad$ $\downarrow$ $\qquad$ $\downarrow$

$dw^{[1]}$ $\qquad$ $dw^{[2]}$ $\qquad$ $dw^{[l]}$

$db^{[1]}$ $\qquad$ $db^{[2]}$ $\quad \cdots$ $\qquad$ $db^{[l]}$

$$w^{[l]} = w^{[l]} - \alpha \, dw^{[l]}$$
$$b^{[l]} = b^{[l]} - \alpha \, db^{[l]}$$

# What are Neural Networks?



$a^{[0]} = x$     $a^{[1]}$     $w^{[1]}, b^{[1]}$
    $[4,3]$    $(4,1)$    → 2 Layer Neural Network.

→ Input Layer

Hidden Layer

↳ Output Layer.

$a_i^{[l]}$ → layer
↳ node in layer

$$z_i^{[1]} = w_i^{[1]T} x + b_i^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

Similarly 4 nodes

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

writing a for loop here is inneficient.

$$\begin{bmatrix} — w_1^{[1]T} — \\ — w_2^{[1]T} — \\ — w_3^{[1]T} — \\ — w_4^{[1]T} — \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4 \end{bmatrix}$$

II

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

⇒ Given input x

$$z^{[1]} = W_{\times}^{[1]} a^{[0]} + b^{[1]} \qquad a^{[1]} = \sigma(z^{[1]})$$

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

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

Multiple ~~exempt~~ training examples.

$$\begin{aligned} X &\longrightarrow a^{[2]} = \hat{y} \\ X^{(1)} &\longrightarrow a^{(2)}(1) = \hat{y}^{(1)} \\ &\vdots \\ X^{(m)} &\longrightarrow a^{[2](m)} = \hat{y}^{(m)} \end{aligned} \quad \Bigg\} \quad m \text{ training examples.}$$

for i = 1 to m

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

$$\Bigg\} \quad \text{add (i) to all training examples.}$$

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \cdots & x_m \\ | & | & & | \end{bmatrix} \quad (n_x, m)$$

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

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$
$$A^{[1]} = \sigma(z^{[1]})$$
$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$
$$A^{[2]} = \sigma(z^{[2]})$$
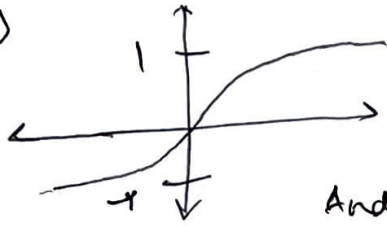
$$A^{[1]} = \begin{bmatrix} | & | & & | \\ a^{[1](1)} & a^{[1](2)} & \cdots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

## Activation function

In $\quad a = \sigma(z) \qquad \sigma(z) = \dfrac{1}{1 + e^{-z}}$

sigmoid here is called the activation function.

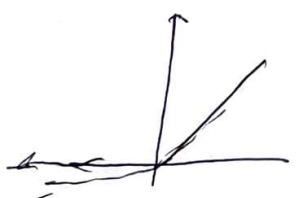## More activation functions.

1)

$a = \tanh(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ : generally Much better So that $CO_2$ mean $= 0$

Andrew loves tanh and says it is mucch more superior except in finale layer as sigmoid fn can be interpreted as a Probability.

2) Rectified linear unit (ReLU function)

$a = \max(0, z)$

→

Generally : In Binary classification systems, sigmoid is preflred. and cxe if you don't know what to do, use ReLU.

} → Leaky Relu ($a = \max(0.01z, z)$ or $(0.001 z, z)$ you get the hint.

# IV

Gradient descent for neural networks:

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

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

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

Cost function $= J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \mathcal{L}(\hat{y}, y)$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \{a^{[2]}$

Gradient descent:

for $i = 1$ to $m$:

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

$\qquad dw^{[1]} = \dfrac{dJ}{dw^{[1]}} \quad , \quad db^{[1]} = \dfrac{dJ}{db^{[1]}} , \; - - -$

$\qquad w^{[1]} = w^{[1]} - \alpha \, dw^{[1]}$

$\qquad b^{[1]} = b^{[1]} - \alpha \, db^{[1]}$

- - - - - - -

Random Initialization

$w$ (weights) needs to be initialized randomly coz lets say

$\underline{w^{[1]} = \text{zeroes}} \quad a_1^{[1]} = a_2^{[1]} \qquad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

In this case $a_1^{[1]} = a_2^{[1]} \quad dz_1^{[1]} = dz_2^{[1]}$

And no matter how much you update, you'll keep computing the exact same function.

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

$\qquad b^{[1]} = $ np. @ zeroes

$\qquad w^{[2]} = $ np. random. $* 0.01$

$\qquad b^{[2]} = $ zeroes

you keep weights small

so sigmoid / tanh doesn't tend to 1 or 0 ~~arose~~