# Week 2 : Basics of NN programming

## Binary classification:
(Exa: Cat (1) or noncat (0))



$64 \times 64 \times 3 = 12288$

$n_x = 12288$

$y = \{1, 0\}$

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

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

$M_{train}$ — no. of training examples

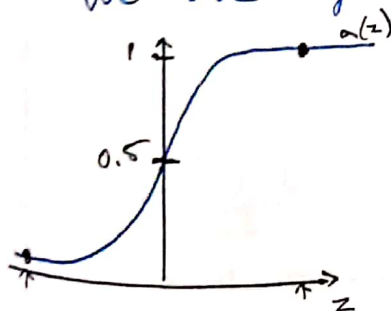$M_{test}$ — no. of test examples

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \ldots & x^{(m)} \\ | & | & & | \end{bmatrix} \qquad X \in \mathbb{R}^{n_x \times m}$$

$$Y = \begin{bmatrix} y^{(1)}, & y^{(2)} \ldots & y^{(m)} \end{bmatrix} \qquad Y \in \mathbb{R}^{1 \times m}$$

## Logistic regression:

Given $x$, we want $\hat{y} = P(y = 1 / x)$ where $0 \leq \hat{y} \leq 1$

∴ We use sigmoid function $\hat{y} = \sigma(w^T x + b)$



$\sigma(z) = \dfrac{1}{1 + e^{-z}}$

If $z$ large, $\sigma(z) \approx \dfrac{1}{1+0} = 1$

If $z$ large negative number,
$\sigma(z) \approx \dfrac{1}{1 + \text{Big no}} \approx 0$

Notation :

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \quad \begin{matrix} \} & b \in \mathbb{R} \\ \\ \} & w \in \mathbb{R}^{n_x} \end{matrix}$$

real number →

↑ $n_x$ dimension vector

Logistic regression cost function:

$\hat{y}^{(i)} = \alpha(\omega^T x^{(i)} + b)$ where $\alpha(z^{(i)}) = \dfrac{1}{1+e^{-z^{(i)}}}$

Given $\{(x^{(1)}, y^{(1)}), \ldots, (x^{(m)}, y^{(m)})\}$, want $\hat{y}^{(i)} \approx y^{(i)}$

prediction ↑ real value ↑

∴ Loss (error) function: - for 1 example

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

If $y=1$, $\mathcal{L}(\hat{y}, y) = -\log \hat{y}$ ← we want $\log \hat{y}$ large ∴ $\hat{y}$ large

$y=0$, $\mathcal{L}(\hat{y}, y) = -\log(1-\hat{y})$ ← we want $\log 1-\hat{y}$ large ∴ $\hat{y}$ small

∴ Cost function: - for m examples

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right]$$
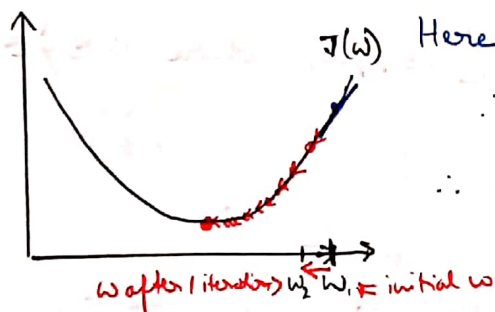
## Gradient Descent:

Repeat (till convergence) {

$\omega := \omega - \alpha \underbrace{\left( \dfrac{dJ(\omega)}{d\omega} \right)}$ ← This can be written as dw while coding and for b, db

↑ learning rate

}



$J(\omega)$ Here Slope is +ve ∴ $\dfrac{dJ(\omega)}{d\omega} > 0$

∴ $\omega := \omega - \alpha \cdot$ something

∴ $\omega$ reduces

$\omega$ after 1 iteration $\omega_2$ $\omega_1$ ← initial $\omega$

Here slope is -ve ∴ $\dfrac{dJ(\omega)}{d\omega} < 0$

$J(\omega)$

∴ $\omega := \omega - \alpha \cdot -$ something

$= \omega := \omega + \alpha \cdot$ something

∴ $\omega$ increases

$\omega_1 \rightarrow \omega_2$

For $J(\omega, b)$,
Repeat {
$\omega := \omega - \alpha \dfrac{\partial J(\omega, b)}{d\omega}$ ←

$b := b - \alpha \dfrac{\partial J(\omega, b)}{db}$
}

we use $\partial$ (partial derivation) instead of d when there are more than 1 variables (Exa: here there is $\omega$ and $b$)

.3

# Derivatives :



$f(a) = 6.0006$ ~~(struck out)~~

when $a = 2$, $f(a) = 6$

$a = 2.0001$, $f(a) = 6.0003$

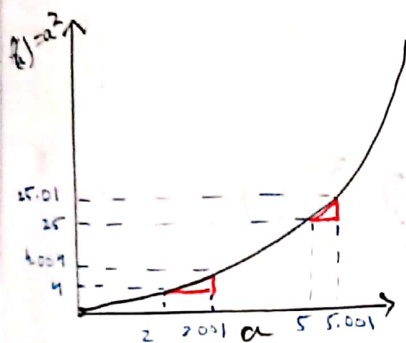Slope/derivative of $f(a)$

$= \dfrac{\text{height}}{\text{width}} = \dfrac{0.003}{0.001} = 3$

$\therefore \dfrac{df(a)}{da}$ or $\boxed{\dfrac{d}{da} f(a) = 3 \text{ constant}}$

$\therefore \dfrac{d f(a)}{da} = 3a$

**Note :** Here slope remains constant (always $\in$ 3)

*remember derivative of $3a$ was 3?*

① when $a = 2$, $f(a) = 4$

$a = 2.001$, $f(a) = 4.004$

$\therefore$ Slope/derivative of $f(a)$ at $a = 2$ is 4

$\dfrac{d}{da} f(a) = 4$ when $a = 2$

② when $a = 5$, $f(a) = 25$

$a = 5.001$, $f(a) = 25.01$

$\therefore$ slope/derivative of $f(a)$ at $a = 5$ is 10

$\dfrac{d}{da} f(a) = 10$ when $a = 5$



$\therefore \dfrac{d}{da} f(a) = \boxed{\dfrac{d}{da} a^2 = 2a}$

## More examples:

$f(a) = a^3 \qquad \dfrac{d}{da} f(a) = 3a^2 \quad \Leftarrow$

$a = 2 \qquad f(a) = 8$

$a = 2.001 \qquad f(a) = 8.012$

$3 \times 2^2 = 12 \qquad \Delta = 12$

$f(a) = \log_e(a) \qquad \dfrac{d}{da} f(a) = \dfrac{1}{a} \quad \Leftarrow$

or

$\ln(a)$

$a = 2 \qquad f(a) = 0.69315$

$a = 2.001 \quad f(a) = 0.69365$

$\dfrac{1}{2} = 0.0005 \qquad \Delta = 0.0005$
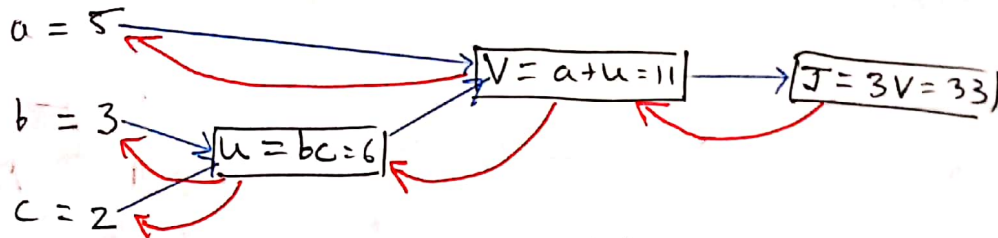
# Computation Graph:

$$J(a,b,c) = 3(\underbrace{a + \underbrace{bc}_{u}}_{v})$$

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

$a = 5$

$b = 3$    $\boxed{u = bc = 6}$

$c = 2$    $\boxed{V = a+u = 11} \longrightarrow \boxed{J = 3V = 33}$

- Left-to-right pass
- Right-to-left pass ← To find derivative

## Finding derivatives:

$a = 5$
$da = 3$

$b = 3$
$db = 6$

$c = 2$
$dc = 9$

$\boxed{u = bc}$   $du = 3$

$\boxed{V = a+u} \longrightarrow \boxed{J = 3V}$    $dv = 3$

(1) $\dfrac{dJ}{dv} = ?$

In coding notation this is dv

$J = 3V$
$V = 11 \rightarrow 11.001$
$J = 33 \rightarrow 33.003$

$\therefore \dfrac{dJ}{dV} = 3$

← Similar to $f(a) = 3a$

$\dfrac{df(a)}{da} = 3$ from previous example

(2) $\dfrac{dJ}{da} = ?$

In coding notation this is da

$a = 5 \rightarrow 5.001$
$V = 11 \rightarrow 11.001$
$J = 33 \rightarrow 33.001$

$\therefore \dfrac{dJ}{da} = 3 = \underbrace{\dfrac{dJ}{dV}}_{3} \cdot \dfrac{dV}{da}$

we found this in step (1)   $1 = 3$

**Chain rule.**

If $a$ affects $V$ affects $J$,

Then $\dfrac{dJ}{da}$ ← effect of $a$ on $J$

$= \dfrac{dJ}{dV} \cdot \dfrac{dV}{da}$ — effect of $a$ on $v$

effect of $V$ on $J$

(3) $\dfrac{dJ}{du} = \dfrac{dJ}{dV} \cdot \dfrac{dV}{du} = 3 \cdot 1 = 3$

$u = 6 \rightarrow 6.001$
$V = 11 \rightarrow 11.001$
$J = 33 \rightarrow 33.003$

(4) $\dfrac{dJ}{db} = \dfrac{dJ}{du} \cdot \dfrac{du}{db} = 3 \cdot 2 = 6$

$b = 3 \rightarrow 3.001$
$u = b \cdot c \rightarrow 6 \rightarrow 6.002$   $c=2$ here only b is ↑

$3.001 \times 2$

(5) $\dfrac{dJ}{dc} = \dfrac{dJ}{du} \cdot \dfrac{du}{dc} = 3 \cdot 3 = 9$

$c = 2 \rightarrow 2.001$
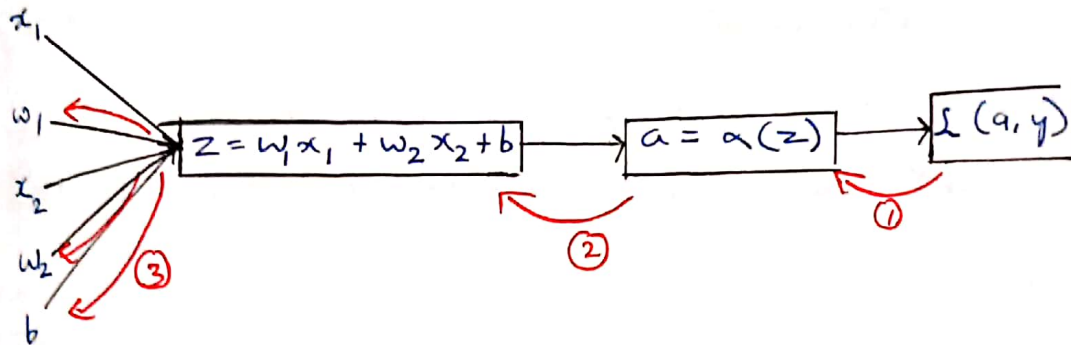$u = b \cdot c = (\rightarrow 6.003 \leftarrow 2.001 \times 3$

# Logistic Regression Gradient Descent:

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

$$\hat{y} = a = a(z)$$

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



① $da = \dfrac{dL(a,y)}{da} = -\dfrac{y}{a} + \dfrac{1-y}{1-a}$ //

② $\underset{\underset{\frac{dL}{dz}}{\wedge}}{dz} = \dfrac{dL}{da} \cdot \dfrac{da}{dz} \quad \overset{a(1-a)}{\swarrow} \quad = a-y$ //

③ $dw_1 = \dfrac{dL}{dw_1} = x_1 \, dz$

$dw_2 = \dfrac{dL}{dw_2} = x_2 \, dz$

$db = \dfrac{dL}{db} = dz$

∴ Gradient descent (one step):

$w_1 := w_1 - \alpha \, dw_1$

$w_2 := w_2 - \alpha \, dw_2$

$b := b - \alpha \, db$

# Gradient Descent for m examples:

The cost function for logistic regression is,

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(a^{(i)}, y)$$

where $a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$

In the previous examples we find $dw_1^{(i)}$, $dw_2^{(i)}$, $db^{(i)}$ for a single training example $(x^{(i)}, y^{(i)})$

$$\frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)})}$$

$$\underbrace{\hspace{3cm}}_{dw_1^{(i)} \text{ for } (x^{(i)}, y^{(i)})}$$

## Algorithm:

$J = 0$; $dw_1 = 0$; $dw_2 = 0$; $db = 0$

For $i = 1$ to $M$

Forward propagation $\begin{cases} z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(z^{(i)}) \\ J \mathrel{+}= -\left[ y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)}) \right] \end{cases}$

Backward propagation $\begin{cases} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw_1 \mathrel{+}= x_1^{(i)} dz^{(i)} \\ dw_2 \mathrel{+}= x_2^{(i)} dz^{(i)} \\ db \mathrel{+}= dz^{(i)} \end{cases}$ $\Big\}$ For $n=2$. otherwise we would find $dw_1, dw_2$ $dw_3, dw_4 \ldots dw_n$

$J \mathrel{/}= m$

$dw_1 \mathrel{/}= m$

$dw_2 \mathrel{/}= m$

$db \mathrel{/}= m$

Using for loops is inefficient, so we use vectorization.

# Vectorization:

We need to compute $z = \omega^T x + b$

where $\omega = \begin{bmatrix} : \\ : \end{bmatrix}$  $x = \begin{bmatrix} : \\ : \end{bmatrix}$   $\omega \in \mathbb{R}^{n_x}$

$x \in \mathbb{R}^{n_x}$

| Non-vectorized | Vectorized |
|---|---|

**Non-vectorized**

```
z = 0
for i in range (n-x):
    z += w[i] * x[i]

z += b
```

**Vectorized**

$z = np.dot(\omega, x) + b$

$\underbrace{\phantom{np.dot(\omega,x)}}$

$\omega^T x$

## Code

```
import numpy as np

a = np.array ([1, 2, 3, 4])   — initialise array

u = np.exp (v)
```
— you can use np functions that'll take effect on all the vector

np.log (v)  np.abs(v)  np.maximum (v,0) values
. . . .

→ Here we will need a loop, instead we can use vectorization

$$dw += x^{(i)} dz^{(i)}$$

However even here, we use a for loop to iterate through all the training examples. Instead we can vectorize

$[z^{(1)} z^{(2)} ... z^{(m)}] \longrightarrow z = np.dot(wT, x) + b$

↑ with broadcasting it'll become

$[x^{(1)} \underset{1}{|} x^{(2)} ... \underset{1}{|} z^{(m)}]$

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

$[b\ b ... b]$

$\mathbb{R}^{1 \times m}$

$A = \alpha(z)$

$[a^{(1)} a^{(2)} ... a^{(m)}]$

$$dz = A - Y \leftarrow [y^{(1)} \dots y^{(m)}]$$

$$[dz^{(1)} \; dz^{(2)} \dots dz^{(m)}] \qquad [a^{(1)} \; a^{(2)} \dots a^{(m)}]$$

$$db = \frac{1}{m} \; np.sum(dz)$$

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

After that we update $w$ and $b$:

$$w := w - \alpha \, dw$$

$$b := b - \alpha \, db$$

Full algorithm:

for iter in range (1000):    $\leftarrow$ running gradient
descent 1000 times
(we can't vectorise this)

$\quad z = np.dot(w.T, X) + b$

$\quad A = \sigma(z)$

$\quad dz = A - Y$

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

$\quad db = \frac{1}{m} np.sum(dz)$

$\quad\quad w := w - \alpha \, dw$

$\quad\quad b := b - \alpha \, db$

# Broadcasting in Python

Python automatically rescales variables to fit the array/vector calculation

① $\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \Rightarrow \begin{bmatrix} 1 \\ 2 \\ 3 \\ n \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$

② $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$

③ $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$

# Common Bugs

① While initialising vectors, declare both the no. of rows & columns.

$a = np.randoms.randu \; (5,1) \quad$ or $(1,5)$

column vector $\quad$ row vector

If you do a = np.random.randn(5), then a won't work and np.dot(a, aᵀ) won't work as expected

② You can use $a = a.reshape((5,1))$.

# Derivation for Logistic Regression Cost Function (Optional)

If $y = 1$, $P(y|x) = \hat{y}$ — ①

If $y = 0$, $P(y|x) = 1 - \hat{y}$ — ②

$$P(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

↖ This equation is verified

If $y=1$, $P(y|x) = \hat{y}^1 (1-\hat{y})^0 = \hat{y}$ — ①

If $y=0$, $P(y|x) = \hat{y}^0 (1-\hat{y})^1 = 1-\hat{y}$ — ②

maximise this
$$\uparrow \log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)}$$
$$= y \log \hat{y} + (1-y) \log (1-\hat{y})$$
$$= -\mathcal{L}(\hat{y}, y) \; \downarrow \; \text{To minimise loss}$$
(↓ since we added −)

$$\boxed{J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})}$$

For m training examples