# I  Binary Classification

* Image is processed in the form of matrices for 3 colours red, green and blue.

* Notations:

$m$ - no. of training vectors
$n_x$ - size of input vector
$n_y$ - size of output vector
$x^{(1)}$ - size of input vector
$y^{(1)}$ - size of output vector

$$x \in R^{n_{(x)}}$$
$$y \in (0, 1)$$

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

$$\underleftrightarrow{\quad m \quad}$$

$$X \in R^{n_x \times m}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots & y^{(m)} \end{bmatrix}$$

$$Y \in R^{1 \times m}$$

## II Logistic Regression -

* Algorithm used to classify algorithm of 2 classes

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

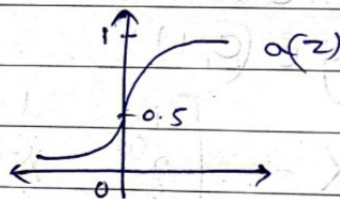* Parameters: $w \in \mathbb{R}^{n_x}$
  $b \in \mathbb{R}$

* $y = wx + b$

* If $x$ is a vector,
  $$y = w(transpose)x + b$$
  $$= w^T x + b$$
  between $(0 \& 1)$

* $\hat{y} = output = \underbrace{a(w^T x + b)}_{z}$



$y \stackrel{?}{=} \frac{1}{\cancel{\infty}}$ : $a(z) = \frac{1}{1+e^{-z}}$

## III) Logistic Regression Cost function

loss function - calculates error for single training example
cost function - average of loss function over entire training set

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

$y = 1$, $L = -\log \hat{y}$ ∴ $\hat{y}$ should be large
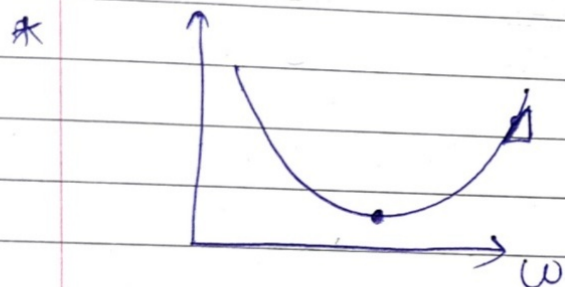not possible as $\hat{y}$ $(0,1)$

$y = 0$, $L = -\log(1-\hat{y})$ ∴ $\log(1-\hat{y}) \Rightarrow$ large
∴ $\hat{y} \Rightarrow$ small

$$J(w, b) = \frac{1}{m} \sum_{i=1}^{M} \left( \hat{y}^{(i)}, y^{(i)} \right) = \frac{-1}{m} \sum_{i=1}^{m} \left( y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)}) \right)$$

## (IV) Gradient Descent

* Goal: To find $w, b$ that minimise $J(w, b)$

* Our cost function is convex & we initialise $w, b$ to 0 (0,0) ⟵ preferable and move steeply downwards

*



$$w = w - \alpha \frac{d J(w)}{dw} \quad \leftarrow \text{'dw'}$$

learning rate

$$w = w - \alpha dw \quad - \text{measure of how much we slope towards } w$$

$$w = w - \alpha \boxed{\frac{d J(w, b)}{dw}} \qquad = w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b = b - \alpha \boxed{\frac{d J(w, b)}{db}} \qquad = w - \alpha \frac{\partial J(w, b)}{\partial b}$$

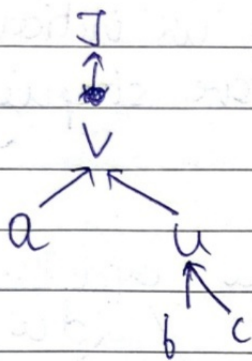## (V) Derivatives ✓

## (VI) More Derivative Examples ✓

## (VII) Computation Graph

It organises computation from left to right.

3

## (VIII) Derivatives with a Computation Graph

dvar — derivative of final output w.r.t intermediate quantities

Computation: R→L

$$x \to y \to z$$

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

## (IX) Logistic regression gradient descent

$$z = w^T x + b$$
$$\hat{y} = a = a(z)$$
$$L(a,y) = -\left[ y \log a + (1-y) \log (1-a) \right]$$

$$z = w_1 x_1 + w_2 x_2 + b \to a = a(z) \to L(a,y)$$

$$dz = \frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} \qquad da = \frac{dL(a,y)}{da}$$

$$= a - y$$

$$= \frac{-y}{a} + \frac{-1-y}{1-a}$$

(X) Logistic regression on m example

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

$$= \frac{1}{m} \sum_{i=1}^{m} \frac{\partial}{\partial \omega_1} \mathcal{L}\left(a^{(i)}, y^{(i)}\right)$$

$x_1$ - feature
$x_2$ - feature
$\omega_1$ - weight of first feature
$\omega_2$ - weight of second feature
$b$ - logistic regression parameter
$M$ - number of training examples
$y(i)$ - expected output of $i$

Algo:

$$J = 0 \quad ; \quad d\omega_1 = 0 \quad ; \quad d\omega_2 = 0 \quad ; \quad db = 0$$

For loop ①:

for $i = 1$ to $m$

$$z^{(i)} = \omega^T x^{(i)} + b$$

$$a^{(i)} = g\left(z^{(i)}\right)$$

$$J += \left[y^{(i)} \log a^{(i)} + 1 - y^{(i)} \log\left(1 - a^{(i)}\right)\right]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

for loop ②
$$\begin{cases} d\omega_1 += x_1^{(i)} dz^{(i)} \\ d\omega_2 += x_2^{(i)} dz^{(i)} \end{cases} \quad n = 2$$

$$db += dz^{(i)}$$

$$J /= m$$

$$d\omega_1 /= m \qquad = \frac{\partial J}{\partial m} \qquad (\text{accumulator})$$

$$w_1 = w_1 - \alpha dw_1$$
$$w_2 = w_2 - \alpha dw_2$$
$$b = b - \alpha db$$

for loops - makes code less efficient - code takes long time

vectorisation - speeds up code - gets rid of for loops

## Vectorisation

* ~~Atomy~~ Numpy library uses vectorisation by default.

* Vectorisation can be done on CPU/ GPU ← faster

Logistic Regression pseudo code

$J = 0;$              $w_1 = 0$

$dw1 = 0$           $w_2 = 0$

$dw_2 = 0$          $b = 0$

$db = 0$

$$dw = np.zeroes(n-x, 1)$$

$w_1 = 0$ ; $w_2 = 0$ ; $b = 0$

forward pass {

for $i = 1$ to m

$\quad z(i) = w1^* x_1(i) + w_2^* x_2(i) + b$

$\quad a(i) = \alpha\, z(i)$

$\quad J \mathrel{+}= (y(i)^* \log a(i) + 1 - y(i)\log(1 - a(i)))$

# backward pass

$dz^{(i)} = np.zeros((n-x, 1))$

$$dz(i) = a(i) - y(i)$$
$$dw_1 \mathrel{+}= dz(i) * x_1(i)$$
$$dw_2 \mathrel{+}= dz(i) * x_2(i)$$
$$db \mathrel{+}= dz(i)$$

$\Bigg\}$  $\boxed{dw \mathrel{+}= x^{(i)} z^{(i)}}$

$$J \mathrel{/}= m$$
$$dw_1 \mathrel{/}= m$$
$$dw_2 \mathrel{/}= m$$
$$db \mathrel{/}= m$$

$\Bigg\} \rightarrow \boxed{dw \mathrel{/}= m}$

# Gradient descent

$$w_1 = w_1 - alpha * dw_1$$
$$w_2 = w_2 - alpha * dw_2$$
$$b = b - alpha * db$$

$\Rightarrow$ Vectorising logistic regression

$z^{(1)} = w^T x^{(1)} + b \qquad z^{(2)} = w^T x^{(2)} + b \qquad z^{(3)} = w^T x^{(3)} + b$

$a^{(1)} = \sigma(z^{(1)}) \qquad a^{(2)} = \sigma(z^{(2)}) \qquad a^{(3)} = \sigma(z^{(3)})$

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

$$w^T \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$Z = \begin{bmatrix} z^{(1)} & z^{(2)} & & z^{(m)} \end{bmatrix}$$
$$= \begin{bmatrix} w^T X + [b \ b \ \cdots \ b] \end{bmatrix}$$
$$= \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(2)} + b & \cdots & w^T x^{(m)} + b \end{bmatrix}$$

7

$$Z = np.dot(\omega^7, x + b)$$

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \dots & a^{(m)} \end{bmatrix} = \boxed{a(z)}$$

$$dz^{(1)} = a^{(1)} - y^{(1)}$$
$$dz^{(2)} = a^{2} - y^{(2)}$$

$$dZ = \begin{bmatrix} dz^{(1)} & dz^{(2)} & \dots & dz^{m} \end{bmatrix}$$

$$A = \begin{bmatrix} a^{(1)} & \dots & a^{(m)} \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix}$$

$$\boxed{dz = A - Y}$$
$$= \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots \end{bmatrix}$$

$\cancel{dw=0}$  $db = 0$
$db + = dz^{(1)}$
$db + = dz^{(2)}$                    $db = \dfrac{1}{m} \sum\limits_{i=1}^{m} dz^{(i)}$
$\quad\quad \vdots$
$db += dz^{(m)}$                     $= \dfrac{1}{m} npsum(dz)$
$db \mathrel{/}=m$

$dw = 0$
$dw + = x^{(1)} dz^{(1)}$                $dw = \dfrac{1}{m} X dz^{7}$
$\quad\quad \vdots$
$dw += x^{m} dz^{(m)}$             $= \dfrac{1}{m} \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ & & & \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$
$dw \mathrel{/}= m$

$$= \dfrac{1}{m} \left( x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)} \right)$$

$$\omega = \omega - \alpha \, dw$$
$$b = b - \alpha \, db$$

8