

m → No. of training examples

* Sigmoid to ReLU

Subject _____ Date _____

structured & unstructured data.

WEEK 2: Binary classification

Training set of 'm' training examples, process entire training set without using for loop.

① Forward propagation

② Backward propagation step / backward pass

LOGISTIC REGRESSION

- Algorithm for binary classification (Input (x) → Output (y))
- Images: 3 separate matrices (R G B)
- Input: Feature vector x

$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \end{bmatrix} \quad \dots \text{Pixel values}$

R {
255
231
⋮

leg:

R · G · B $\Rightarrow 64 \times 64$

G {
255
134
⋮

$x \Rightarrow 64 \times 64 \times 3$

B {
255
134
⋮

(Total no. of pixels)

93 {
134
⋮

$D = D_x$
= Dimension of input feature vector

<u>INPUT</u>	<u>BINARY CLASS.</u>	<u>OUTPUT</u>
feature vector x		0 or 1
<u>→ NOTATION</u>		
1) (x, y)	$\downarrow \leftrightarrow 0 \text{ or } 1$	
	$x \in \mathbb{R}^{n_x}$	
	(n dimensional feature vector)	
2) " m " training examples $\rightarrow m_{\text{train}}$	$\rightarrow (x_1, y_1), (x_2, y_2) \dots (x_m, y_m)$	
	$m_{\text{test}} = \text{no. of test examples}$	
3) $X = \underline{x_1} \quad \underline{x_2} \quad \underline{x_3} \dots$	$\underline{\underline{x_1}}$	
$X = \begin{bmatrix} & & & & \\ x^1 & x^2 & x^3 & \dots & x^m \\ & & & & \end{bmatrix} \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \end{matrix}$	$\leftarrow m \rightarrow$	
$X = \mathbb{R}^{n_x \times m}$		
Python: $x.\text{shape} = (n_x, m)$		
(n_x, m) dimensional matrix		
rows \searrow columns		

$$4) \mathbf{y} = [y^{(1)} \ y^{(2)} \dots \ y^{(m)}] : \mathbf{y} \in \mathbb{R}^{L \times m}$$

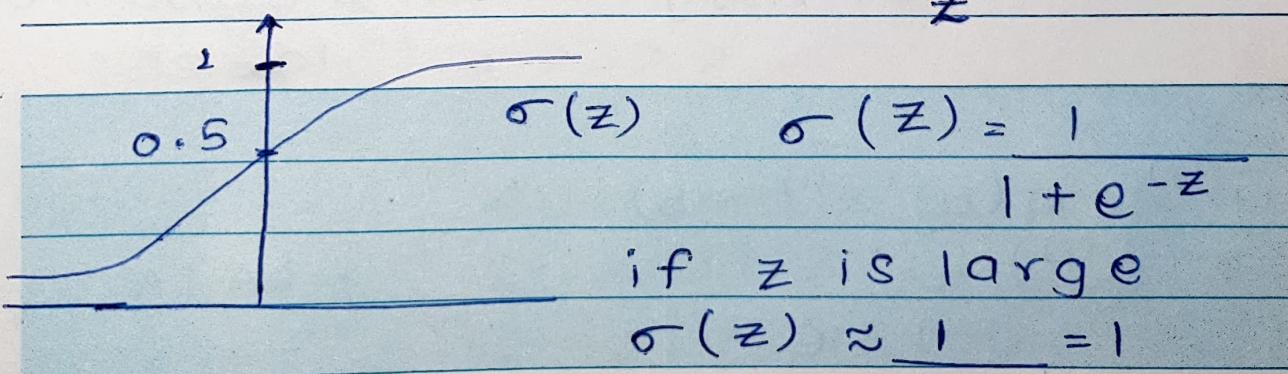
Y. shape = (L, m)

Subject _____ row \leftarrow column \rightarrow (x or y)
Date _____

IMP: take data associated with different training examples and stack them in different columns.

LOGISTIC REGRESSION

- used for binary classification problems (output label $y = 0/1$)
- Given x , want $\hat{y} = P(y=1|x)$
 $x \in \mathbb{R}^{D_x}$ $0 \leq \hat{y} \leq 1$
- Logistic Regression:
 $w \in \mathbb{R}^{D_x} \rightarrow$ vector (n -dimension)
- $b \rightarrow$ real number
- output: $\hat{y} = w^T x + b$
NOT good for bin. class.
 $w^T x + b$ can be -ve or $>> 1$
- Actually: $\hat{y} = \sigma(\underbrace{w^T x + b}_z)$



if z is small or a very large -ve no,

TATA CONSULTANCY SERVICES $\sigma(z) = \frac{1}{1 + e^{-z}} \approx 0$



TATA

Alternative: $x_0 = 1 \in \mathbb{R}^{Dx+1}$

$$\hat{y} = \sigma(\theta^T x)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{Dx} \end{bmatrix} \quad \left. \begin{array}{l} \{ b \\ \omega \} \\ \text{significance of } 'b' \& 'w' \end{array} \right.$$

NOTE: parameters of logistic regression \rightarrow "w" ~~are~~ $[Dx]$
 = dimensional vector
 "b" \rightarrow a real number

COST FUNCTION:

$$\hat{y} = \sigma(\omega^T x + b)$$

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

Training set: m parameters
 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$
 output expected $\rightarrow \hat{y}^{(i)} \approx y^{(i)}$
 training set \leftrightarrow ground truth
 label

$$\hat{y}^{(i)} = \sigma(\omega^T x^{(i)} + b) \quad x^{(i)} \quad \left. \begin{array}{l} x^{(i)} \\ y^{(i)} \\ z^{(i)} \end{array} \right\} \text{DATA}$$

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

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

Experience certainty. IT Services Business Solutions Consulting

data associated with i^{th} example
 (superscript (i) in parenthesis)

Loss(error) function

$$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

as small as possible

Subject _____

Date _____

How good output \hat{y} is when true label is y

→ squared error doesn't make gradient descent work well

→ loss function in logistic regression makes optimisation easy → $L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$
as small as possible

eg: If $y=1$; $L(\hat{y}, y) = -\log \hat{y}$
want $\log \hat{y}$ to be >>> &
 $-\log \hat{y}$ <<<

⇒ want \hat{y} >>> (close to 1
because $0 \leq \hat{y} \leq 1$)

IF $y=0$: $L(\hat{y}, y) = -\log(1-\hat{y})$

FOR LOSS function to be small,

$-\log(1-\hat{y}) \Rightarrow <<<$

$\therefore \log(1-\hat{y}) \Rightarrow >>>$

$(1-\hat{y}) \Rightarrow >>>$

$\Rightarrow \hat{y} <<<$ (as small as possible)

LOSS fx? ⇒ single training ex.
cost fx? ⇒ entire training set



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

In training: find ω & b that minimize J

NOTE: Logistic Regression is a v.v. small neural network

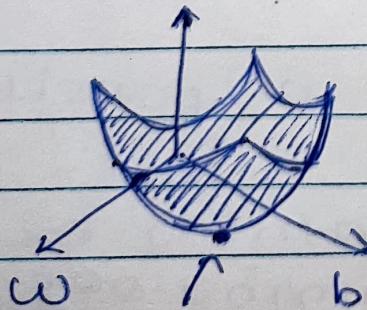
GRADIENT DESCENT

$$\hat{y} = \sigma(\omega^T x + b)$$

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

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

Find ω, b to minimise cost function $J(\omega, b) \rightarrow$ convex fx?

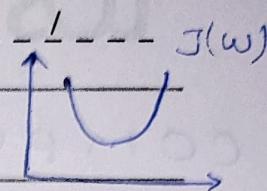


- ① Initialize ω & b
- ② Gradient descent
start at pt 1 & move downwards on the steepest path by 1 step.

Experience certainty.

IT Services
Business Solutions
Consulting

Least value of $J(\omega, b)$

NEURAL NETWORK

- Gradient descent

Repeatedly update w'

Repeat { \rightarrow Learning rate

$$w := w - \alpha d J(w)$$



Updating w

$$d w$$

derivative, update/
change ($d w$)

Learning rate controls how big
of a step is taken on each
iteration of gradient descent

$$\text{CODE: } w := w - \alpha d w$$

For $J(w, b)$

$$w := w - \alpha \frac{d J(w, b)}{d w} \quad \left. \begin{array}{l} \partial J(w, b) \\ \partial w \end{array} \right\}$$

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

$$b := b - \alpha d b \quad \& \quad w := w - \alpha d w$$

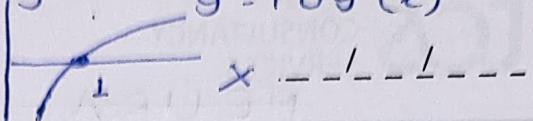
\Rightarrow convex fx? only has 1 optima.

CALCULUS & DERIVATIVES!

① st. line \Rightarrow slope = constant

tcs TATA
CONSULTANCY
SERVICES

LOG: $y = \log(x)$

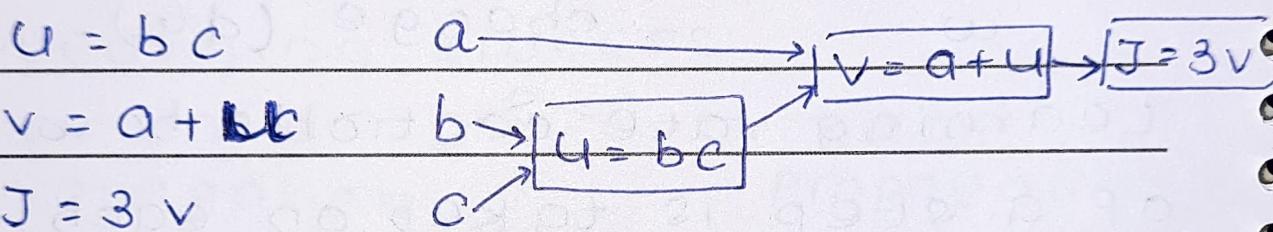


COMPUTATION GRAPH:

- TO explain forward & back propagation

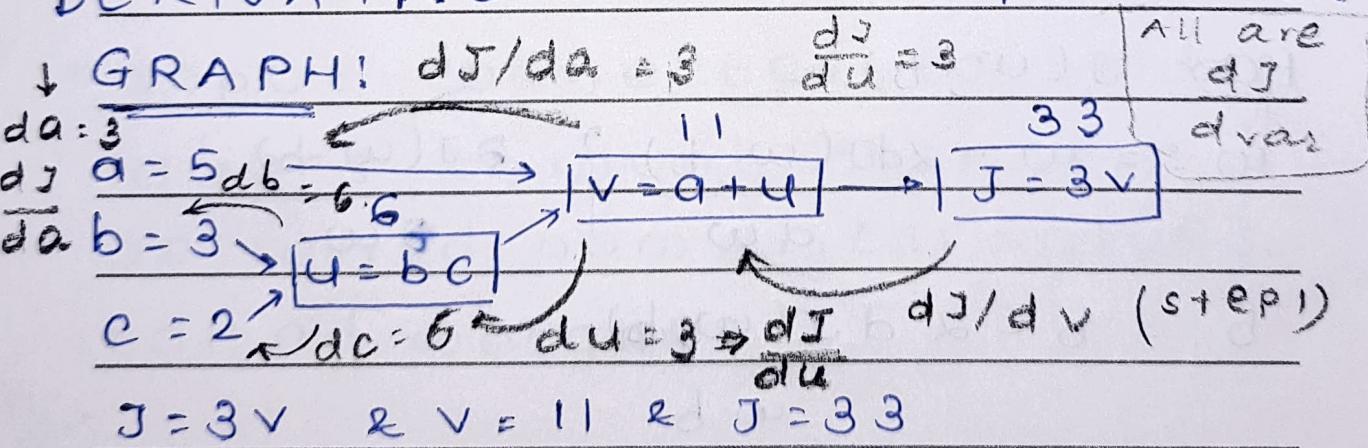
→ Left to right computation
(forward pass)

$$\text{eg: } J(a, b, c) = 3(a + bc)$$



NOTE: "Backward" propagation yields final output variable derivative

DERIVATIVES WITH COMPUTATION



$$\therefore \frac{dJ}{dv} = 3 \quad \frac{dJ}{da} = 3v = 3(a+u) = 3$$

Building on belief Or $\frac{dJ}{da} = \frac{dJ}{dv} \times \frac{dv}{da} = 3 \times 1 = 3$

CODE: d output final var = dJvar = dvar



TATA
CONSULTANCY
SERVICES

Forward \Rightarrow cost fix? (minimize)

Backward \Rightarrow derivatives

GRADIENT DESCENT FOR LOGISTIC

REGRESSION:

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

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

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

$$L(a, y) = -(y \log(a) + (1-y) \log(1-a)) \Rightarrow [a = \hat{y}] \Rightarrow \text{LOSS fx?}$$

computation graph:

$$\begin{aligned}
 & x_1 \xrightarrow{\omega_1} d\omega_1 = \frac{dL}{d\omega_1} = x_1 dz ; \quad d\omega_2 = \frac{dL}{d\omega_2} = x_2 dz ; \quad db = dz \\
 & z = \omega_1 x_1 + \omega_2 x_2 + b \xrightarrow{a = \sigma(z)} L(a, y) \quad \frac{dL}{da} = \frac{-y}{a} + \frac{1-y}{1-a} \\
 & \omega_2 \xrightarrow{b} dz = \frac{dL}{dz} = \frac{da}{da} \quad da = \frac{dL}{da} \quad \frac{da}{dz} = a(1-a) \\
 & \quad = a - y \quad = -\frac{y}{a} + \frac{1-y}{1-a} \\
 & d\omega_1 = \frac{dL}{da} \cdot \frac{da}{dz}
 \end{aligned}$$

$$\begin{aligned}
 \text{updates} \rightarrow \omega_1 &:= \omega_1 - \alpha d\omega_1 & b &:= b - \alpha db \\
 \omega_2 &:= \omega_2 - \alpha d\omega_2
 \end{aligned}$$

LOGISTIC REGRESSION ON 3 EXAMPLES

$$J(\omega, b) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(a^{(i)}) + (1-y^{(i)}) \log(1-a^{(i)}) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

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

$$\frac{\partial J(\omega, b)}{\partial \omega_1} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_1} L(a^{(i)}, y^{(i)})}_{\partial \omega_1 \dots (x^{(i)}, y^{(i)})} \quad dw = np \cdot \text{zeros}$$

initialize: $J=0, d\omega_1=0, d\omega_2=0, db=0$ ($(n+1, 1)$)
 for $i=1$ to m : for loop #1

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

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

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

$$\begin{aligned}
 dz^{(i)} &= a^{(i)} - y^{(i)} && \text{FOR loop #2} \\
 d\omega_1 &+= x_1^{(i)} dz^{(i)} && \downarrow \quad \left. \begin{array}{l} \text{FOR } j = 1 \text{ to } n_x \\ d\omega_j \neq \dots \end{array} \right\} \eta = 2 \\
 d\omega_2 &+= x_2^{(i)} dz^{(i)} \\
 db &+= dz^{(i)} \quad dw += x^{(i)} dz^{(i)}
 \end{aligned}$$

- 2 for loops
 ① for loop over
 'm' training
 examples
 ② for loop over
 features

$$J/m \quad J = J/m$$

$$\begin{aligned}
 d\omega_1 / &= m & d\omega_1 &= d\omega_1 / m \\
 d\omega_2 / &= m & d\omega_2 &= d\omega_2 / m \\
 db / &= m & db &= b / m
 \end{aligned}$$

\downarrow
 accumulators

(sum over entire training set)

$$\begin{aligned}
 \omega_1 &:= \omega_1 - \alpha d\omega_1 \\
 \omega_2 &:= \omega_2 - \alpha d\omega_2 \\
 b &:= b - \alpha db
 \end{aligned}
 \quad \text{] UPDATES}$$

* vectorization
 helps us get
 rid of for loops

$$\sigma(z) = a \leftarrow \text{activation}$$

\uparrow
 sigmoid
 activation
 function

* DERIVATIVES:

$$\frac{da}{dz} = \frac{d}{dz} \sigma(z) = \sigma(z) \times (1 - \sigma(z)) = a(1-a); \quad \frac{dL}{da} = \frac{a-y}{a(1-a)} \Rightarrow \frac{dL}{dz} = a-y$$

DERIVATIVE OF SIGMOID.

* VECTORIZATION : Getting rid of for loops

$$z = \omega^T x + b$$

non-vectorized

$$z = 0$$

for i = 1 to (n_x)
 (i in range (n_x)):

$$z += \omega[i] * x[i]$$

$$z += b$$

vectorized

\Rightarrow VECTORISED CODE RUNS MUCH
 FASTER

GPU } SIMD
 CPU }

single instruction
 & multiple data
 stream
 (PARALLELISM)

$$z = \underbrace{n_p \cdot \text{dot}(\omega, x)}_{\omega^T x} + b$$

• VECTORIZATION

$$u = \underset{\substack{\uparrow \\ \text{vector}}}{A} \underset{\substack{\searrow \\ \text{matrix}}}{v}$$

$$u_i = \sum_j A_{ij} v_j$$

$$u = np.zeros(10, 1)$$

for i ...

for j ...

$$u[i] += A[i][j] * v[j]$$

ALTERNATIVE

$$u = np.dot(A, v)$$

$$+ np.log(v) \leftarrow \text{logarithm}$$

$$+ np.abs(v) \leftarrow \text{absolute}$$

$$+ np.maximum(v, 0) \leftarrow \text{max.}$$

$$+ v ** 2 \leftarrow \text{square}$$

$$+ \frac{1}{v} \leftarrow \text{inverse}$$

ex: $v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$ to have $u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$

$$u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = \text{math.exp}(v[i])$$

ALTERNATIVE

import numpy as np

$$u = np.exp(v)$$

• VECTORIZING LOGISTIC REGRESSION :

$$z^{(1)} = \omega^T x^{(1)} + b \quad \& \quad a^{(1)} = \sigma(z^{(1)}) \quad \dots \quad z^{(m)} = \omega^T x^{(m)} + b$$

$$\& a^{(m)} = \underbrace{\sigma(z^{(1)})}_{\text{ACTIVATION}}$$

$$x = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(n)} \\ | & | & | & | & | \end{bmatrix} \Rightarrow (n_x, m) \text{ matrix}$$

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

$$\begin{bmatrix} z^{(1)} & z^{(2)} & z^{(3)} & \dots & z^{(m)} \end{bmatrix} = \omega^T x + \underbrace{\begin{bmatrix} b & b & \dots & b \end{bmatrix}}_{(1 \times m)}$$

$$= \begin{bmatrix} z_1 & & & & \\ \omega^T x^{(1)} + b & \omega^T x^{(2)} + b & \dots & \omega^T x^{(n)} + b \end{bmatrix} \xrightarrow{L \times m \quad z^n}$$

command: $z = np.dot(\omega.T, x) + b$ \leftarrow real no. (1×1 matrix)
BROADCASTING \leftarrow to calculate all 'z' at same time.

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(z)$$

• VECTORIZING LOGISTIC REGRESSION:

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)} \quad \dots \quad dz^{(m)} = a^{(m)} - y^{(m)}$$

$$\therefore d\mathbf{z} = [dz^{(1)} \quad dz^{(2)} \quad \dots \quad dz^{(m)}]$$

$$\mathbf{A} = [a^{(1)} \quad a^{(2)} \quad \dots \quad a^{(m)}] \quad ; \quad \mathbf{y} = [y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(m)}]$$

$$\therefore d\mathbf{z} = \mathbf{A} - \mathbf{Y} = [a^{(1)} - y^{(1)} \quad \dots \quad a^{(m)} - y^{(m)}]$$

$$d\omega = 0$$

$$d\omega_1 + = x^{(1)} dz^{(1)}$$

$$d\omega_2 + = x^{(2)} dz^{(2)}$$

:

$$d\omega_m + = x^{(m)} dz^{(m)}$$

$$d\omega = d\omega/m$$

$$\Rightarrow d\omega/m = m$$

$$db = 0$$

$$db_1 + = dz^{(1)}$$

$$db_2 + = dz^{(2)}$$

:

$$db_m + = dz^{(m)}$$

$$db = db/m$$

$$\Rightarrow db/m = m$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

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

$$d\omega = \frac{1}{m} \times dz^T$$

$$d\omega = \frac{1}{m} \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$\begin{bmatrix} dz^{(1)} \\ dz^{(m)} \end{bmatrix}$$

$$d\omega = \frac{1}{m} \begin{bmatrix} x^{(1)} dz^{(1)} & \dots & x^{(m)} dz^{(m)} \end{bmatrix}$$

• IMPLEMENTING LOGISTIC REGRESSION:

Traditional approach:

$$\circ J = 0; d\omega_1 = 0, d\omega_2 = 0, db = 0$$

• for $i = 1$ to m :

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

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

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

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

$$d\omega_1 + = x_1^{(i)} dz^{(i)}$$

$$d\omega_2 + = x_2^{(i)} dz^{(i)}$$

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

$$\circ J = J/m, d\omega_1 = d\omega_1/m, d\omega_2 = d\omega_2/m, db = db/m$$

loop

code implementation:

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

$$= np.dot(w.T, x) + b$$

$$A = \sigma(z)$$

$$dz = A - y$$

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

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

$$\omega := \omega - \alpha dw$$

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

• BROADCASTING:-

$$cal = A.sum(axis=0)$$

* axis=0 → vertical

axis=1 → Horizontal

$$\text{percentage} = 100 * A / (cal.reshape(1, 4))$$

A → (3, 4) & (1, 4) vector divisor

e.g. of broadcasting → fit to dimensions

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

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

$$\oplus (m, n) \quad \left\{ \begin{array}{l} + \rightarrow (1, n) \rightsquigarrow (m, n) \\ , \rightarrow (m, 1) \rightsquigarrow (m, n) \end{array} \right.$$

- JUPYTER • RANK 1 ARRAY = 1 square bracket
- vector → 2 square bracket

import numpy as np

$a = np.random.randn(5)$ → creates 5 random Gaussian variables stored in array a

print(a.shape) = (5,) → rank 1 vector (neither row nor column)

print(np.dot(a, a.T)) → a multiplied by a^T , gives a number

NOT EFFICIENT (rank 1 vectors) array

USE COLUMN{(n, 1) arrays} OR ROW{(1, n) array} vector instead:

```
import numpy as np
a.shape = (5, 1)
a = np.random.randn(5, 1) ← 5x1 vector (column)
```

print(a.T) → row vector

assertions → documentation for code

assert (a.shape == (5, 1))