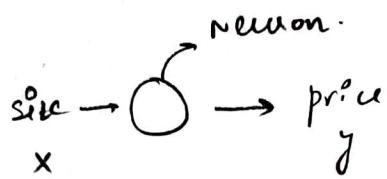
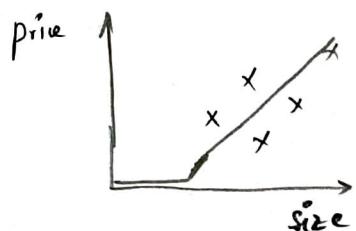
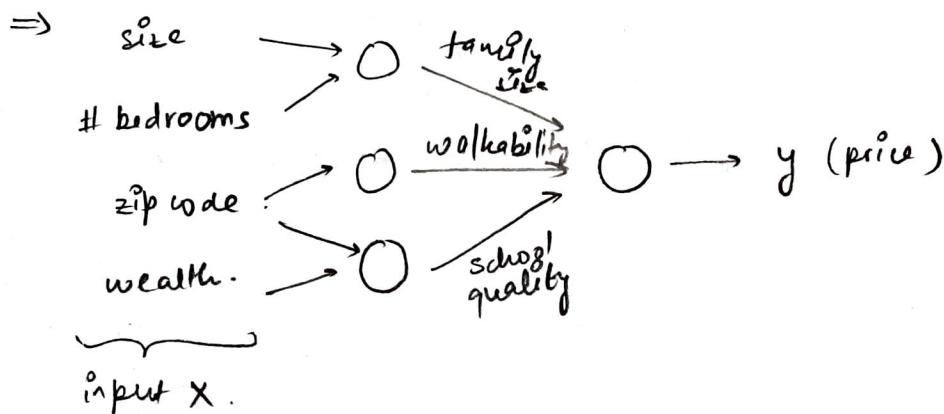


Neural Networks and Deep Learning

Housing price prediction



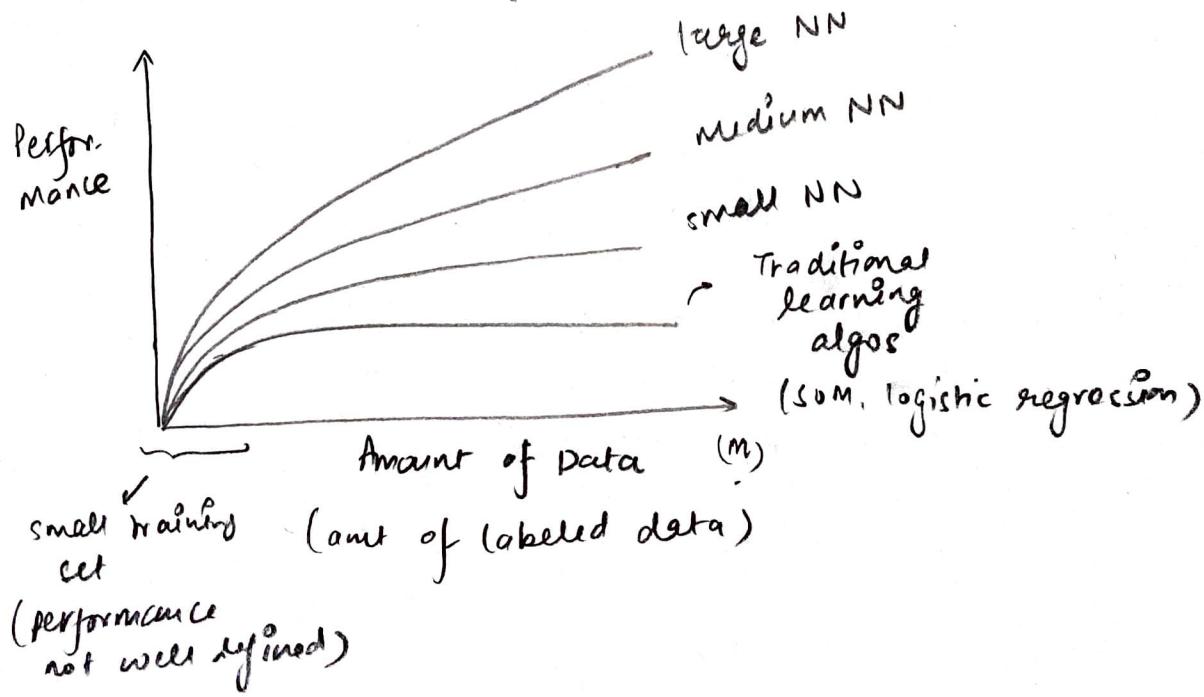
ReLU
Rectified linear unit.



structured data → proper tabular data, for ex. real estate data.

unstructured data → audio, images, text

Scale drives deep learning progress



\Rightarrow Binary classification

(input image) \rightarrow 1 (cat) vs (0)
 ↴
 64×64 pixels. $\underbrace{y}_{\text{non cat}}$

RGB $\rightarrow 64 \times 64$ matrices

$$x = \begin{bmatrix} 255 \\ 231 \\ 21 \\ \vdots \\ 255 \\ 124 \\ 124 \\ 124 \\ 126 \end{bmatrix} \hookrightarrow 64 \times 64 \times 3 \\ = 12288 \\ n = n_x = 12288$$

$x \rightarrow y$
 (input) $\hookrightarrow 0 \text{ or } 1$

$$m = m_{\text{train}}$$

$$m_{\text{test}} = \# \text{ test examples.}$$

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

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

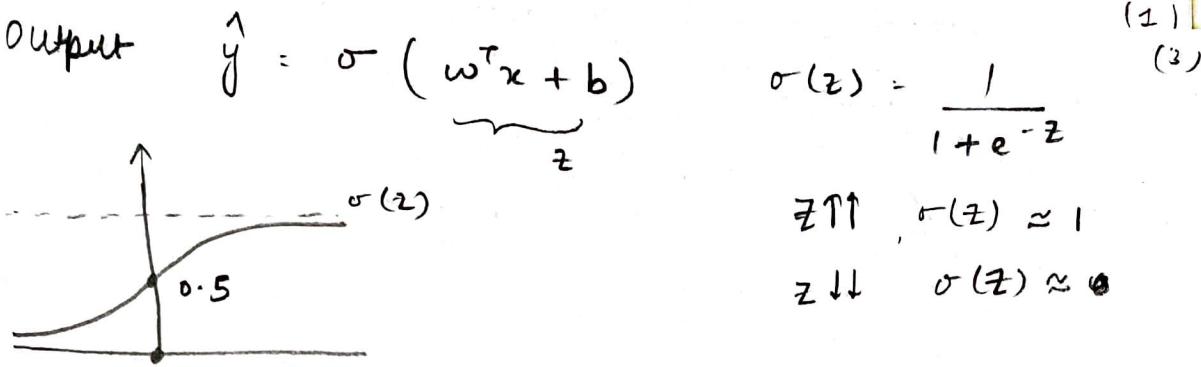
$$X = \left[\begin{array}{c|c|c|c} & & & \\ \hline x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ \hline & & & \end{array} \right] \quad \begin{matrix} \uparrow \\ n_x \\ \downarrow \\ \text{m} \end{matrix} \quad \begin{matrix} x \in \mathbb{R}^{n_x \times m} \\ x. \text{shape} \rightarrow (n_x, m) \end{matrix}$$

$$y = \left[\begin{array}{c} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{array} \right] \quad y \in \mathbb{R}^{1 \times m} \quad y. \text{shape} = (1, m)$$

\Rightarrow Logistic Regression \rightarrow used when $y \in \{0, 1\}$

Given x , want $\hat{y} : P(y=1|x) \rightarrow \text{probability of } y=1 \text{ given } x \in \mathbb{R}^{n_x}$

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



Logistic Regression Cost function

$$\hat{y}^{(i)} = \sigma(\omega^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}}$$

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

Loss function:

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

If $y=1$: $L(\hat{y}, y) = -\log \hat{y}$, for minimum $L(\hat{y}, y)$,
we want \hat{y} to be large. ($\hat{y} \rightarrow 1$)

If $y=0$: $L(\hat{y}, y) = -(1-\hat{y}) \log (1-\hat{y})$ $\rightarrow \hat{y}$ as small as possible. ($\hat{y} \rightarrow 0$)

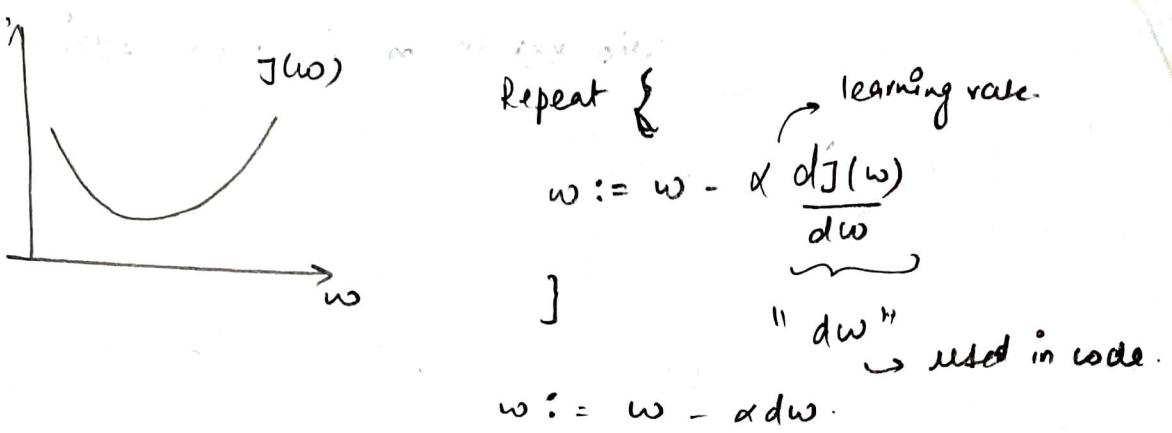
Cost function: $J(\omega, b)$

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

$$-\frac{1}{m} \sum_{i=1}^m [\hat{y}^{(i)} \log \hat{y}^{(i)} + (1-\hat{y}^{(i)}) \log (1-\hat{y}^{(i)})]$$

Gradient Descent

→ we want to find ω, b so that $J(\omega, b)$ is minimum

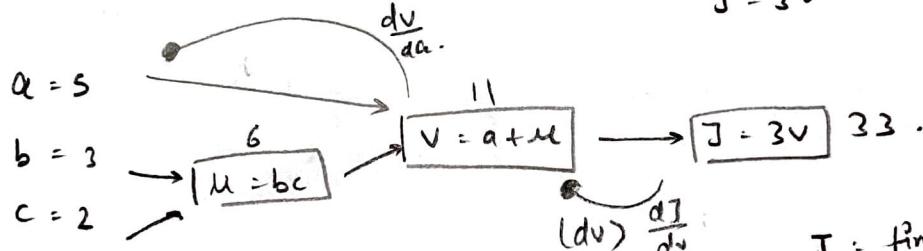


$$\left. \begin{array}{l} J(w, b) \\ w := w - \frac{\partial J(w, b)}{\partial w} \\ b := b - \alpha \frac{\partial J(w, b)}{\partial b} \end{array} \right\} \text{Actual implementation.}$$

\Rightarrow Derivatives

\Rightarrow computation graph

$$J(a, b, c) = 3(a + bc) \quad u = bc \quad v = a + u \quad J = 3v$$



J : final output variable.

$$\frac{dJ}{dv} = 3 \quad \frac{dJ}{da} = 3 = \frac{dv}{da} \times \frac{dJ}{dv}$$

$\rightarrow \frac{d(\text{final output var})}{d\text{var.}}$ } \rightarrow "dvar" \rightarrow in python code.

$$\rightarrow "dv" = 3, "da" = 3 \rightarrow \frac{dJ}{db} = \frac{dJ}{dv} \times \frac{dv}{du} \times \frac{du}{db}$$

$$\frac{dJ}{du} = 3, du = 3 \quad = 3 \times 1 \times c \quad = 3c = db$$

$$\rightarrow \frac{dJ}{dc} = dc = 3b.$$

$$db|_{at c=2} = 6$$

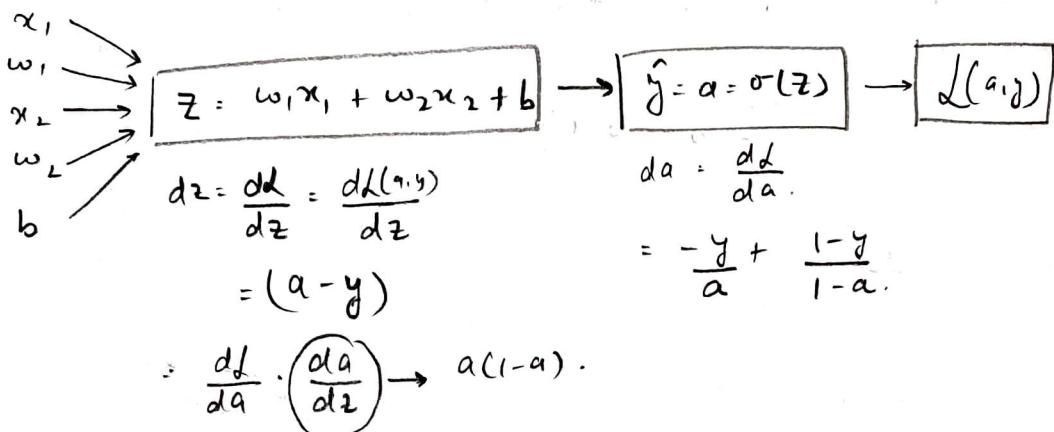
$$dc|_{b=3} = 9$$

Gradient Descent for logistic regression (one example)

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

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

Let's write as a computation graph. we have two features x_1 and x_2 .



$$\frac{dL}{dw_1} = "dw_1" = x_1 \cdot dz \quad "dw_2" = x_2 \cdot dz \quad db = dz.$$

$$\begin{cases} w_1 := w_1 - \alpha dw_1 \\ w_2 := w_2 - \alpha dw_2 \\ b := b - \alpha db \end{cases}$$

\Rightarrow Gradient Descent on m examples

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

$$\rightarrow a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^T x^{(i)} + b) \quad dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

$$\frac{\partial J(\omega, b)}{\partial \omega_1} = \alpha \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{d}{d \omega_1} L(a^{(i)}, y^{(i)})}_{d \omega_1^{(i)}} \quad ,$$

Now, calculate as from previous section.

logistic regression on m examples

$$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0 \quad \} \text{ Initialize}$$

for $i = 1$ to m :

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

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

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

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

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

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

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

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

$$b = b - \alpha db.$$

$$J / = m$$

$$dw_1 = dw_1 / m$$

$$dw_2 = dw_2 / m.$$

\Rightarrow Vectorization

$$z = w^T x + b \quad w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w \in \mathbb{R}^{n_m} \quad x \in \mathbb{R}^{n_m}$$

Non vectorised.

$$z \rightarrow$$

for i in range n_m .

$$z += w[i] * x[i]$$

vectorised

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

$$z += b$$

$$\rightarrow v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \quad u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$$u = np \cdot \exp(v)$$

$$u = np \cdot \log(v)$$

$$u = np \cdot \text{abs}(v)$$

\Rightarrow Vectorizing Logistic Regression

$$z^{(1)} = \omega^T x^{(1)} + b \quad z^{(2)} = \omega^T x^{(2)} + b$$

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

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

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

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

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \omega^T X + [b, b, \dots, b]_{(1 \times m)}$$

$$\omega^T = [\omega_0, \omega_1, \dots, \omega_m] \quad \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ | & | & | \end{bmatrix} + [b, b, \dots, b]$$

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}]$$

$$Z = np.dot(\omega, X) + b \xrightarrow{(1,1)} \text{but python automatically makes } b \text{ a row vector.}$$

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

\Rightarrow Vectorizing Gradient Descent in Logistic Regression

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

$$\Delta Z = [dz^{(1)}, dz^{(2)}, \dots, dz^{(m)}]_{1 \times m}$$

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

$$\Delta Z = A - Y = [a^{(1)} - y^{(1)}, a^{(2)} - y^{(2)}, \dots]$$

$$d\omega = 0$$

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

$$d\omega + = y^{(2)} dz^{(2)}$$

$$db = 0$$

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

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

$$d\omega/m$$

$$db/m$$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)} : \frac{1}{m} \cdot np \cdot \text{sum}(dz)$$

$$d\omega = \frac{1}{m} X (dz)^T$$

$$= \frac{1}{m} \left[\begin{array}{c|c|c|c} 1 & 1 & \cdots & 1 \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \\ \hline 1 & 1 & \cdots & 1 \end{array} \right]_{m \times m} \left[\begin{array}{c} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \end{array} \right]_{m \times 1}$$

$$= \frac{1}{m} \left[x^{(1)} dz^{(1)} + x^{(2)} dz^{(2)} + \cdots + x^{(m)} dz^{(m)} \right]$$

→ Implementing Logistic Regression

$$Z = w^T x + b \quad : np \cdot \text{dot}(w, x) + b$$

$$A = \sigma(Z)$$

$$dz = A - Y$$

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

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

$$\underline{\text{q1}} \quad w = w - \alpha dw$$

$$b = b - \alpha db$$

thus, this is ~~is~~ a single iteration of Gradient Descent.

\Rightarrow Broadcasting in Python

calc: A sum (axis=0)

axis=0 \rightarrow sum downwards
 axis=1 sum horizontally

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

General principle

(m, n)	$\begin{matrix} + \\ - \\ \times \\ / \end{matrix}$	$(1, n) \rightsquigarrow (m, n)$
matrix		$(m, 1) \rightsquigarrow (m, n)$

a = np.random.randn(5)
 a.shape = (5,)
 "rank 1 array"

} don't use

a = np.random.randn(5, 1) \rightarrow a.shape=(5, 1)

column vector

a = np.random.randn(1, 5) \rightarrow a.shape=(1, 5)

row vector

\Rightarrow Explanation of logistic regression cost function

$$\hat{y} = \sigma(w^T x + b) \quad \text{where} \quad \sigma(z) = \frac{1}{1+e^{-z}}$$

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

$y = 1 : p(y|x) = \hat{y}$ (probability of $y=1$ is \hat{y})

$y = 0 : p(y|x) = 1 - \hat{y}$ (probability of $y=0$ is $1 - \hat{y}$)

$$\text{Let } p(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

$$y=1 : p(y|x) = \hat{y}^1 (1-\hat{y})^0 = \hat{y}$$

$$y=0 : p(y|x) = \hat{y}^0 (1-\hat{y})^1 = 1 - \hat{y}$$

Cost on m examples

$$\log p(\text{labels in training set}) = \log \prod_{i=1}^m p(y^{(i)} | x^{(i)})$$

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

$$\text{Cost (minimise)} = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

→ np.exp(z) → exponential function to all elements of matrix z.

$$\rightarrow \alpha = \sigma(z) = \frac{1}{1 + e^{-z}}, \sigma'(z) = \alpha(1 - \alpha).$$

$$\rightarrow x = \begin{bmatrix} 0 & 3 & 4 \\ 2 & 6 & 4 \end{bmatrix} \|x\|_{\text{row}} = \text{np.linalg.norm}(x, \text{axis}=1, \text{keepdims} = \text{True})$$

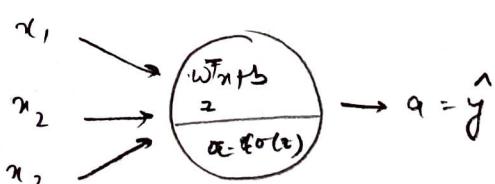
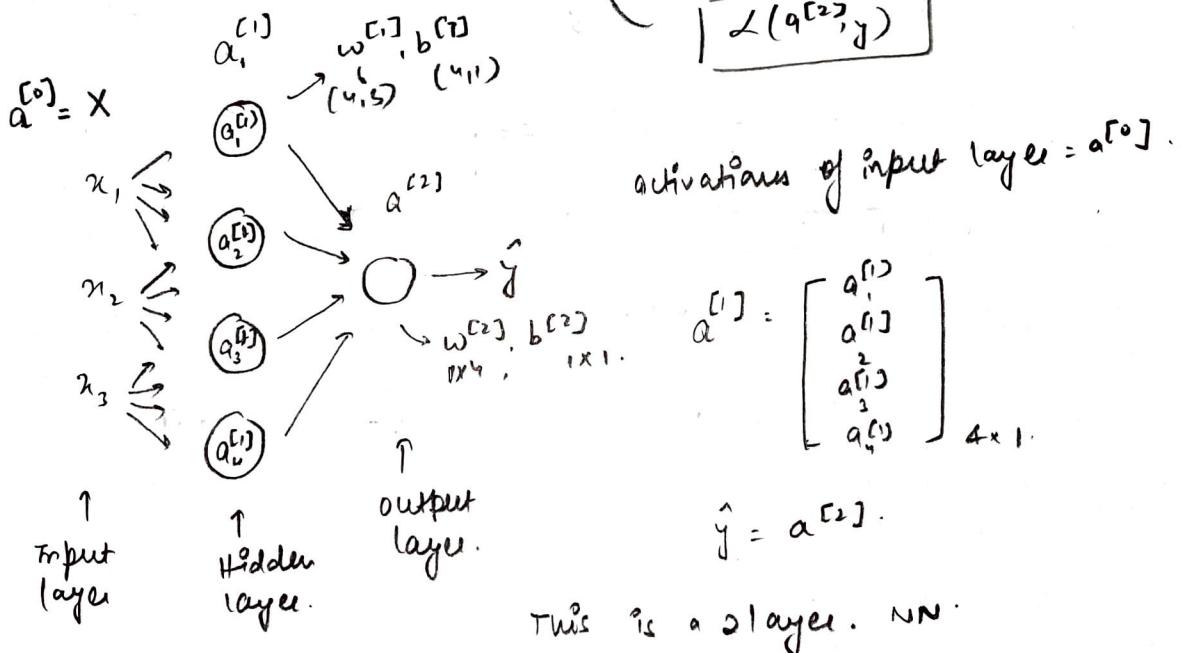
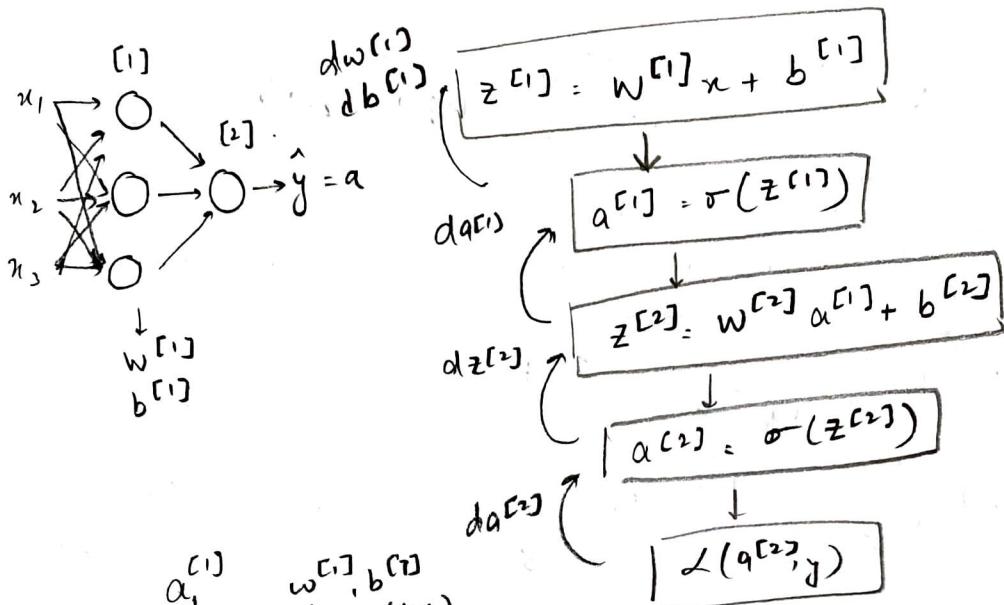
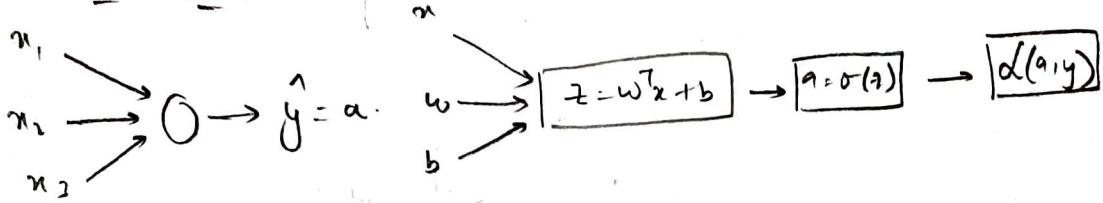
$$\|x\| = \begin{bmatrix} 0 & \frac{3}{\sqrt{58}} & \frac{4}{\sqrt{58}} \\ \frac{2}{\sqrt{58}} & \frac{6}{\sqrt{58}} & \frac{4}{\sqrt{58}} \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{58}}{\sqrt{58}} \end{bmatrix}$$

→ keepdims → keepdimensions.

→ x.sum = np.sum(x, axis=1, keepdims=True).

Week - 2 Shallow Neural Networks

Neural network



$$\left. \begin{aligned} z_1^{[1]} &= w_1^{[1]}\mathbf{x} + b_1^{[1]} \\ a_1^{[1]} &= \sigma(z_1^{[1]}) \end{aligned} \right\} \quad \begin{aligned} z_2^{[1]} &= w_2^{[1]}\mathbf{x} + b_2^{[1]} \\ a_2^{[1]} &= \sigma(z_2^{[1]}) \end{aligned} \quad \begin{aligned} z_3^{[1]} &= w_3^{[1]}\mathbf{x} + b_3^{[1]} \\ a_3^{[1]} &= \sigma(z_3^{[1]}) \end{aligned}$$

$a_i^{[1]}$ → node in layer.

$$\left. \begin{aligned} z_4^{[1]} &= w_4^{[1]}\mathbf{x} + b_4^{[1]} \\ a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned} \right\} \quad \begin{aligned} z_1^{[1]} &= w_1^{[1]}\mathbf{x} + b_1^{[1]} \\ a_1^{[1]} &= \sigma(z_1^{[1]}) \end{aligned}$$

Now, vectorise.

$$w^{[1]T} = \begin{bmatrix} | & | & | & | \\ w_1^{[1]} & w_2^{[1]} & w_3^{[1]} & w_4^{[1]} \\ | & | & | & | \end{bmatrix}_{3 \times 4}$$

$$w^{[1]T} = \begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} =$$

$$\begin{bmatrix} w_1^{[1]T}\mathbf{x} + b_1^{[1]} \\ w_2^{[1]T}\mathbf{x} + b_2^{[1]} \\ w_3^{[1]T}\mathbf{x} + b_3^{[1]} \\ w_4^{[1]T}\mathbf{x} + b_4^{[1]} \end{bmatrix} = z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = \begin{bmatrix} \sigma(z_1^{[1]}) \\ \vdots \\ \sigma(z_4^{[1]}) \end{bmatrix}$$

Given input $x = a^{[0]}$.

$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]}$$

(4,1) (4x3) (3x1) (4,1)

w is a vector obtained by stacking up transposes of consecutive w 's.

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

(4,1) (4,1)

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

(1,1) (1x4) (4,1) (1,1)

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

(1,1) (1,1)

All of this was for one training example, so we will vectorize for m training examples.

Single training example.

$$x \rightarrow a = \hat{y}$$

$$\left\{ \begin{array}{l} z^{[1]} = w^{[1]} x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

for m examples.

$$\begin{aligned} x^{(1)} &\rightarrow \hat{y}^{(1)} = a^{[2](1)} \\ x^{(2)} &\rightarrow \hat{y}^{(2)} = a^{2} \\ x^{(m)} &\rightarrow \hat{y}^{(m)} = a^{[2](m)} \end{aligned}$$

$a^{[2](i)}$ \rightarrow example no. i
 \downarrow layer 2.

so, using loop

for $i = 1$ to m :

$$\begin{aligned} z^{(1)(i)} &= w^{[1]} x^{(i)} + b^{[1]} \\ a^{(1)(i)} &= \sigma(z^{(1)(i)}) \\ z^{(2)(i)} &= w^{[2]} a^{(1)(i)} + b^{[2]} \\ a^{(2)(i)} &= \sigma(z^{(2)(i)}) \end{aligned}$$

using a loop, but
we want to
vectorize

$$X = \begin{bmatrix} | & | & | & | \\ X^{(1)} & X^{(2)} & X^{(3)} & \dots & X^{(m)} \\ | & | & | & | \end{bmatrix} \rightarrow \mathbb{R}^{n_x \times m}$$

(n_x, m)

$$Z^{[1]} = w^{[1]} X + b^{[1]} \quad (\text{all capital cases})$$

where $Z^{[1]} = \begin{bmatrix} | & | & | & | \\ Z^{1} & Z^{[1](2)} & Z^{[1](3)} & \dots & Z^{[1](m)} \\ | & | & | & | \end{bmatrix}$

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

$\xrightarrow{\text{m examples}}$
 corresponds to
 1 node

\Rightarrow Explanation of vectorised equations

$$z^{(1)(1)} = w^{[1]} X^{(1)} + b^{[1]}, \quad z^{(1)(2)} = w^{[1]} X^{(2)} + b^{[1]},$$

lets say $b^{[1]} = 0$

$$z^{(1)(3)} = w^{[1]} X^{(3)} + b^{[1]}$$

$$w^{[1]} X^{(1)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \xrightarrow{\text{column vector}}$$

$$w^{[1]} X^{(2)} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \xrightarrow{\text{column vector}}$$

$$Y = X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \dots \\ | & | & | \end{bmatrix}$$

$$\begin{aligned} w^{[1]} X &= \begin{bmatrix} | & | & | \\ Z^{(1)(1)} & Z^{(1)(2)} & \dots \\ | & | & | \end{bmatrix} \\ &= Z^{[1]} \end{aligned}$$

So, for m training examples:-

$$X = A^{[0]}$$

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

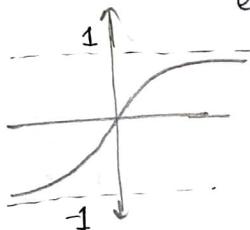
$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

Activation Functions

Till now, we have used sigmoid activation function. So generalising, we would represent $g(z)$.

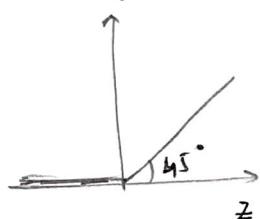
eg. $a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$



$\tanh(z)$ almost always works better as activation function in hidden layers;

slope is nearly zero for $|z| \uparrow$, so gradient descent works slowly.

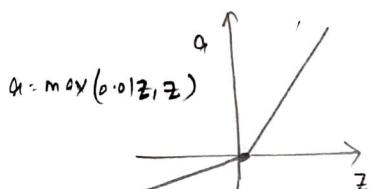
$$a = \max(0, z)$$



Rectified Linear Unit

(ReLU) \rightarrow use for hidden layers.

In practice, ReLU is used.



leaky ReLU

(small slope for negative z)

Sigmoid (z) \rightarrow new use except for last layer in binary classification.

ReLU $\rightarrow a = \max(0, z)$

⇒ Why do we need non linear activation functions?

Given x :

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

$$a^{(1)} = g(z^{(1)}) = w^{(1)}x + b^{(1)}$$

$$a^{(1)} = z^{(1)}$$

$$a^{(2)} = w^{(2)}(w^{(1)}x + b^{(1)}) + b^{(2)}$$

$$z^{(2)} = w^{(2)}a^{(1)} + b^{(2)}$$

↓
So, we will end up with a linear map only.

We can use linear activation function in output layers.
(in some cases).

⇒ Derivatives of activation functions

(i) Sigmoid

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

$$= g(z)(1-g(z))$$

$$g'(0) = \frac{1}{4} \quad = g(1-g) \quad \text{as } g(z)=a.$$

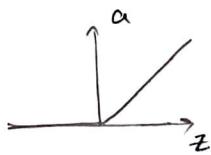
(ii) tanh activation function

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

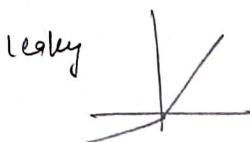
$$\text{if } g(z) = a$$

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

(iii) ReLU and Leaky ReLU



$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & z \geq 0 \end{cases}$$



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

Gradient Descent for Neural Network

Parameters : $W^{[2]}, b^{[1]}, W^{[2]}, b^{[2]}$ $n_m = n^{[0]}$
 $(n^{[1]}, n^{[0]}) \quad (n^{[0]}, 1) \quad (n^{[2]}, n^{[1]}) \quad (n^{[2]}, 1)$ $n^{[1]}, n^{[2]} = 1$.

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

After initializing parameters.

Gradient Descent:

Repeat :

predict ($\hat{y}^{(i)}$ for $i = 1, \dots, m$)

$$dW^{[1]} = \frac{\partial J}{\partial w^{[1]}}, \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}} \quad \dots$$

$$w^{[1]} = w^{[1]} - \alpha dW^{[1]}$$

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

$$w^{[2]} = w^{[2]} - \alpha dW^{[2]}$$

$$b^{[2]} = b^{[2]} - \alpha db^{[2]}$$

}

Formula for derivatives

$$z^{[1]} = w^{[1]}x + b^{[1]} \quad A^{[1]} = g^{[1]}(z^{[1]}), \quad z^{[2]} = w^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} = g^{[2]}(z^{[2]}) = \sigma(z^{[2]}).$$

Back propagation

$$dz^{[2]} = A^{[2]} - y$$

} Assuming binary classification with sigmoid function.

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

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

$$dz^{[1]} = \underbrace{w^{[2] T} dz^{[2]}}_{(n^{[1]}, m)} \times \underbrace{g^{[1]}'(z^{[1]})}_{(n^{[1]}, m)}$$

element wise
product

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

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

$\underbrace{\quad}_{(n^{[1]}, 1)}$ to avoid rank 1 arrays.

\Rightarrow ④ x and dx have same matrix dimensions

Summary of Gradient Descent

For single training example.

$$dZ^{[2]} = \alpha^{[2]} - y$$

$$dW^{[2]} = dZ^{[2]} \alpha^{[1] T}$$

$$db^{[2]} = dZ^{[2]}$$

$$(n^{[1]}, 1) \quad dZ^{[1]} = W^{[2] T} dZ^{[2]} \times g^{[1]}'(z^{[1]})$$

$$dW^{[1]} = dZ^{[1]} x^T$$

$$db^{[1]} = dZ^{[1]}$$

For m training examples

$$dZ^{[2]} = A^{[2]} - Y$$

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

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

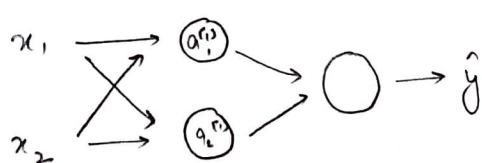
$$(n^{[1]}, m) \quad dZ^{[1]} = W^{[2] T} dZ^{[2]} \times \underbrace{g^{[1]}'(z^{[1]})}_{\text{elementwise multiplication}}$$

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

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

Random Initialization

Initialising to 0 won't work



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

$a_1^{[1]}$ and $a_2^{[1]}$ will be equal:
thus, when we compute back-propagation.

$$dz_1^{(1)} = dz_2^{(2)}$$

dw becomes $\begin{bmatrix} u & v \\ u & v \end{bmatrix}$ $w^{(1)} = w^{(1)} - \alpha \cdot dw$
 $1^{\text{st}} \text{ row} = 2^{\text{nd}} \text{ row}.$

So, initialise randomly.

$$w^{(1)} = \text{np.random.randn}(2, 2) * 0.01$$

$$b^{(1)} = \text{np.zeros}(2, 1)$$

$$w^{(2)} = \text{np.random.randn}(1, 2) * 0.01$$

$$b^{(2)} = 0$$

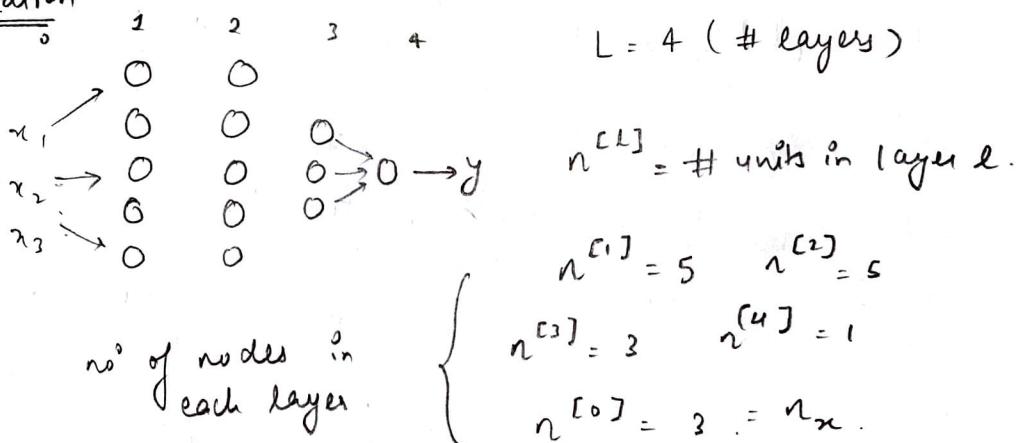
→ we want weights
to be small
initially we don't
want to end up in flat
part of say $\tanh(z)$
while backpropagation.

⇒ Week-4

Deep 1-layer neural network

Deep → means more hidden layers

Notation



$a^{(l)}$: activations in layer l . $X = a^{(0)}$ $a^{(L)} = \hat{y}$

$a^{(l)} = g^{(l)}(z^{(l)})$, $w^{(l)}$: weights for $z^{(l)}$
 $b^{(l)}$.

→ forward propagation in deep network

Consider network drawn in previous page

$a^{[0]} = x \rightarrow$ single training example layer 2

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

$$\underbrace{a^{[1]} = g^{[1]}(z^{[1]})}_{\text{layer 1}} \quad a^{[2]} = g^{[2]}(z^{[2]})$$

layer 1

layer 4

$$z \text{ and } a \text{ small.} \quad z^{[4]} = w^{[4]}a^{[3]} + b^{[4]}$$

$$a^{[4]} = g^{[4]}(z^{[4]})$$

$$\text{So, } z^{[\ell]} = w^{[\ell]}a^{[\ell-1]} + b^{[\ell]} \quad \left. \begin{array}{l} \\ a^{[\ell]} = g^{[\ell]}(z^{[\ell]}) \end{array} \right\} \text{general rule}$$

Vectorize → For m training examples.

$$z^{[1]} = w^{[1]}X + b^{[1]} \quad X = A^{[0]}$$

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

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

$$A^{[2]} = g^{[2]}(z^{[2]})$$

$$z^{[4]} = w^{[4]}A^{[3]} + b^{[4]}$$

$$A^{[4]} = g^{[4]}(z^{[4]})$$

for loop · to calculate

\Rightarrow getting matrix dimensions right

$$\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 \\ x_1 & 0 & 0 & 0 & 0 & 0 \rightarrow y \\ x_2 & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

$$L = 5$$

$$n^{[1]} = 3$$

$$n^{[2]} = 5$$

$$n^{[3]} = 4$$

$$n^{[4]} = 2$$

$$n^{[5]} = 1$$

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

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

$$z^{[1]} \rightarrow (3,1) \rightarrow (n^{[1]}, 1), x \rightarrow (n^{[0]}, 1)$$

$$\text{so } w^{[1]} = (n^{[1]}, n^{[0]})$$

$$\text{generally, } w^{[\ell]} = (n^{[\ell]}, n^{[\ell-1]})$$

$$w^{[2]} = (n^{[2]}, n^{[1]}) = (5, 3)$$

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

$$b^{[\ell]} = (n^{[\ell]}, 1)$$

\rightarrow dimensions of $dw = w$

$$dw^{[\ell]} = (n^{[\ell]}, n^{[\ell-1]})$$

$$db^{[\ell]} = (n^{[\ell]}, 1)$$

In vectorised implementation

w, b, dw, db dimensions would not change.

$$\underbrace{z^{[1]} = w^{[1]}x + b^{[1]}}_{(1^1, m)} \rightarrow \text{dimension is still } n^{[1]} \times 1, \quad \text{it is broadcasted and then added to } w \cdot x.$$

$$z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

$$dz^{[l]}, dA^{[l]} : (n^{[l]}, m)$$

Why Deep Representations?

Building Blocks of Deep Neural Net

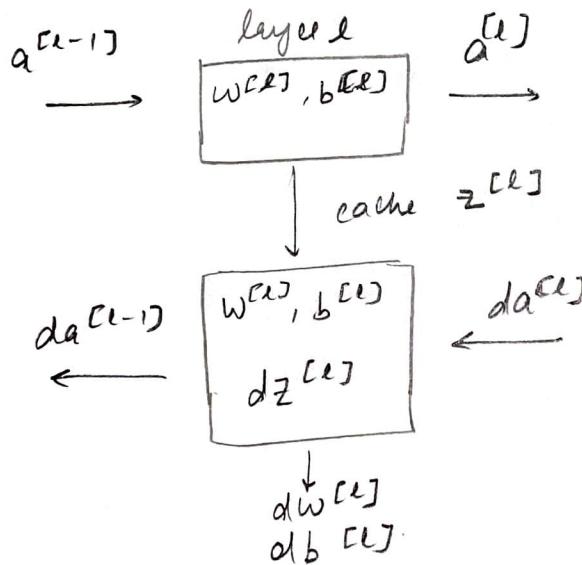
layer l : $W^{[l]}, b^{[l]}$

forward prop: $a^{[l-1]}$, output $\rightarrow a^{[l]}$

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \quad \text{cache } [z^{[l]}]$$

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

Back prop. \rightarrow input $da^{[l]}$, output $da^{[l-1]}$
cache $(z^{[l]})$



\Rightarrow Forward and Backward propagation

FP for layer l

Input $a^{[l-1]}$, output $a^{[l]}$, cache $z^{[l]}$

$$\left. \begin{array}{l} z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} = g(z^{[l]}) \end{array} \right\} \quad \begin{array}{l} z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = g(z^{[l]}) \end{array}$$

Backward propagation for layer l

Input $da^{[l]}$, output $da^{[l-1]}, dw^{[l]}, db^{[l]}$

$$dz^{[l]} = da^{[l]} * g^{[l]}'(z^{[l]})$$

$$dw^{[l]} = dz^{[l]} * a^{[l-1]^\top}$$

$$db^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = w^{[l]^\top} dz^{[l]}$$

vectorised $dz^{[l]} = dA^{[l]} * g^{[l]}'(z^{[l]})$

$$dw^{[l]} = \frac{1}{m} dz^{[l]} \cdot A^{[l-1]^\top}$$

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

$$dA^{[l-1]} = W^{[l]^\top} dz^{[l]}.$$

\Rightarrow Parameters v/s Hyperparameters

w, b

\hookrightarrow learning rate

no. of hidden layers L

no. of hidden units $n^{(1)}, n^{(2)}, \dots$

activation functions.