

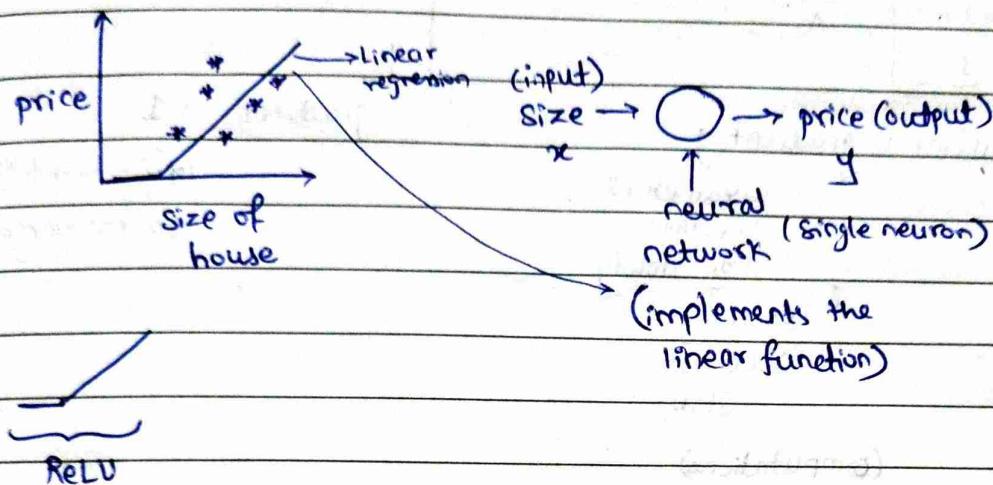
# \* Deep Learning And Neural Networks \*



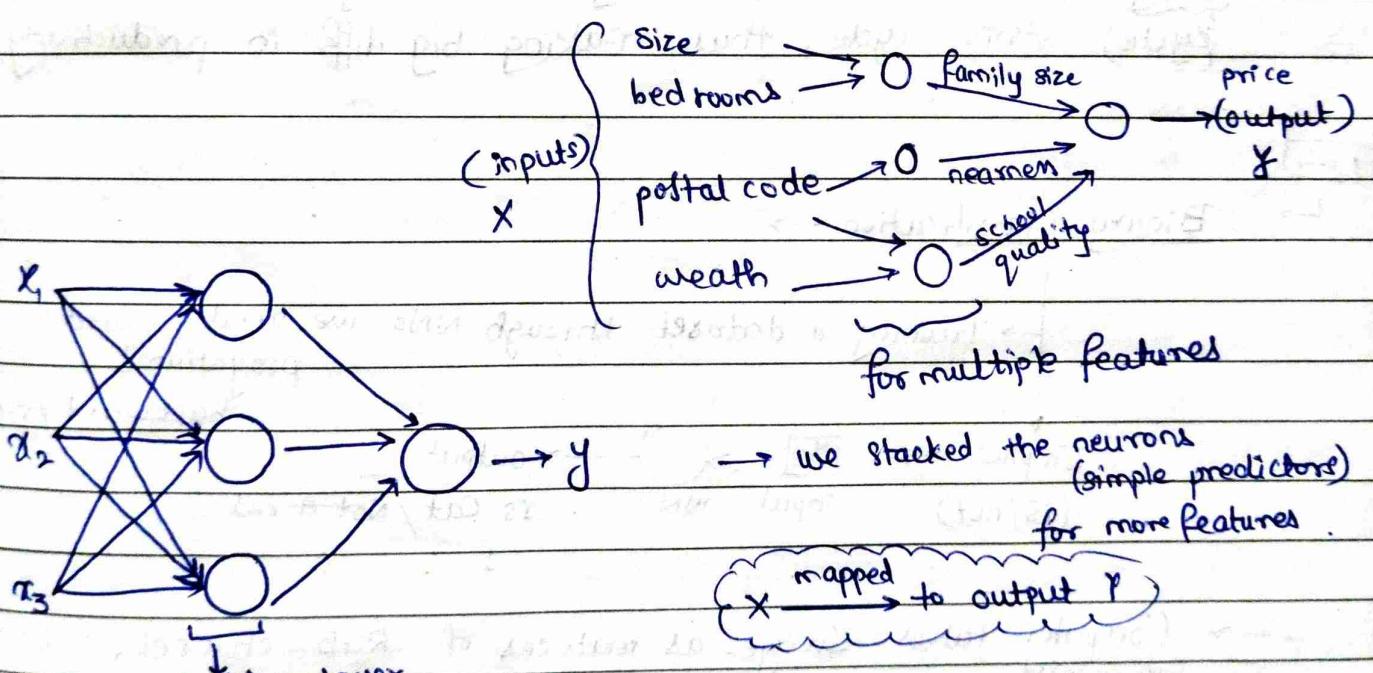
L=1

→ Neural Networks ⇒

↳ Let take an ex. of Housing price Prediction



↳ Rectified Linear Unit → Multiple parameters also included having non-linear fun<sup>n</sup>



(Every feature is connected to every neuron)

L=2

- Structured Data → We can have databases having data in form of rows & columns, well-defined features.

- Unstructured Data → Mainly Audio files, images or texts  
↳ NNs made it easy to recognise unstructured data.

L=31  
For High performance :  
 1) We need a neural network that is large enough to take the advantage of large amounts of data.  
 2) We also need a lot of data for the training.

→ Also algorithms making the computation faster.

↓  
Sigmoid func → to → ReLU func (faster)

 almost 0 gradient  
parameters change very slowly  
learning becomes slow  
  
 gradient → 1  
most likely to converge.

(computation)  
→ We need faster algorithms, as training N.N.s is basically a loopy structure and if training takes longer time so does (cycle) this cycle, thus making big diff. in productivity.

L=4

↪ Binary Classification →

→ Training a dataset through NNs we need forward propagation & backward propaga.  
Implies →  →  → output  
(is/not) input NN IS Cat / Not a cat  
1 0

→ Computer takes Images as matrices of RGB channel,  
(we know this)  
RGB

We need to turn these pixel intensity values into a feature vector by unfolding these pixels into an input  $X \rightarrow X: \text{dimensions} \times m \times n \times 3$   
 $\downarrow$   
vector  $n = m \times n \times 3$

→ A single training sample is denoted by  $(x, y)$  = pair

$\rightarrow m$ : training sample  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$   
no. of training samples (examples)  
 $x \in \mathbb{R}^{n \times 1}$  &  $y \rightarrow \{0, 1\}$   
(feature vector)  
 $\rightarrow$  training set

$$X = \begin{bmatrix} | & | & | & | & | \\ X^{(1)} & X^{(2)} & X^{(3)} & \dots & X^{(m)} \\ | & | & | & & | \end{bmatrix} \quad n_x = \text{rows} \quad x \in \mathbb{R}^{n_x \times m}$$

$\longleftrightarrow$  columns

$n_x \times m$  dimensional matrix  
 $X \cdot \text{shape} = (n_x, m)$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \quad y \in \mathbb{R}^{1 \times m}$$

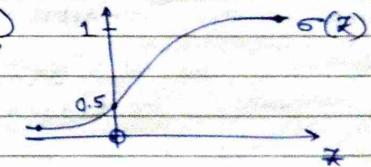
$Y \cdot \text{shape} = (1, m)$

L=5 → Logistic Algorithm = "Used when the output label  $y$  is 0/1."  
Regression  
↳ Binary Classification problem

→ Given:  $X$ ,  $\hat{y} = p(y=1|x)$   
feature input vector  
 $0 \leq \hat{y} \leq 1$  → probability of  $y$  being 1

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

Output:  $\hat{y} = \sigma(w^T x + b)$   
Sigmoid function to be applied  
to ensure range of  $\hat{y}$  bw (0 → 1)



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

If  $z \rightarrow \text{large (+ve)}$   $\sigma(z) \approx 1$   
If  $z \rightarrow \text{large (-ve)}$   $\sigma(z) \approx 0$

L=6 → Logistic regression cost func -

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

Given:  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , we want  $\hat{y}^{(i)} \approx y^{(i)}$   
predicted  $\downarrow$  correct sample

→ Loss function: measures how large the error is b/w the predicted value of output ' $\hat{y}$ ' & the true value  $y$  (provided).

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

not convex (optimized)

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

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

we need to minimize this error  $[\max \hat{y} \leq 1]$

$$\text{If } y=1: L(\hat{y}, y) = -\log \hat{y} \leftarrow \text{want } \log \hat{y} \uparrow, \text{ want } \hat{y} \uparrow \text{ s.e. } \hat{y}=1$$

$$\text{If } y=0: L(\hat{y}, y) = -\log(1-\hat{y}) \leftarrow \text{want } \log(1-\hat{y}) \uparrow, \text{ want } \hat{y} \downarrow$$

$$\rightarrow \text{Cost fun.} \rightarrow J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

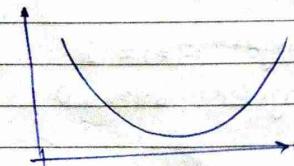
measures how well the param.  $w$  &  $b$  predict the training set

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

→ Loss function is applied to single training examples & cost function acts like avg of parameters.

We need to minimize Cost fun. (Overall)  
To find  $w, b$ )

[L=7] → Gradient Descent →



Repeat  $\ell$

$$w := w - \alpha \frac{dJ(w)}{dw}$$

learning rate

" $dJ$ "  
white implementation  
we take

$$J(w, b)$$

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

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

L=8 Computation Graph →

$$\text{Let } J(a, b, c) = 3(a+b+c)$$

$$U = b \cdot c$$

$$V = a+b \cdot c$$

$$J = 3V$$

$$\frac{dJ}{dU} = 3$$

$$L=9$$

$$\frac{dJ}{dU} = \frac{dJ}{dV} \cdot \frac{dV}{dU} = 3 \cdot \frac{dV}{dU}$$

$$3 \cdot 2 = 6$$

$$3 \cdot 2 = 6$$

$$dJ = dJ \cdot \frac{dU}{db} = 6$$

$$dJ = 6$$

$$\frac{dJ}{dc} = 9$$

$$3 \cdot 3 = 9$$

$$dJ = dJ \cdot \frac{dU}{dc} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{db} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

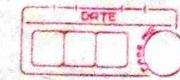
$$dJ = 9$$

$$dJ = dJ \cdot \frac{dU}{da} = 9$$

L=11 Logistic regression on iris example →  
↳ Actual algorithm

$\Rightarrow \text{d}w = np.zeros(n\_x, 2)$  instead of  $\text{d}w_1 \& \text{d}w_2 = 0$

L=14 → Vectorising Logistic Regression →



$$\begin{aligned} \rightarrow z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b & z^{(3)} &= w^T x^{(3)} + b \\ \rightarrow a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

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

$$w^T X = \underbrace{\begin{bmatrix} \cdot & \cdot & \cdot & \dots & \cdot \end{bmatrix}}_{w^T} \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{bmatrix}$$

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

$\bar{z} = \text{np.dot}(w^T, X) + b$  → real op.  
converted to  
 $\bar{z}$  calculated Broadcasting in python.  
in one-line.

L=12 → Vectorisation →

→ We need to compute:  $z = w^T x + b$

$$w = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \quad x = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}$$

Column vector

Non-vectorised:  
 $z = 0$

for i in range(n\_x):  
 $z += w[i] * x[i]$   
 $z += b$

Vectorised:  
 $z = \text{np.dot}(w, x) + b$        $w \in R^{n_x}$   
 $w^T x$   
much faster  $\therefore$  ✓

L=15

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

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

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

$$dz = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & a^{(3)} - y^{(3)} & \dots \end{bmatrix}$$

L=13 → PTR :-

→ Whenever possible, avoid explicit for-loops.

Numpy → build-in functions use ✓

## Vectorisation in Logistic Regression

$$J=0, dw_1=0, dw_2=0, db=0$$

If for  $i$  in range  $1, m$ :

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

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

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

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

vector value replaced

$$\left. \begin{array}{l} \text{loop for} \\ \text{more} \\ \text{features} \end{array} \right\} \begin{array}{l} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{array} \rightarrow \text{Efficient} \rightarrow dw += x^{(i)} * dz^{(i)}$$

$$\tilde{x} = w^T x + b$$

$$\left. \begin{array}{l} \tilde{x} = np.\text{dot}(w^T, x) + b \\ A = \sigma(\tilde{x}) \\ dz = A - y \end{array} \right\} \text{compute for all the values of } (i)$$

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

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

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

↓ Just step for gradient descent  
 ↤ A for loop will be used for multiple steps.  
 ↤ axis=0 (vertical) & (horizontal)

L=16 → Broadcasting example →

$$\text{eg 1: } \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 100 \xrightarrow{\substack{\text{python takes} \\ \text{numerical constant}}} \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} \Rightarrow \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

(4,1) vector

$$\text{eg 2: } \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 400 & 500 & 600 \end{bmatrix} \xrightarrow{\substack{(1,n) \text{ vector} \\ \text{python copies} \\ \text{this vector } m \\ \text{times to make it } m \times n \text{ vector}}} \begin{bmatrix} 101 & 202 & 303 \\ 404 & 505 & 606 \end{bmatrix}$$

L=17

Recap →

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

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

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

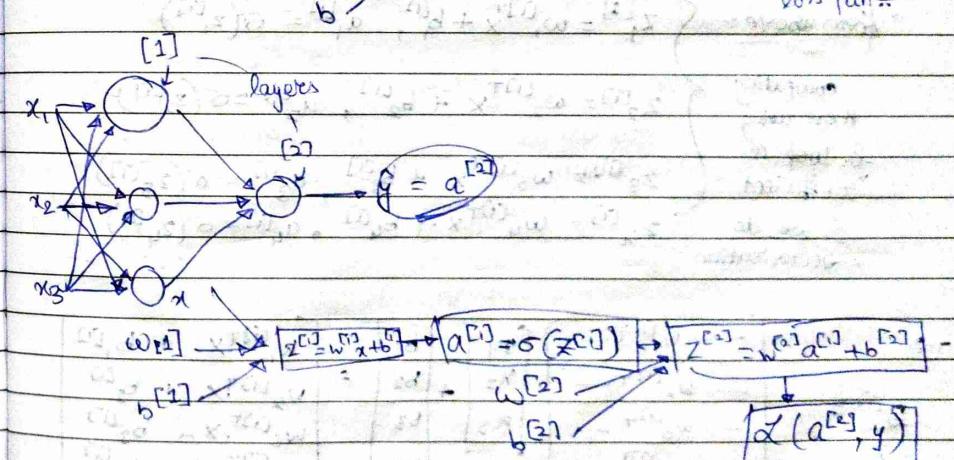
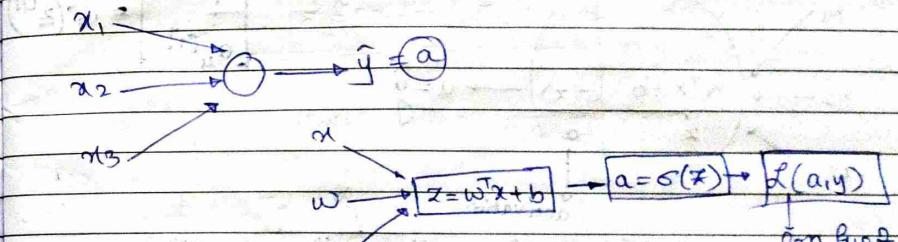
$$\log P(y|x) = \log \hat{y}^y (1-\hat{y})^{1-y}$$

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

\* last on in' exples  $\Rightarrow y$

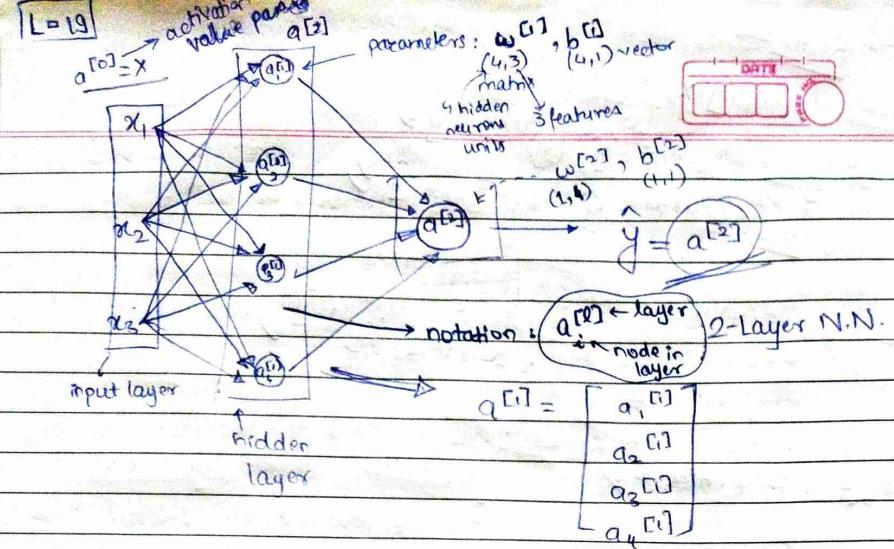
L=18

What is Neural Network →

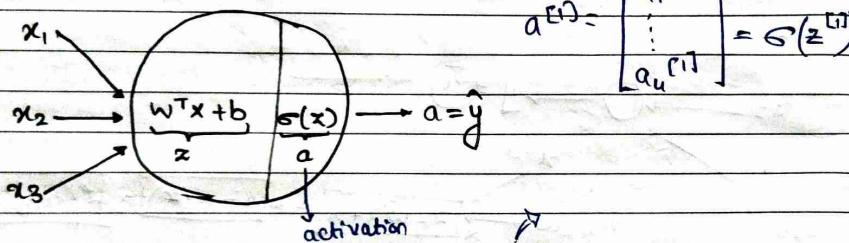


$$\begin{aligned} & z^{[1]} = w^{[1]} x^{[1]} + b^{[1]} \xrightarrow{a^{[1]} = \sigma(z^{[1]})} z^{[2]} = w^{[2]} a^{[1]} + b^{[2]} \xrightarrow{a^{[2]} = \sigma(z^{[2]})} \dots \\ & z^{[L-1]} = w^{[L-1]} x^{[L-1]} + b^{[L-1]} \xrightarrow{a^{[L]} = \sigma(z^{[L-1]})} z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \xrightarrow{a^{[L]} = \sigma(z^{[L]})} \text{out} \end{aligned}$$

$$f(x(a^{[L]}, y))$$



→ Representation →



→ From above:-

$$z_1^{[1]} = \omega_1^{[1] T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = \omega_2^{[1] T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = \omega_3^{[1] T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = \omega_4^{[1] T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

computing these using for loop is inefficient  
so we do vectorisation

→

$$z^{[1]} = \begin{bmatrix} \omega_1^{[1] T} \\ \omega_2^{[1] T} \\ \omega_3^{[1] T} \\ \omega_4^{[1] T} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} = \begin{bmatrix} \omega_1^{[1] T} x + b_1^{[1]} \\ \omega_2^{[1] T} x + b_2^{[1]} \\ \omega_3^{[1] T} x + b_3^{[1]} \\ \omega_4^{[1] T} x + b_4^{[1]} \end{bmatrix}$$

$\underbrace{\omega^{[1]}}_{W^{[1]}}$

$\underbrace{b^{[1]}}_b$

$\underbrace{\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}}_{z^{[1]}}$

Given input  $X: a^{[0]}$

$z^{[1]} = W_1^{[1]} x + b^{[1]}$

$(4,1) \quad (4,3) \quad (3,1) \quad (4,1)$

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

$\Rightarrow z^{[2]} = W_2^{[1]} a^{[1]} + b^{[2]}$

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

L = 20

Vectorisation across multiple variables/example →

→ Non-vectorised

for  $i=1$  to  $m$ ,  $\rightarrow$  inefficient

$z^{[1](i)} = W_1^{[1]} x^{(i)} + b^{[1]}$

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

$z^{[2](i)} = W_2^{[1]} a^{[1](i)} + b^{[2]}$

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

$i^{th}$  training sample

Vectorisation →

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

$\underbrace{x^{(1)}}_{\text{1st row}}, \underbrace{x^{(2)}}_{\text{2nd row}}, \dots, \underbrace{x^{(m)}}_{\text{mth row}}$

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

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

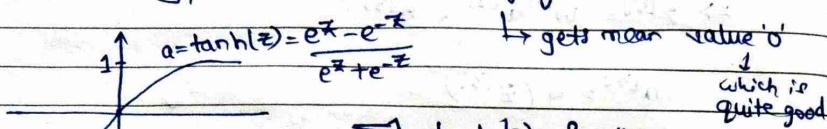
# Implementation →

① Forward Propagation examples =

L-2) Activation Functions →

↳ Sigmoid

↳  $\tanh(z)$  is always superior & better than sigmoid  
 ↳ tanh-hyperbolic ↳ ranging  $\in [-1, 1]$



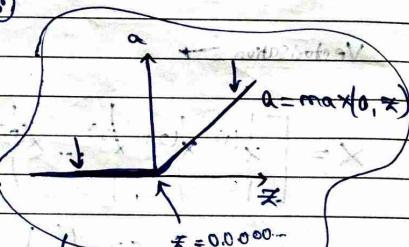
→  $\tanh(z)$  function for hidden layer.

→ Sigmoid  $g(z)$  function for output layer.

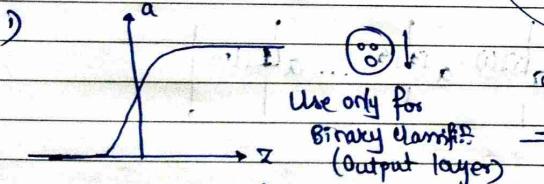
→ A common disadvantage in both these fun<sup>t</sup>s is that, when  $z$  is very large/small, the derivatives approach zero.

gradient descent can become smaller

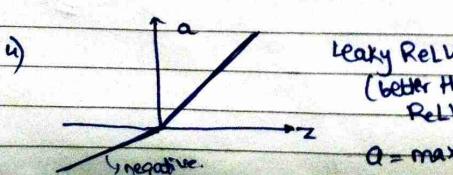
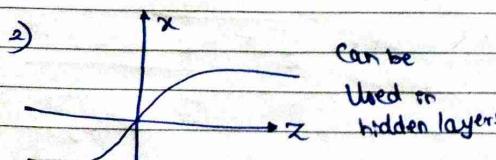
ReLU → choice!!!



\* Pros & Cons of activation fun<sup>t</sup> →



$$\text{sigmoid: } a = \frac{1}{1+e^{-z}}$$

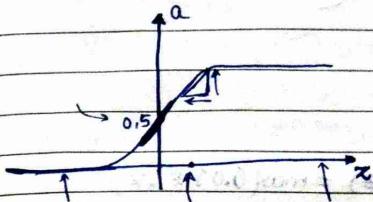


most efficient

$$a = \max(0.01z, z)$$

L-2) → When we need to implement back propagation, we need to compute derivative of the activation fun<sup>t</sup>:

1) Sigmoid Activation fun<sup>t</sup>:



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

$$z=10, g(z) \approx 1$$

$$\frac{d}{dz} g(z) \approx 1(1-1) \approx 0$$

$$z=-10, g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0(1-0) \approx 0$$

$$z=0, g(z) \approx 0.5$$

$$\frac{d}{dz} g(z) \approx 0.5(0.5) = 0.25$$

2) Tanh activation fun<sup>t</sup> →

$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

$$z=10, \tanh(z) \approx 1$$

$$\therefore g'(z) = 1 - (1)^2 = 0$$

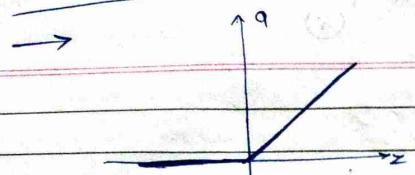
$$z=-10, \tanh(z) \approx -1$$

$$\therefore g'(z) = 1 - (-1)^2 = 0$$

$$z=0, \tanh(z) = 0$$

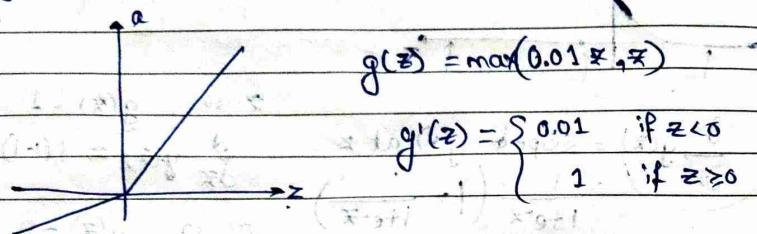
$$\therefore g'(z) = 1$$

### 3) ReLU & Leaky ReLU



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

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



$$g(z) = \max(0.01 \cdot z, z)$$

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

### [L=23] Gradient Descent for NN.

parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$   
 $(n^{[1]}, n^{[0]})$

$$\text{Cost fun} \hat{y} = J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}_i, y_i)$$

Gradient descent:

Repeat {

→ compute predicts  $(\hat{y}^{(i)}, i=1, \dots, m)$

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

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

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

### \* forward Propagation →

$$\tilde{z}^{[0]} = W^{[0]} x + b^{[0]}$$

$$A^{[1]} = g^{[1]}(\tilde{z}^{[0]})$$

$$\tilde{z}^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(\tilde{z}^{[2]}) = \sigma(\tilde{z}^{[2]}) \approx$$

### \* Back propagation →

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

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

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

"To not output funny 1-D array"

$$dZ^{[1]} = \underbrace{W^{[1]T} dz^{[2]}}_{(n^{[1]}, m)} * \underbrace{g^{[1]}'(z^{[1]})}_{\text{product (row, m)}} \quad \text{*(element-wise)}$$

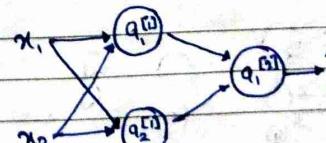
$$dW^{[1]} = \frac{1}{m} dZ^{[1]} x^T$$

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

### [L=24]

#### → Random Initialisation =

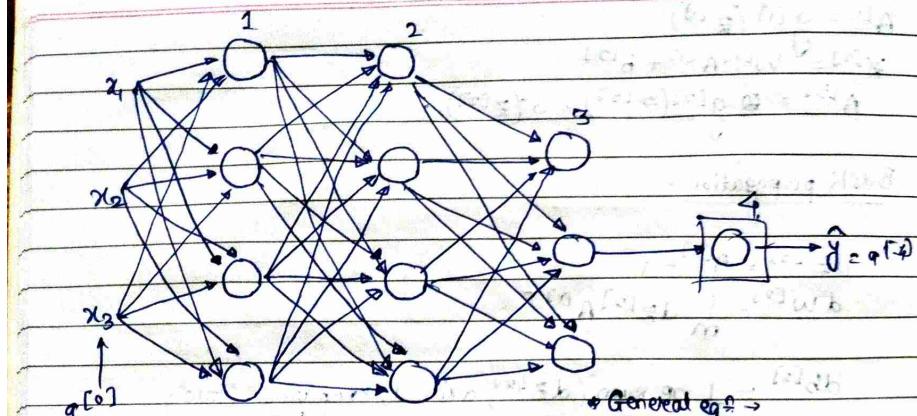
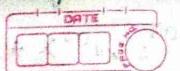
→ Initiating  $W=0, b=0$  won't work as multiple iterations gives exact value so no good outcome of training the model (NN)



Randomly →

$W^{[1]} = \text{np.random.randn}((2, 2)) * 0.01$   
 $b^{[1]} = \text{np.zeros}((2, 1))$   
 ↓  
 → Key to initialize to 0  
 $W^{[2]} = \text{np.random.randn}((1, 2)) * 0.01$   
 $b^{[2]} = 0.0$   
 there may  
 be better  
 constants  
 than 0.01  
 (but relatively  
 small no.)

L-26

Forward Network in Deep Network →  
propagation

$$x : z^l = W^l x + b^l$$

feature vector (input)

$$\begin{aligned} z^l &= W^l a^{l-1} + b^l \\ a^l &= g^l(z^l) \end{aligned}$$

\* Vectorize  $x = a^0$

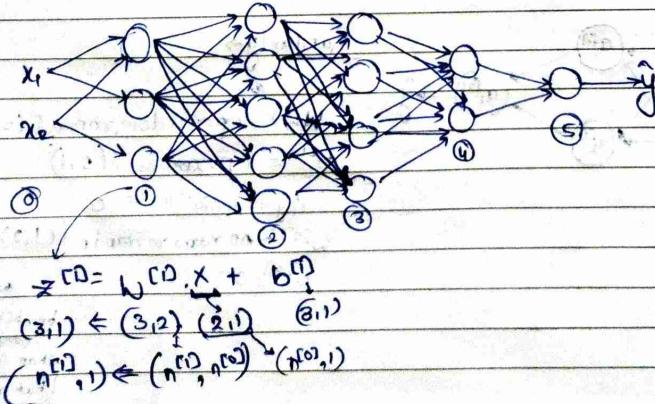
$$z^2 = W^2 a^1 + b^2 \quad \rightarrow z^1 = W^1 x + b^1$$

For layer 2:  $a^2 = g^2(z^1)$        $A^1 = g^1(z^1)$

$$\begin{aligned} z^3 &= W^3 a^2 + b^3 \\ a^3 &= g^3(z^2) \end{aligned}$$

$A^2 = g^2(z^2)$

→ Stacking

L-26 → Dimensions for parameters  $W^l$  &  $b^l$  → later part

$$W^l : (n^l, n^0)$$

$$W^2 : (5,3) \quad (n^2, n^1)$$

$$(z^2) = W^2 \cdot a^1 + b^2$$

$$(5,3) \quad (5,3) \quad (3,1) \quad (1,1) \Rightarrow W^2 : (n^2, n^{l-1})$$

$$\rightarrow b^2 : (n^2, 1)$$

$$W^3 : (4,5) \quad W^4 : (1,2)$$

$$W^0 : (2,4)$$

$$dW^2 : (n^2, n^{l-1})$$

$$db^2 : (n^2, 1)$$

As in training samples

$$(z^1) \rightarrow : (n^1, m)$$

↳  $z^1 \rightarrow$  pythone broadcasting  
↳  $b \rightarrow$  pythone broadcasting

→ Deep Representation →

simple funt's like

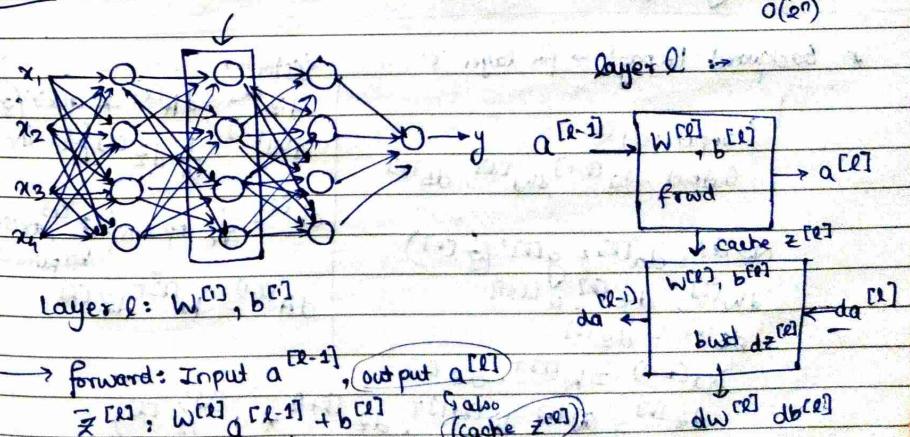
→ Early layers can detect the edges of the image, and combining them together in later layers for learning more complex funt's

→ Extracting features ↴

→ When 'not more' hidden layers are needed, then we need to have exponentially large no. of hidden units.

$$\sim 2^{(n-1)} \\ O(2^n)$$

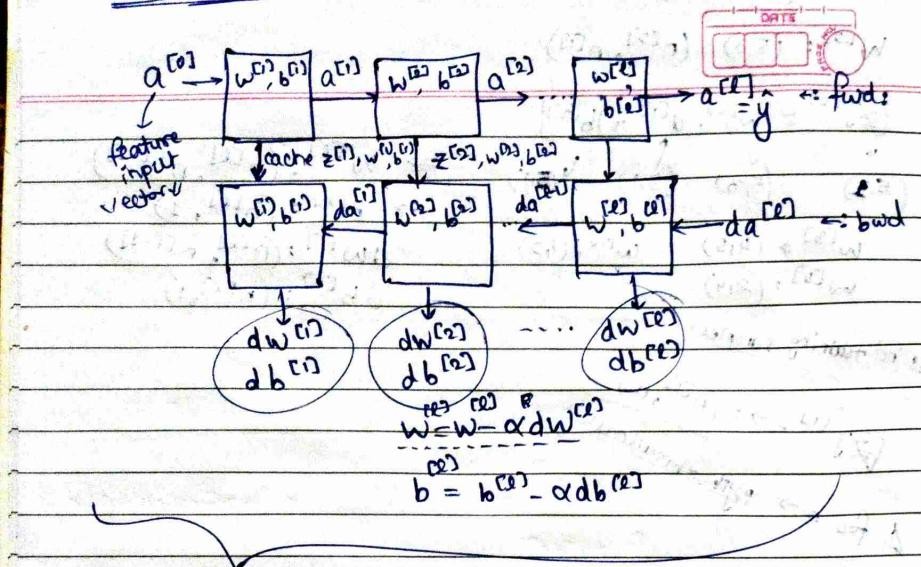
L-27



→ forward: Input  $a^{l-1}$ , output  $a^l$   
 $\bar{z}^l, W^l, a^{l-1} + b^l$  also (cache  $z^l$ )  
 $a^l : g^l(z^l)$

→ backward: Input  $da^l$ , output  $da^{l-1}$   
cache  $z^l$   
 $dw^l, db^l$

## → Basic Mechanism →



One iteration of gradient descent

\* forward Propag. → for layer 'l'

Input: \$a^{[l-1]}

Output: \$a^{[l]}\$, cache (\$z^{[l]}\$)

Vector P. →

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

\* Backward Propag. → for layer 'l'

Input: \$da^{[l]}

Output: \$da^{[l-1]}, dw^{[l]}, db^{[l]}

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

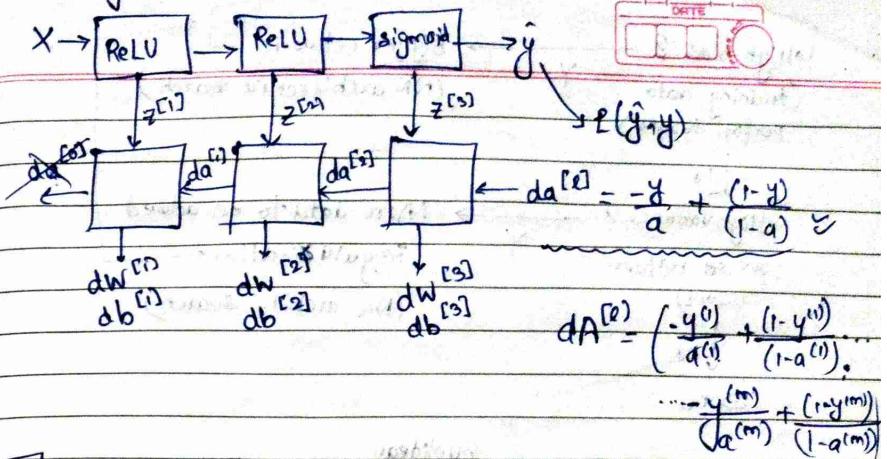
$$dw^{[l]} = dz^{[l]} \cdot a^{[l-1]}$$

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

$$da^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

$$dz^{[l]} = (W^{[l+1]})^T \cdot dz^{[l+1]} * g^{[l]}'(z^{[l]})$$

## → Summary:



$L=28$

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

→ Hyperparameters: Learning rate ( $\alpha$ )

# iterations

(# : no. of) control

# hidden layers

# hidden units:  $n^{[1]}, n^{[2]}, \dots$

Choice of activ. fun<sup>2</sup>

Course = 2 →

→ Machine learning model

Training set

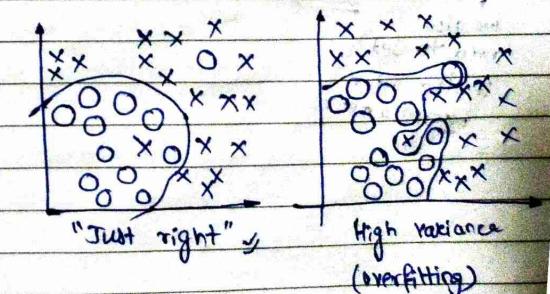
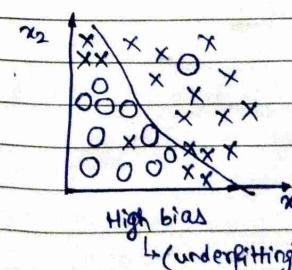
→ Hold-out / Cross-Valid. / Development (Dev) set

Test set

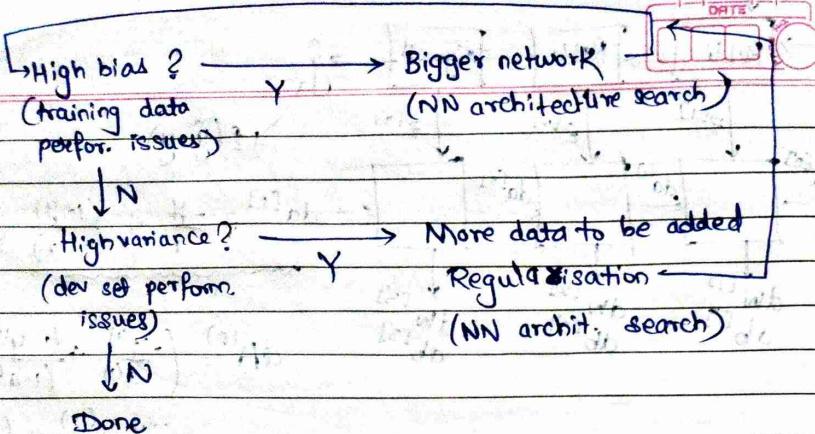
→ Make sure dev and test set come from same distribution

freedom

→ Bias & Variance:



## → Basic Recipe:



## → Neural Network

$$J(w^{[1]}, w^{[2]}, w^{[3]}, \dots, w^{[L]}, b^{[1]}, b^{[2]}, b^{[3]}, \dots, b^{[L]})$$

$$= \frac{1}{m} \sum_{i=1}^m d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l-1]}} (w_{ij}^{[l]})^2 \quad w \in (n^{[L-1]}, n^{[L]})$$

"Frobenius norm"

$$\begin{aligned} dW^{[l]} &= (\text{from backprop.}) + \frac{\lambda}{m} w^{[l]} \\ w^{[l]} &:= w^{[l]} - \alpha dW \end{aligned}$$

"Weight decay"

$$w^{[l]} := w^{[l]} - \alpha [(\text{from backprop.}) + \frac{\lambda}{m} w^{[l]}]$$

$$\text{using step 3 trick: } (1 - \frac{\alpha \lambda}{m}) w^{[l]} = [w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]}] - \alpha (\text{from backprop.})$$

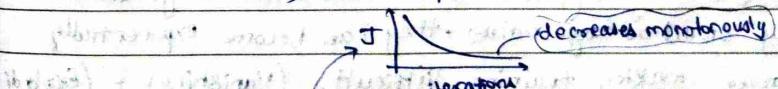
→ Regularization reduces overfitting !!

large 'λ' can have negative effects on our model if too large.

- If λ is too large,  
 $w^{[l]} \downarrow \rightarrow$  losing weights  
(losing relevant features)

→ To debug the gradient descent →

Create a plot of the cost function J to the no. of iterations



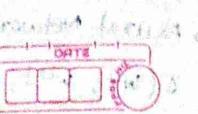
$$\Rightarrow \boxed{J(\dots) = \sum_i d(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{[l]}\|_F^2}$$

## \* Dropout Regularization

→ Based on a "specific probability" some nodes from the hidden layers are dropped out & with the few left nodes we train our model. (May seem vague but works! ☺)

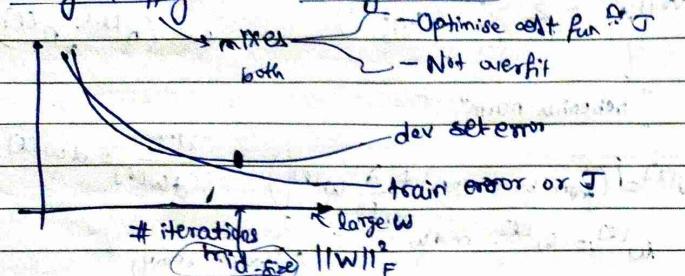
→ Data Augmentation (Regularization)

- A few changes in data
- If an image is provided, flip it/rotate it, we get new data to train & be efficient.



→ Early Stopping →

\* Orthogonalisation:



→ Variance:

$$\text{Var}(w_i) = \frac{\sigma}{n} \rightarrow \text{for ReLU}$$

→  $W^{[L]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{[L-1]}}\right)$

→ for tanh:

$$* \tanh = \sqrt{\frac{1}{n^{[L-1]}}}$$

Xavier's random initialization

$$\sqrt{\frac{2}{n^{[L-1]} + n^{[L]}}}$$

# Derivative/Gradient Computation:

\* Normalizing training data → (Done so that we get a symmetric structure easy for gradient descent to take steps).

→ i) Subtract mean: descent to take steps).

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu$$

→ Normalize variance:

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} * x^{(i)T}$$

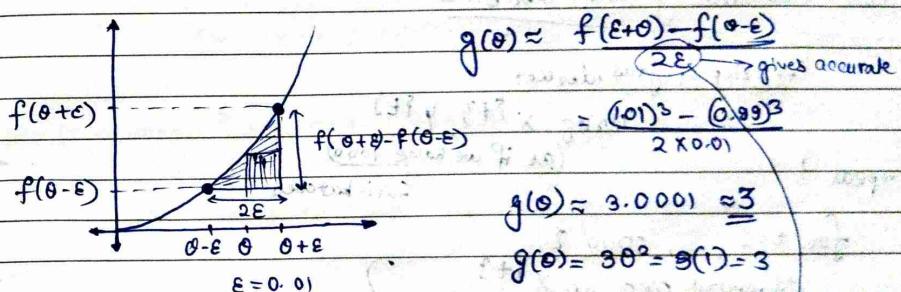
→ The model takes a very long time if the derivatives or gradients take a very small or large values, they can become exponentially smaller/larger making training difficult. (Vanishing) + (Exploding)

$w^{[L]} >$  Identity matrix (large-very)

$w^{[L]} <$  Identity matrix (very small)

$$z = w_0 x_0 + w_1 x_1 + \dots + w_n x_n + b$$

large  $n \rightarrow$  smaller  $w_i$   
no. of input features



→ We take  $w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}$  & reshape it into a big vector 'o'

$$J(o) = J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]})$$

→ We take  $dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}$  and reshape into a big vector 'dθ'.

→ for each 'i',

$$\rightarrow d\theta_{\text{approx.}[i]} = J(o, o_2, \dots, o_{i+\epsilon}, \dots) - J(o, o_2, \dots, o_{i-\epsilon}, \dots)$$

$$\approx d\theta[i] = \frac{\partial J}{\partial o_i} \quad | \quad d\theta_{\text{approx.}} \approx d\theta$$

Check,  $\frac{\|d\theta_{\text{approx.}} - d\theta\|_1}{\|d\theta_{\text{approx.}}\|_2 + \|d\theta\|_2} \approx 10^{-2}$  or smaller → then great! approximation

$$\epsilon = 10^{-2}$$

\* Tip →

- Don't use Grad. Checking in training - to debug
- If algo. fails gradcheck, look at other components to identify the bug  
i.e. if  $d\theta$  approx. is very far from  $d\theta$ , then for diff  $i$ , which values give very diff.

→ Remember Regularization

→ Doesn't work with Dropout

→ If we have a very large nos. of training samples we divide them in few groups (called batches).  
(mini-)

$$X^{[t]}, y^{[t]}$$

→ Mini-batch gradient descent =

→ 1 step of grad. descent

Using  $X^{[t]}, y^{[t]}$

(as if we have 1000)

repeat {

for  $t=1, \dots, 5000$  }

forward prop. on  $X^{[t]}$

$$\hat{y}^{[t]} = W^{[t]} X^{[t]} + b^{[t]}$$

vector implementation

(1000 examples)

$$A^{[t]} = g(\hat{y}^{[t]})$$

Now we can process 1000 mini-batches of training sets than large amounts together.

→ Compute cost fun.:

$$J^{[t]} = \frac{1}{1000} \sum_{i=1}^1 L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum \|W^{[t]}\|_F^2$$

→ backprop. to compute grad. over  $J^{[t]}$  (using  $X^{[t]}, y^{[t]}$ )

$$W^{[t+1]} = W^{[t]} - \alpha dW^{[t]}, b^{[t+1]} = b^{[t]} - \alpha db^{[t]}$$

3

this is 1 pass → "1 epoch"

→ Mini-size batch size must be between (1, m); not too small or too large.  
→ If size = 1 → (Stochastic grad. descent)  
↳ loses speed up from vectorisation  
(never hits minimum can lead to wrong direction)

→ If size = m → Batch grad. descent

↳ Too long time per iteration

• Guidelines →

→ for small training set: Use batch grad. descent.  
( $m \leq 2000$ )

→ for large — : Use minibatch size

$m \in (64, 128, 256, 512)$

→ Make sure mini-batches  $X^{[t]}, y^{[t]}$  fits in CPU/GPU.  
 $1024$   
 $2^P \rightarrow$  efficient

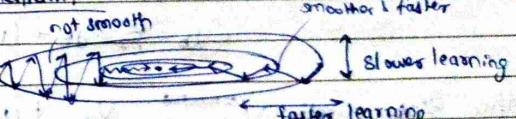
→ To implement Exponentially Weighted Averages =

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$V_t$  as approximately avg. over:

$$\frac{1}{1-\beta}$$

→ Gradient descent momentum →



→ Momentum:

On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

$$VdW = \beta VdW + (1-\beta)dW$$

$$Vdb = \beta Vdb + (1-\beta)db$$

Friction  $\Rightarrow$   $w = w - \alpha VdW$ ,  $b = b - \alpha Vdb$

velocity  $\downarrow$  as  $\alpha$  omitted  
 $dW$  may be in few cases

make smoother

Hyper-parameters  $\rightarrow \alpha, \beta$

$\beta = 0.9 \downarrow$  (pretty well)

→ Root mean square prop →  $\beta_2$ : hyperparameters  
(RMS)

element-wise:



$$\textcircled{1} \rightarrow S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \rightarrow \text{small}$$

$$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \leftarrow \text{large}$$

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

$\sqrt{S_{dw}}$

$\sqrt{S_{db}}$

delete if  $\approx 0$ , we can have large  $\alpha$  without much deviation  
in our oscillations.

### Adam Optimization algorithm →

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On Iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw, \quad \} \text{momentum}$$

$$V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

combined  $\times$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2 \quad \} \text{RMS}$$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$$\text{After Bias correction: } V_{dw}^{\text{corrected}} = V_{dw} / (1-\beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1-\beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1-\beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1-\beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}, \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

Hyperparameters:

- Very effective ↴
- Imp.  $\alpha$ : needs to be tuned oftenly
- $\beta_1$ : 0.9 (dw)
- $\beta_2$ : 0.999 (dw<sup>2</sup>)
- $\epsilon$ :  $10^{-8}$  (recommended)

→ Learning Rate Decay

1 epoch = 1 pass

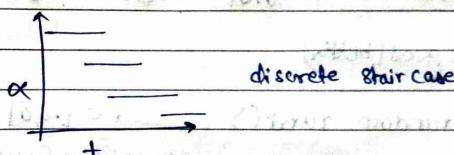


$$\alpha = \frac{\alpha_0}{1 + \text{decay rate} \times \text{epoch\_num}}$$

→ Other formulas:

$$\alpha = 0.95^{\text{epoch\_num}} \cdot \alpha_0 \rightarrow \text{exponentially decay rate}$$

$$\alpha = \frac{K(\text{some constant})}{\sqrt{\text{epoch\_num}}} \cdot \alpha_0 \quad \text{or} \quad \alpha = \frac{K}{\sqrt{t}} \cdot \alpha_0$$



### \* Tuning Process →

→ Priority for hyperparameters to be tuned:

- ①  $\alpha \rightarrow$  most imp.
- ②  $\beta$  [default] = 0.9  
 $\beta_{prop}$ :  $\beta_1 \sim 0.99, \beta_2 \sim 10^{-8}$
- ③ no. of hidden units
- ④ mini-batch size
- ⑤ no. of layers
- ⑥ learning rate decay

→ Chopping random values rather than grid allows us to explore diff. possible values for most imp. hyperparameters.

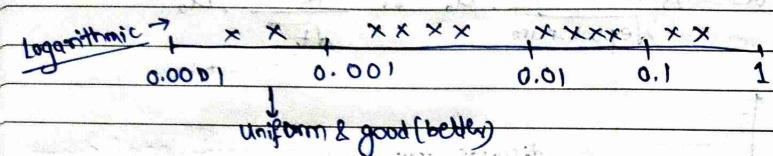
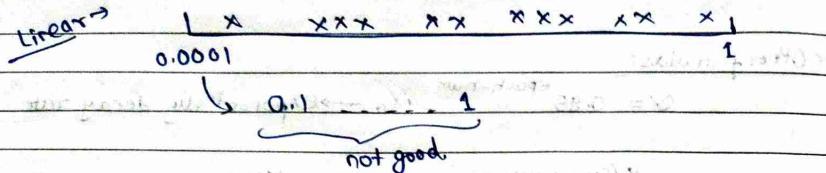
→ In Coarse to fine approach from a bigger area we get a smaller region which have better values and get our hyperparameters (random)

→ Logarithmic scale makes more sense & works better than a linear scale.



→ Appropriate scale for hyperparameters →

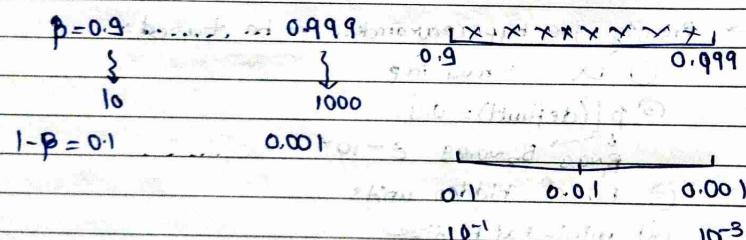
let say  $\alpha \in [0.0001, \dots, 1]$



$$\tau = -4 * \text{np.random.rand}() \quad \leftarrow \tau \in [-4, 0]$$

$$\alpha = 10^\tau \quad \leftarrow \tau \in [+10^{-4}, 10^0]$$

→ for exp. weightage averages -



$$\tau \in [-3, -1]$$

$$\beta = 1 - 10^\tau$$

- \* Two approaches →
  - 1) Panda approach (babysitting the model)
    - ↳ when we have large models & diff. to train
  - 2) Carrier approach (training models parallelly)
    - ↳ pick the best one
    - ↳ enough computers to train models

→ Normalizing input features can speed up learning.

$$\rightarrow U = \frac{1}{m} \sum_i x^{(i)}$$

$$x = x - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i (x^{(i)} - \mu)^2$$

$$x = x / \sigma$$

→ But activations are also included, Can we normalize  $a^{[2]}$  so as to train  $w^{[2]}, b^{[2]}$  faster?

→ Given some intermediate values in NN. → R<sup>th</sup> layer

$$U = \frac{1}{m} \sum_i z^{(i)} \rightarrow \bar{z}^{(m)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \bar{z}^{(m)})^2$$

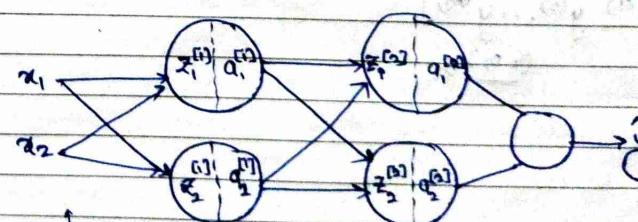
$$\boxed{\begin{aligned} z^{(i)}_{\text{norm.}} &= \frac{z^{(i)} - \bar{z}^{(m)}}{\sqrt{\sigma^2 + \epsilon}} & \Rightarrow \bar{z}^{(i)} &= \gamma z^{(i)}_{\text{norm.}} + \beta \\ \text{don't want mean as 0} & & \text{learnable parameters of the model} & \\ \text{and } \sigma^2 & & \text{(allows set mean other than 0/1)} & \\ \text{variance} & & & \end{aligned}}$$

$$\boxed{\begin{aligned} \gamma &= \sqrt{\sigma^2 + \epsilon} \\ \text{true? : } \bar{z}^{(i)} &= \bar{z}^{(m)} \\ \beta &= \bar{U} \end{aligned}}$$

!!

→ Use  $\tilde{z}^{(i)}$  instead of  $z^{(i)}$  ↳ !!

\* Fitting →

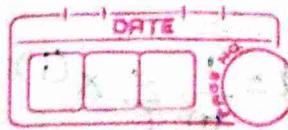


Without norm:

$$\boxed{\begin{aligned} \text{With norm: } X &\xrightarrow{w^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\text{Batch (BN) norm}} \tilde{z}^{[1]} \xrightarrow{w^{[2]}, b^{[2]}} a^{[2]} = g^{[2]}(\tilde{z}^{[2]}) \\ &\dots = \hat{y} \end{aligned}}$$

## → Softmax

Suppose we have 4 classes,



$$\Omega = 4$$

(4,1) :-

$$\text{let } z^{[1]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

Softmax  
activ. fun.  
 $\rightarrow g^{(1)}(\cdot)$

$$\text{S.A.F.} \rightarrow g^{(1)}(z^{[1]}) = \begin{bmatrix} e^5 / \sum_{j=1}^4 t_i \\ e^2 / \sum_{j=1}^4 t_i \\ e^{-1} / \sum_{j=1}^4 t_i \\ e^3 / \sum_{j=1}^4 t_i \end{bmatrix} = \begin{bmatrix} 0.0842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix} \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

## → Loss function :-

$$\hat{y}^{(0)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ cat} \quad y_1 = 1 \quad y_2 = 0$$

$$a^{[2]} = \hat{y}^{(1)} = \begin{bmatrix} 0.9 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad (\text{c} = 4)$$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^n y_j \log \hat{y}_j$$

$$= -y_2 \log \hat{y}_2$$

$$= -\log \hat{y}_2$$

(make  $\hat{y}_2 \rightarrow \text{large}$ )

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

Zero gradient  $\rightarrow$  Saddle point ✓  
(derivative)

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

$(\underbrace{4, m})_m$