

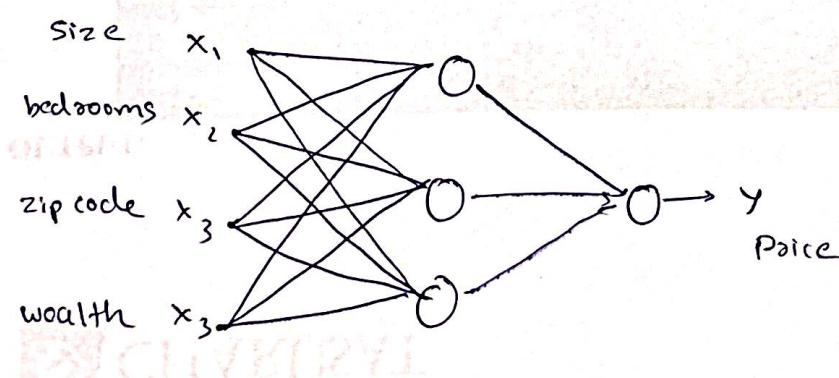
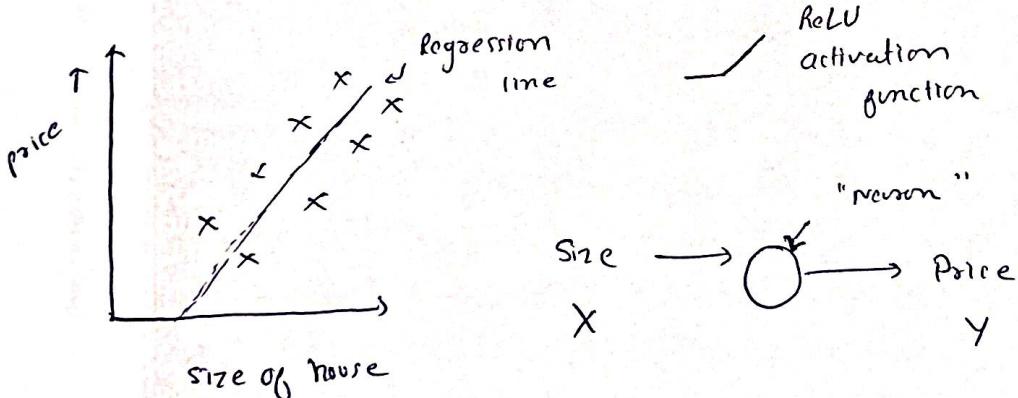
# Deep Learning

Course 1

## \* Learning outcomes

1. Neural networks and deep learning
2. Improve Deep learning networks. Hyperparameter tuning, Regularization and optimization
3. Structuring your ml project [ end-to-end DL ]
4. Convolutional Neural Networks
5. Natural Language Processing : Building sequence models [ RNN, LSTM ]

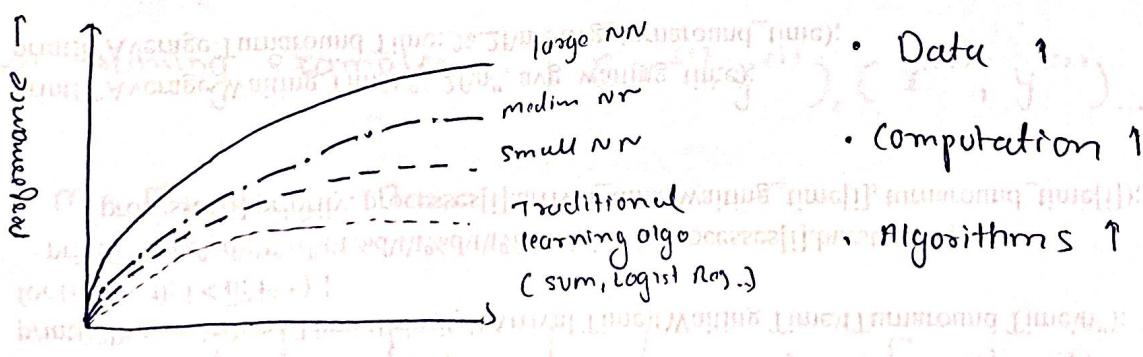
### 1. Neural Networks & DL



## → Supervised Learning [ labeled data ] \*\*

Input (Features)	Output (Label)	Application
Home feat.	Price	Real Estate
Ad. user info	Click on ad?	Online Ad.
Image	Object (1...1000)	Photo tagging
Audio	Text transcript	Speech recog.
English	Chinese	Machine Translation
Image, Radar info	Position of cars	Autonomous driving
		Custom / complex Hybrid NN

## ⇒ Scale and deep learning progress

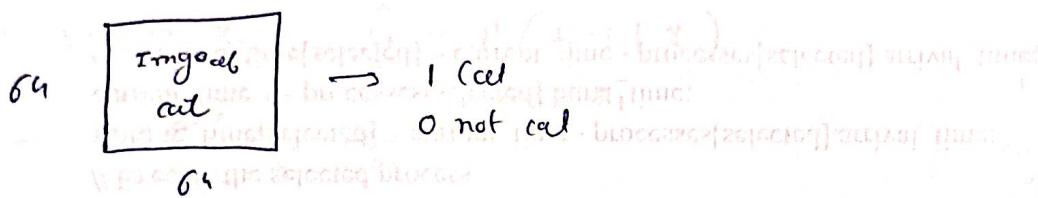


Amount of data →  
 (labeled data)  
 $(X, y)$  for training

## CHYBNOV

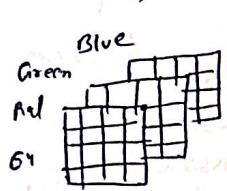
## \* Binary classification

DEFINITION



DEFINITION

-  $\text{Representation} \rightarrow \text{feature vector } x \in \mathbb{R}^{n_x}$



- feature vector  
vector

$$x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(m)} \end{bmatrix}$$

no. of feature

$$= 64 \times 64 \times 3$$

$$n_x = 12288$$

$$\text{or } n = 12288$$

$\Rightarrow$  Notation

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

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

$$- M = M_{\text{train}}; M_{\text{test}}$$

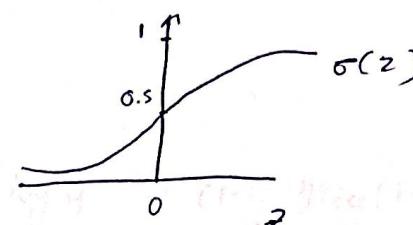
(number of training samples; number of test samples)

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

$\underbrace{\hspace{1cm}}_m \quad \underbrace{\hspace{1cm}}_m$

$$X \in \mathbb{R}^{n_x \times m} \quad x.\text{shape} = (n_x, m) \quad y.\text{shape} = (1, m)$$

## \* logistic Regression

- DID YOU  
GET IT?
- Given  $x$ ,  $\hat{y} = P(y=1|x)$   
 $0 \leq \hat{y} \leq 1$
  - $x \in \mathbb{R}^{n_x}$   $\rightarrow$  feature vector
  - Parameters:  $w \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$
  - Output  $\hat{y} = \sigma(w^T x + b)$
  - We have to learn  $w$  &  $b$  parameters so, that we can get  $\hat{y}$  as the probability  $P(y=1|x)$ .
- $\sigma(z) = \frac{1}{1+e^{-z}}$
- 

## Cost function

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

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

- loss function: must be as small as possible.

$L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$ , But, we don't use this because gradient descent doesn't work properly

$\Rightarrow L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log(1-\hat{y}))$  single training example

If  $y=1$ :  $L(\hat{y}, y) = -\log \hat{y}$   $\leftarrow$  want  $\log \hat{y}$  large, want  $\hat{y}$  large  $\approx 1$

If  $y=0$ :  $L(\hat{y}, y) = -\log(1-\hat{y})$   $\leftarrow$  want  $\log(1-\hat{y})$  large, want  $\hat{y}$  small  $\approx 0$

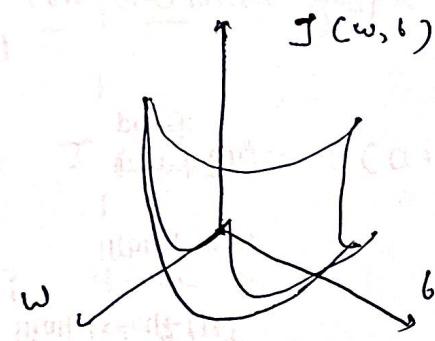
- cost function:

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

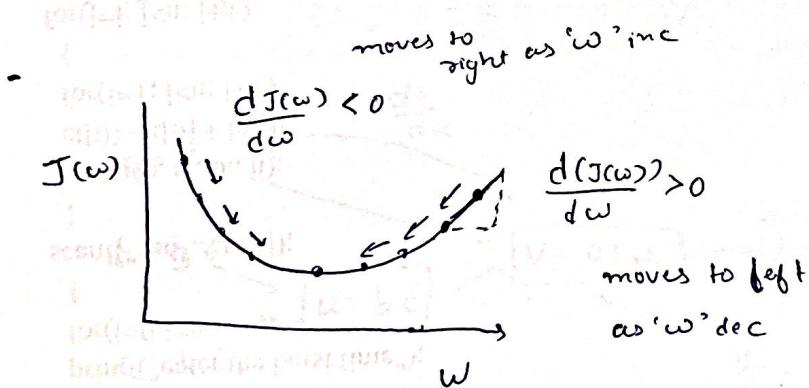
see how  $w \in \mathbb{R}^n$  are working on other dataset.

## \* Gradient Descent

- $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$
- $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})$
- Want to find optimal  $w$  &  $b$  that minimize  $J(w, b)$



-  $J(w, b)$  dimension is always perpendicular to all the parameter dimension.



$$w := w - \alpha \left[ \frac{d J(w)}{d w} \right]$$

$$\therefore w := w - \alpha \frac{d J(w)}{d w} \quad \text{or} \quad \frac{\partial J(w, b)}{\partial w}$$

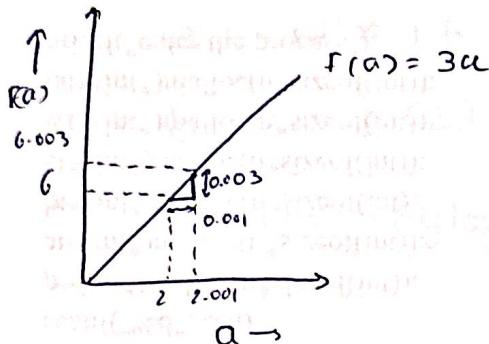
$d \rightarrow$  for single param.  
 $\partial$  - for multiple param.

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

RECAP

$\Rightarrow$  Intuition about derivatives

Intuition



Change of output

$$\text{slope of } f(a) = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.001} = 3$$

$$\boxed{\text{Derivative of } f(a)} = \frac{df(a)}{da} \rightarrow \frac{\text{change in } f(a)}{\text{change in } a}$$

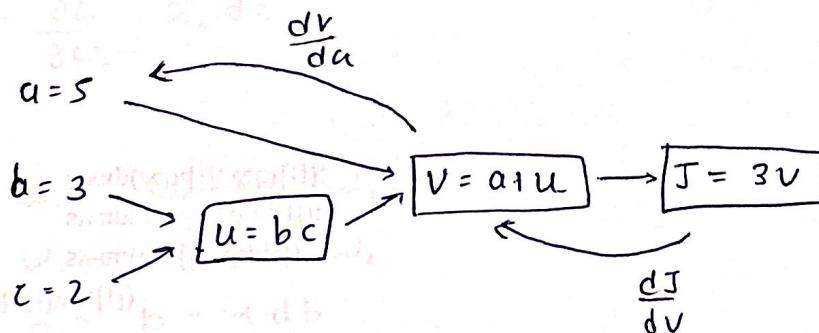
$\Rightarrow$  computational graph

$$J(a, b, c) = 3(a + bc)$$

$$u = bc$$

$$v = a + u$$

$$J = 3v$$



Computational graph

$$\therefore \frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

$$\boxed{\frac{d \text{Final output var}}{d \text{var}}}$$

$$\therefore \frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du}$$

$$\therefore \frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db}$$

## $\Rightarrow$ logistic regression Gradient Descent

With Data Class (output plane)

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

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

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

zero cost function for 0 or 1

for 0.5

</div

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

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

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

$$J/m = m$$

$$d\omega_1/m; d\omega_2/m; dl/m,$$

Möller - Möller et al. Moller

Brillouin - Brattain - Brattain

$$\omega_1 := \omega_1 - \alpha d\omega_1$$

$$\omega_2 := \omega_2 - \alpha d\omega_2$$

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

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

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

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

## Vectorization

- Used to remove explicit for loops so computational time can be reduced.

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

Non-vectorized

$$z = 0$$

for i in range(n):

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

$$z += b$$

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

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

for i in range(n):

$$u[i] = math.exp(v[i])$$

$$\omega = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Vectorized

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

$$z += b$$

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

(9)

## Vectorization of logistic regression

$$J = 0, \{ \cancel{dw_1 = 0}, \cancel{dw_2 = 0} \}, db = 0$$

$$dw = \text{np. zeros}((n_x, 1))$$

↓  
vector

for  $i = 1 \dots m$ :

$$z^{(i)} = w^T x^{(i)} + b \quad [x^{(i)}]$$

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

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

$$dz^{(i)} = a^{(i)}(1-a^{(i)}) \cdot y^{(i)} \quad [y^{(i)}]$$

for  $j = 1 \dots n_x$ :

$$\left\{ \begin{array}{l} dw_1 = x_1^{(i)} dz^{(i)} \\ dw_2 = x_2^{(i)} dz^{(i)} \end{array} \right\} \quad n_x = 2 \quad dw = x^{(i)} dz^{(i)}$$

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

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

$$dw/m$$

$$\Rightarrow X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad (n_x \times m)$$

$$w^T = [\cancel{w_1} \ \cancel{w_2} \ \dots]$$

$$w^T = [ \dots ]$$

$$z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ \dots \ b]$$

$$= [w^T x^{(1)} + b, w^T x^{(2)} + b, \dots, w^T x^{(m)} + b]$$

Syntax:  $z = np. dot(w^T, x) + b$

$A = [0^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = \sigma(z)$

↓  
Real no converted  
to  $n \times 1$  vector

## $\Rightarrow$ Vectorizing logistic Regression

Input Data

$$\text{d}z^{(1)} = \sigma^{(1)} - y^{(1)} \quad \text{d}z^{(2)} = \sigma^{(2)} - y^{(2)}$$

$$\Rightarrow \text{d}z = [\text{d}z^{(1)} \ \text{d}z^{(2)} \dots \ \text{d}z^{(m)}]_{1 \times m}$$

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

$$\Rightarrow \text{d}z = A - Y = [\sigma^{(1)} - y^{(1)} \ \sigma^{(2)} - y^{(2)} \dots]$$

$$\text{d}w = 0$$

$$\text{d}w^1 = x^{(1)} \text{d}z^{(1)}$$

$$\text{d}w^2 = x^{(2)} \text{d}z^{(2)}$$

:

$$\text{d}w^1 = m$$

$$\text{d}b = 0$$

$$\text{d}b^1 = \text{d}z^{(1)}$$

$$\text{d}b^2 = \text{d}z^{(2)}$$

:

$$\text{d}b^1 = \text{d}z^{(m)}$$

$$\text{d}b^1 = m$$

Vectorized Logistic reg.  $*+*$ 

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

$$= \text{np. dot}(w.T, x) + b$$

$$A = \sigma(Z)$$

$$\text{d}Z = A - Y$$

$$\text{d}w = \frac{1}{m} x \text{d}Z^T$$

$$\text{d}b = \frac{1}{m} \text{np. sum}(\text{d}Z)$$

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

$$b := b - \alpha \text{d}b$$

$$\text{d}w = \frac{1}{m} x \text{d}Z^T$$

$$= \frac{1}{m} \left[ x^1 \dots x^m \right] \begin{bmatrix} \text{d}Z^1 \\ \vdots \\ \text{d}Z^m \end{bmatrix}$$

$$= \frac{1}{m} \left[ x^1 \text{d}Z^1 + \dots + x^m \text{d}Z^m \right]_n$$

$\Rightarrow$  Broadcasting

## broadcasting

$$(m, n) \xrightarrow{\text{matrix}} (1, n) \rightsquigarrow (m, n)$$

Eq

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

⇒ Python - Numpy Vectors

-  $a = \text{np. random} - \text{randn}(5)$

CODE

a.shape = (5,) { don't use

"thank I carry"

-  $a = np.random.rand(5, 1) \rightarrow a.shape = (5, 1)$

$a = \text{np. random. randn}(1, 5) \rightarrow a. \text{shape} = (1, 5)$  now vector ✓

- assert Ca.shape == (5,1))

$a = a \cdot \text{reshape}((5, 1))$

## \* logistic regression cost function

QUESTION

Other observation

13

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

- Interpretation  $\hat{y} = p(y=1|x)$

$$\text{If } y=1 : p(y|x) = \hat{y}$$

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

CODE:

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

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

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

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

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

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

$\Rightarrow$  cost of m examples

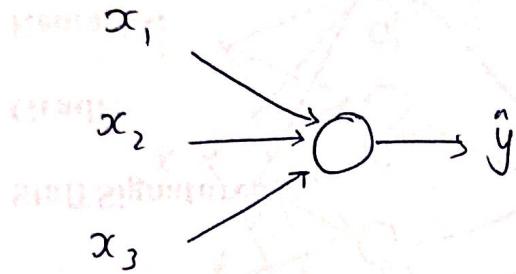
$$\log P(\text{labels in training set}) = \log \prod_{i=1}^m p(y^{(i)}|x^{(i)})$$

$$\log P(\dots) = \sum_{i=1}^m \log p(\hat{y}^{(i)}|x^{(i)})$$

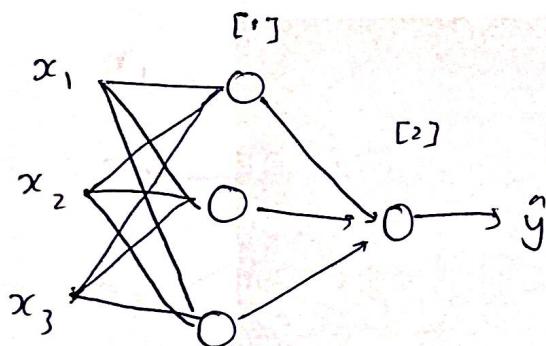
$$= - \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

cost :  $J(\omega, b) = \frac{1}{m} \sum_{j=1}^m L(\hat{y}^{(j)}, y^{(j)})$   
(minimize)

## \* What is Neural Network

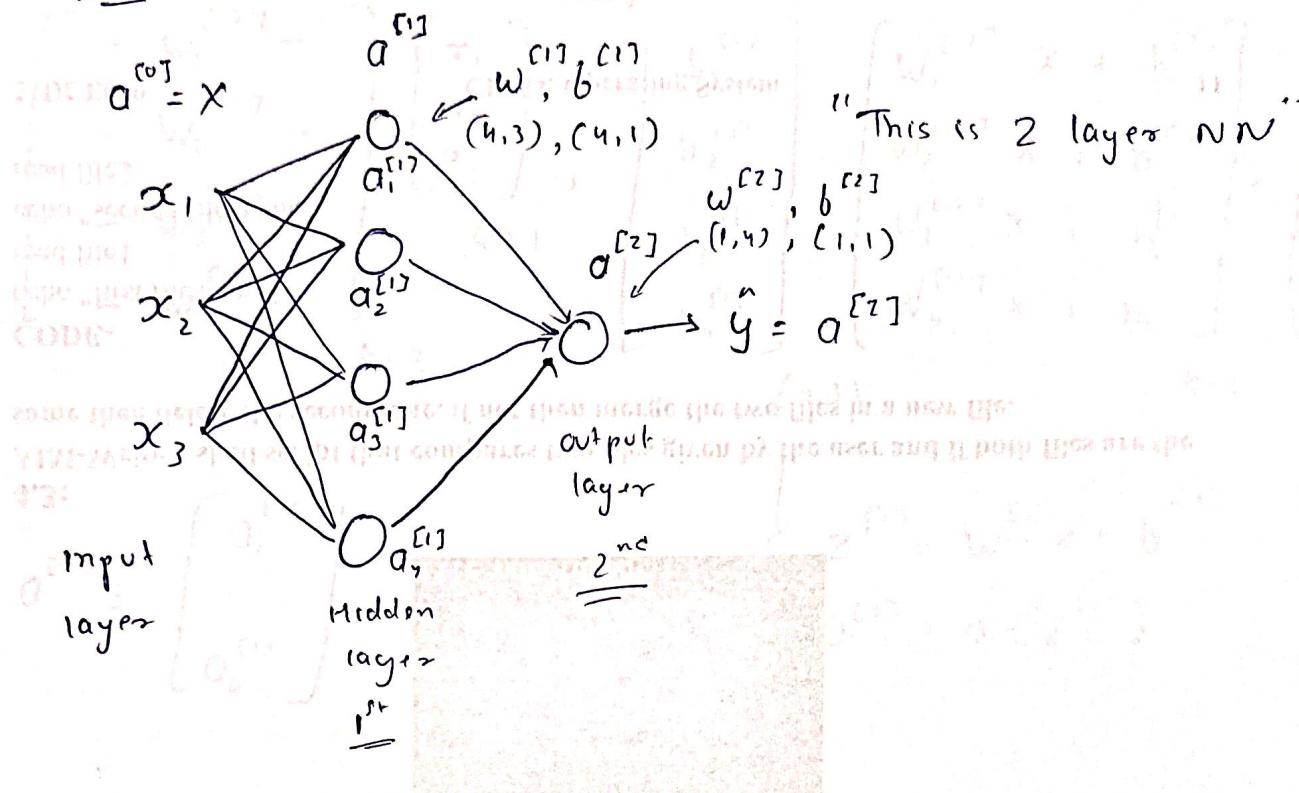


$$\begin{aligned} \text{Input } x &\rightarrow \boxed{Z = W^T x + b} \\ \omega &\rightarrow \boxed{a = \sigma(Z)} \rightarrow \boxed{L(a, y)} \\ b & \end{aligned}$$

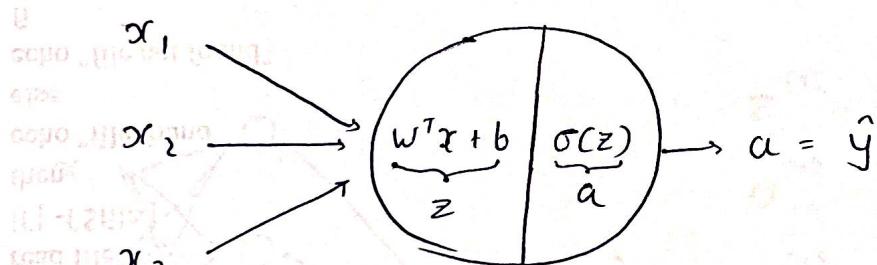


$$\begin{aligned} \frac{dW^{[1]}}{db^{[1]}} & \leftarrow \frac{da^{[1]}}{dz^{[2]}} \leftarrow \frac{dZ^{[2]}}{dZ^{[2]}} \\ \begin{matrix} x \\ w^{[1]} \\ b^{[1]} \end{matrix} & \rightarrow \boxed{Z^{[1]} = W^{[1]}x + b^{[1]}} \rightarrow \boxed{a^{[1]} = \sigma(Z^{[1]})} \rightarrow \boxed{Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}} \\ & \rightarrow \boxed{L(a^{[2]}, y)} \leftarrow \boxed{a^{[2]} = \sigma(Z^{[2]})} \rightarrow \frac{da^{[2]}}{dZ^{[2]}} \end{aligned}$$

## $\Rightarrow$ Neural Network representation

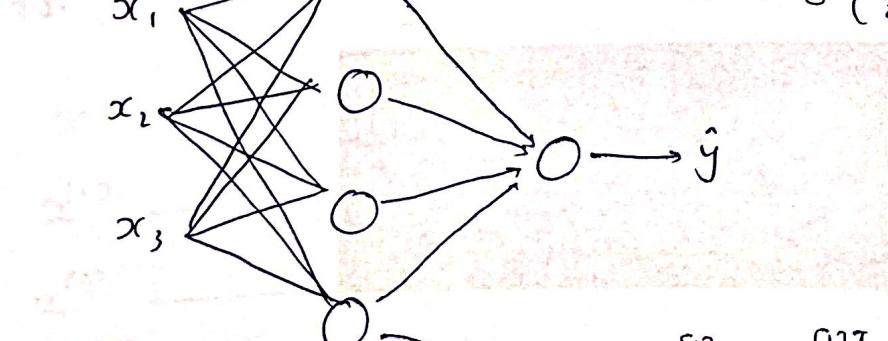


$\rightarrow$  Obtaining Activation Function



$$z_1^{[1]} = w_1^{[1] T} x + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$



$$z_1^{[1]} = w_1^{[1] T} x + b_1^{[1]}$$

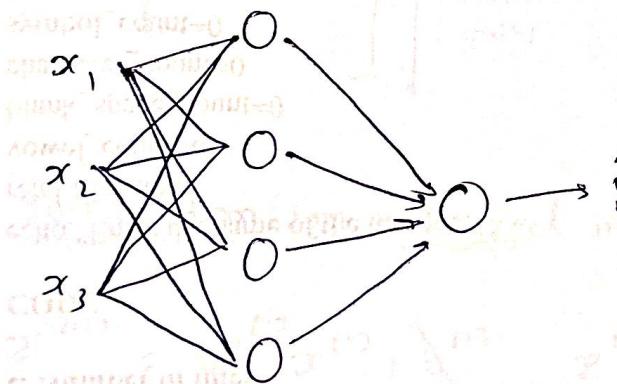
$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z^{[1]} = \underbrace{\begin{bmatrix} -w_1^{[1]T} \\ -w_2^{[1]T} \\ -w_3^{[1]T} \\ -w_4^{[1]T} \end{bmatrix}}_{4 \times 3} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{3 \times 1} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}}_{4 \times 1} = \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix}_{4 \times 1} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}_{4 \times 1}$$

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

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

$\Rightarrow$  Vectorizing across multiple examples



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

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

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

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

$$x \longrightarrow a^{[1]} = y$$

$$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)}$$

$$x^{(2)} \longrightarrow a^{[2](2)} = \hat{y}^{(2)}$$

$$\vdots$$

$$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)}$$

$a^{[2](1)}$  examples  
layers

- For  $i=1$  do  $m$ :

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

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

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

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

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

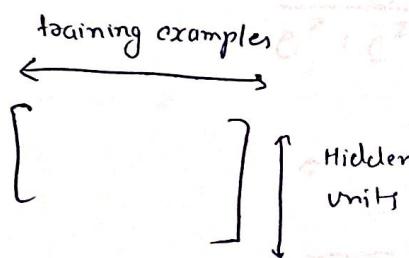
$$a^{[1]} = \sigma(z^{[1]})$$

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

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

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

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



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

⇒ Justification for vectorized implementation

$$z^{[1](1)} = w^{[1]} x^{(1)} + b^{[1]}, \quad z^{[1](2)} = w^{[1]} x^{(2)} + b^{[1]}, \quad z^{[1](3)} = w^{[1]} x^{(3)} + b^{[1]}$$

$$w^{[1]} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \quad w^{[1]} x^{(1)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{[1]} x^{(2)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \quad w^{[1]} x^{(3)} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

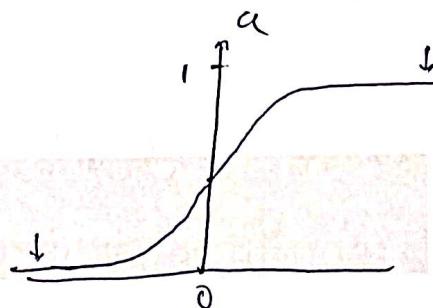
$$w^{[1]} \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ z^{[1](1)} & z^{[1](2)} & z^{[1](3)} \\ 1 & 1 & 1 \end{bmatrix} = Z^{[1]}$$

$$w^{[1]} x^{(i)} = z^{[1](i)}$$

## \* Activation function

sigmoid function (output layer)

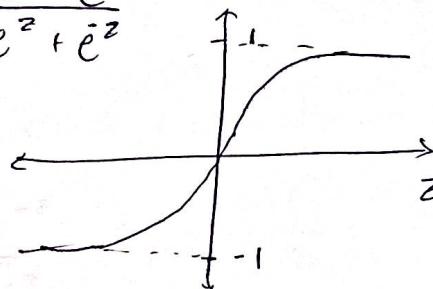
$$a = \frac{1}{1 + e^{-z}}$$



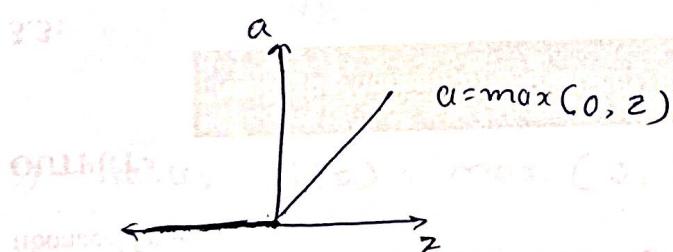
When  $z$  is very large or very small, gradient becomes very small or 0 which slows down GD called vanishing GD.

tanh function (Better than sigmoid) { hidden layers }

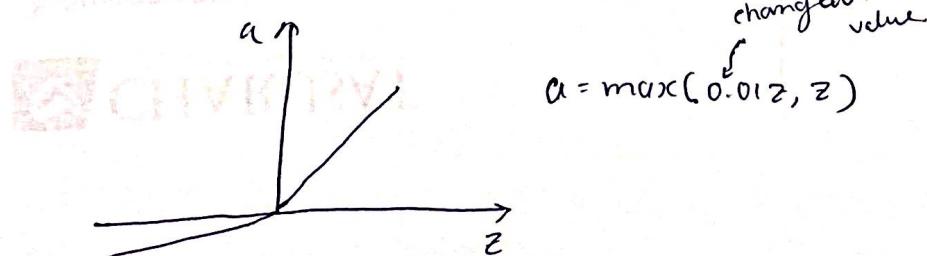
$$a = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



ReLU function (choose as a Default function)



Leaky ReLU function



## $\Rightarrow$ Derivatives of Activation function

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

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

$$\frac{d \left( \frac{e^z}{e^z + 1} \right)}{dz} = \frac{d \left( 1 - \frac{1}{e^z + 1} \right)}{dz} = \underline{\underline{g(z)(1-g(z))}}$$

$$= a(1-a)$$

$$2) \quad g(z) = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = \frac{d g(z)}{dz} = 1 - (\tanh(z))^2$$

CODE:

(Leaky ReLU)

$$3) \quad \text{ReLU}, \quad g(z) = \max(0, z)$$

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

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

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



20

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

$$d\omega^{[2]} = dZ^{[2]} \alpha^{[1]_H}$$

$$Hb^{[2]} = dz^{[2]}$$

$$dz^{(1)} = w^{(2)T} dz^{(2)} + g^{(1)T} (z^{(1)})$$

לְמִזְבֵּחַ וְלְכָבֵד כְּבָדָה וְלְמִזְבֵּחַ וְלְכָבֵד כְּבָדָה

$$d\omega^{1,1} = dz \wedge x^7$$

$$d\mathfrak{b}^{(1)} = dz^{(1)}$$

\* ~~vectorization~~ ~~vectorization~~

$$dZ^{(z)} = A^{(z)} - A$$

$$dW^{(2)} = \sum_{i,j} A^{(2)}_{ij} dz_i$$

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

[17] *Luz* [27] *T* *luz* *luz*

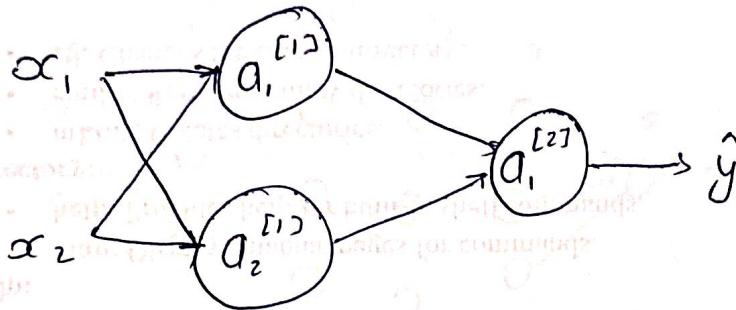
$$dz^{[1]} = w^{[2]'} dz^{[2]} + g^{[1]'}(z^{[1]})$$

$$\omega^{(n)} = \frac{1}{m} dz^{[ij]} x^T$$

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

## Random initialization

Input layer      Hidden layer      Output layer



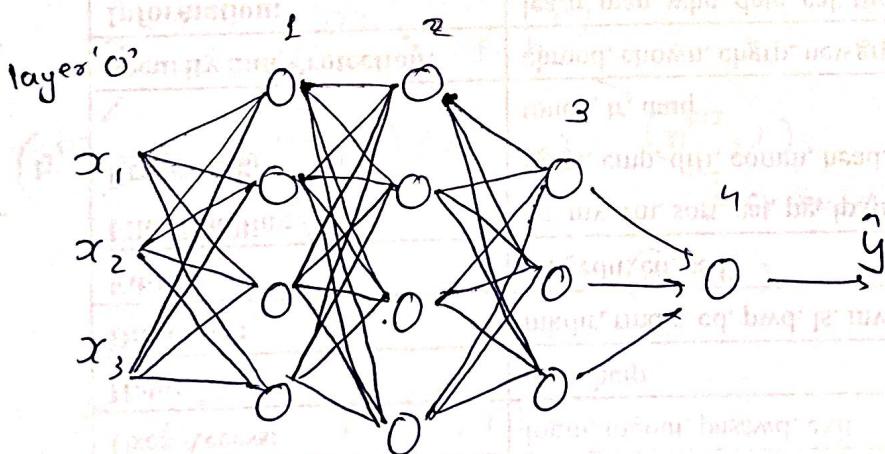
$$w = \text{np.random.randn}(2, 2) * 0.01 \quad \text{if this is large no. we get } z \text{ as large, that will saturate the activation function.}$$

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

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

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

## \* Deep neural networks notation



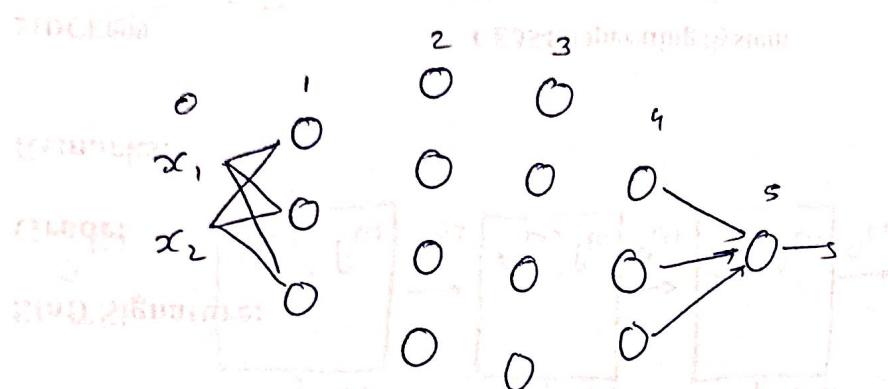
$$L = 4 \quad \# \text{ layers}$$

$$n^{(0)} = \# \text{ units in layer 0} \quad (n^{(0)} = 3, n^{(1)} = 4, n^{(2)} = 4, n^{(3)} = 3, n^{(4)} = 1)$$

$$a^{(L)} = \text{activations in layer } L$$

$$a^{(L)} = g^{(L)}(z^{(L)}) \quad w^{(L)} = \text{weights for } z^{(L)}$$

## $\Rightarrow$ Matrix Dimensions



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

$w^{[1]}: (n^{c1}, n^{co1})$
$b^{[1]}: (n^{c1}, 1)$

Dimensions of the vectors and matrices involved in the computation:

- $x: (n^{c1}, 1)$
- $w^{[1]}: (n^{c1}, n^{co1})$
- $b^{[1]}: (n^{c1}, 1)$
- $z^{[1]}: (n^{co1}, 1)$

### - Vectorized implementation

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

Dimensions of the vectors and matrices involved in the computation:

- $x: (n^{c1}, m)$
- $w^{[1]}: (n^{c1}, n^{co1})$
- $b^{[1]}: (n^{co1}, m)$
- $z^{[1]}: (n^{co1}, m)$

Note:  $b^{[1]}$  is broadcasted across all columns of  $x$ .

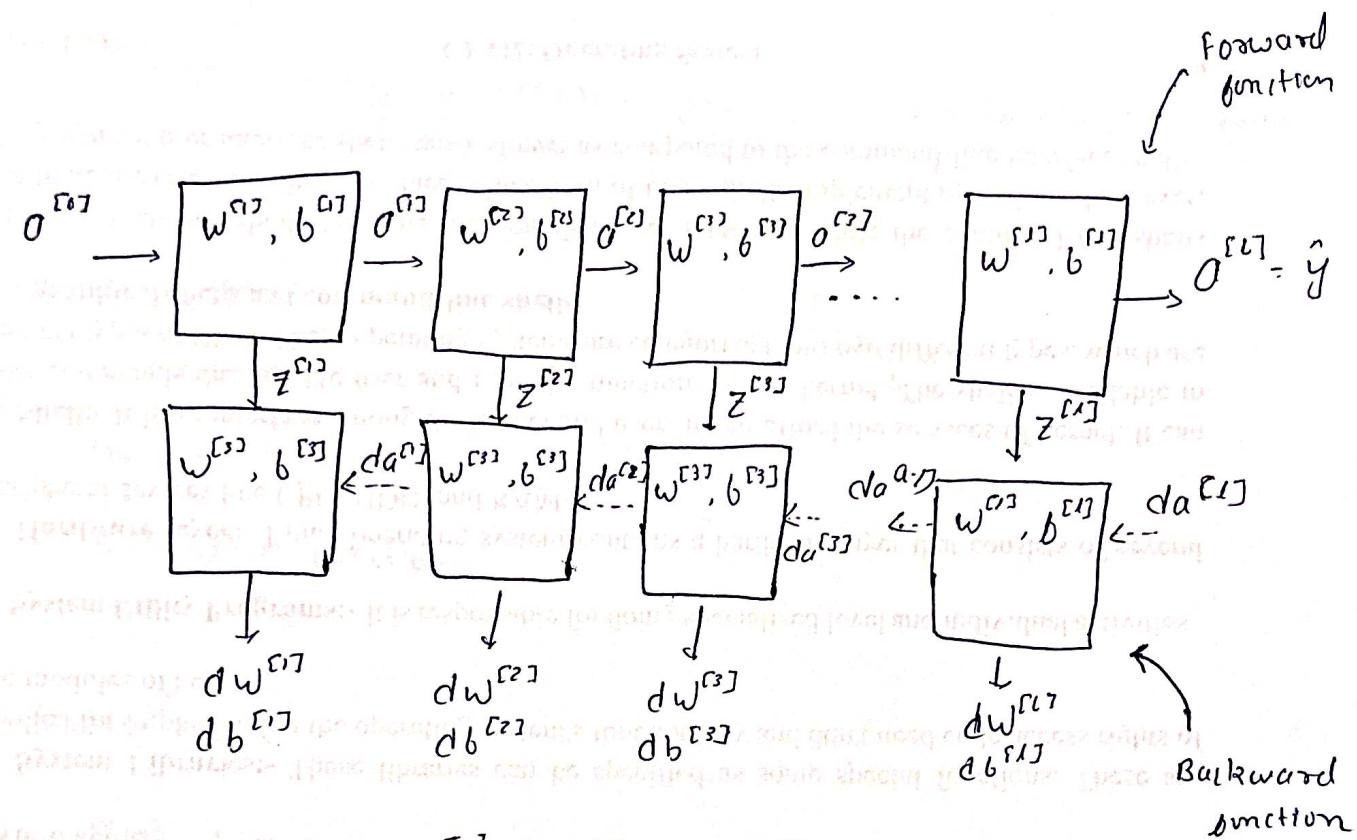
$$z^{[1]}, a^{[1]}: (n^{co1}, 1)$$

$$z^{[1]}, A^{[1]}: (n^{co1}, m)$$

$$\lambda=0 \quad A^{[0]} = X = (n^{co1}, m)$$

$$dz^{[1]}, dA^{[1]}: (n^{co1}, m)$$

## ⇒ forward and backward functions



$$\begin{aligned} w^{[l]} &= w^{[l]} - \alpha dW^{[l]} \\ b^{[l]} &= b^{[l]} - \alpha db^{[l]} \end{aligned}$$

### • forward :

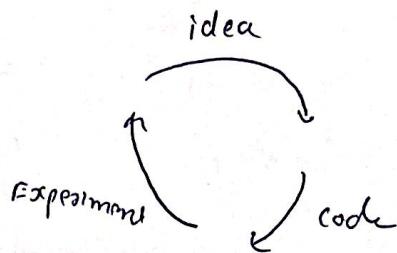
- input  $a^{[L-1]}$ , output  $a^{[L]}$  & cache ( $z^{[L]}$ )
- $z^{[L]} = w^{[L]} \cdot a^{[L-1]} + b^{[L]}$
- $a^{[L]} = g^{[L]}(z^{[L]})$

### • backward :

- input  $da^{[L]}$ , output  $da^{[L-1]}$  cache ( $z^{[L]}$ )
- $dw^{[L]}$
- $db^{[L]}$

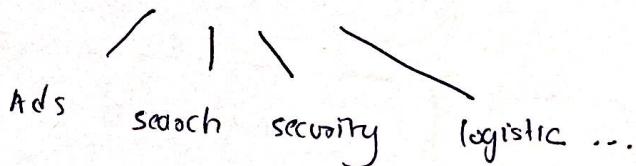
- Applied ML is highly iterative process, we have to fine tune hyperparameters.

- # layers
- # hidden units
- learning rate
- activation functions.
- ...



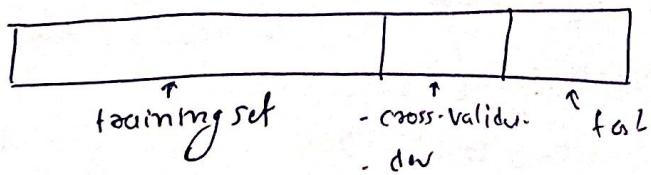
- Applications :

NLP, Vision, speech, structural data



=> Train / dev / test sets

- Data



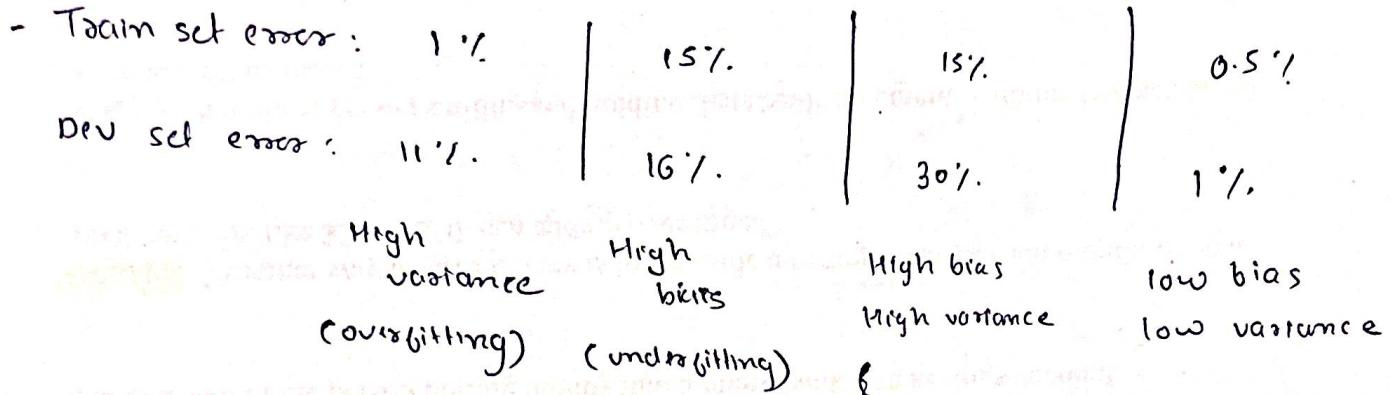
70 / 30 %. or 60 / 20 / 20 %.

- Big Data : 1,000,000

98 / 1 / 1 % for large dataset.

- make sure dev and test comes from same distribution.

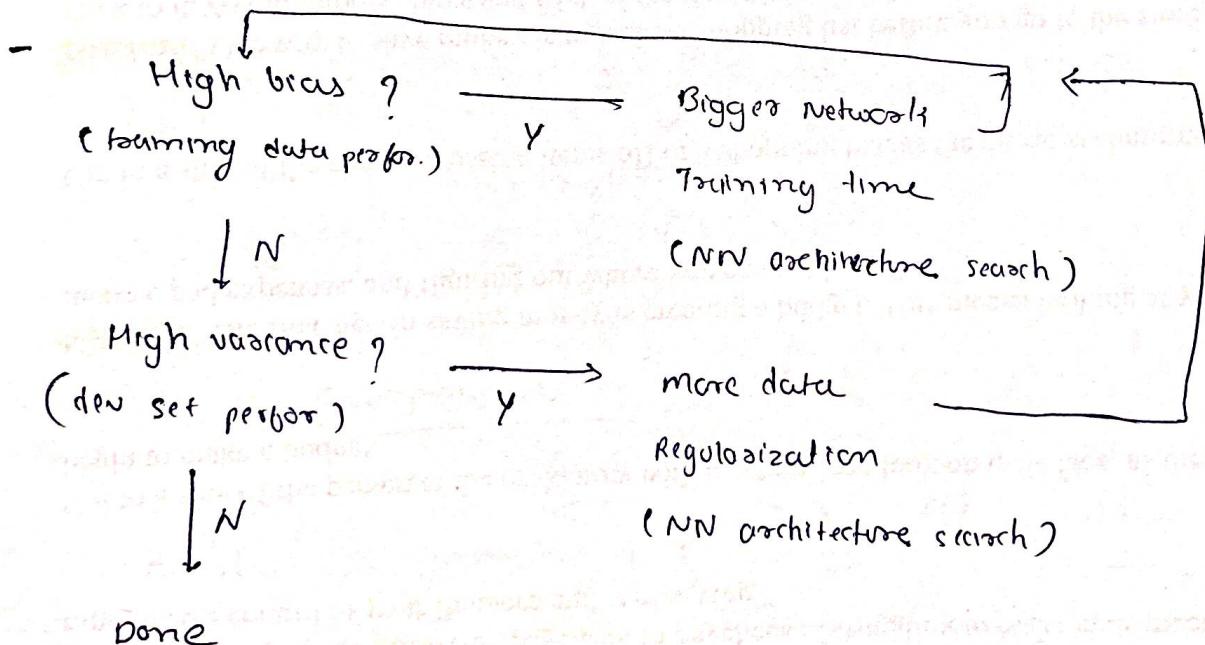
## $\Rightarrow$ Bias and variance



- High variance  $\rightarrow$  complex model (To fix use regularization)  
low variance  $\rightarrow$  simple model

High bias  $\rightarrow$  more errors

low bias  $\rightarrow$  low errors



## \* Regularization

- In logistic regression

$$\min_{w,b} J(w,b) \quad , \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

regularization parameter

$$J(w,b) = \frac{1}{m} \sum_{j=1}^m L(y^{(j)}, \hat{y}^{(j)}) + \frac{\lambda}{2m} \|w\|_2^2$$

-  $L_2$  regularization  $\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$

-  $L_1$  regularization  $\frac{\lambda}{2m} \sum_{j=1}^m |w_j| = \frac{\lambda}{2m} \|w\|_1$

- In neural Network

$$J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{j=1}^m L(\hat{y}^{(j)}, y^{(j)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|^2$$

$$\|w^{(l)}\|^2 = \sum_{i=1}^{n^{(l-1)}} \sum_{j=1}^{n^{(l)}} (w_{ij}^{(l)})^2$$

w:  $(n^{(1)}, n^{(L)})$

Frobenius norm  $\| \cdot \|_F^2 \leftrightarrow \| \cdot \|_F^2$

-  $d w^{(l)} = \text{from backprop.} + \frac{\lambda}{m} w^{(l)}$   $\frac{\partial J}{\partial w^{(l)}}$

$$w^{(l)} := w^{(l)} - \alpha d w^{(l)}$$

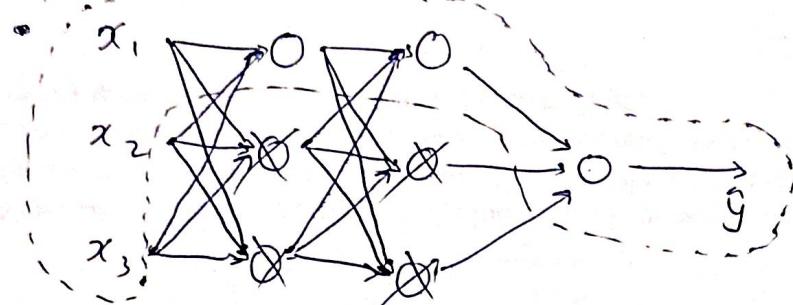
- $L_2$  Regularization or weight decay

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

$$= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{from backprop})$$

⇒ How does regularization prevent overfitting?

(9)

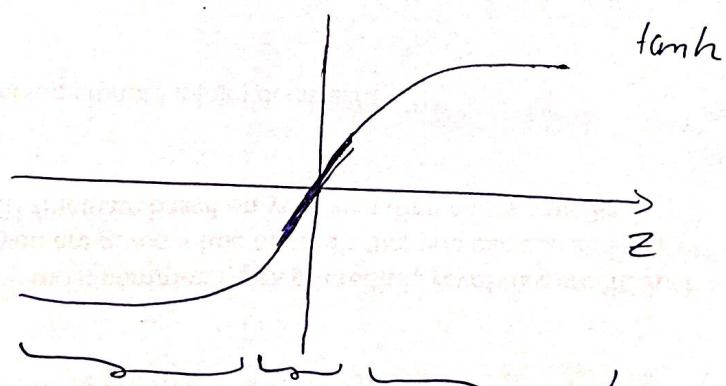


very high

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

$\omega^{[1]} = 0$

- It will cancell out some hidden units by setting it's  $\omega$  to zero. Thus NN becomes simple and overfitting is solved

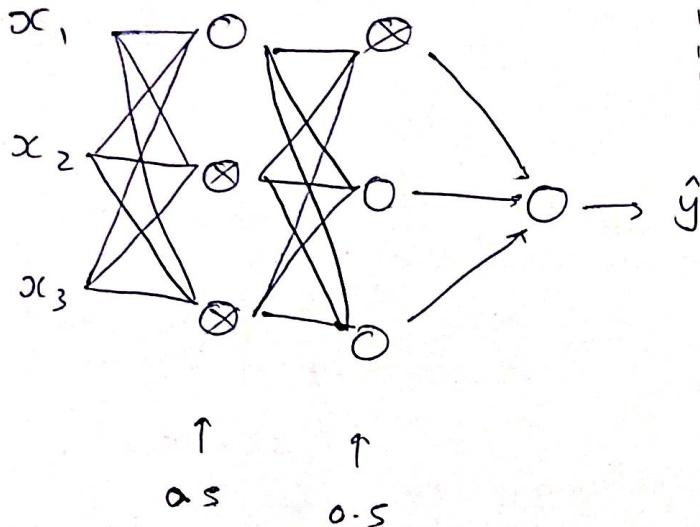


$$g(z) = \tanh(z)$$

$$\lambda \uparrow \rightarrow \omega^{[1]} \downarrow \quad z^{[1]}_l = \omega^{[1]} a^{[l-1]} + b^{[1]} \quad (\text{mostly simple nn or linear})$$

## ⇒ Dropout Regularization

(5)



\* Other ways of regularization

- Data Augmentation (inc data)
- early stopping

This means 50% of hidden units will be dropped to making nn & more simple networks. i.e 50% probability that a unit will be considered or not.

- For  $l=3$  layers NN     $\text{keep-prob} = 0.8$  ( 80% hidden unit will be kept and 20% will be drop)

$$d_3 = \text{np. random. rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep-prob}$$

$$a_3 = \text{np. multiply}(a_3, d_3) \quad \# \quad a_3 \neq d_3$$

$$a_3 / = \text{keep-prob.}$$

- So units  $\rightarrow$  10 unit shot off (20%)

$$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$$

$\downarrow$  reduced by 20%,  
 $/ = 0.8$ .

- At  $a^{[0]}$   $\rightarrow$  No dropout (input layer)

$\Rightarrow \underline{\text{Normalizing inputs}}$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

1. subtract mean:

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

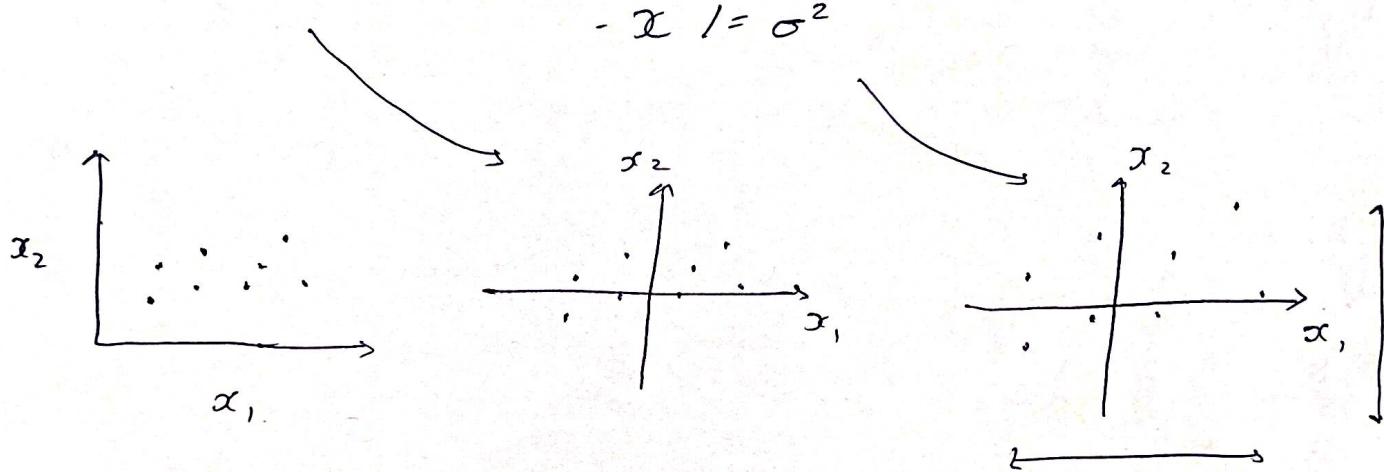
$$\mathbf{x} := \mathbf{x} - \mu$$

2. Normalize variance

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

elementwise

$$\mathbf{x} / \sigma^2$$



- Use the same  $\mu$  &  $\sigma$  to normalize test size.

$\Rightarrow \underline{\text{Vanishing Gradient}} \neq \underline{\text{Exploding gradients}}$

- In deep NN

If  $W^{[L]} > 1$  (Identity matrix) then it result exploding

If  $W^{[L]} < 1$  then it result vanishing GD.

- This can be solved by initializing the weight

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

$$g^{(l)}(z) = \text{ReLU}(z)$$