

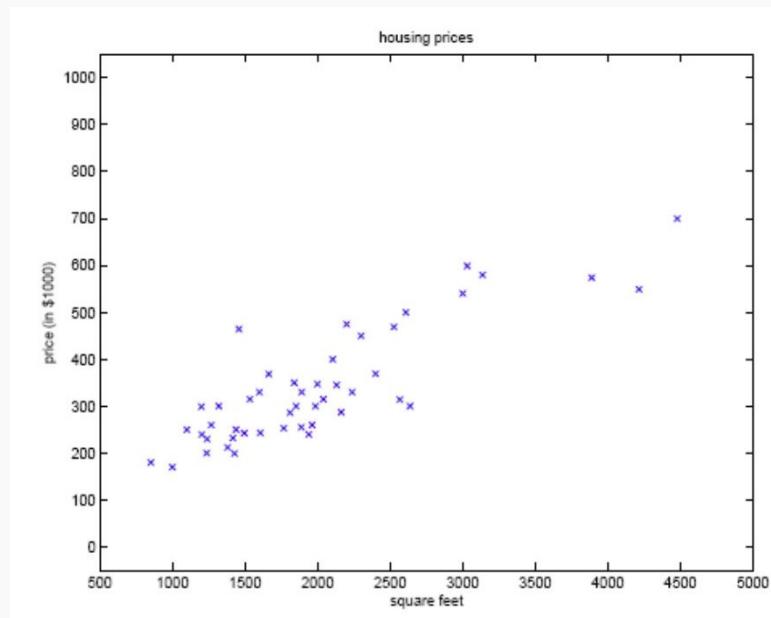
# **Applied Deep Learning**

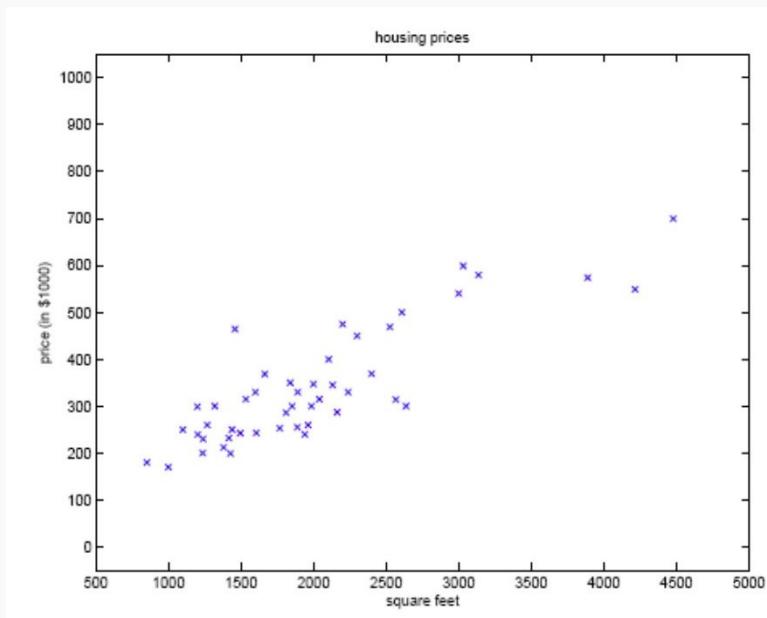
## Lecture 2

---

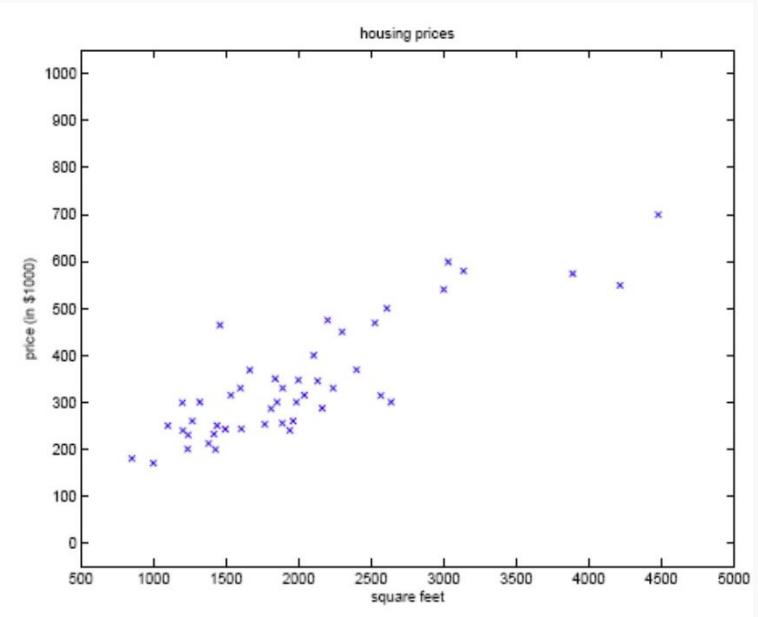
Ganapathy Krishnamurthi

IIT-Madras

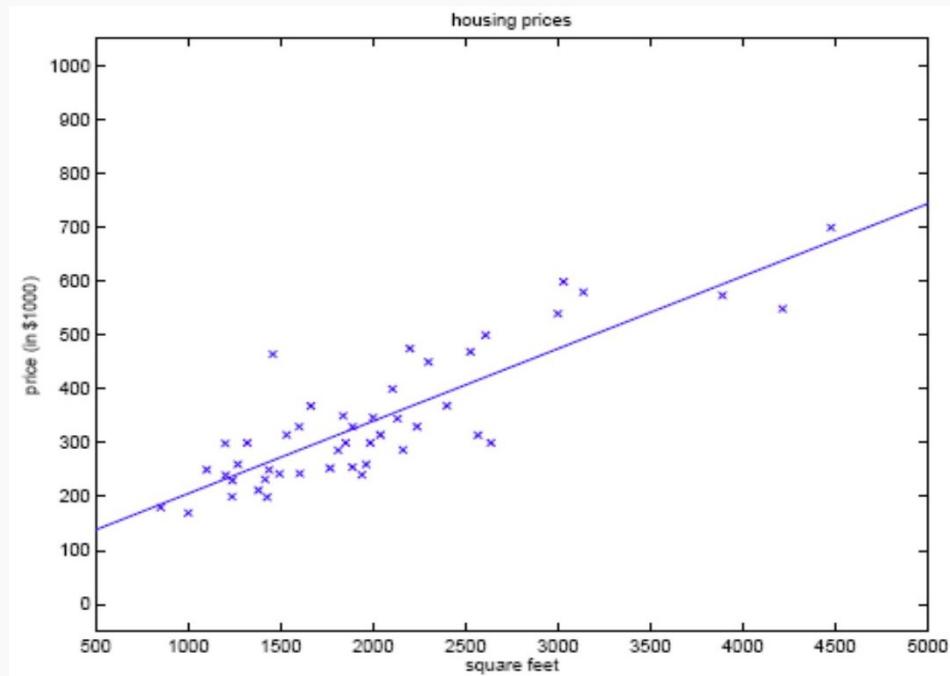




Living area (feet <sup>2</sup> )	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮



Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:



Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮

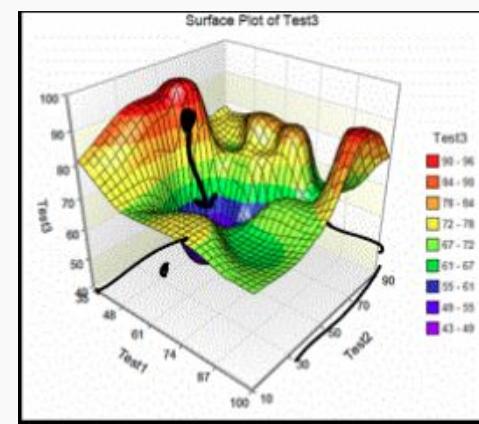
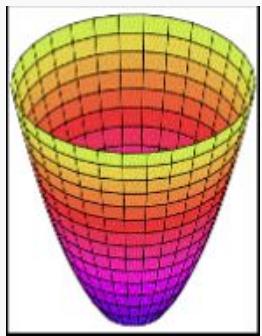
Living area (feet <sup>2</sup> )	Price (1000\$s)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮

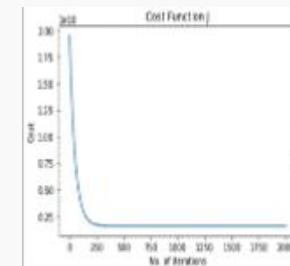
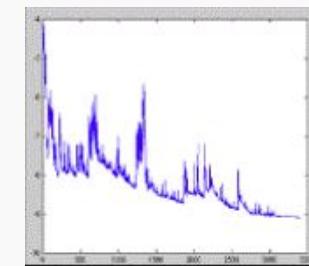
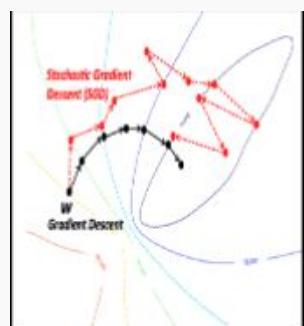
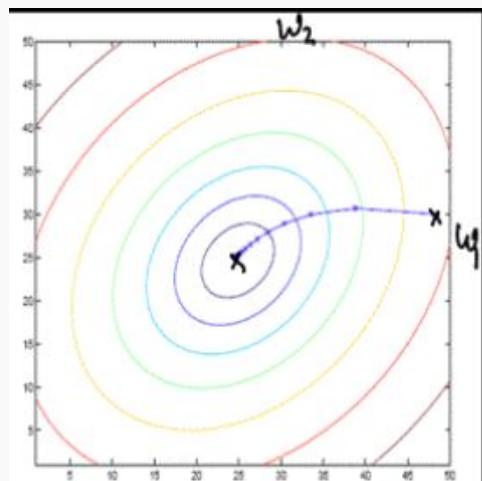
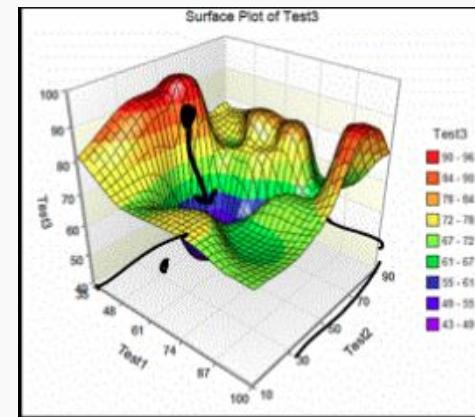
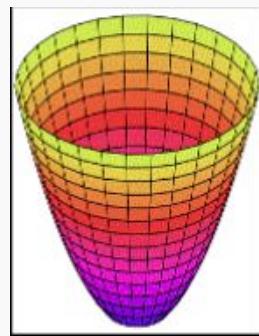
$$J = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (y^{(i)} - X_w)^2$$

$X_w$  is the input matrix.

Solve for  $\nabla_w J = 0 \rightarrow$  Gives linear system of equations for  $w$ , which is normal equation.





# Recap - Linear Regression

## kaggle → housing price dataset

# GrLivArea	# BedroomA...	# SalePrice
1710	3	208500
1262	3	181500
1786	3	223500
1717	3	140000
2198	4	250000
1362	1	143000
1694	3	307000
2090	3	200000
1774	2	129900
1077	2	118000

Living area and bedroom in table Input

and Price is the output ( $y$ )

and Price is the output ( $y$ )

$$\vec{x} = [x_1^{(1)}, x_2^{(0)}]$$

index of the  
sample

component index  
feature index

$m \rightarrow$  samples or no of data points

$n$  - number of feature

# Recap - Linear Regression

Models: Linear, ANN, Logistic, CNN

$$\frac{\partial J}{\partial w_2} = w_2 - \alpha \frac{\partial J}{\partial w_2} \quad \hat{w}_1 = w_1 - \alpha \frac{\partial J}{\partial w_1}, \text{ start with guess}$$

# GrLivArea	# BedroomA...	# SalePrice
1710	$x_1^{(1)}$	3 $x_2^{(1)}$
1262	$x_1^{(2)}$	181500 $y^{(2)}$
1786	3	223500 $y^{(3)}$
1717	3	140000
2198	4	250000
1362	1	143000
1694	3	307000
2090	3	200000
1774	2	129900
1877	2 $x_1^{(m)}$	118000 $y^{(m)}$

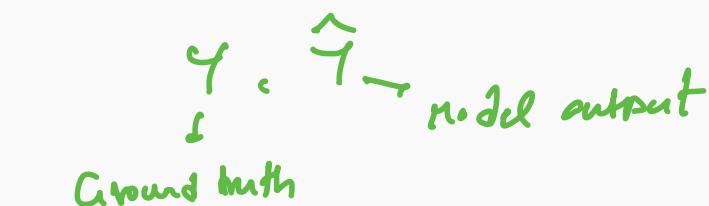
Living area and bedroom in table Input  $\vec{x}$  and Price is the output ( $y$ )

$$\vec{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}) \rightarrow y^{(i)}$$

Input

Target output

$$\begin{aligned}\hat{y}^{(1)} &= w_0 + w_1 x_1^{(1)} + w_2 x_2^{(1)} \\ \hat{y}^{(2)} &= w_0 + w_1 x_1^{(2)} + w_2 x_2^{(2)} \\ \hat{y}^{(3)} &= w_0 + w_1 x_1^{(3)} + w_2 x_2^{(3)}\end{aligned}$$



Input	Output	Model
$1, x_1^{(1)}, x_2^{(1)}, x_3^{(1)}, \dots, x_n^{(1)}$	$y^{(1)}$	$\hat{y}^{(1)}$
.	.	.
$1, x_1^{(m)}, x_2^{(m)}, x_3^{(m)}, \dots, x_n^{(m)}$	$y^{(m)}$	$\hat{y}^{(m)}$
$\vec{X} : m \times n$	Ground Truth ( $\vec{y}$ )	$\hat{\vec{y}}$

Where m is the examples and n is the features

Example  $\rightarrow [\vec{x}^{(i)}, y^{(i)}]$

$$\hat{y} = b + wx \rightarrow \text{Single feature} \rightarrow [w_0 + w_1 x]$$

Living Area

$b \rightarrow w_0$ , Which is bias unit

$$\hat{y}^{(i)} = w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} \dots + w_n x_n^{(i)} \rightarrow \text{General case } n \text{ features}$$

$w_0 x_0$

$$x_0 = 1$$

$$By \rightarrow 2 \rightarrow 4$$

$$By \rightarrow 4 \rightarrow 16$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \quad \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \quad \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

n+1 - columns      (n+1) x 1

$$\vec{\hat{y}} = X w$$

The cost function  $J = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$ . This is mean square error (MSE).  
 Find  $w$  that minimizes  $J$ .

↳ Least Squares

$$\hat{J}^{(i)} \rightarrow \hat{J}^{(i)}(w) = \text{minimizing } J(w)$$

Set  $\nabla_w J = 0$

$$\left\{ \begin{array}{l} \frac{\partial}{\partial w} \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 = 0 \\ \frac{\partial}{\partial w} \left[ \frac{1}{m} \left\| \vec{y} - \vec{\hat{y}} \right\|_2^2 \right] = 0 \\ \frac{\partial}{\partial w} \left[ [y - X_w]^T [y - X_w] \right] = 0 \\ \frac{\partial}{\partial w} \left[ y^T y - w^T x^T y - y^T X_w + w^T X^T X_w \right] = 0 \end{array} \right\}$$

$\left\{ \begin{array}{l} \|w\|_2^2 \\ |w| \rightarrow \text{Sparse} \\ w \rightarrow 0 \end{array} \right.$

$-2X^T y + 2X^T X_w = 0$

$$Xw \quad w^T x = y \quad Ax = b$$

$$[Xw = y]$$

$$w = (X^T X)^{-1} X^T y$$

$\Rightarrow$   $\uparrow$

Normal Equations

$$\frac{\partial J(\hat{J}(w))}{\partial w} = \begin{bmatrix} \frac{\partial J}{\partial \hat{y}} & \frac{\partial J}{\partial w} \end{bmatrix}$$

$$\begin{aligned} & \left\| \vec{y} - \vec{\hat{y}} \right\|^2 \\ &= [\vec{y} - \vec{\hat{y}}]^T [\vec{y} - \vec{\hat{y}}] \end{aligned}$$

$$\begin{aligned} \frac{\partial (V^T a)}{\partial V} &= \vec{a} \\ \frac{\partial (V^T A_v)}{\partial V} &= 2A_v \end{aligned}$$

Matrix Calculus

Normal  $\sim N(0, 1)$

$$\rightarrow (\vec{a}, \hat{\vec{y}})$$

$$|\vec{a}| = |\vec{y}| \cos \theta$$

$$\cos \theta = 1$$

$$Y \sim \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

$$\hat{Y} \sim \begin{bmatrix} : \\ : \\ : \end{bmatrix} \sim N(0, I)$$

Better to use Gradient Descent

Algorithm

1. Initialize (guess)  $\vec{w}$ , set  $\alpha$ , Which is learning rate.

2. Iterate  $\vec{W} = \vec{W} - \alpha \nabla_{\vec{w}} J$

3. Check for convergence / stopping criteria.

If not met, go to 2

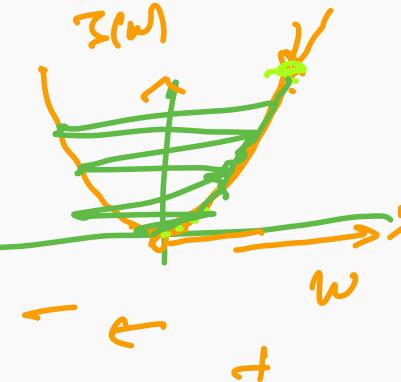


Learning rate schedule

$(\sigma^3)$

$$y \sim \begin{bmatrix} \text{soft}^{(1)} \\ \text{soft}^{(2)} \\ \vdots \\ \text{soft}^{(n)} \end{bmatrix} =$$

$$(\vec{y} - \hat{\vec{y}})$$



$$J(w) \approx w^2$$

$$\frac{dJ}{dw} = w$$

$$-w$$

$$\frac{d}{dt} \frac{d}{dw} \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y}_i)^2$$

Update Val deha set

For  $\nabla_w J$

$$[w_j := v_j - \alpha \frac{\partial J}{\partial w_j}]$$

$$\frac{\partial J}{\partial q} \quad \frac{\partial \hat{y}}{\partial w}$$

$$[x^{(1)} \rightarrow]$$

$$x^{(2)}$$

$$\vdots (m)$$

SGB  $\rightarrow$  1 pt  
Batch  $\rightarrow$  All  
mini-Batch  $\rightarrow$  subset

This is the batch- gradient descent.

$$[J = \frac{1}{m} \sum_i (y^{(i)} - \hat{y}^{(i)})^2]$$

$$\leftarrow \frac{\partial J}{\partial w} = \frac{1}{m} \sum_i -2(y^{(i)} - \hat{y}^{(i)}) \frac{\partial \hat{y}^{(i)}}{\partial w}.$$

$$\left( \frac{\partial J}{\partial w_j} \right) = \frac{1}{m} \sum_{i=1}^m -2(y^{(i)} - \hat{y}^{(i)}) \frac{\partial \hat{y}^{(i)}}{\partial w_j}.$$

$$\hat{y}^{(i)} = w_0 + w_1 x_1^{(i)} + \dots + w_n x_n^{(i)}$$

$$\left[ \frac{\partial \hat{y}^{(i)}}{\partial w_j} = x_j^{(i)} \right] \rightarrow j^{th} \text{ feature of } i^{th} \text{ example}$$

$$\left[ \frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m -2(y^{(i)} - \hat{y}^{(i)}) x_j^{(i)} \right]$$

$$\delta_i = \text{Error}^{(i)} \text{ Input}$$

$$\left[ \frac{\partial J}{\partial \vec{w}} = \frac{1}{m} \sum_{i=1}^m -2\delta^{(i)} \vec{x}^{(i)} \right]$$

$$= \frac{1}{m} \text{Sum} [-2\text{Error}^{(i)} * \text{Feature}^{(i)}]$$

hyper parameter optimization

$$\alpha = \frac{\alpha}{100}$$

every N iterations

$$\frac{\partial J}{\partial w}$$

$$\left[ \begin{array}{c} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{array} \right]$$

# Mini-Batch Gradient Descent

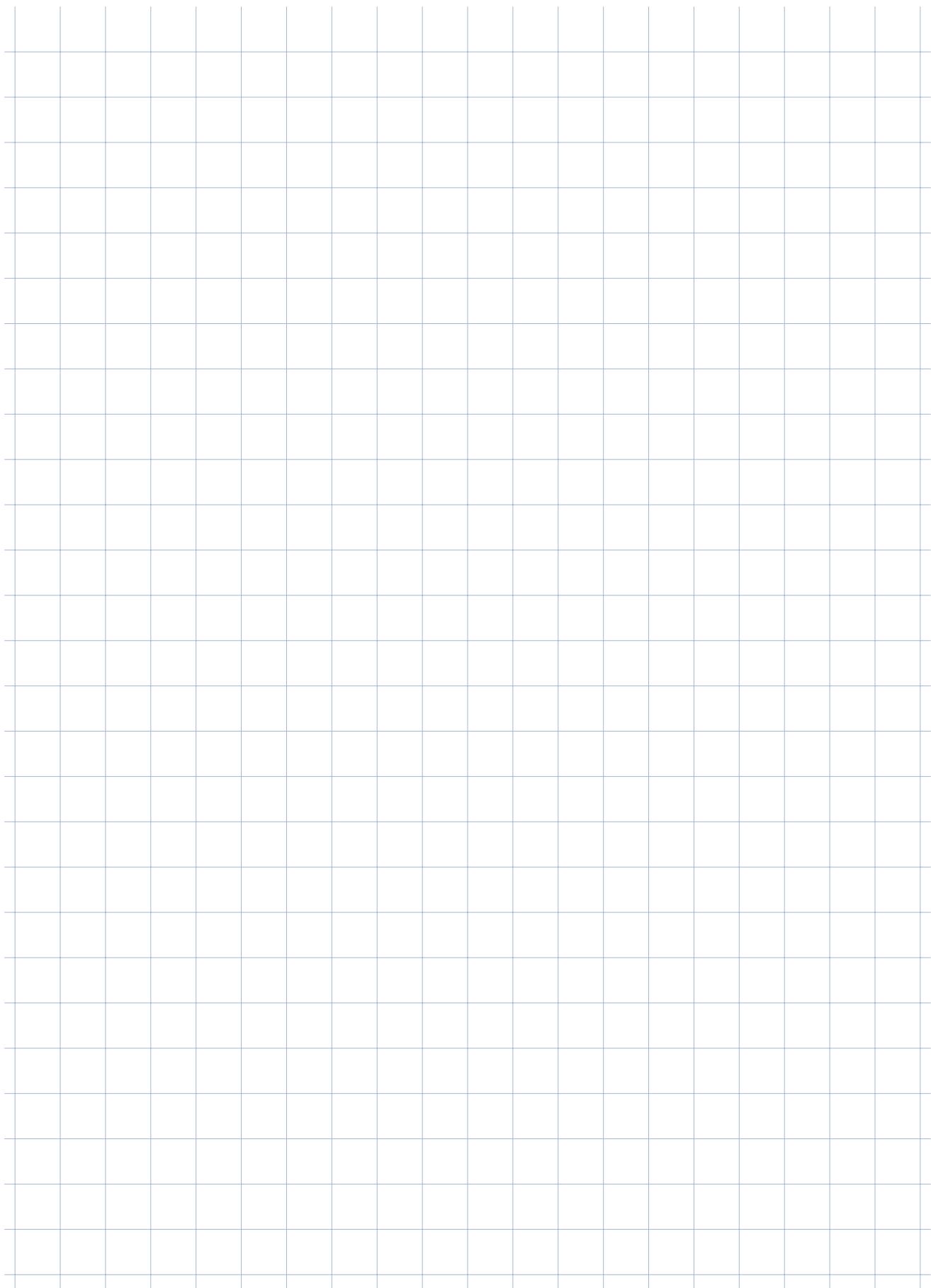
The diagram shows a green circle containing a function  $f$  with a derivative  $\frac{\partial f}{\partial w_j}$ . A green arrow points from this circle to a green bracket labeled  $\sum_{mini-batch} [-2Error * Feature]$ . Another green arrow points from this bracket to a green circle containing a vector  $w$ . Inside this circle, a green arrow points to the right, labeled  $= w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_n x_n^{(i)}$ . Below this, a green arrow points down to the update equation:  $= w_j - \frac{\partial f}{\partial w_j}$ .

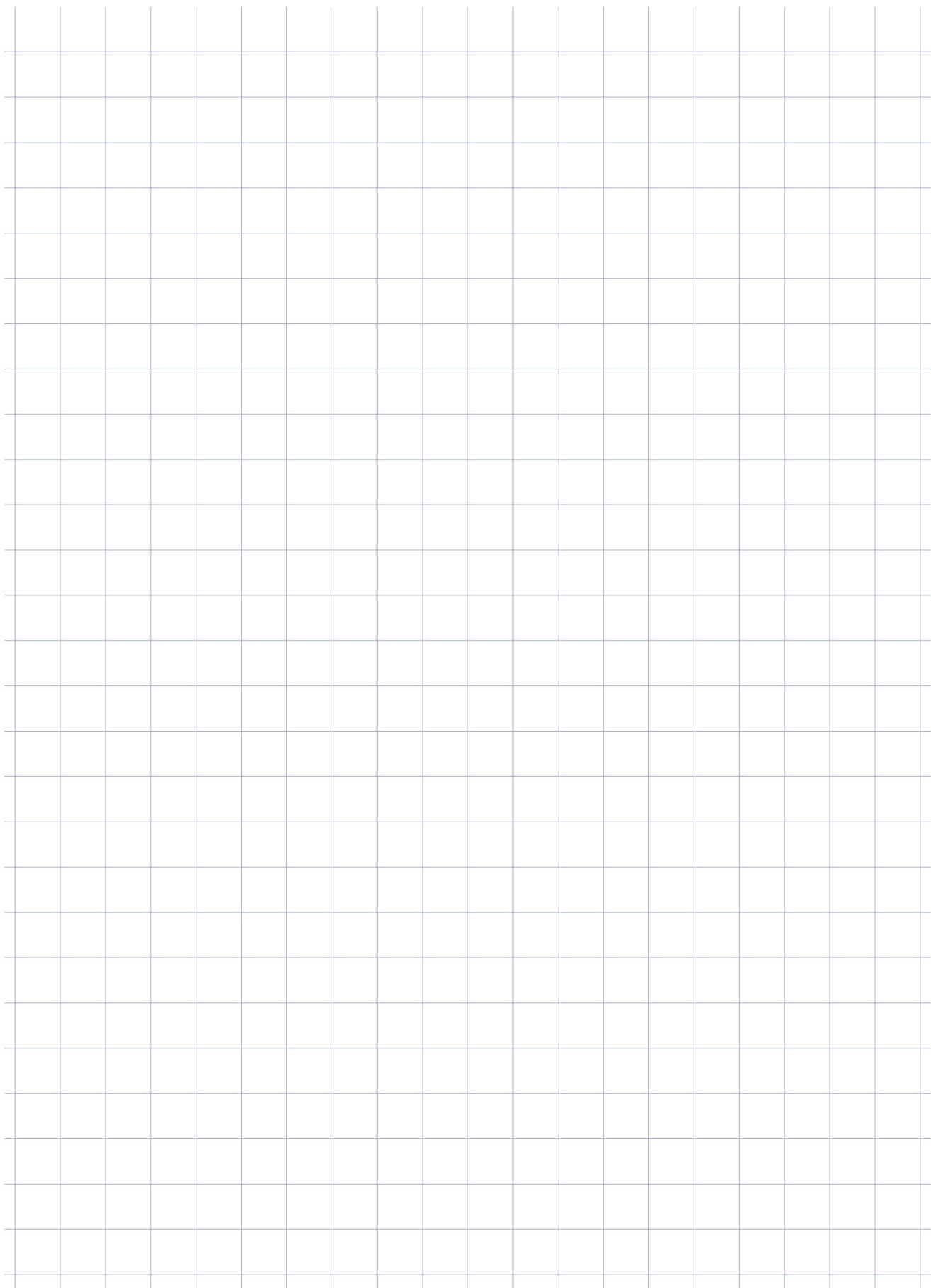
Stochastic Gradient Descent (SGD)  $\rightarrow$  Mini-batch size = 1.

General Algorithm

1. Initialize  $\vec{w}, \alpha$
2. Split data into mini-batches, Randomly shuffle.
3.  $\vec{w} = \vec{w} - \sum_{mini-batch} \left( \frac{2Error * Feature}{m} \right)$  - Once you see the whole data  $\rightarrow$  epoch.
4. Until stopping

[  
  | Epoch  
  : 100s Epoch.  
]





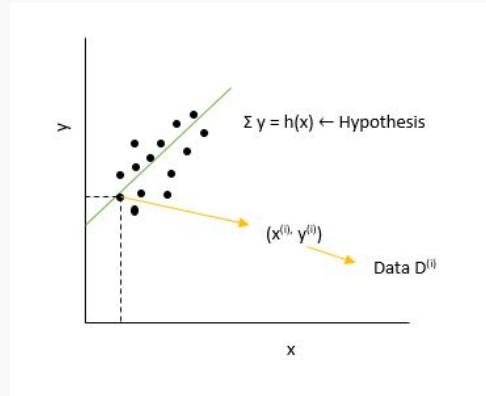
# Topics for today

1. Bias/variance Trade off
2. Regularization
3. Probabilistic Interpretation → MAP/MLE
4. Logistic Regression

$\hat{y} = h(x)$  to Hypothesis function.

$\hat{y} = w^T x \rightarrow$  Linear regression.

Least Square Cost function:  $\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}))^2 = J$



$$y^{(i)} = h(x^{(i)}) + \delta^{(i)}$$

Where  $h(x^{(i)})$  is  $\hat{y}$  and  $\delta^{(i)}$  is  $P(y|x)$ .

# Probabilistic Interpretation

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Where  $P(h|D)$  is the Posterior Probability which is the probability of hypothesis h Given data is D.

$P(D|h)$  is the likelihood (How likely is Data D for given hypothesis h).

$P(h)$  is the prior probability (Before we know anything about the data)

# MAP

MAP: Maximum a posteriori probability ( $h \in H$  (Space/ set of all hypothesis)) says the best hypothesis is the one which maximizes  $P(h|D)$  mathematically,

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

From Bayes' Theorem,  $h_{MAP} = \operatorname{argmax}_{h \in H} \frac{P(D|h)P(h)}{P(D)}$ .

$P(D)$  is independent of  $h$ .

If  $P(h_i)$  are all equal (that is, equally likely priors), then

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(D|h)$$

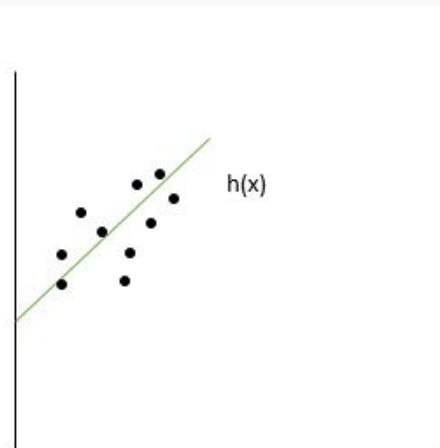
i.e., For uniform priors,

# MAP

$$h_{MAP} = h_{max-likelihood}$$

Now, consider a general regression problem:

$$\begin{array}{ccc} y^{(i)} & = & h(x^{(i)}) + \delta^{(i)} \\ Data & Hypothesis & Error \end{array}$$



# MAP

Assume:

- Independent  $\delta^{(i)}$ . That is, all training examples, are independent.
- Normally distributed  $\delta^{(i)} : \delta \sim N(0, \sigma^2)$  (Mean = 0, Variance =  $\sigma^2$ )

$$P(\delta^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{\delta^2}{2\sigma^2}\right\}$$

Now, the regression problem

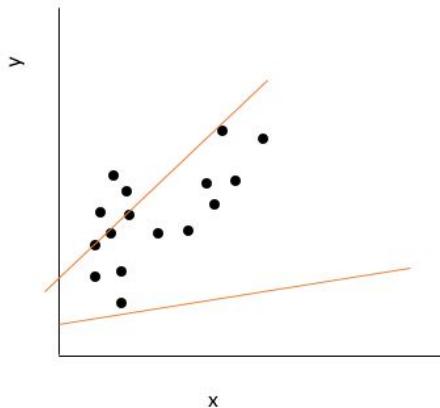
Hypothesis Space ( $H$ ) - Set of functions that the learning algorithm is allowed to choose from for a solution.

Ex: Linear Regression has  $H = \text{Set of functions linear in input feaures.}$

$$h \in H : S \cdot t \underset{h \in H}{\operatorname{argmin}} \|y - \hat{y}\|^2$$

*Least square formulation*

# Least squares from Max-Likelihood



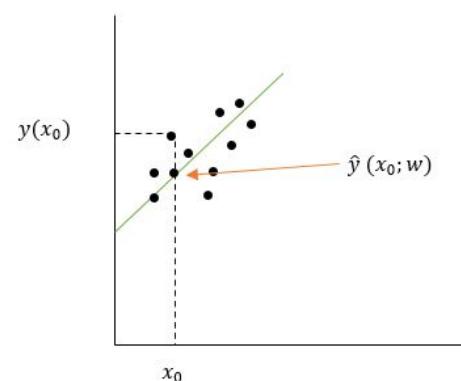
$$y = h(x) + \delta$$

Assumptions:

- I.I.D  $\delta$
- $\delta$  is normally distributed.  
 $\delta \sim N(0, \sigma^2)$

Question: Given a particular fit, how likely is the data?

$$\begin{aligned} & P(\vec{y} | \vec{x}; w, \sigma^2) \\ &= P(y^{(1)}, y^{(2)}, \dots, y^{(m)} | \vec{x}; w) \\ &\prod_{i=1}^m P(y^{(i)} | x^{(i)}; w, \sigma^2) \end{aligned}$$



$$\therefore Likelihood = \prod_{i=1}^m P(y^{(i)} | x^{(i)}; w, \sigma^2)$$

$$\log Likelihood = \sum_{i=1}^m \log(P(y^{(i)} | x^{(i)}; w))$$

$$N(y^{(i)}; \hat{y}(x^{(i)}), \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma^2} \exp \left[ \frac{-(y^{(i)} - \hat{y}^{(i)})^2}{2\sigma^2} \right]$$

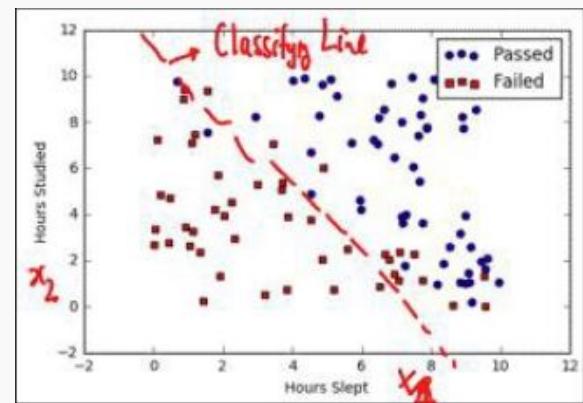
$$= \frac{-m}{2} \log 2\pi - m \log \sigma - \sum_{i=1}^m \frac{(y^{(i)} - \hat{y}^{(i)})^2}{2\sigma^2}$$

$$\therefore \underset{w}{\operatorname{argmax}} \log(P(D|h)) = \underset{w}{\operatorname{argmin}} \sum_{i=1}^m \frac{(y^{(i)} - \hat{y}^{(i)})^2}{2\sigma^2}$$

$$W_{MLE} = W_{Least-square}$$

We can also obtain  $W_{\text{Regularized least squares}}$  by finding  $W_{MAP}$  by using an appropriate prior.

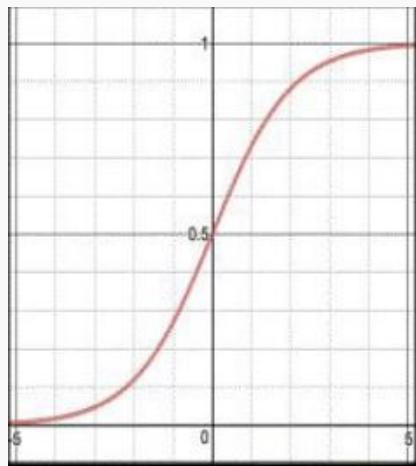
$x_A$	$x_B$	
Studied	Slept	Passed
4.85	9.63	1
8.62	3.23	0
5.43	8.23	1
9.21	6.34	0



Classify Haplers in two steps.

1. Finding classifying line (boundary)
2. Assigning a label (value) to each point

## Sigmoid function



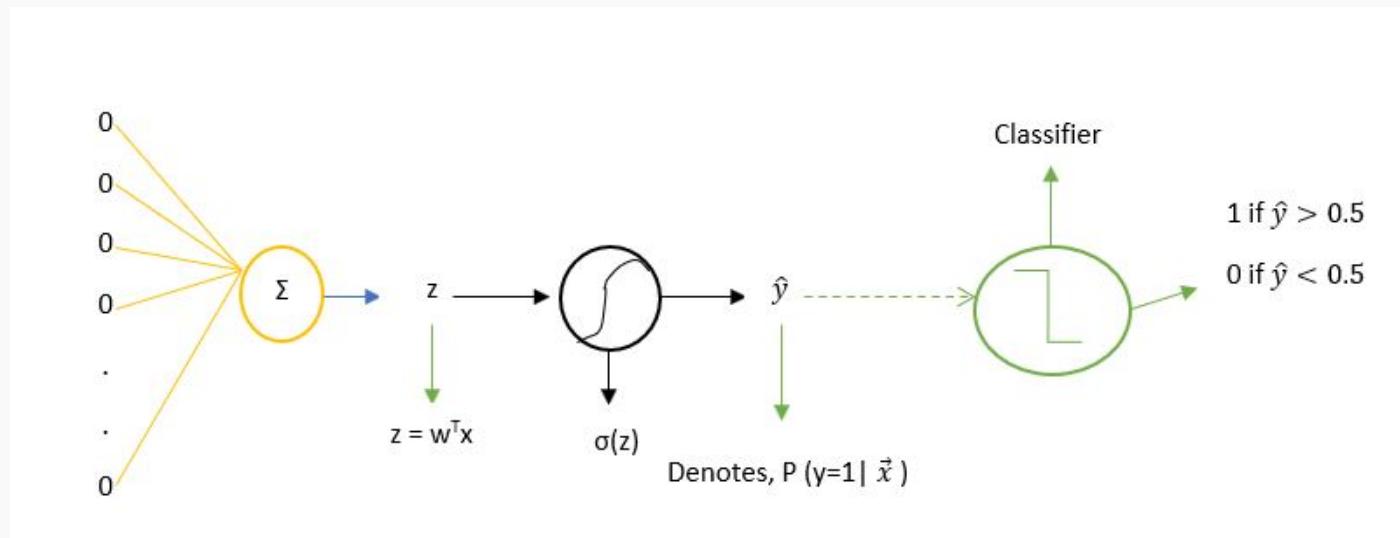
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$z \rightarrow -\infty, \sigma(z) \rightarrow 0$$

$$z \rightarrow \infty, \sigma(z) \rightarrow 1$$

$$z = 0, \sigma(z) = 0.5$$

## Linear regression model



# Summarize Linear Regression

$$\vec{X} \rightarrow (x_1^{(r)}, x_2, \dots, x_n) \quad \text{Input features}$$

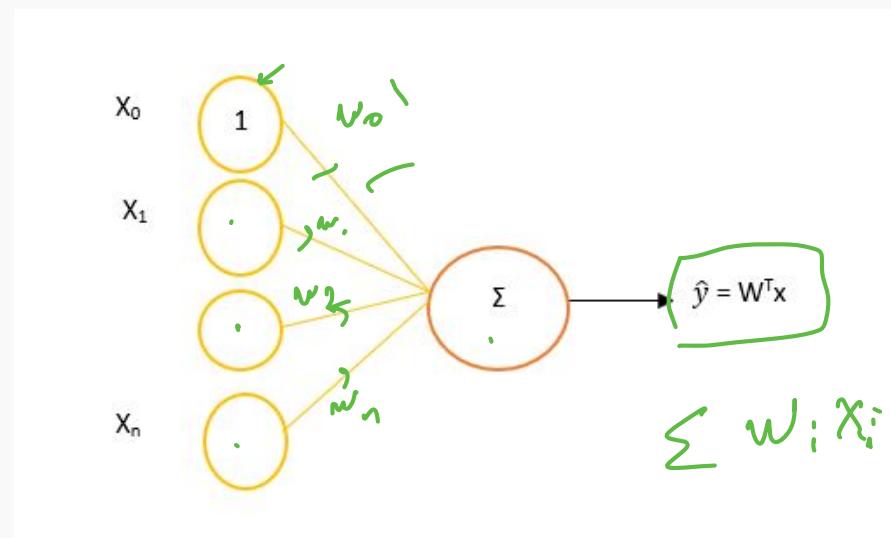
$$\hat{y} = W^T x + b$$

such that

Find  $\hat{y}$  such that  $\frac{1}{m} \sum \|\vec{y} - \hat{y}\|_2^2$  is minimized.

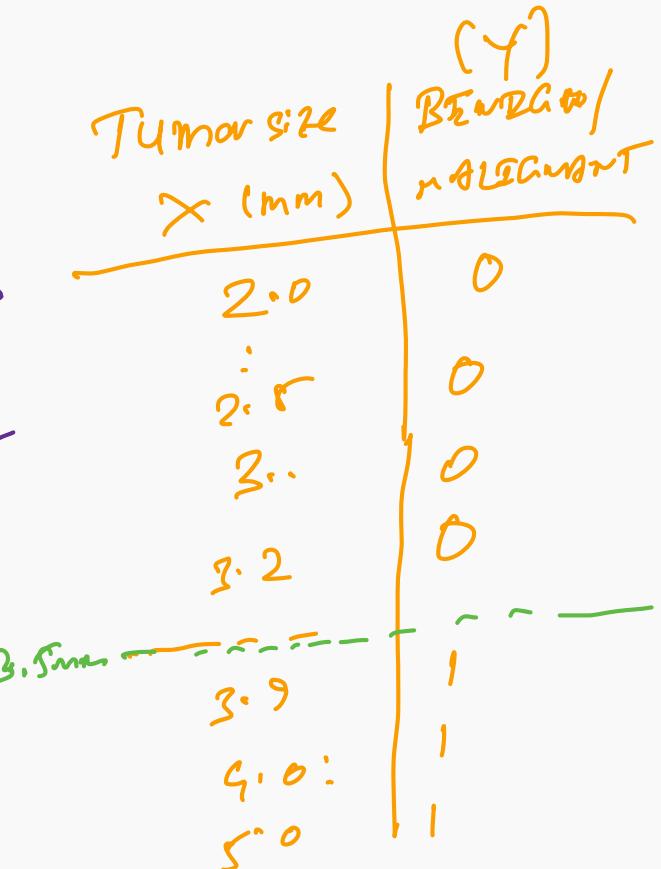
$$\begin{aligned} \frac{\partial J(\vec{w})}{\partial w_j} &= \frac{1}{m} \sum_{i=1}^m -2(y^{(i)} - \hat{y}^{(i)}) x_j^{(i)} \\ w_j &= w_j - \frac{\partial J(\vec{w})}{\partial w_j} \\ &\quad (y^{(i)} - \hat{y}^{(i)}) x_j^{(i)} \\ g^{(i)}(\vec{w}) & \end{aligned}$$

$g(\vec{w})$



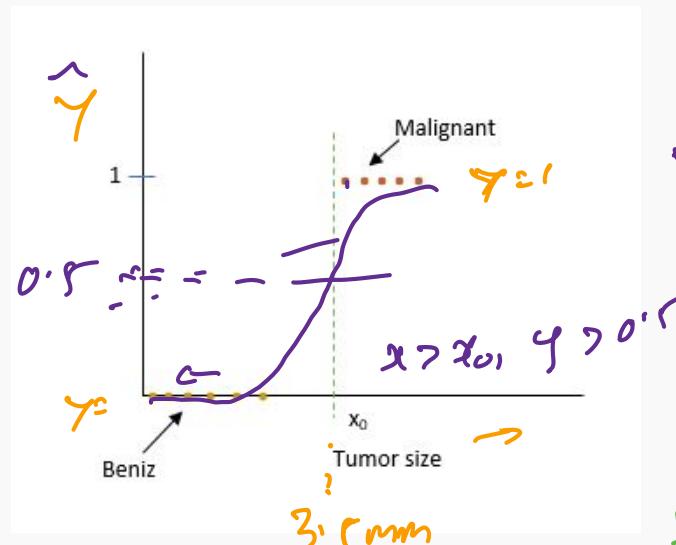
n- features

# Patient Data



Classification:

One feature



$$y = \begin{cases} 1 & x > x_0 \\ 0 & x < x_0 \end{cases}$$

# Recap

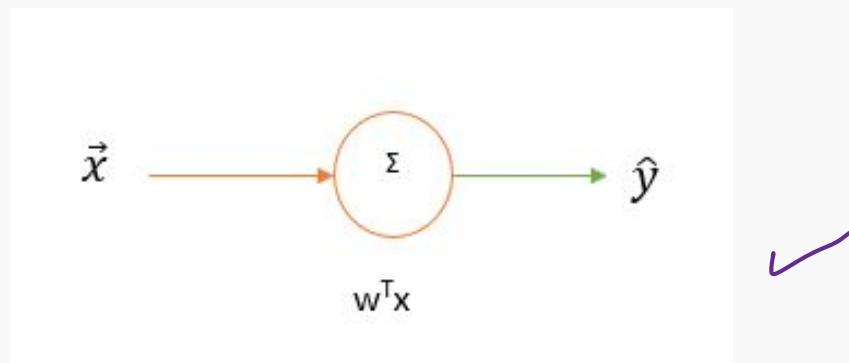
Linear Regression → For continuous data

Maps:  $\vec{x} \rightarrow y$ , where  $y$  is continuous.

# Recap

Linear Regression → For continuous data

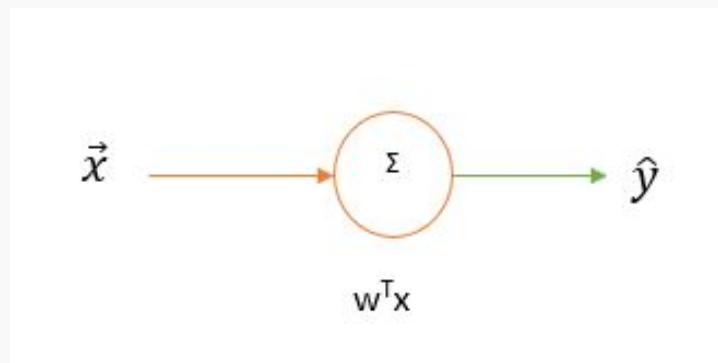
Maps:  $\vec{x} \rightarrow y$ , where  $y$  is continuous.



# Recap

Linear Regression → For continuous data

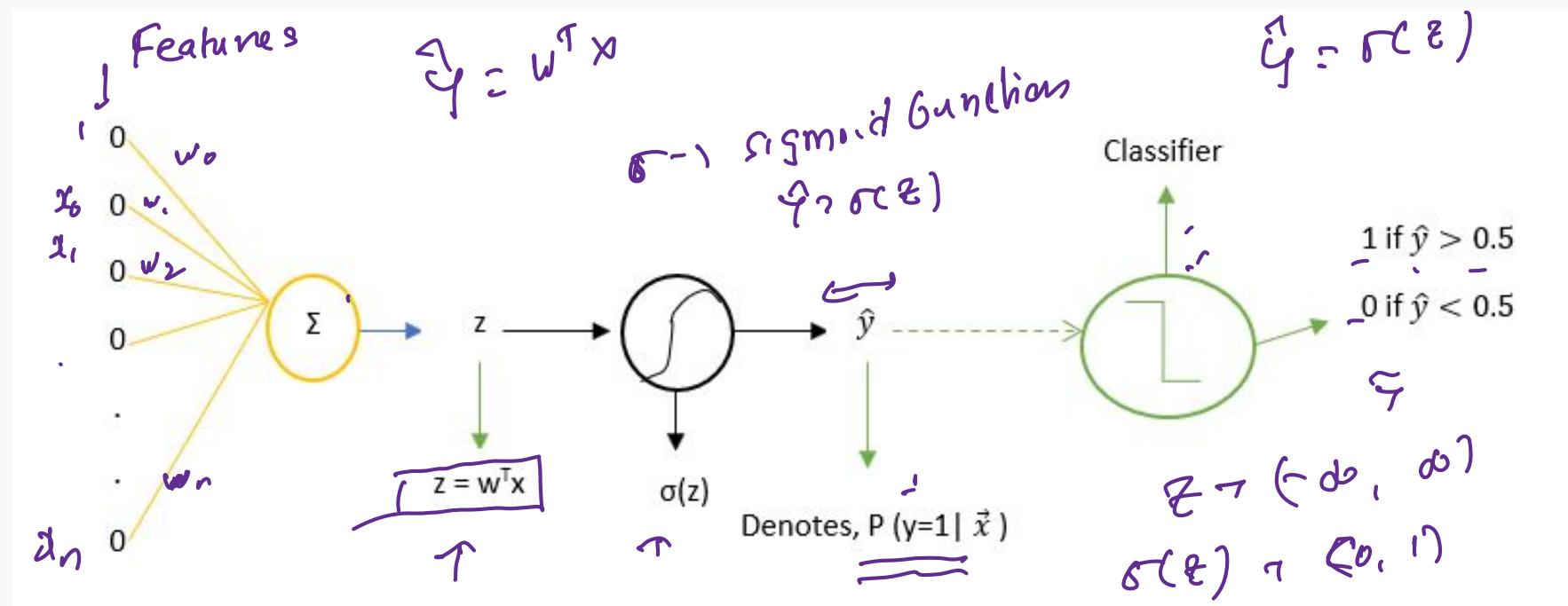
Maps:  $\vec{x} \rightarrow y$ , where  $y$  is continuous.



$$L = \sum_{i=1}^m (y^{(i)}_{True} - \hat{y}^{(i)}_{Predicted})^2 \quad \checkmark$$

Logistic Regression → For discrete data (Classification)

Maps:  $\vec{x} \rightarrow y$ , where  $y$  is discrete which is either 0 or 1.



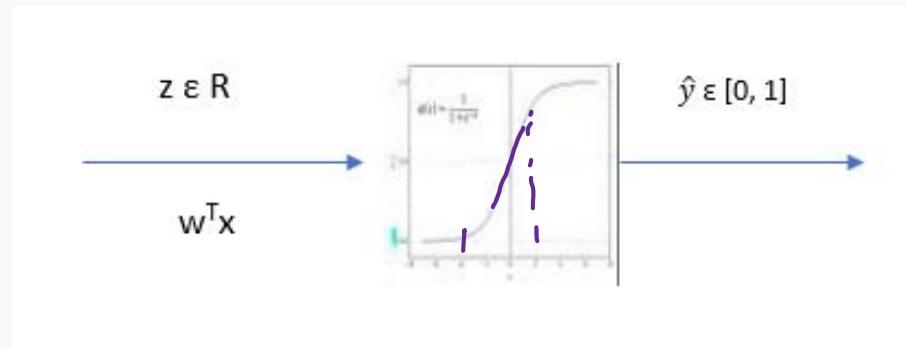
This  $z$  is continuous and unbounded. Sigmoid function maps to a continuous number between 0 and 1. This represents probability that  $P(y=1 | \vec{x})$ .

$$\underline{\underline{P(Y=1 | \vec{x})}}$$

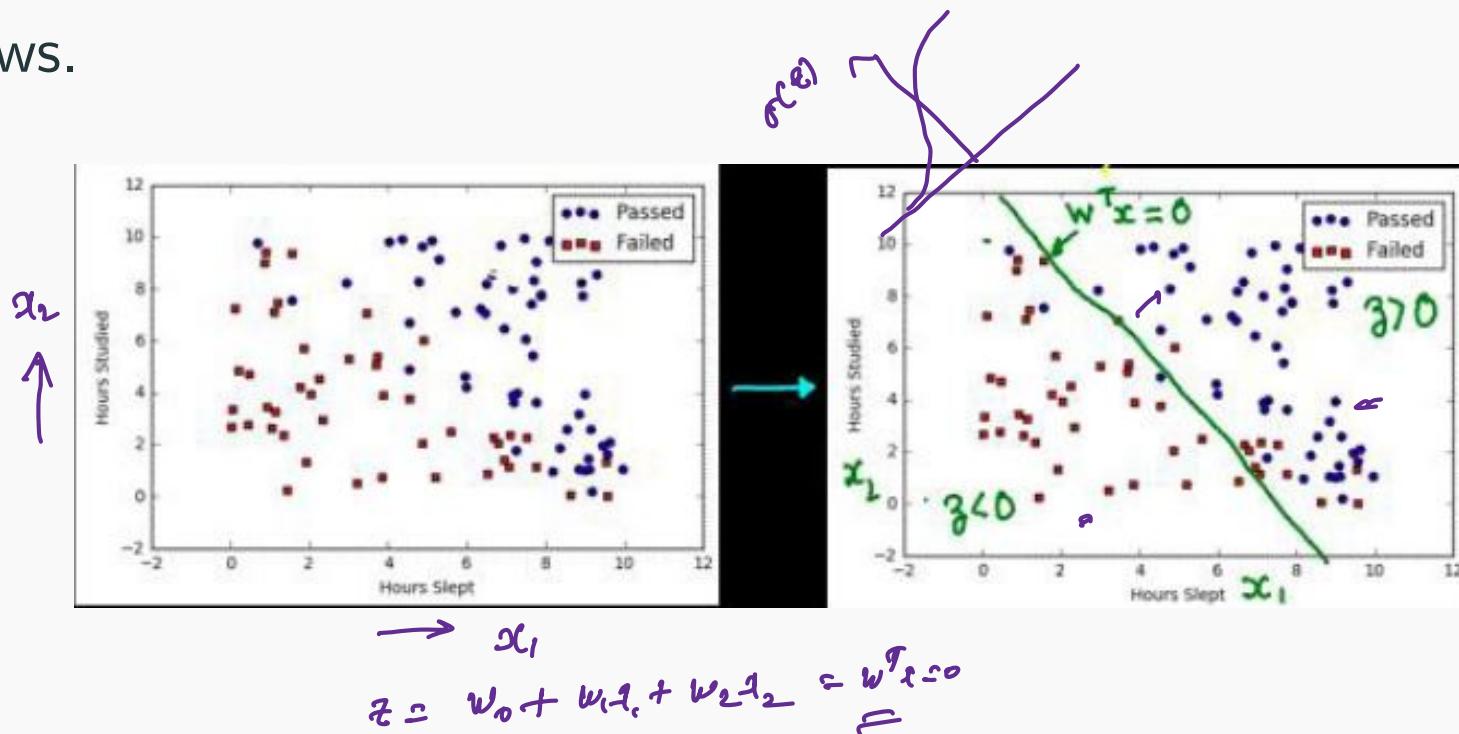
Sigmoid:  $\sigma(z) = \frac{1}{1+e^{-z}}$

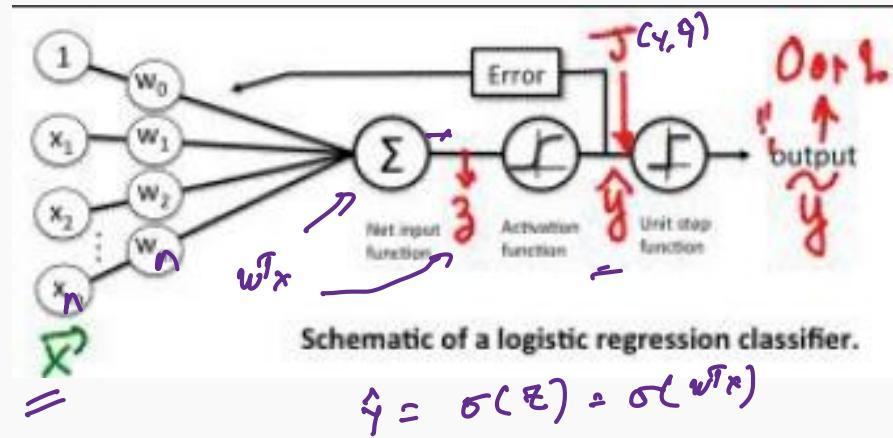
$$z \rightarrow -\infty \rightarrow \sigma(z) \rightarrow 0$$

$$z \rightarrow \infty \rightarrow \sigma(z) \rightarrow 1$$



Geometrically, the net effect of the logistic regression process is as follows.





Example: OR Gate

$[0, 1]$

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

$$z = 0$$

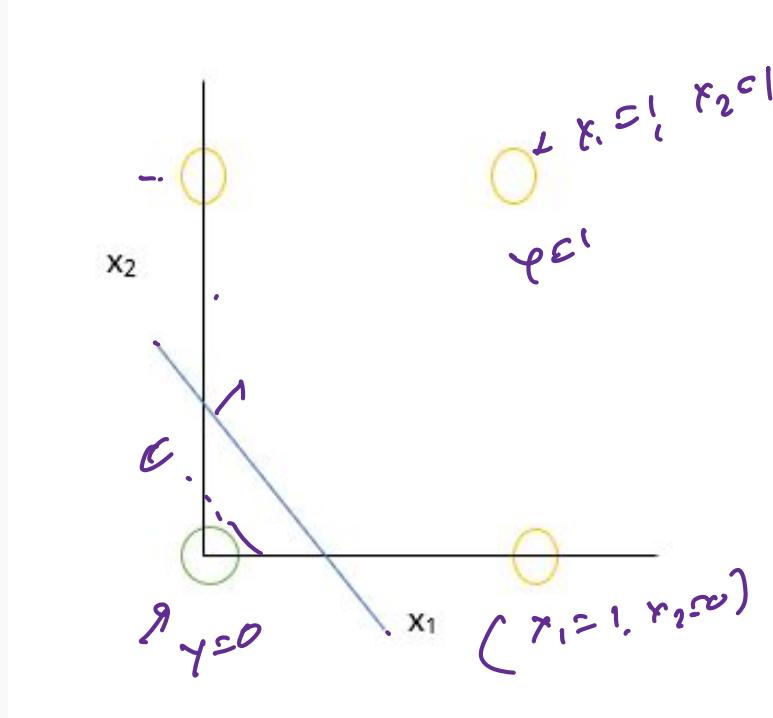
$x_1$	$x_2$	$y$	$z =$	$\hat{y} = \sigma(z)$
0	0	0	< 0 (Need bias)	
0	1	1	> 0	
1	0	1	> 0	
1	1	1	> 0	

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

$$\hat{y}_{\text{OR}} = 1 - e^{-z}$$

$$\hat{y}_{\text{OR}} = 0 + e^{-z}$$

$$z = w_0 + w_1 x_1 + w_2 x_2$$



Find  $\vec{w} \cdot S \cdot t$  we get the OR gate.

$$z = w_0 + w_1 x_1 + w_2 x_2$$

$$w_0 = -1, w_1 = 2, w_2 = 2$$

$$\Rightarrow -1 + 2x_1 + 2x_2 = 0$$

$$\Rightarrow \underline{x_1 + x_2 = 1/2}$$

Loss Function? What about LSE?

$z \rightarrow \sigma(z)$   
4 threshold

$$L = \sum (y^{(i)} - \hat{y}^{(i)})^2$$

$\hat{y} \geq 0$   
 $\hat{y} < 1$

$y$	$\hat{y}$	LSE
0	0	0
1	1	0
0	1	1
1	0	1

Cross-Entropy cost function  $\rightarrow$  Binary Cross Entropy

$$\begin{aligned} \text{Partial cost function } J^{(i)} &= -\{ y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \} \\ &\stackrel{\text{TARGET / GROUND TRUTH}}{=} \sum_{i=1}^m J^{(i)} \end{aligned}$$

$$\begin{aligned} \hat{y} &= \sigma(z) \\ &= \sigma(w^T x) \\ z > 0 &\Rightarrow w^T x - \text{hyper plane} \end{aligned}$$

Gradient Descent:  $w = w - \alpha \nabla_w J$

$\rightarrow$  For cross-entropy  $J \rightarrow \frac{\partial J^{(i)}}{\partial w_j} = -(y^{(i)} - \hat{y}^{(i)})x_j^{(i)}$   
 $-(y^{(i)} - \hat{y}^{(i)})x_j^{(i)}$  Similar to the least squares  $\nabla_w J$

$y = \sigma(z)$   
z - even  
large  
number

$$-[y \log y + (1-y) \log(1-y)]$$

$$\rightarrow y=1, \quad \hat{y} \approx 1$$

$$y \cdot 1 \log 1 + (1-1) \cdot \log(1-\hat{y}) \quad \left| \begin{array}{l} y=0 \\ y=1 \end{array} \right.$$

$$y=1, \quad \hat{y}=0$$

$$-[0+0] \Rightarrow 0$$

## Revisiting OR gate : Gradient Descent

$$z = w_0 + w_1 x_1 + w_2 x_2$$

$$x_0 = 1, \quad z \geq 0,$$

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

Example	$x_0$	$x_1$	$x_2$	$y$	$\tilde{y}$
(1)	1	0	0	0.	1
(2)	1	0	1	1	1
(3)	1	1	0	1	1
(4)	1	1	1	1	1

$$w \leftarrow w - 2 \left[ \frac{\partial J}{\partial w} \right]$$

$$w_j \leftarrow w_j - \frac{\partial J}{\partial w_j}, \quad j=0, 1, 2$$

Do not do this  $\tilde{y}$ . In practice use  $\hat{y}$ .

$$\rightarrow \alpha = 1$$

$$\text{Say } w = [2 \ 2 \ 2]^T$$

$$[\nabla_w J] = \sum_{i=1}^m -(y^{(i)} - \hat{y}^{(i)}) \vec{x}^{(i)}$$

$$= (1) \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + (0) + (0) + (0)$$

$$= [1 \ 0 \ 0]$$

$$\left( \begin{array}{c} \frac{\partial J}{\partial w_0} \\ \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \end{array} \right)$$

$$f(x_1, x_2)$$

$$\nabla f \rightarrow \left( \begin{array}{c} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{array} \right)$$

$$\frac{\partial f}{\partial x}$$

$$\frac{\partial J}{\partial w_0} = (0-1) \cdot 1$$

$$\frac{\partial J^{(i)}}{\partial w_j} = -(\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$$\left( \frac{\partial J}{\partial w_0} \right) = - (0-1) x_0^{(i)} = 1$$

$$\left( \frac{\partial J}{\partial w_1} \right) = - (0-1) x_1^{(i)} = 0$$

$$\left( \frac{\partial J}{\partial w_2} \right) = - (0-1) x_2^{(i)} = 0$$

$$w \leftarrow w - 1 \cdot \frac{\partial J}{\partial w}$$

$$\begin{bmatrix} \end{bmatrix} \leftarrow \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 , y$$

$$J = (y - \hat{y})^2 \text{ L, one data point}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial w_1}$$

$$= 2(y - \hat{y}) \cdot \frac{\partial \hat{y}}{\partial w_1}$$

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

$$\frac{\partial \hat{y}}{\partial w_1} = 0 + x_1 + 0 = x_1$$

$$\frac{\partial J}{\partial w_1} = 2(y - \hat{y}) \cdot x_1 \quad \left. \begin{array}{l} \text{for one} \\ \text{Data Point} \end{array} \right\}$$

$$\frac{\partial J}{\partial w_2} = 2(y - \hat{y}) \cdot x_2$$

$$\frac{\partial J}{\partial w_1} = 2(y^{(i)} - \hat{y}^{(i)}) \cdot x_1$$

$$\frac{\partial J}{\partial w_j} = 2(y^{(i)} - \hat{y}^{(i)}) \cdot x_j \quad \begin{array}{l} \text{gradient} \\ \text{combination} \\ \text{from one} \\ \text{data point} \\ (i) \end{array}$$

m - data points

$$\frac{\partial J}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m 2(y^{(i)} - \hat{y}^{(i)}) \cdot x_j$$

↳ Do this for all j

$$w_j \leftarrow [w_j - 2 \frac{\partial J}{\partial w_j}] \rightarrow \begin{array}{l} \text{Linear Regression} \\ \text{Logistic Regression} \end{array}$$

$$\Rightarrow w = w - [1 \ 0 \ 0] = [1 \ 2 \ 2]$$

↓

$$[0 \ 2 \ 2]$$

↓

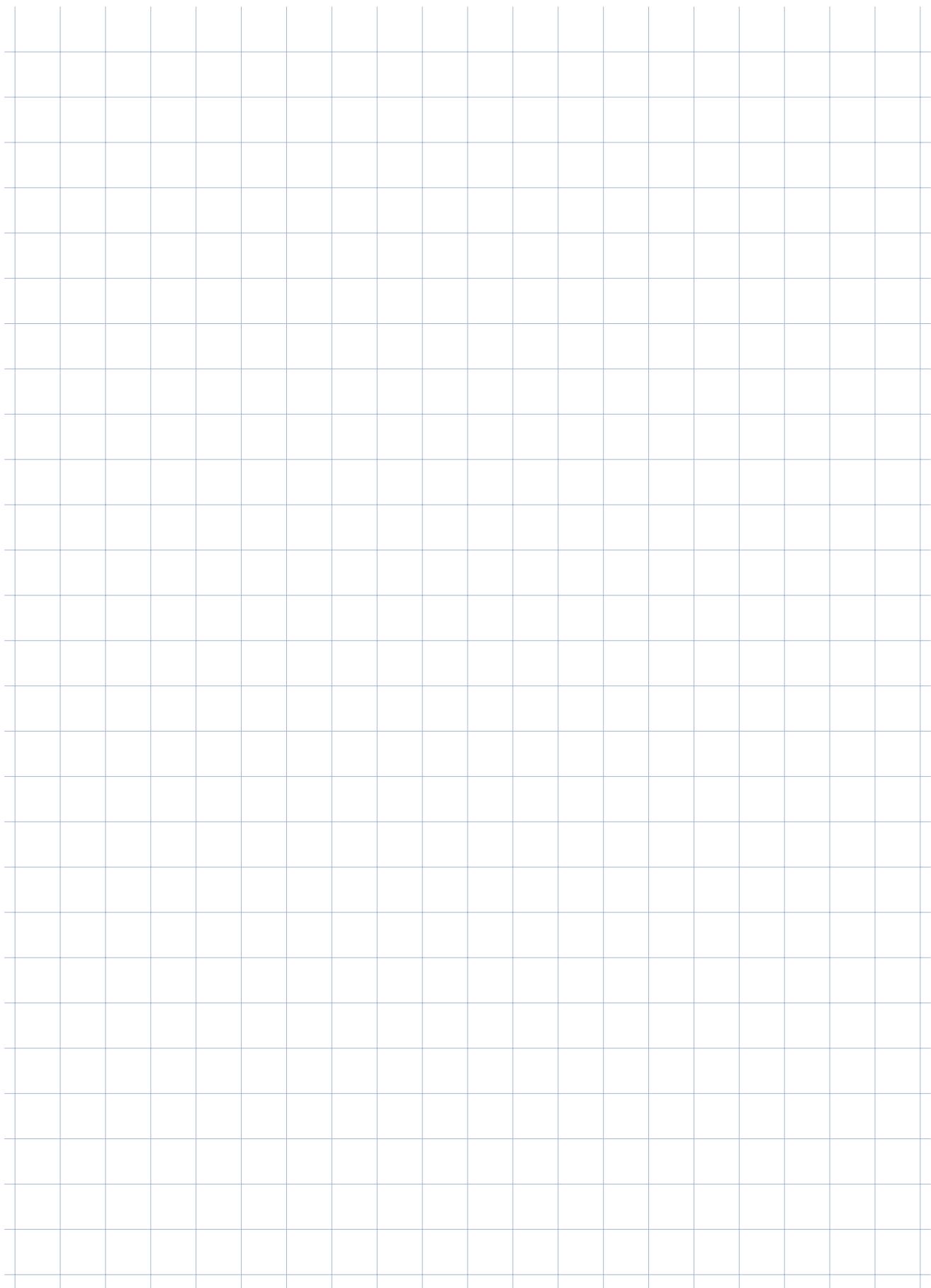
$$[-1 \ 2 \ 2]$$

$$\tilde{y} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} = y$$

$$\nabla_w J = 0$$

When number of classes > 2, k=2 Binary classification

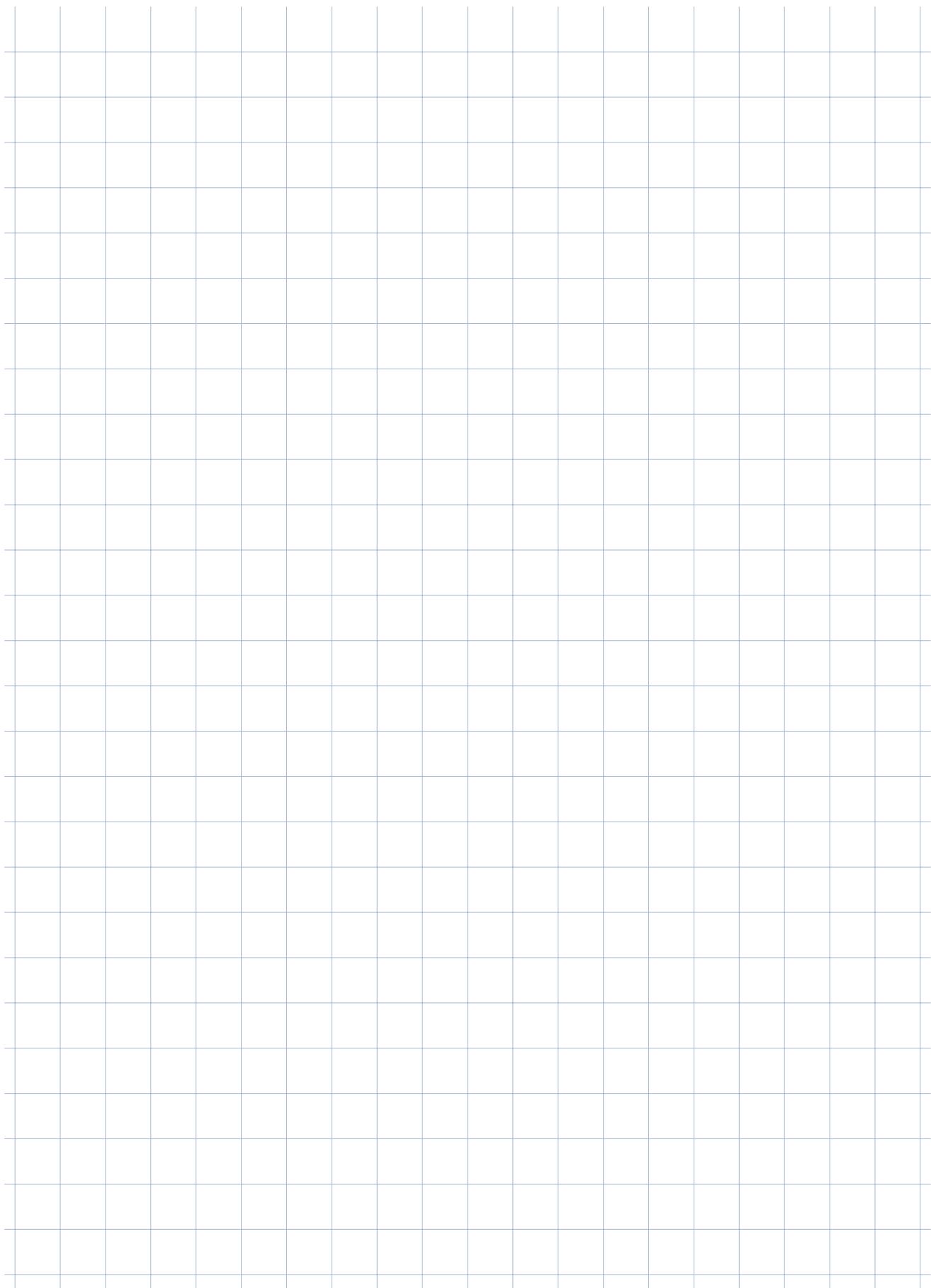
1. Some representation for  $y$ .
2. Loss function.
3. Non-linearity



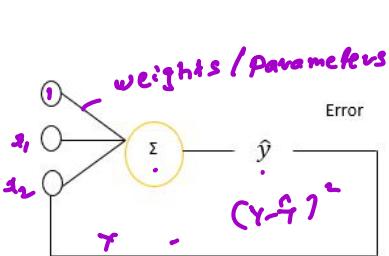
Linear Regression

Logistic Regression

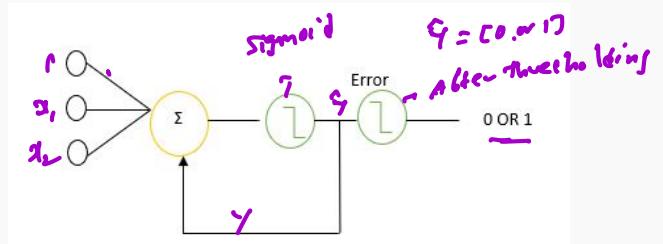




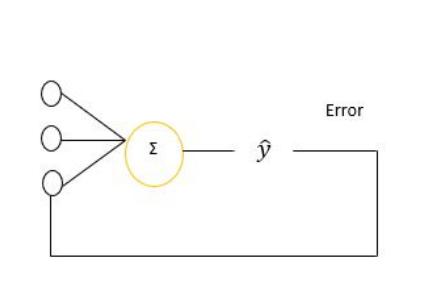
# Linear Regression



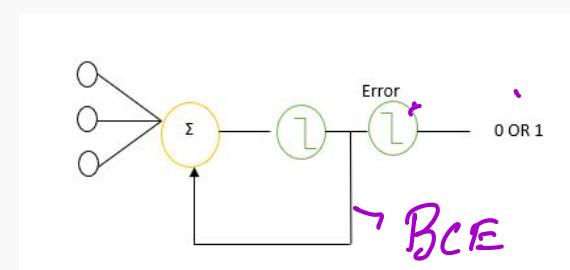
# Logistic Regression



## Linear Regression



## Logistic Regression



Feed forward process: i.e. Finding output: Given (Inputs, Parameters)

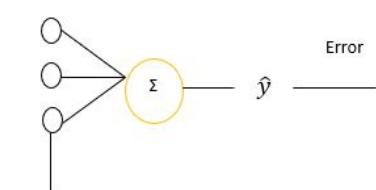
$$\text{Loss fn} \Rightarrow L(y, \hat{y}) = -(y - \hat{y})^2$$

Actual price      Estimated price

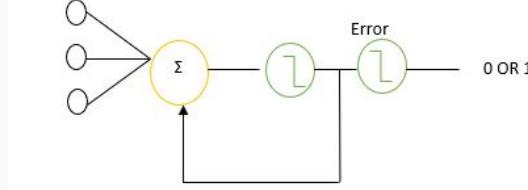
$$\text{Logistic } y = 0, \infty$$
$$\hat{y} = \text{prob}(y=1|x)$$

$$\hat{y} = \frac{1}{1 + e^{-y}}$$

## Linear Regression



## Logistic Regression



Feed forward process: i.e. Finding output: Given (Inputs, Parameters)

$$\hat{y} = \vec{w}^T \vec{x} = \vec{w} \cdot \vec{x}$$

i.e.  $\hat{y}^{(i)} = \sum_{j=0}^n w_j x_j^{(i)}$ ;  $x_0^{(i)} = 1$

$$\hat{y} = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}} = \text{sigmoid}(\vec{w} \cdot \vec{x})$$

A quantized classifier is applied after this.

Loss function:  $L = \sum_{i=1}^m L^{(i)}$

$\nwarrow$  Over the Samples  
or  
Training Examples

Loss function:  $L = \sum_{i=1}^m L^{(i)}$

$$L^{(i)} = \frac{1}{2} \left[ \underbrace{y^{(i)}_{\text{Groundtruth}} - \hat{y}^{(i)}_{\text{Prediction}}} \right]^2$$

$$\begin{aligned} L^{(i)} &= - \left\{ \underbrace{y^{(i)} \ln \hat{y}^{(i)}}_{\text{Cross-entropy}} + \underbrace{(1 - y^{(i)}) \ln (1 - \hat{y}^{(i)})}_{\text{Costfunction}} \right\} \\ &= - \sum_{k=0}^1 y_k^{(i)} \ln \hat{y}_k^{(i)} \\ &\quad \hookrightarrow \text{over the categories} \end{aligned}$$

---

Loss function:  $L = \sum_{i=1}^m L^{(i)}$

$$L^{(i)} = \frac{1}{2} \left[ \begin{matrix} y^{(i)} \\ \text{Groundtruth} \end{matrix} - \begin{matrix} \hat{y}^{(i)} \\ \text{Prediction} \end{matrix} \right]^2$$

$$\begin{aligned} L^{(i)} &= - \left\{ \begin{matrix} y^{(i)} \ln \hat{y}^{(i)} \\ \text{Cross-entropy} \end{matrix} + \begin{matrix} (1 - y^{(i)}) \ln (1 - \hat{y}^{(i)}) \\ \text{Costfunction} \end{matrix} \right\} \\ &= - \sum_{k=0}^1 y_k^{(i)} \ln \hat{y}_k^{(i)} \end{aligned}$$

---

Gradient:

→ EXPECTATION

Loss function:  $L = \frac{1}{m} \sum_{i=1}^m L^{(i)}$

$$L^{(i)} = \frac{1}{2} \left[ \underset{\text{Groundtruth}}{y^{(i)}} - \underset{\text{Prediction}}{\hat{y}^{(i)}} \right]^2$$

$$\begin{aligned} L^{(i)} &= - \left\{ \begin{array}{l} y^{(i)} \ln \hat{y}^{(i)} \\ \text{Cross-entropy} \end{array} + (1 - y^{(i)}) \ln (1 - \hat{y}^{(i)}) \right\} \\ &\quad \text{Costfunction} \\ &= - \sum_{k=0}^1 y_k^{(i)} \ln \hat{y}_k^{(i)} \end{aligned}$$

Gradient:

$$\frac{\partial L^{(i)}}{\partial w_j} = \{\hat{y}^{(i)} - y^{(i)}\} x_j^{(i)}$$

$\overbrace{n\text{-weights}}$  ) parameter

$$\frac{\partial L^{(i)}}{\partial w_j} = \underbrace{\{\hat{y}_1^{(i)} - y_1^{(i)}\} x_j^{(i)}}_{w_j = w_j - \alpha \frac{\partial L^{(i)}}{\partial w_j}}$$

## General Algorithm for finding weights for the feedforward network

1. Initialize weights via some guess  $\overrightarrow{w_{init}}$

Epoch  $\rightarrow$  ALL TRAINING DATA

2. Use feedforward network to find prediction  $\hat{y}$ . Typically this will be different from ground truth  $y$ .

3. Improve weights through gradient descent.

$$\frac{\delta L}{L} = \frac{|L_{\text{Prev}} - L_{\text{current}}|}{L} < 10^{-3}$$
$$w = w - \alpha \frac{\partial L}{\partial w}$$

Gradient Descent Step

Repeat step 2 and 3 til convergence.

Multinomial classification of  $K > 2$  classes requires

- Representation for  $y \rightarrow$

- Nonlinearity  $\rightarrow$  Sigmoid

- Cost/ Loss function - Categorical Cross Entropy

Representation for  $y$

One-hot vector: say 3 classes

one hot vector representation

1 of  $k$

Bad predictions

$$\hat{y} \quad y$$

$$\begin{bmatrix} 0.6 \\ 0.1 \\ 0.05 \\ \dots \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ OR } \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ OR } \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

Class1                    Class2                    Class3

Good predictions

$$\rightarrow \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 0.01 \\ 0.2 \\ 0.9 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

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

$$\begin{bmatrix} P(y_1 = 1|x) \\ P(y_2 = 1|x) \\ \vdots \\ \vdots \\ P(y_k = 1|x) \end{bmatrix}$$

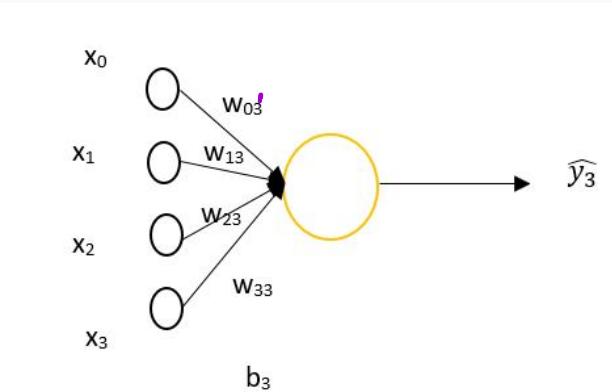
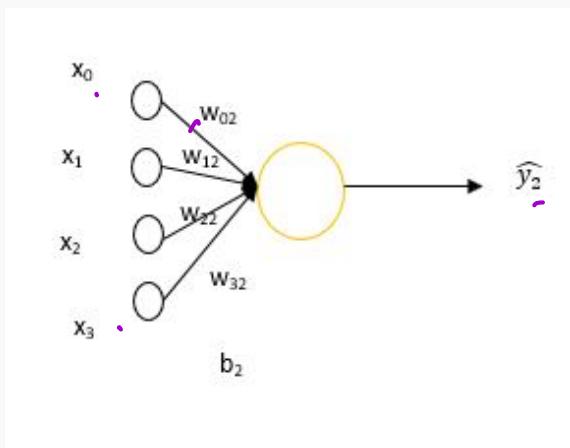
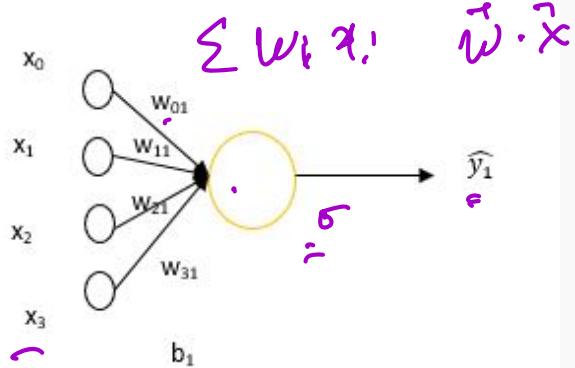
1 of  $k$   
classes

In general,  $\hat{y} =$

( $\hat{y}_1$ ,  $-\hat{y}_1$ )

Schematic:  $K=3$  classes

3. 7



Separate  $\hat{y}$  for each class  
using  $k$  different ( $k=3$ ) sets of weights / neurons

$\hat{y}_1 = g(w_{01}x_0 + w_{11}x_1 + w_{21}x_2 + w_{31}x_3)$   
 $W_{ij}$ , where i is the feature and j is the class

$$\vec{y} = g(\underbrace{W^T x}_{z_i} + \vec{b})$$

$$z_i = \underline{w^T x} + \vec{b}$$

What is g?

Suppose we chose  $g = \text{Sigmoid function}$

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

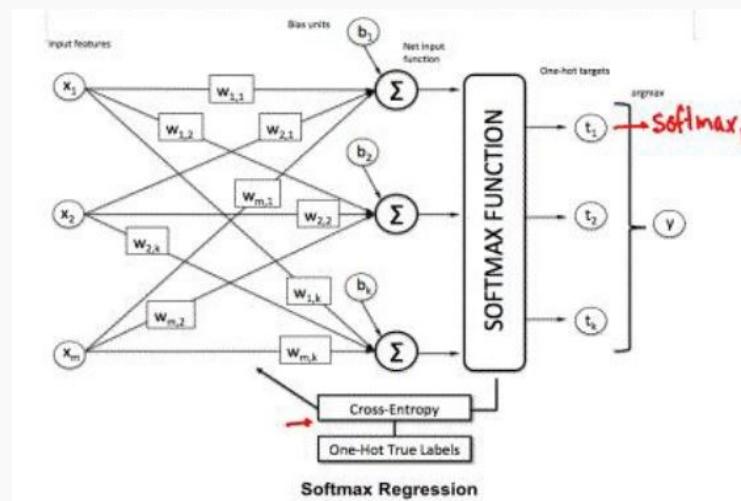
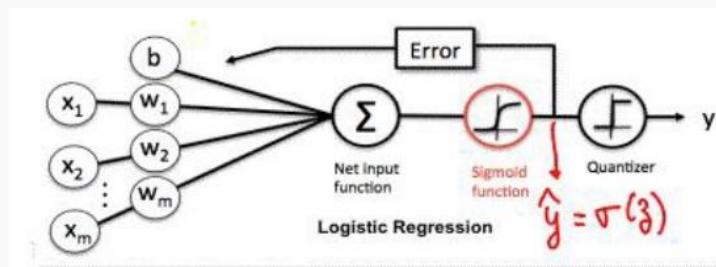
$$z_1 \rightarrow \circ \rightarrow \hat{y}_1 = g(z_1)$$

$$g = \sigma$$

$$z_2 \rightarrow \circ \rightarrow \hat{y}_2 = g(z_2)$$

$$z_3 \rightarrow \circ \rightarrow \hat{y}_3 = g(z_3)$$

In general, if we set  $g = \text{sigmoid}$ ,  $\sum_{k=1}^k \hat{y}_k \neq 1$ . So, we use softmax.  $\text{softmax}_j = \frac{\exp(z_j)}{\sum_{k=1}^k \exp(z_k)}$



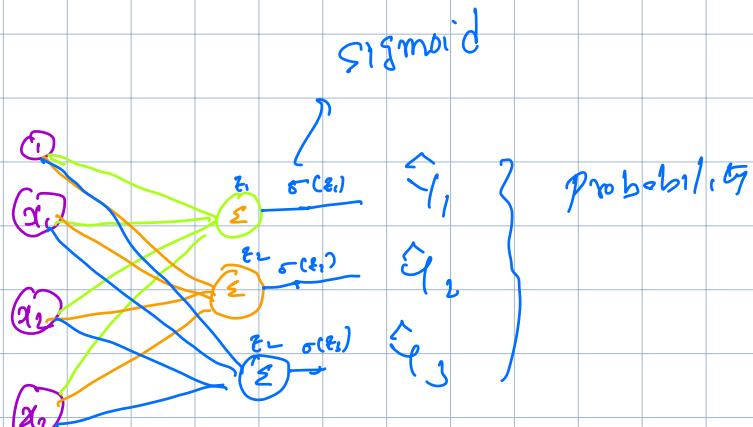
Cost function (Cross-Entropy) :

$$L^{(i)} = - \sum_{k=1}^K y_k^{(i)} \ln \hat{y}_k^{(i)}$$

$y_1 \quad y_2 \quad y_3$   
 $[0 \ 1 \ 0]$

$\hat{y}_1 \quad \hat{y}_2 \quad \hat{y}_3$   
 $[0.05 \ 0.8 \ 0.15]$

↴ for one data point



Independent

$$y_1 + y_2 + y_3 \neq 1$$

$$\left\{ \begin{array}{l} \hat{y}_1 = p_0(y_1|x) \\ \hat{y}_2 = p_1(y_2|x) \\ \hat{y}_3 = p_2(y_3|x) \end{array} \right.$$

Softmax output

$$\hat{y}_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

softmax  
output

Compare this with binary cross-entropy, which is  $-\sum_{k=1}^2 y_k^{(i)} \ln \hat{y}_k^{(i)}$

$$L = \sum_{i=1}^m L^{(i)} = - \sum_{i=1}^m \left[ \sum_{k=1}^K y_k^{(i)} \ln \hat{y}_k^{(i)} \right]$$

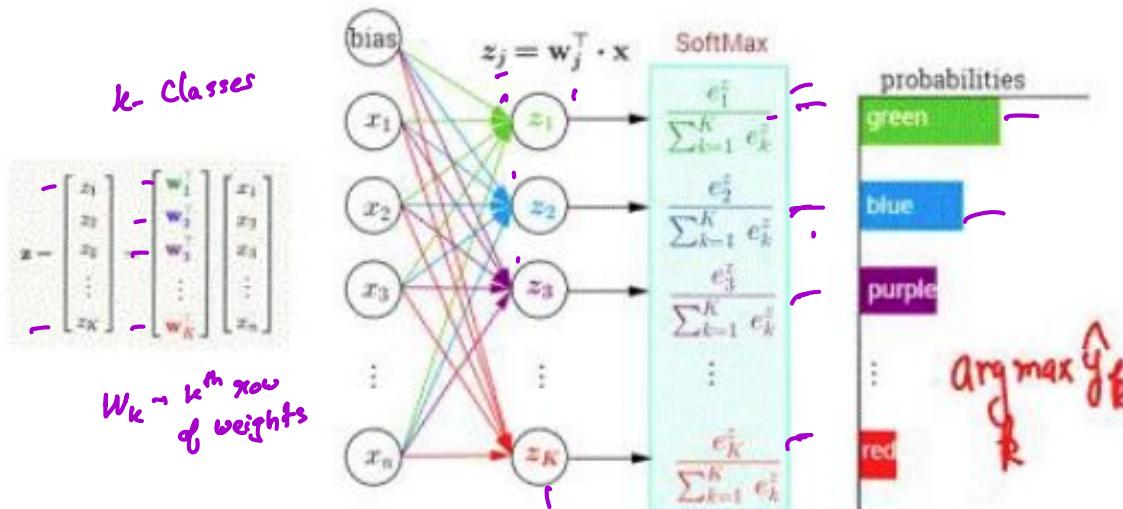
*GrT*      *probability out put*

# Summary: Binary vs Multinomial Classification

Binary	Multinomial
<p>If there are two classes</p> <ul style="list-style-type: none"> <li>→ <math>\hat{y} = \sigma(\vec{w} \cdot \vec{x})</math></li> <li>    <math>\hat{y}</math>: Scalar <math>\epsilon [0, 1]</math></li> <li>→ <math>w</math> is a vector: <math>(n + 1) \times 1</math></li> <li>→ <math>L^{(i)} = - \sum_{k=1}^2 y_k^{(i)} \ln \hat{y}_k^{(i)}</math></li> </ul> <p><i>n - No of feature</i></p> <p><i>+ - Bias unit</i></p>	<p>For more than 2 classes</p> $\hat{y}_j = \frac{e^{w_j \cdot \vec{x}}}{\sum_{i=1}^K e^{w_i \cdot \vec{x}}} \quad j \sim \text{class index}$ <p><math>\hat{y}</math>: One hot vector</p> <p><math>w</math> is a matrix: <math>(n + 1) \times K</math></p> $L^{(i)} = - \sum_{k=1}^K y_k^{(i)} \ln \hat{y}_k^{(i)}$ <p><i>x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub></i></p> $\begin{pmatrix} w_{01} & w_{02} & \dots & w_{03} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \dots & w_{n3} \end{pmatrix}$

# Examples

## Multi-Class Classification with NN and SoftMax Function



$$\sum_k y_k \ln \hat{y}_k$$

$$\hat{y}_k = \frac{e^{z_k}}{\sum_{k=1}^K e^{z_k}}$$

$$\ln(\hat{y}_k)$$
 ~~$z_k = \dots$~~ 

$$z_k \leq z_k$$

$$\hat{Y}_1 = \frac{e^{z_1}}{\sum_{k=1}^K e^{z_k}}$$

$$\ln \hat{Y}_1 = \ln(e^{z_1}) - \ln\left(\sum_{k=1}^K e^{z_k}\right)$$

$$= z_1 - \ln\left(\sum_{k=1}^K e^{z_k}\right)$$

$L = - \sum_{k=1}^K y_k \ln \hat{y}_k$

$$z_i = w_i^T x$$

↳ class index

$$\frac{\partial L}{\partial w}$$

$$\rightarrow \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix}$$

features,  $n=3$

$$\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \xrightarrow{\text{one hot encoding}} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \xrightarrow{\text{softmax}} \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{pmatrix}$$

$$q = \left( \begin{matrix} q_1 \\ q_2 \\ q_3 \end{matrix} \right)$$

$$G \cdot T \quad \gamma_2 \left( \begin{array}{c} 0 \xrightarrow{\gamma_1} \\ 1 \xrightarrow{\gamma_2} \\ 0 \xrightarrow{\gamma_3} \end{array} \right)$$

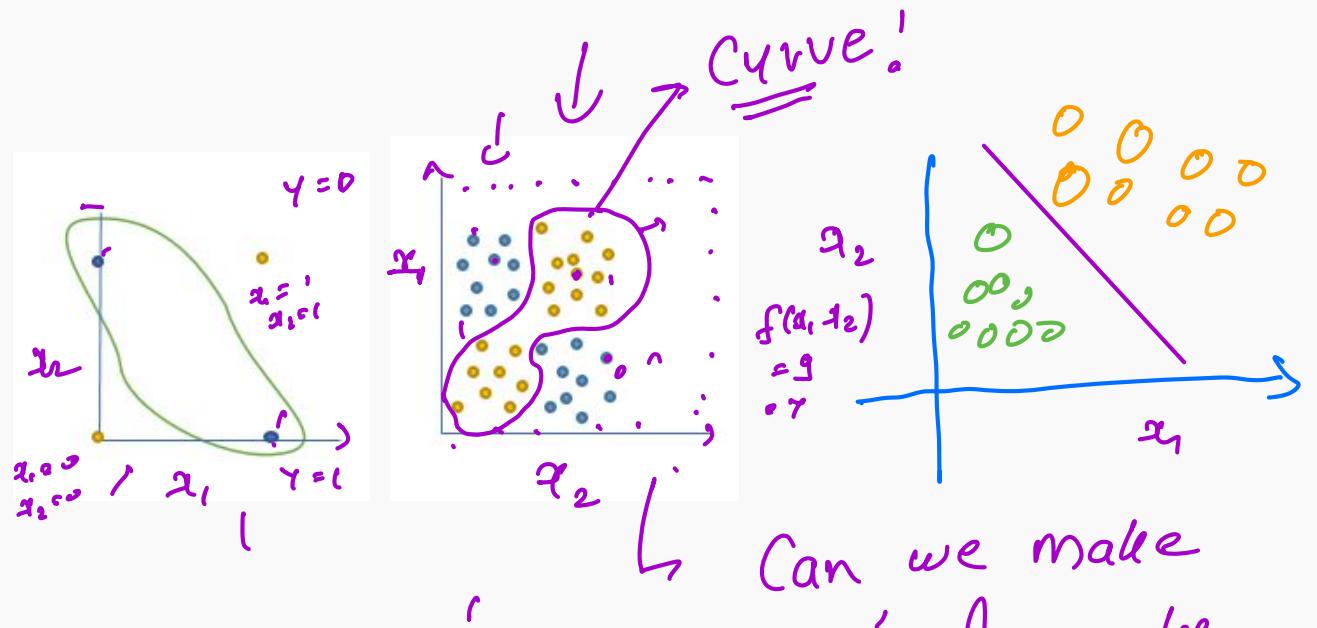
$$G \cdot T : Y = \begin{pmatrix} Y_1 \\ Y_2 \\ Y_3 \end{pmatrix}$$

one date part

$$\ln \hat{q} = \begin{pmatrix} \ln q_1 \\ \ln q_2 \\ \ln q_3 \end{pmatrix} \quad \begin{pmatrix} q_1 \\ q_2 \\ q_3 \end{pmatrix}$$

$$\sum_k \gamma_k \ln \gamma_k = (\gamma_1, \gamma_2, \gamma_3) \begin{pmatrix} \ln \gamma_1 \\ \ln \gamma_2 \\ \ln \gamma_3 \end{pmatrix}^{-1}$$

# Why Deep networks?



Can we make  
this linearly  
separable?

Not linearly separable: XOR

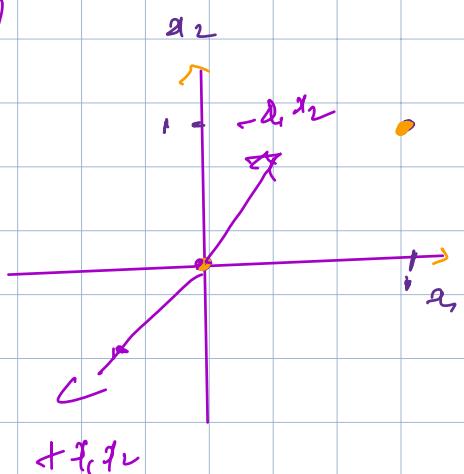
		$x_1$	$x_2$	$y$
		A	B	XOR
0	0	0		0
0	1		1	1
1	0		1	1
1	1		0	0

→ A classifier that  
will mimic this  
logic gate?

$$x_1 \in [0, 1] \\ x_2 \in [0, 1]$$

$$Y = x_1 + x_2 - 2x_1 x_2$$

$x_1$	$x_2$	$Y$
0	0	0
0	1	+1
1	0	+1
1	1	0



$$Y = x_1 + x_2 - 2x_1 x_2 \leftarrow \text{logic gate}$$

$\hookrightarrow$  non-linear term

$$x_1 x_2 = x_3$$

$$Y = [x_1 + x_2 - 2x_3]$$

$\hookrightarrow$  hyperplane

# Why Deep networks?

Solution:

Use nonlinear features (Kernel Trick) can lead to combinational explosion (cures of dimensionality)  $\rightarrow X_1^2, \underline{X_1 X_2}, \underline{\underline{X_2^2}}$ .

Trade dimension for complexity



→ Neural Networks with "hidden" layers.

"learn"

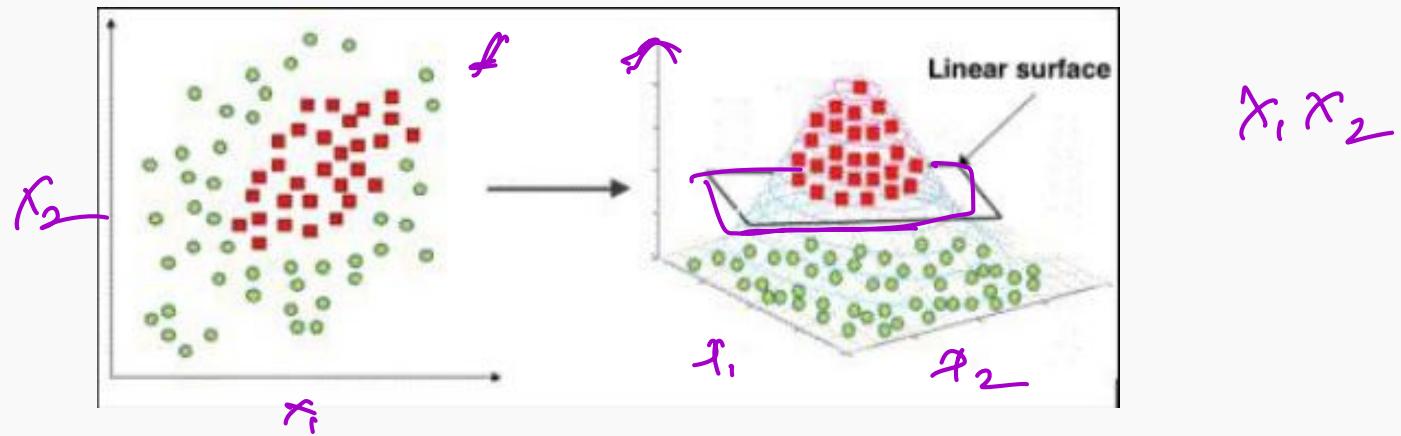
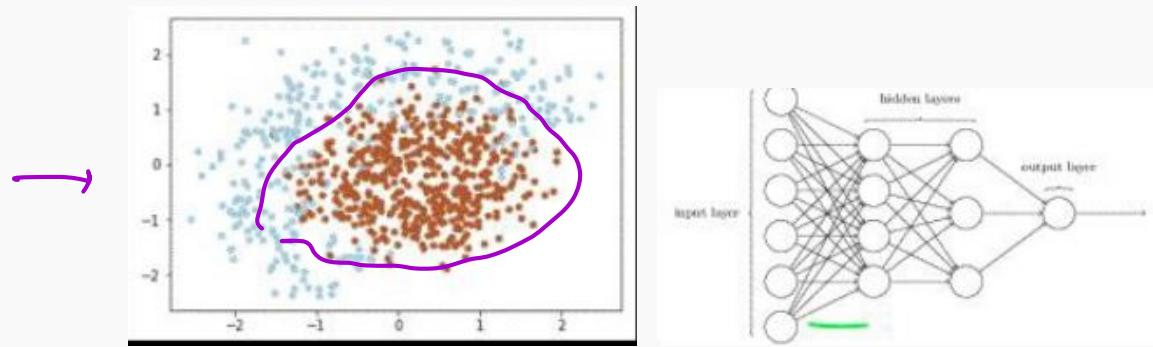


Universal Approximation Theorem

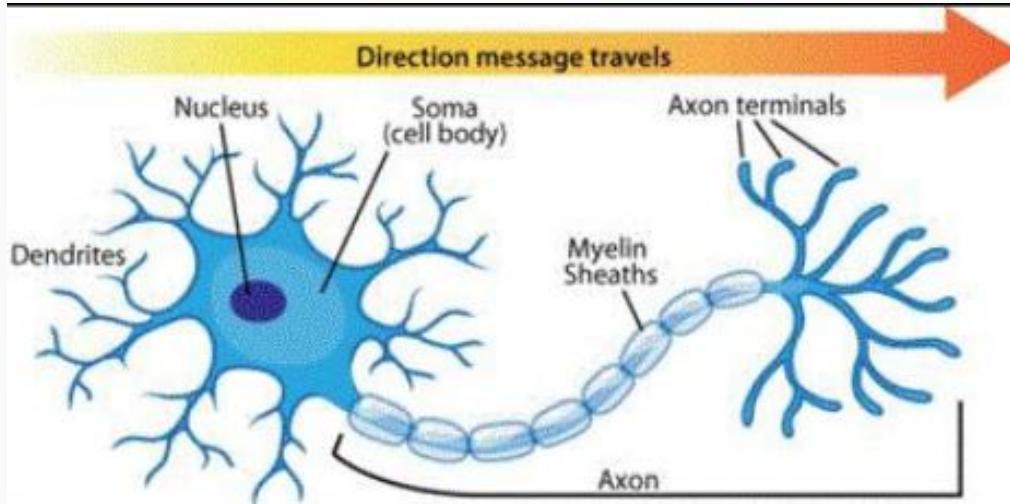


Any function can be approximated arbitrarily close by a Neural Network with a single hidden layer.

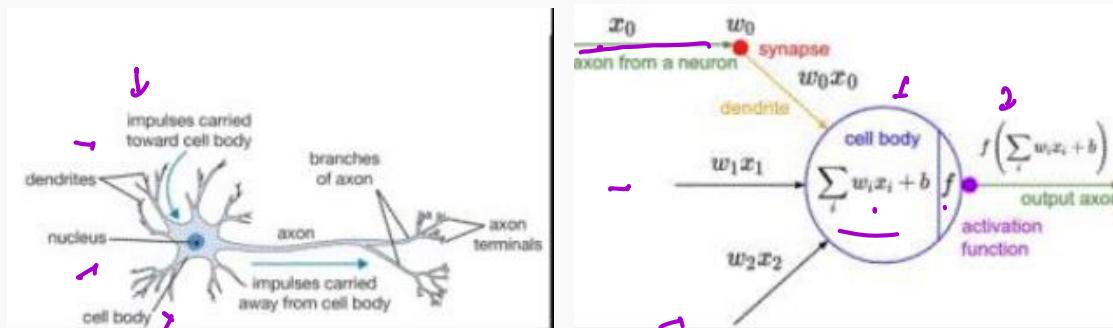
# Deep Networks



# Biological Analogy



→



# How does our Brain Work?

- A neuron is connected to other neurons via its input and output links  
*↑ signal from a neuron*
- Each incoming neuron has an activation value and each connection has a weight associated with it
- The neuron sums the incoming weighted values and this value is input to an activation function (Sigmoid)
- The output of the activation function is the output from the neuron

# Neural Networks

---

# Neural Networks

Linear Model:

---

# Neural Networks

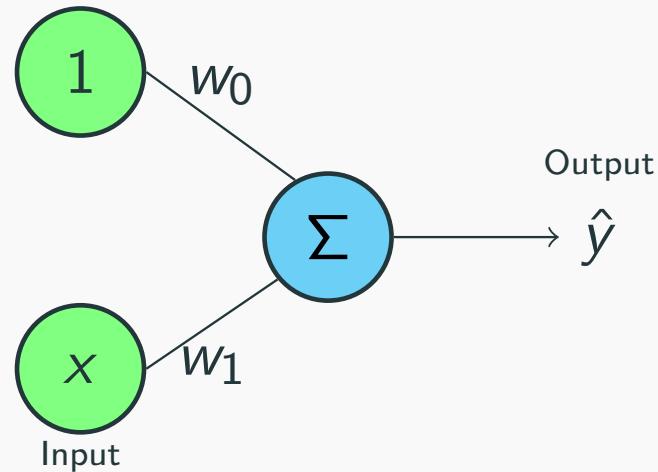
Linear Model:

$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

---

# Neural Networks

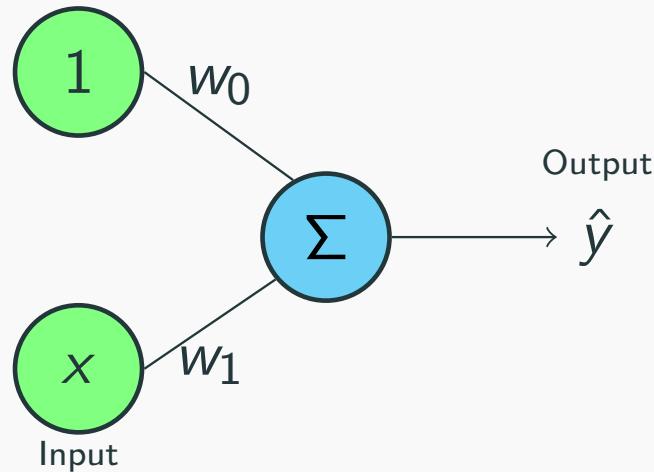
## Linear Model:



$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

# Neural Networks

Linear Model:



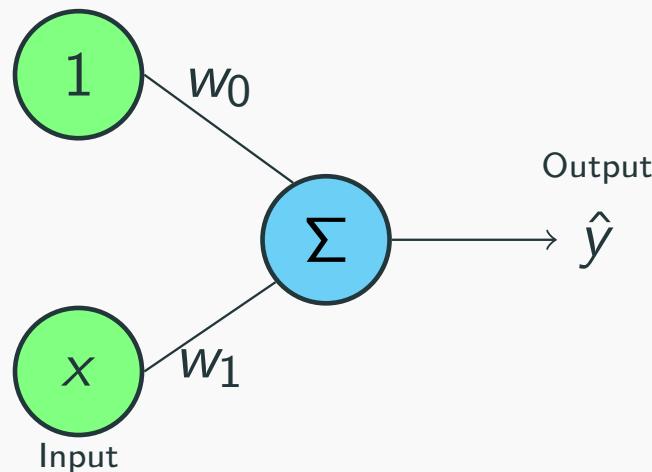
Neural network model:

$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

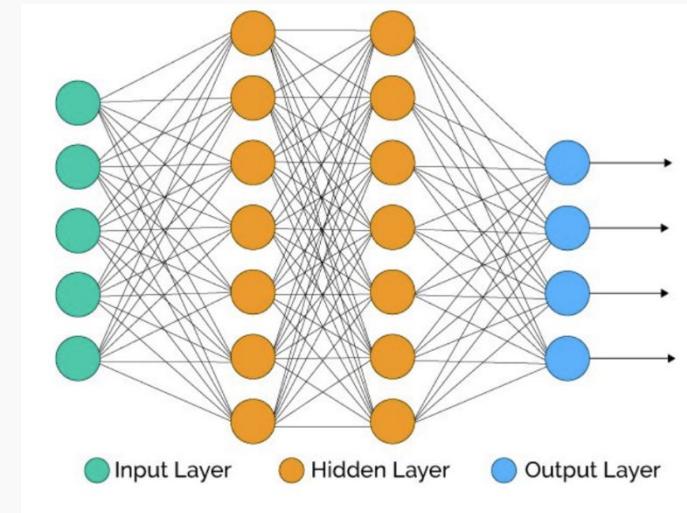
# Neural Networks

Neural network model:

Linear Model:



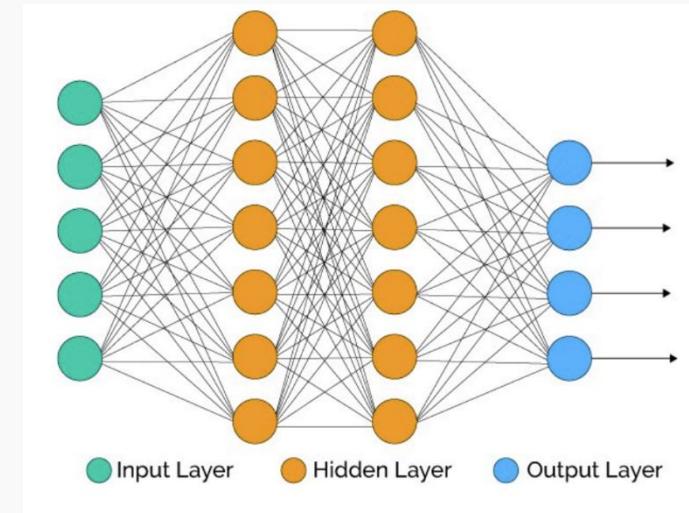
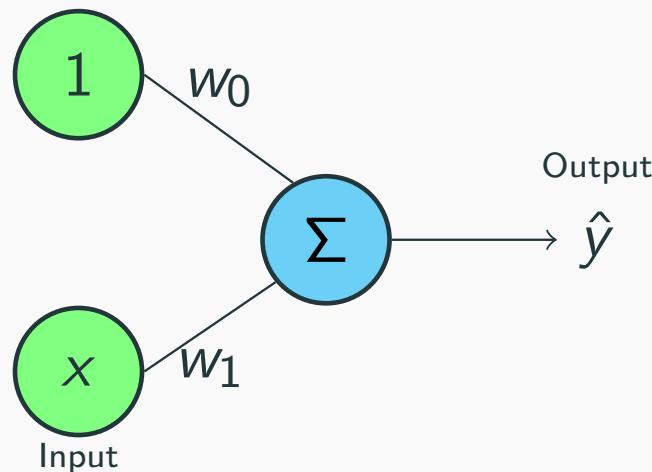
$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$



# Neural Networks

Neural network model:

Linear Model:



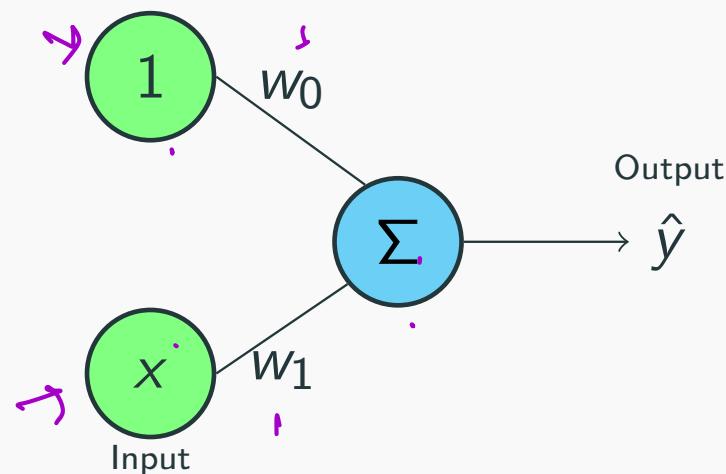
$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

A diagram showing the activation function  $g$ . It consists of three circles. The leftmost circle is green and contains the letter 'a'. The middle circle is light blue and contains the Greek letter  $\Sigma$ . An arrow points from the  $\Sigma$  circle to the rightmost circle, which is also light blue and contains the letter 'g'. This represents the function  $a = \Sigma \rightarrow g$ .

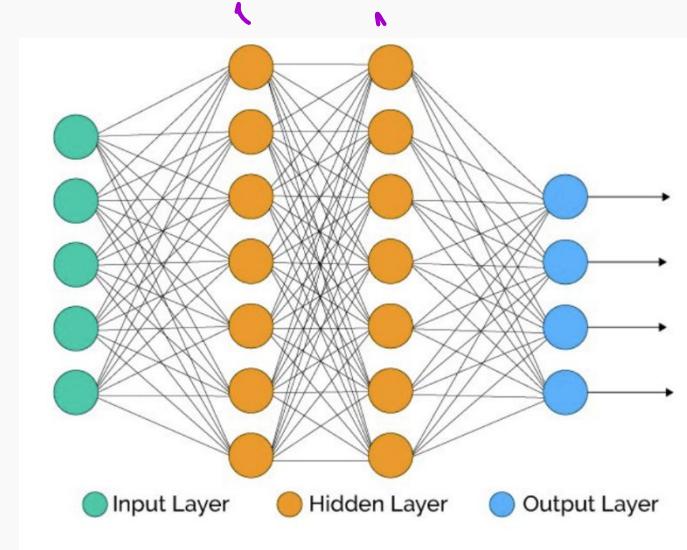
# Neural Networks

Neural network model:

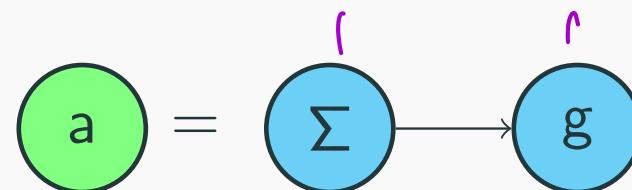
Linear Model:



$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$



$$a = g\left(\sum w_i x_i\right)$$

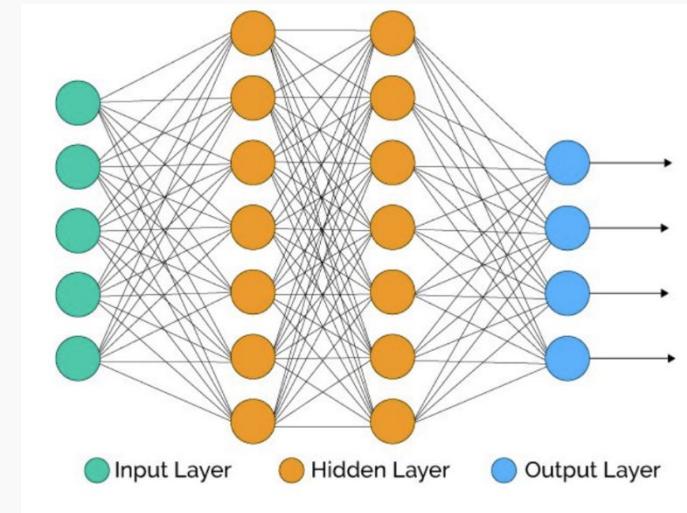
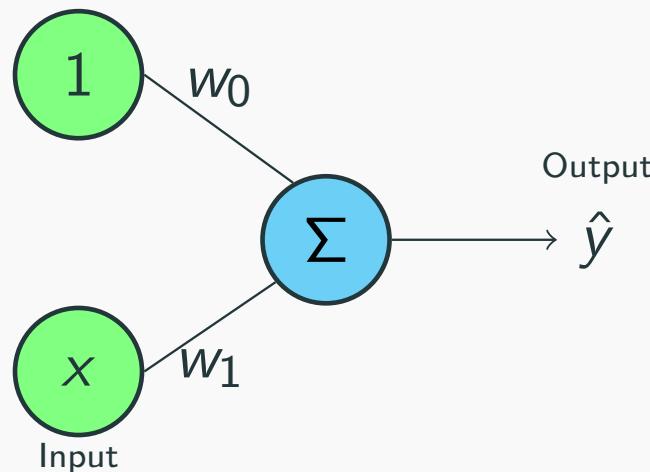


A Neuron has a linear and a nonlinear operation.

# Neural Networks

Neural network model:

Linear Model:



$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$



A Neuron has a linear and a nonlinear operation.

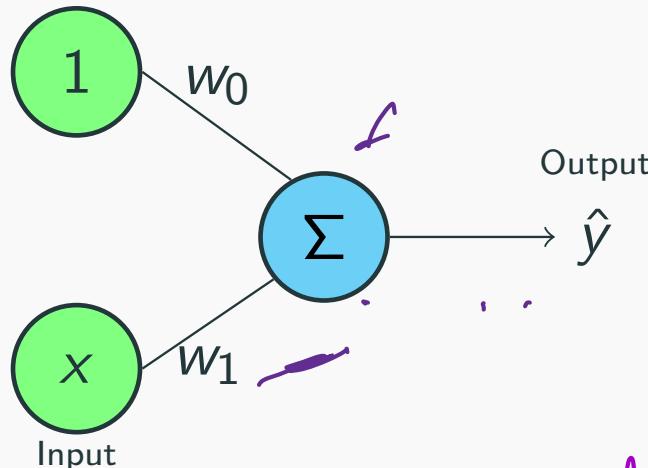
---

Example of non-linear function is  $Sigmoid(x) = \frac{1}{1+exp(-x)}$

# Neural Networks

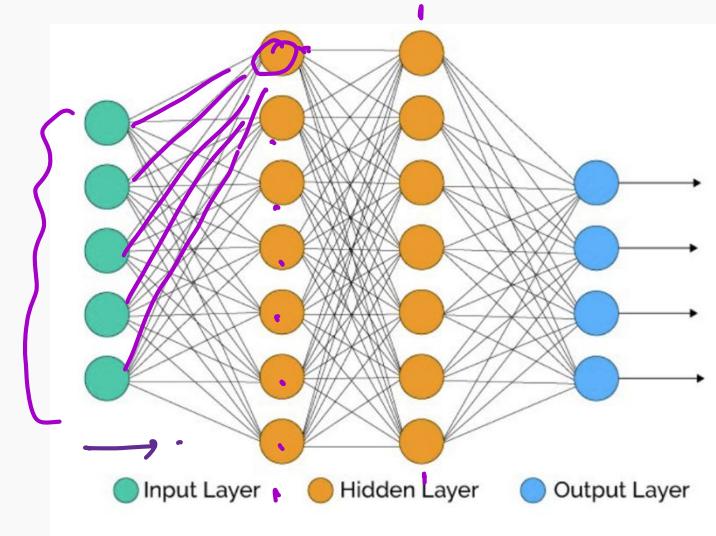
Neural network model:

Linear Model:



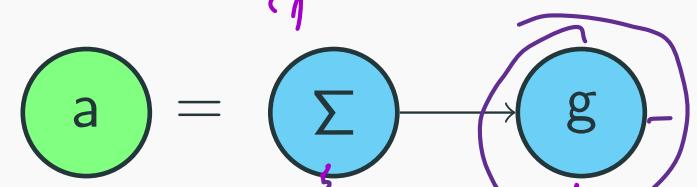
At least a hidden layer

$$\hat{y} = w_0 + w_1 x = \sum w_i x_i$$

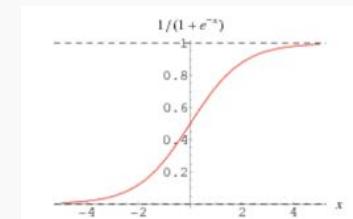


$$a = g(\sum w_i x_i)$$

weighted sum  $\rightarrow$  non-linearity

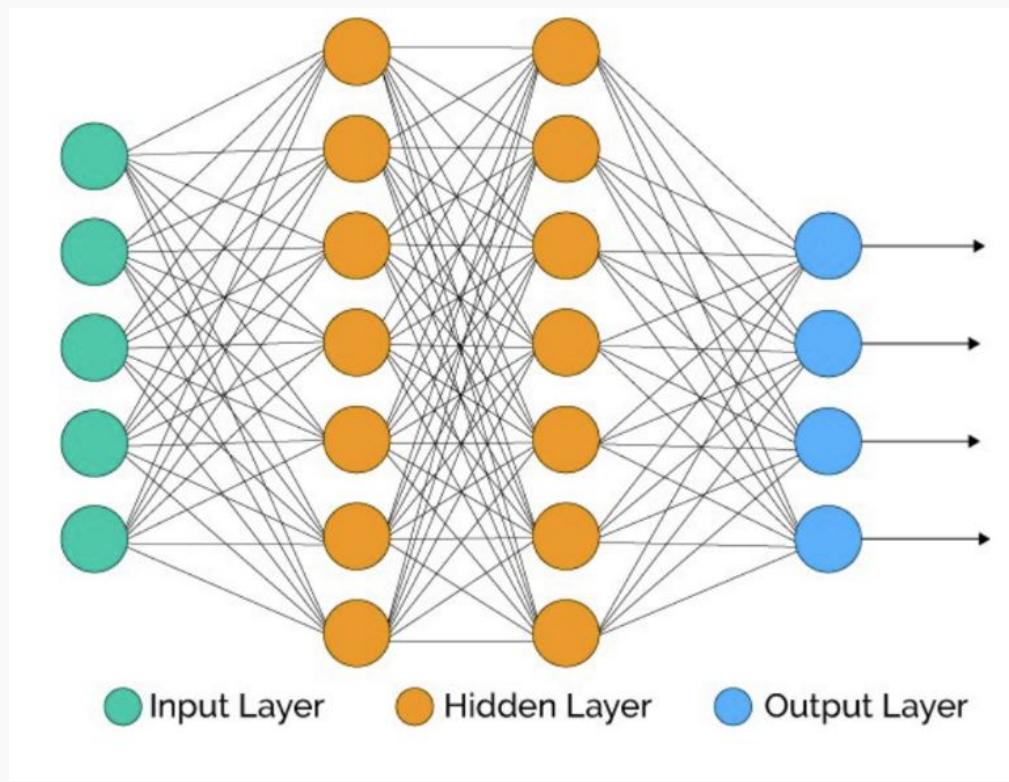


A Neuron has a linear and a nonlinear operation.  
Creates non-linear features

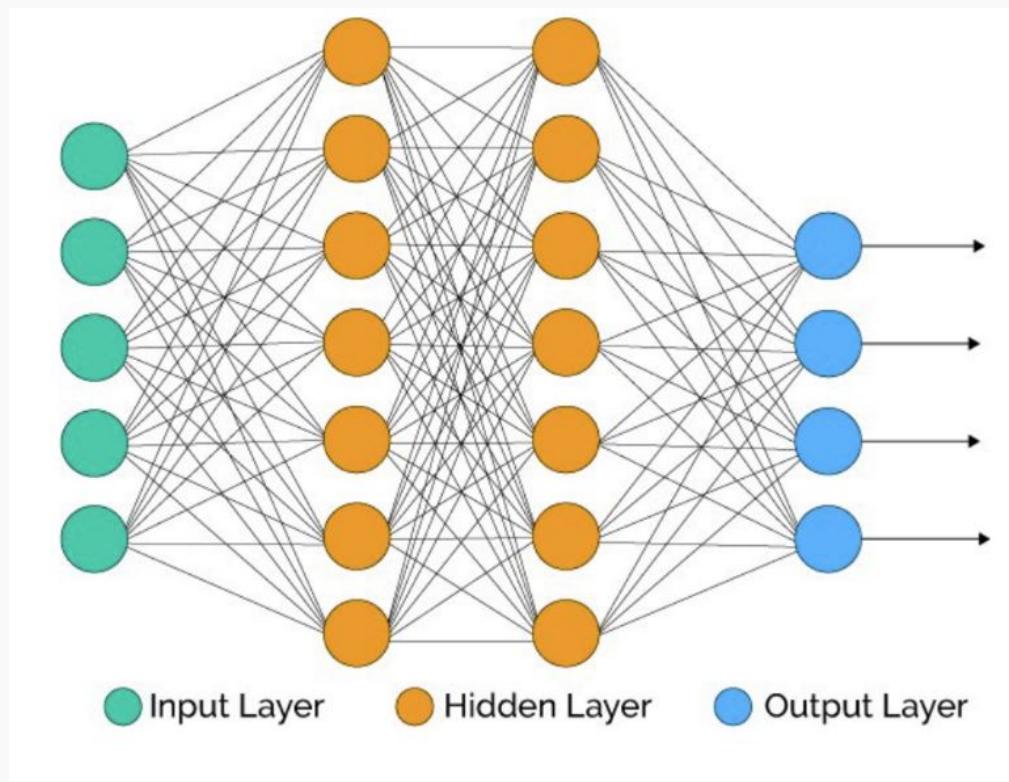


Example of non-linear function is  $Sigmoid(x) = \frac{1}{1+\exp(-x)}$

# Universal Approximation Theorem



# Universal Approximation Theorem

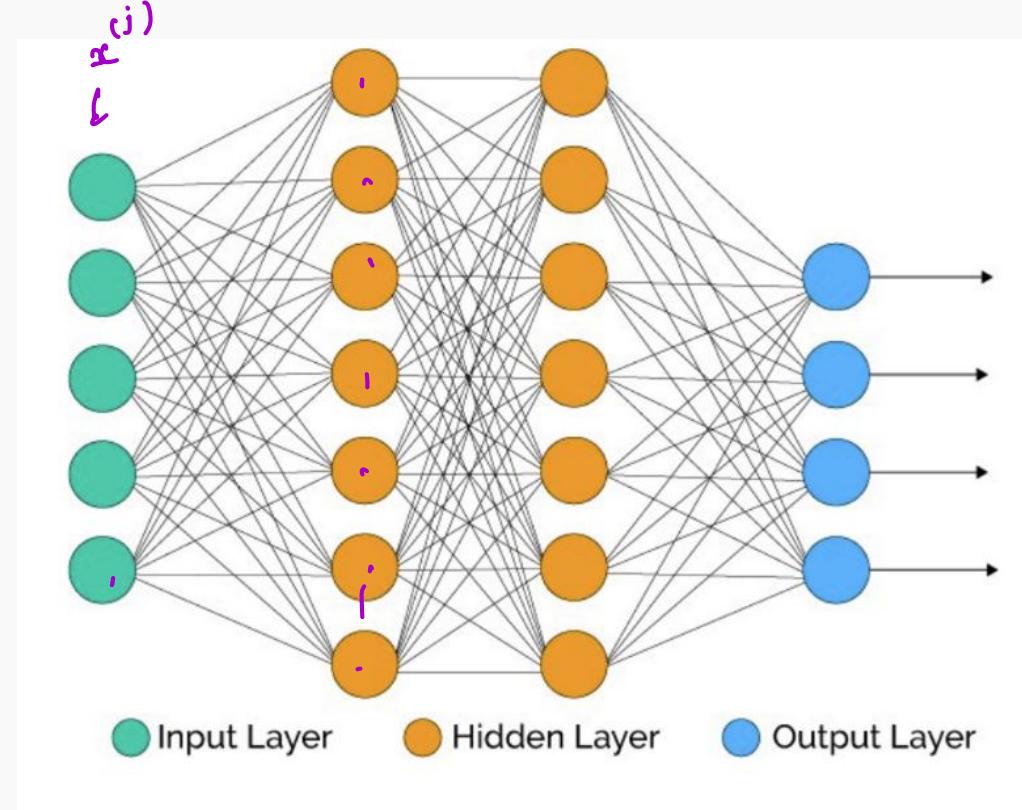


# Universal Approximation Theorem

$$x_1^{(1)} \ x_2^{(1)} \dots x_n^{(1)} \rightarrow y^{(1)}$$

↓  
Let's do  
data!

$$x_1^{(m)} \ x_2^{(m)} \dots x_n^{(m)} \rightarrow y^{(m)}$$



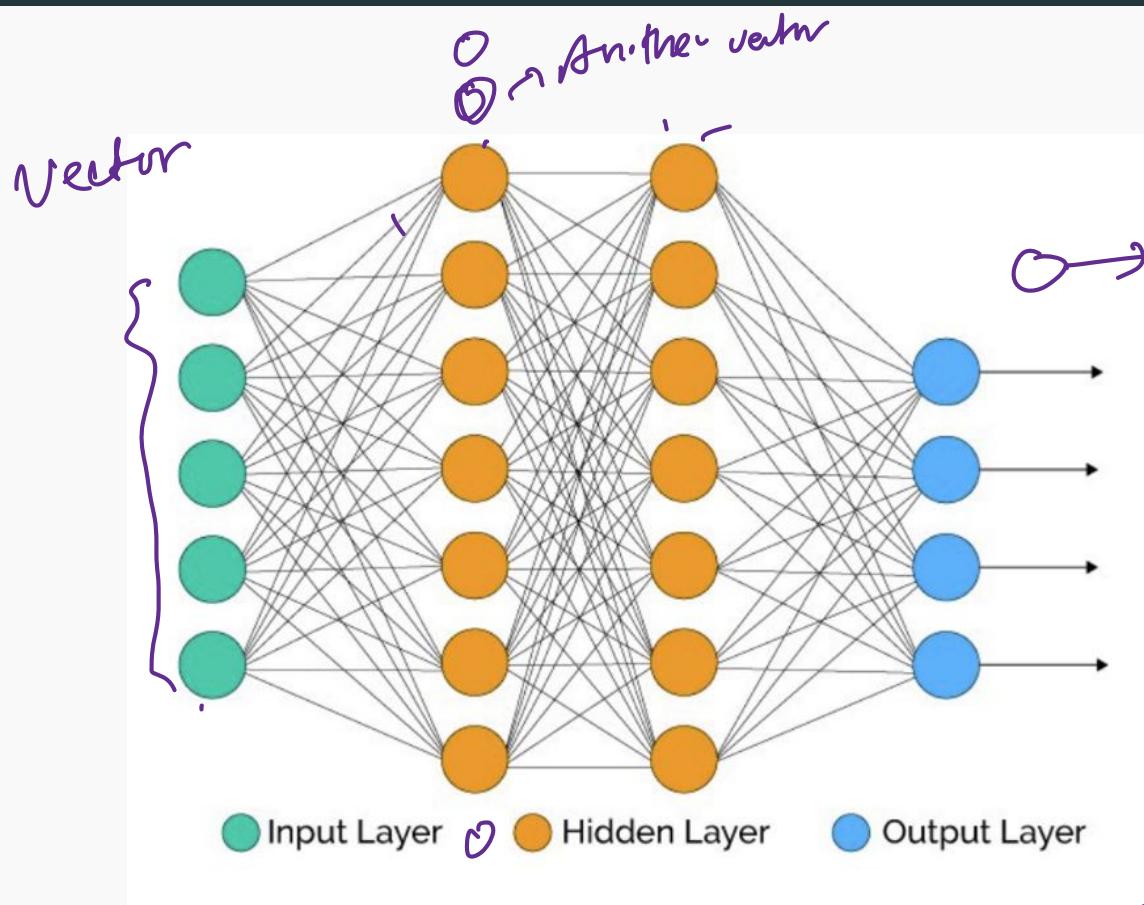
**Theorem**

$$x_1 \cdot x_2 \cdot x_n \in \{0, 1\}^n \rightarrow \mathbb{R}$$

$f(x_1, x_2, \dots, x_n) \rightarrow \emptyset \rightarrow$  training data

*Given sufficient data and neurons, a Neural Network can approximate any function to any desired accuracy.*

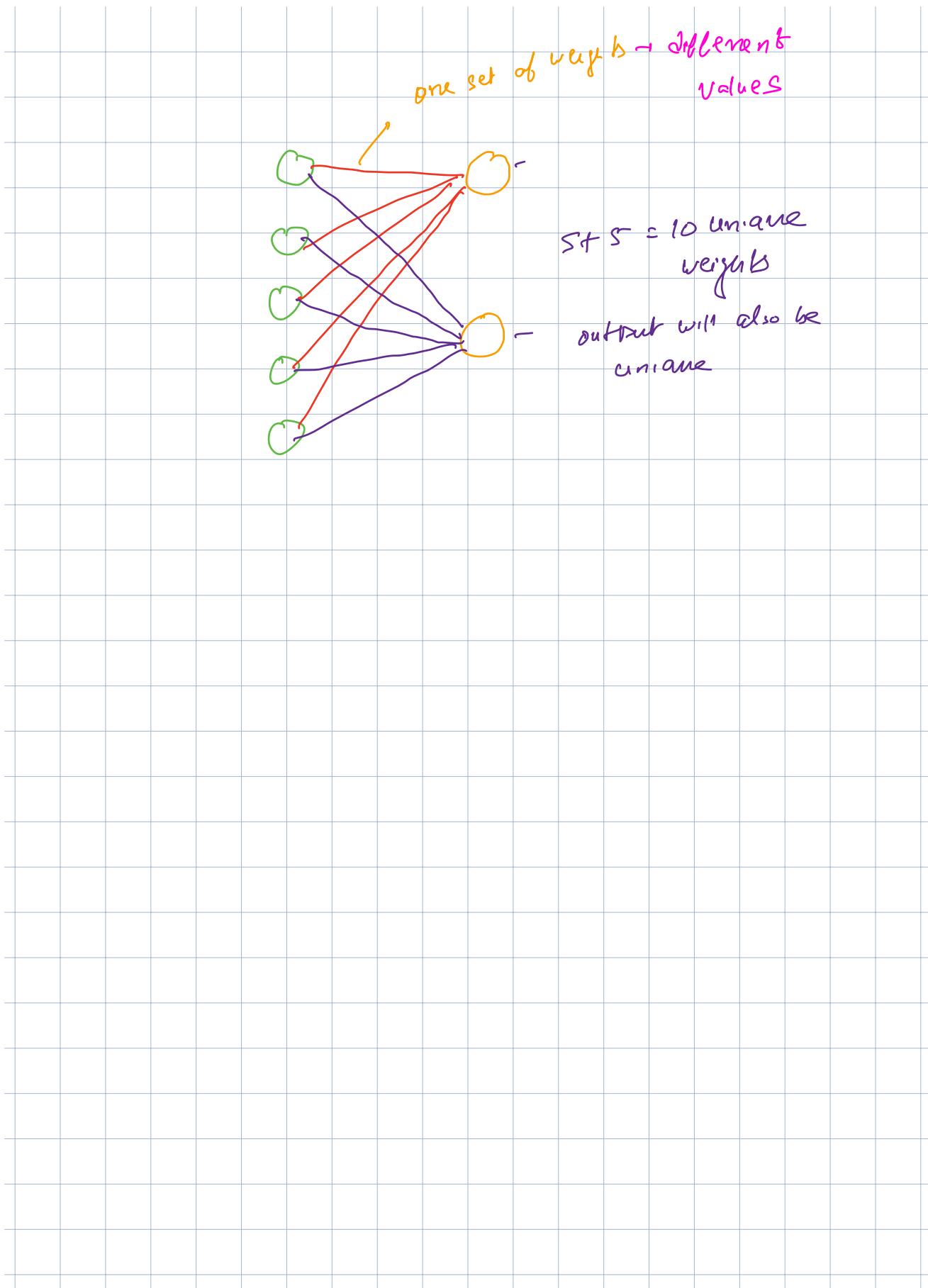
# Universal Approximation Theorem



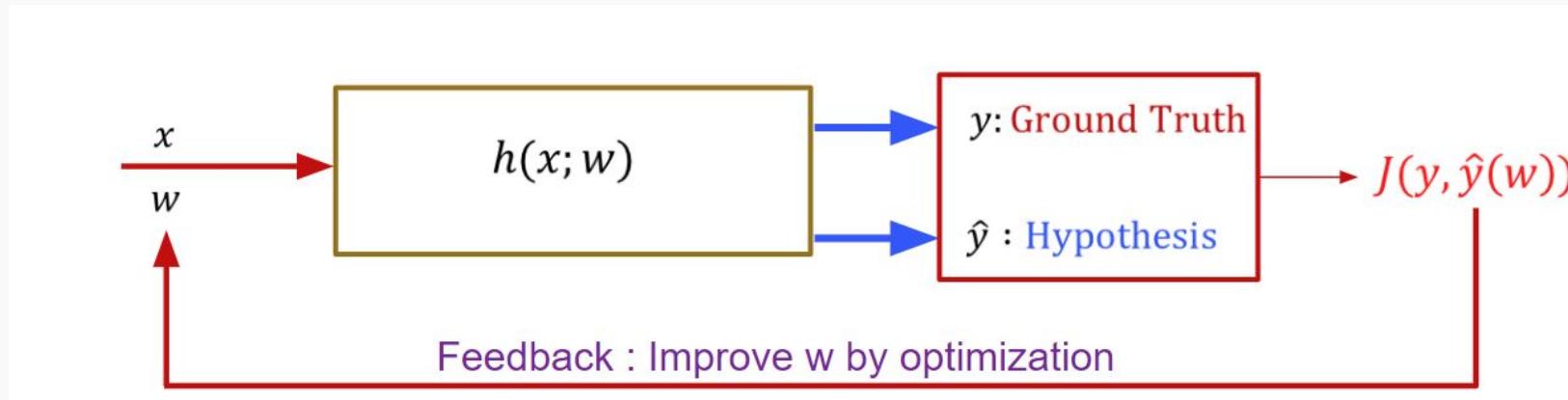
$$f(x_1, x_2, x_3, x_4, x_5) = \begin{cases} 0 & \\ 1 & \end{cases}$$

## Theorem

*Given sufficient data and neurons, a Neural Network can approximate any function to any desired accuracy.*

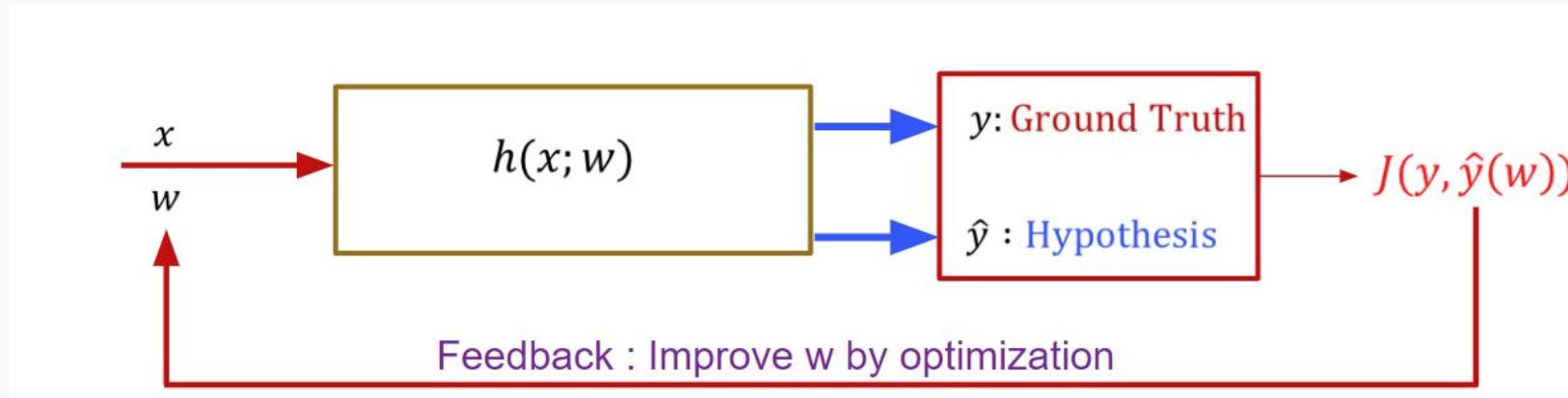


# Recall : Learning the parameters via feedback



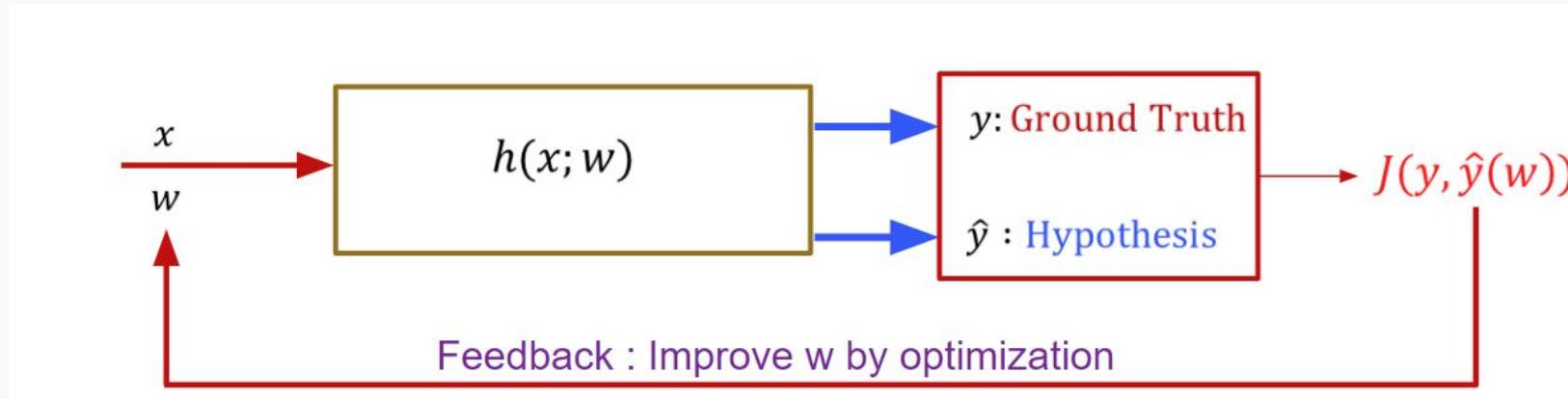
- To learn the parameters, we follow this paradigm
- Guess for the form of the **hypothesis function**  $h(x; w)$
- For an arbitrary guess for  $w$ 
  - We will get some  $\hat{y} = h(x; w)$  which will not match the ground truth  $y$ .
  - Define a cost function  $J(y, \hat{y}(w))$  depending on the difference

# Recall : Learning the parameters via feedback



- To learn the parameters, we follow this paradigm
- Guess for the form of the **hypothesis function**  $h(x; w)$
- For an arbitrary guess for  $w$ 
  - We will get some  $\hat{y} = h(x; w)$  which will not match the ground truth  $y$ .
  - Define a cost function  $J(y, \hat{y}(w))$  depending on the difference
- Find optimal  $w$  by minimizing  $J(w) \rightarrow$  **Feedback Process**
  - Requires Backpropagation for  $\frac{\partial J}{\partial w}$

# Recall : Learning the parameters via feedback



- To learn the parameters, we follow this paradigm
- Guess for the form of the **hypothesis function**  $h(x; w)$
- For an arbitrary guess for  $w$ 
  - We will get some  $\hat{y} = h(x; w)$  which will not match the ground truth  $y$ .
  - Define a cost function  $J(y, \hat{y}(w))$  depending on the difference
- Find optimal  $w$  by minimizing  $J(w) \rightarrow$  **Feedback Process**
  - Requires Backpropagation for  $\frac{\partial J}{\partial w}$

Even complicated Neural Networks work in exactly the same way!

# Calculation in a single neuron

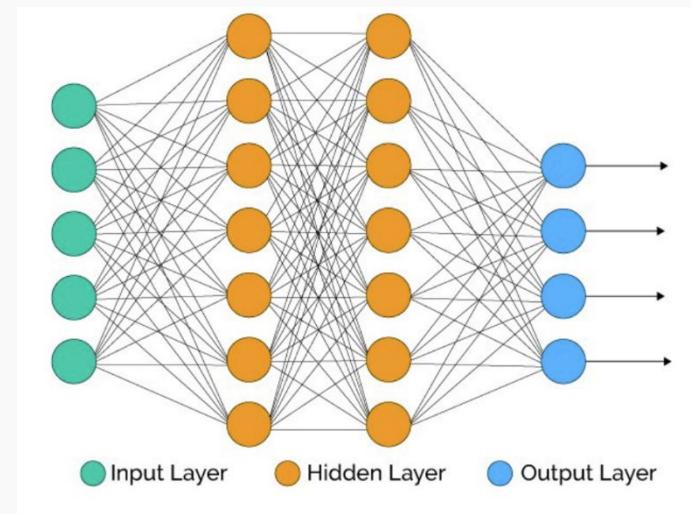
# Recall : Neural Networks

# Recall : Neural Networks

Neural Network model:

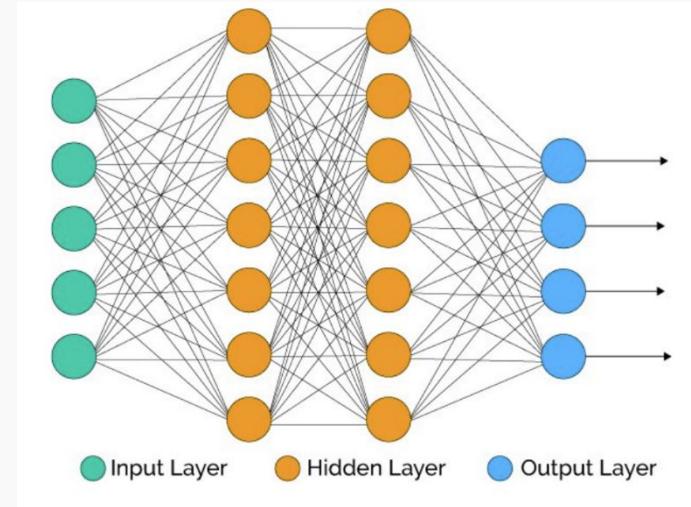
# Recall : Neural Networks

Neural Network model:



# Recall : Neural Networks

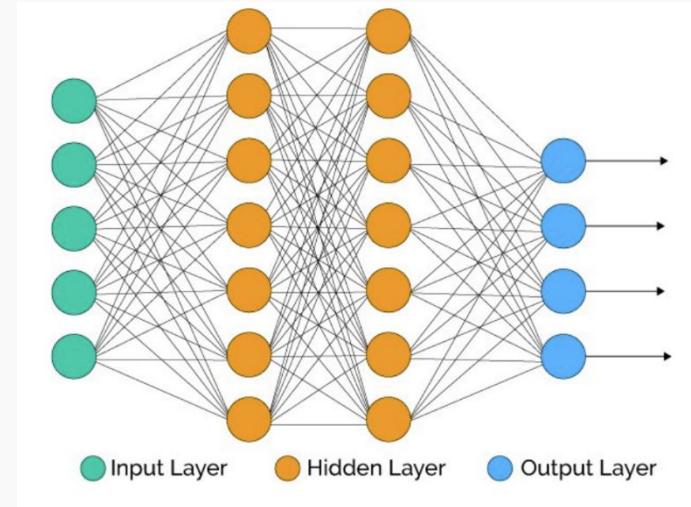
Neural Network model:



$$a = \sum \rightarrow g$$

# Recall : Neural Networks

Neural Network model:

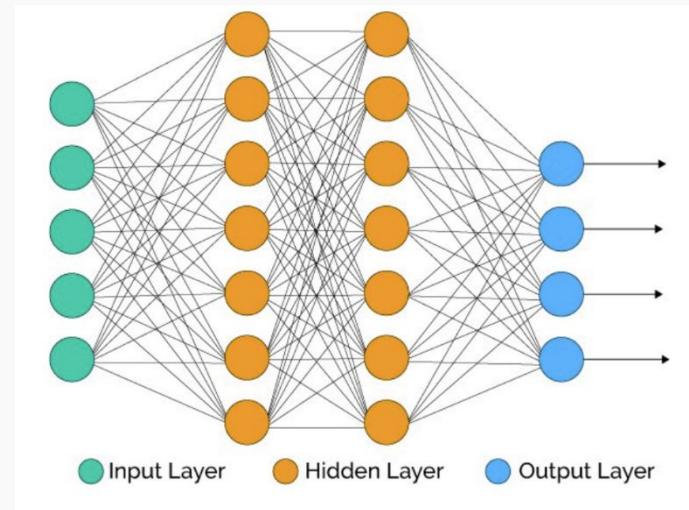


$$a = g\left(\sum w_i x_i\right)$$



# Recall : Neural Networks

Neural Network model:



$$a = g(\sum w_i x_i)$$



A Neuron has a linear and a nonlinear operation.

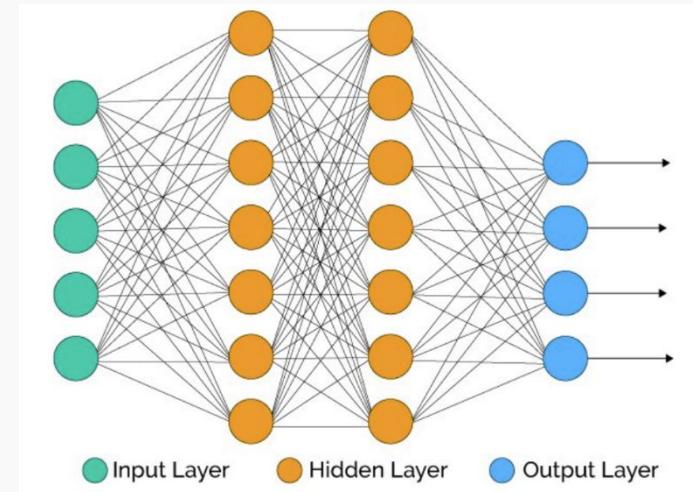
# Recall : Neural Networks

Neural Network model:

Common activation Function:

$$\text{Sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

$$a = g(\sum w_i x_i)$$



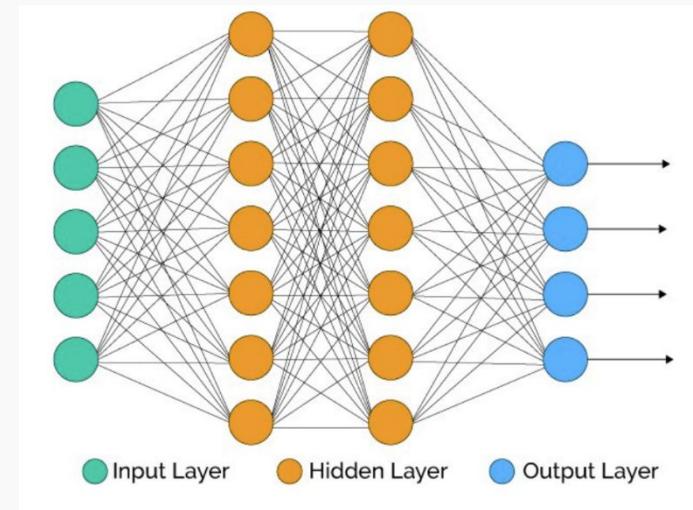
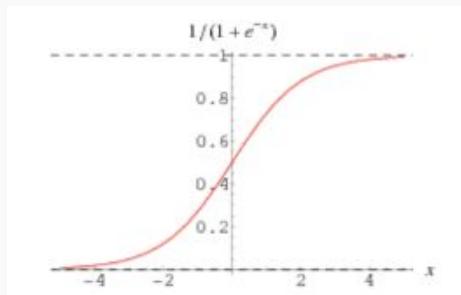
A Neuron has a linear and a nonlinear operation.

# Recall : Neural Networks

Neural Network model:

Common activation Function:

$$\text{Sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$



$$a = g(\sum w_i x_i)$$

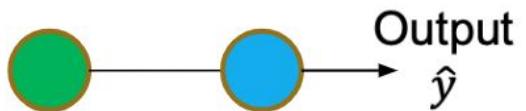


A Neuron has a linear and a nonlinear operation.

# Weights and biases

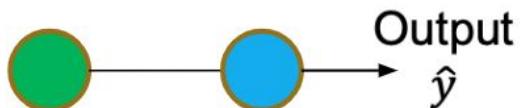
# Weights and biases

What is shown

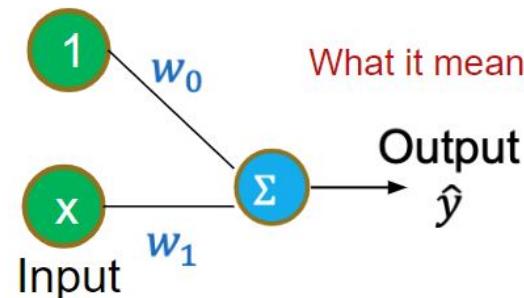


# Weights and biases

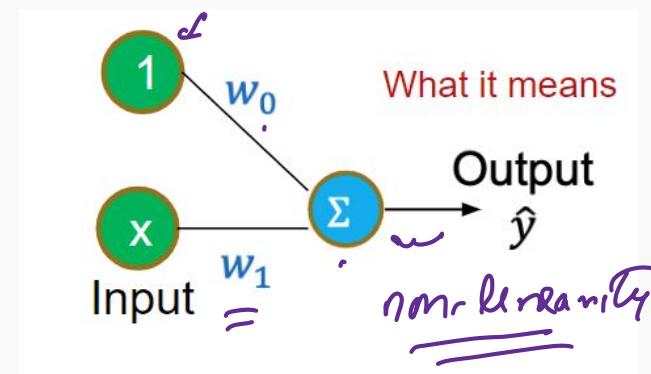
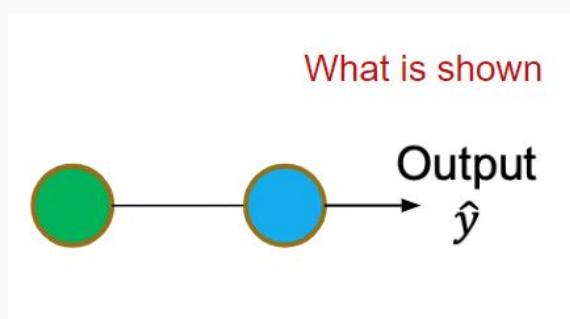
What is shown



What it means

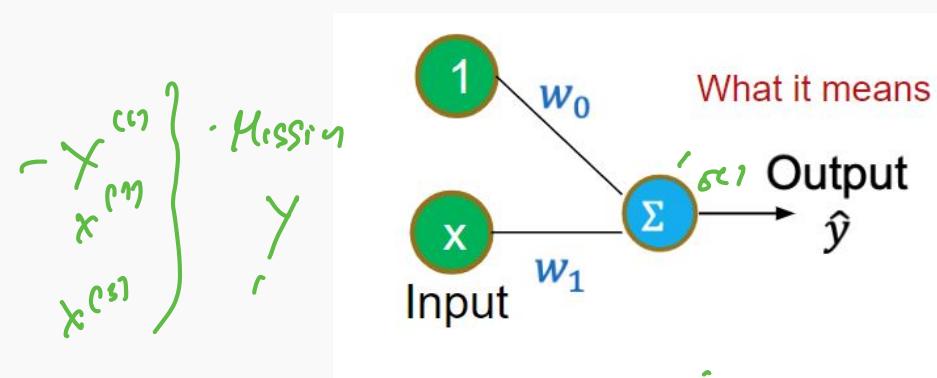
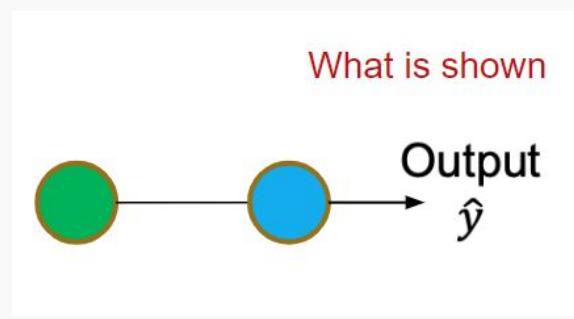


# Weights and biases



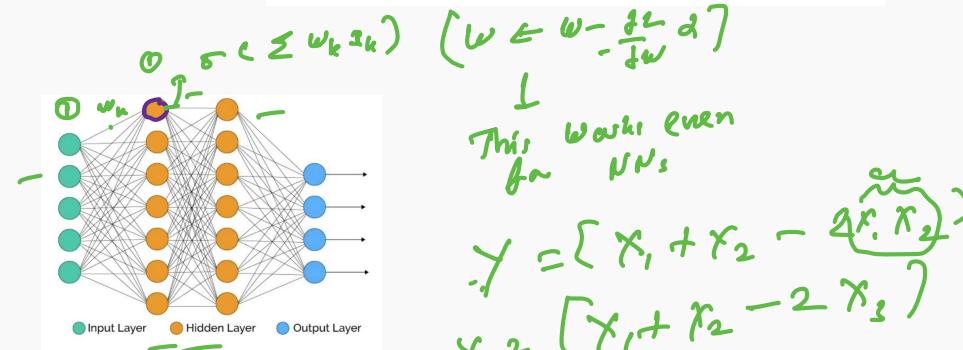
- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

# Weights and biases



DATA  
NORMALIZATION

PCA



- Constant neurons representing a “feature” with 1 are usually not depicted in the pictorial representation.
- The parameters going from the constant neuron is called a bias and the neuron is called a “bias unit”
- The other parameters are called “weights”

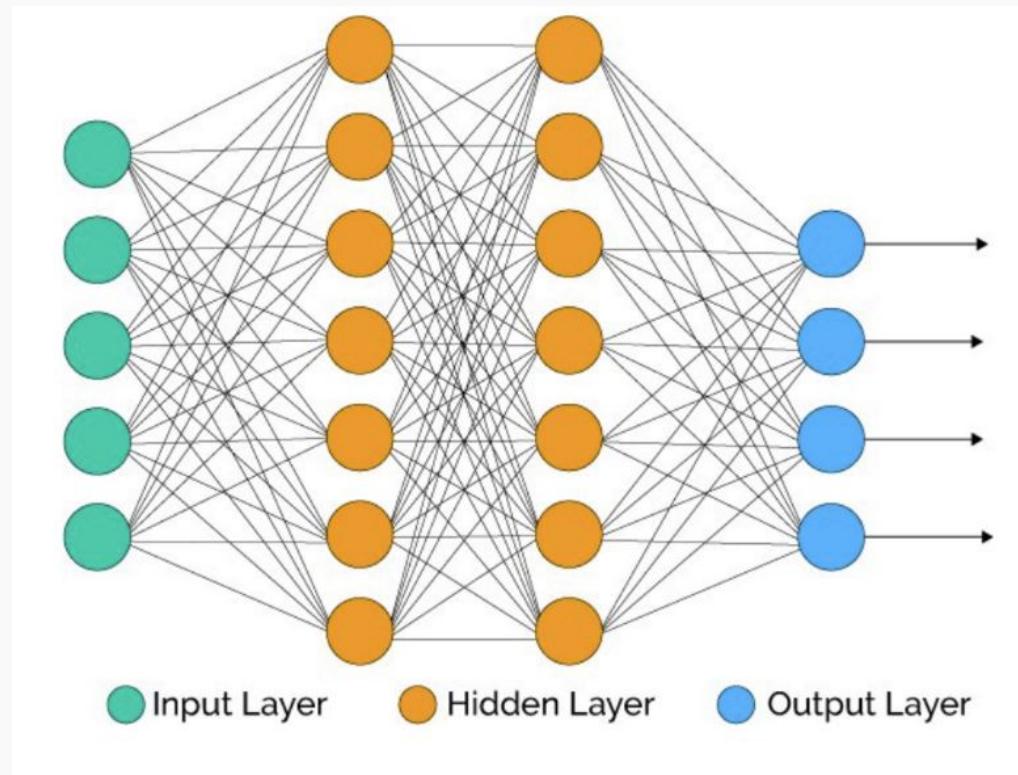
$$y = [x_1 + x_2 - 2(x_1 \cdot x_2)]$$

$$y_2 = [x_1 + x_2 - 2x_1]$$

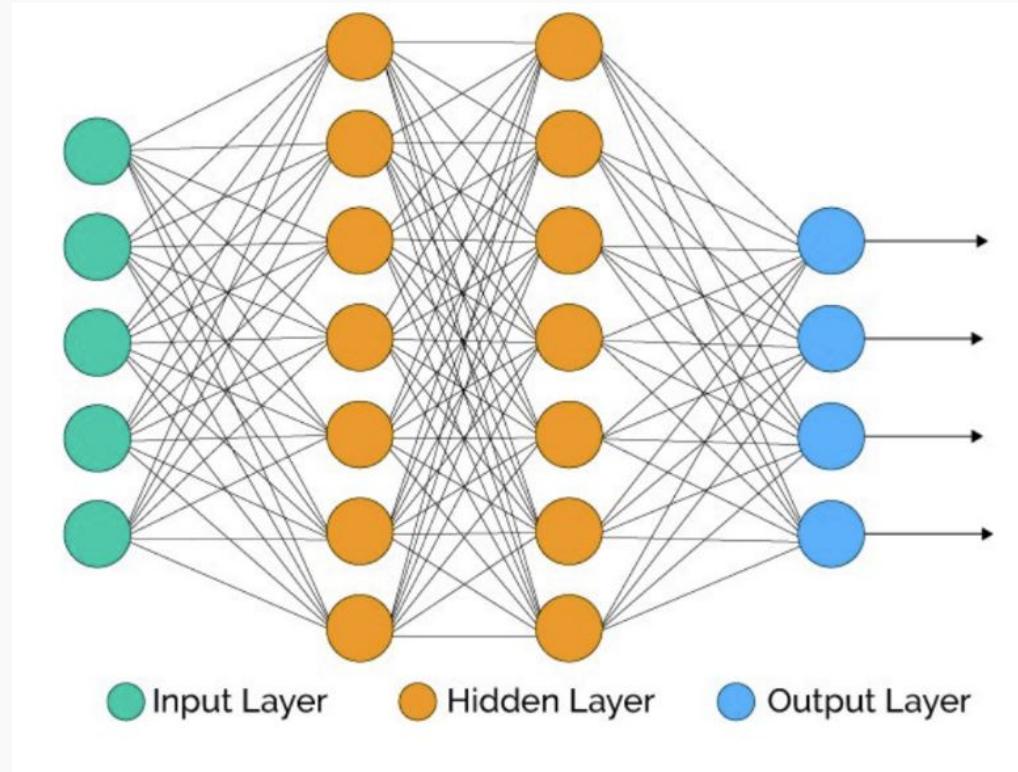
$x_0 - b_1$   
bias

# The output of a single neuron

# The output of a single neuron

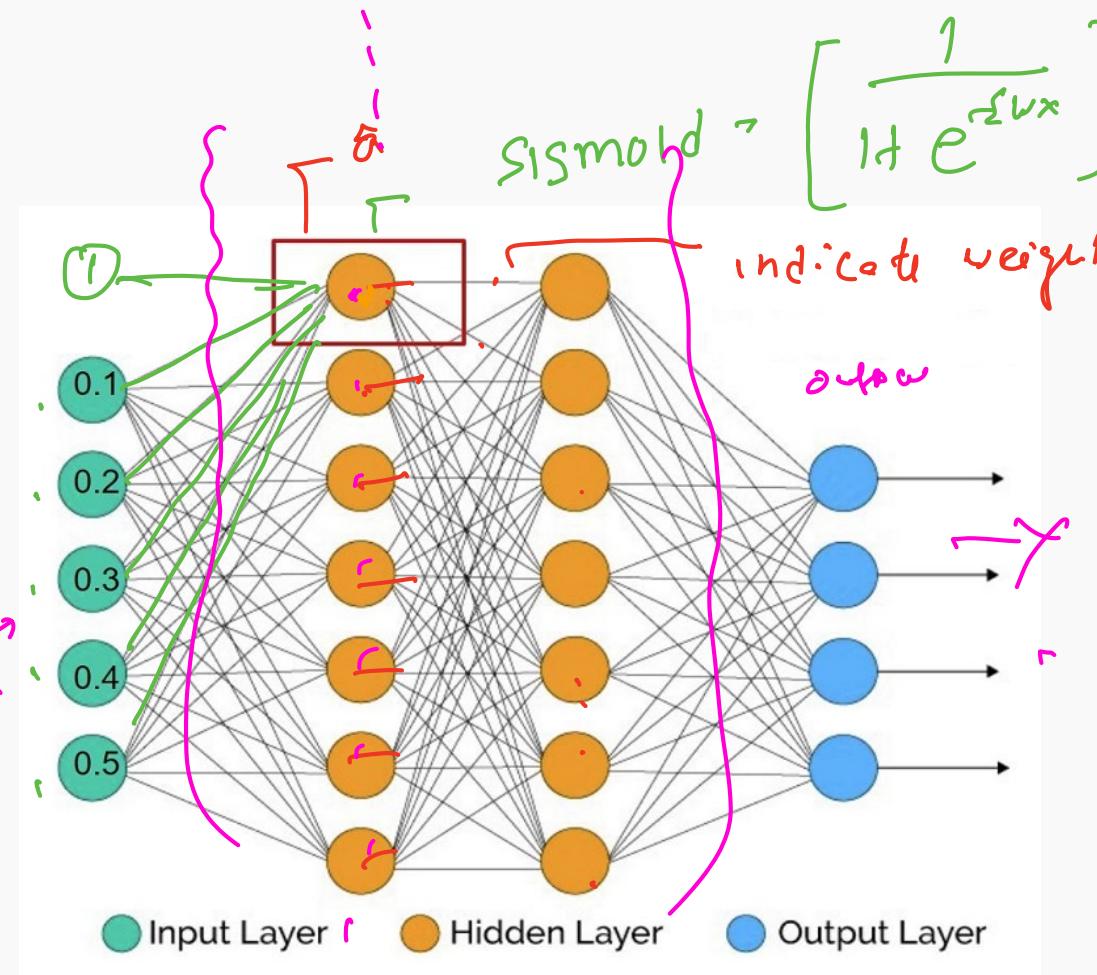


# The output of a single neuron



**Question :** If all weights are 1, then what is the output of the shown neuron?

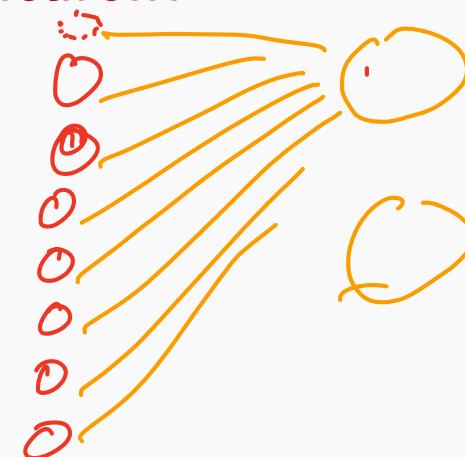
# The output of a single neuron



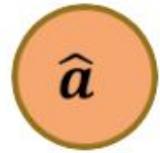
$$y = f(x)$$

that multiply the output

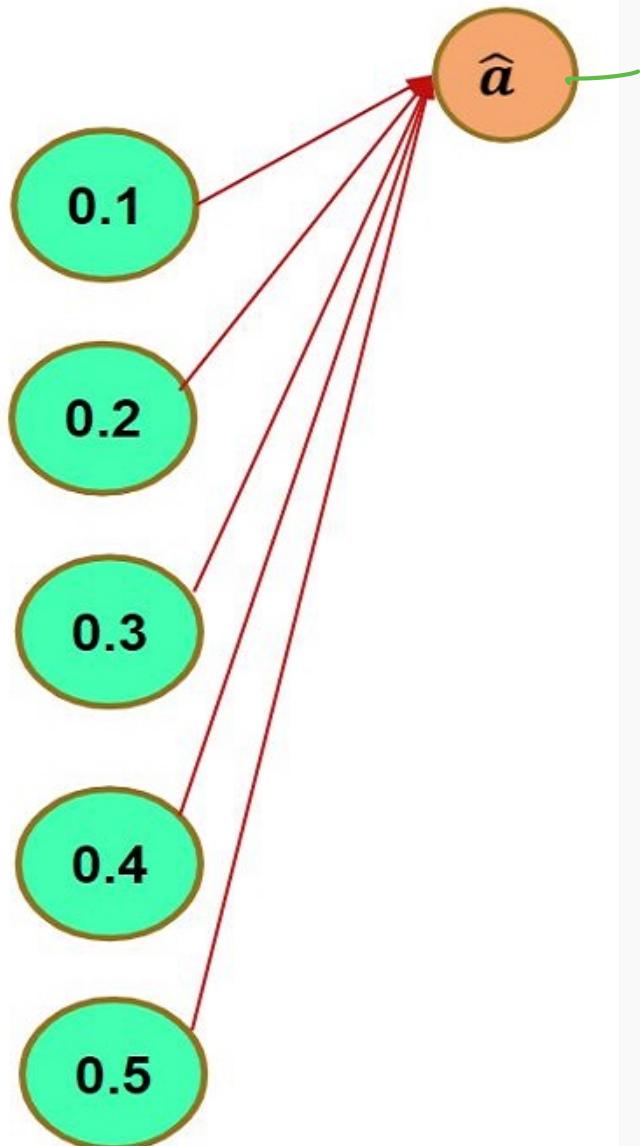
**Question :** If all weights are 1, then what is the output of the shown neuron?



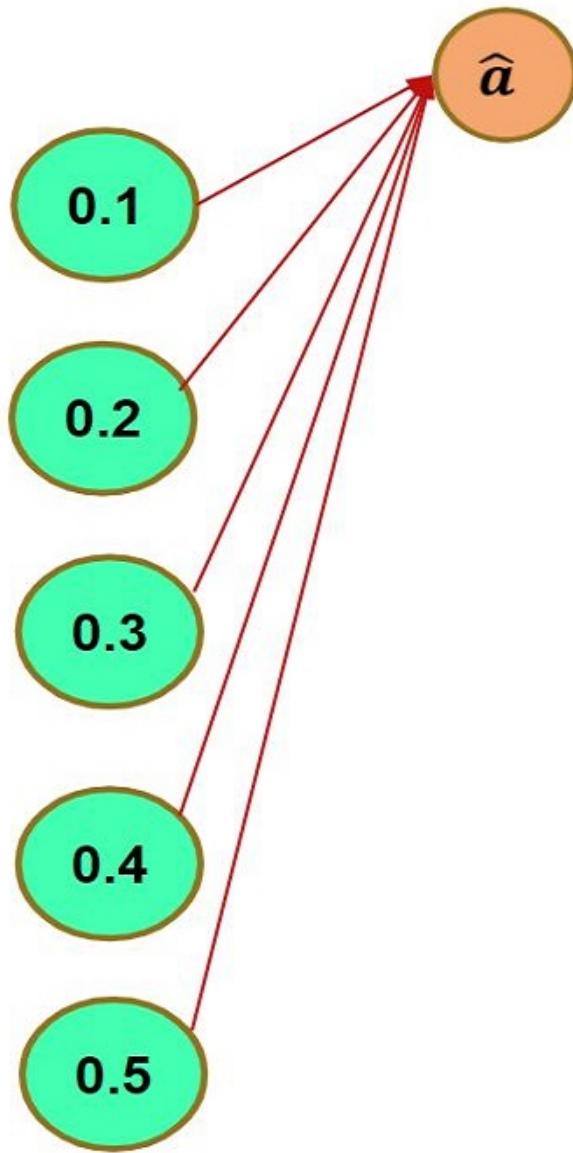
$$\begin{aligned} x^{(i)} &\sim a^{(i)} \\ f &= \dots \end{aligned}$$



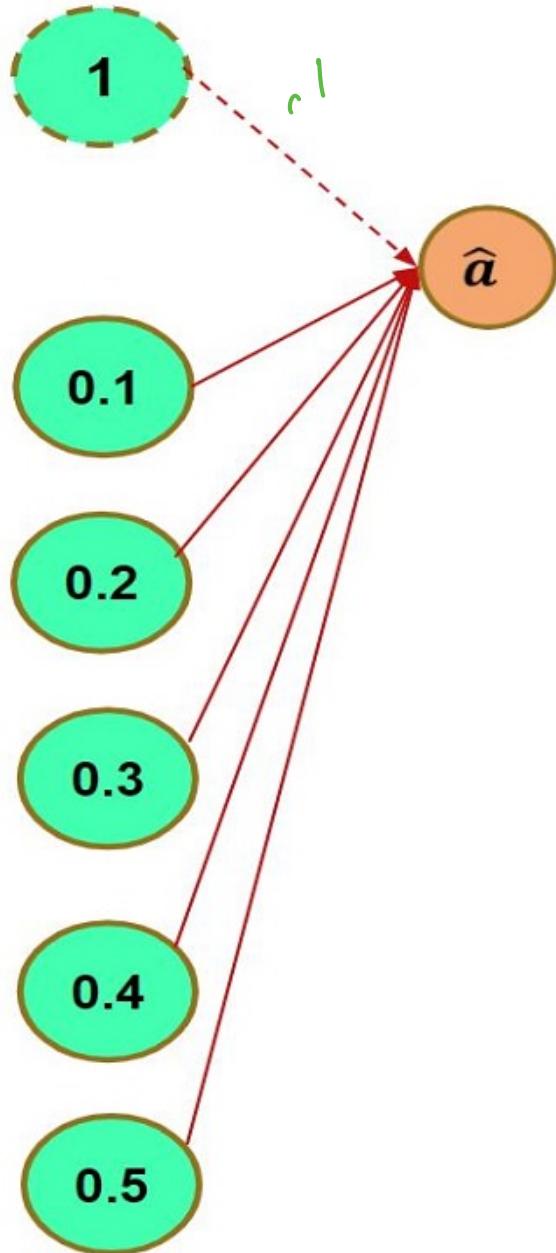
Find the output of the neuron  $\hat{a}$   
All weights are equal to 1



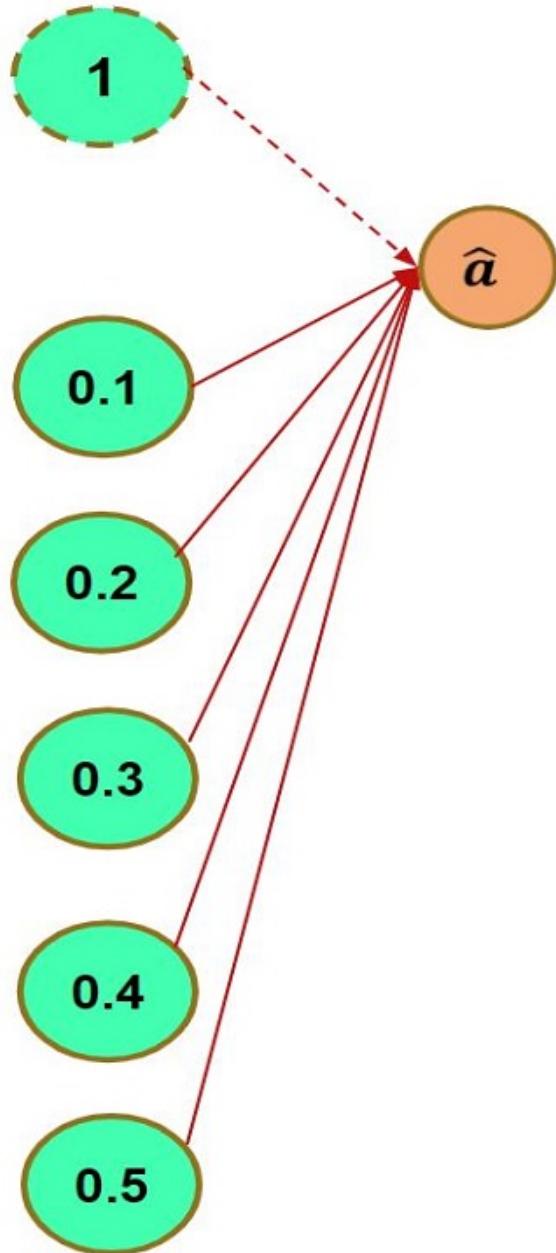
Find the output of the neuron  $\hat{a}$   
All weights are equal to 1



Find the output of the neuron  $\hat{a}$   
All weights are equal to 1



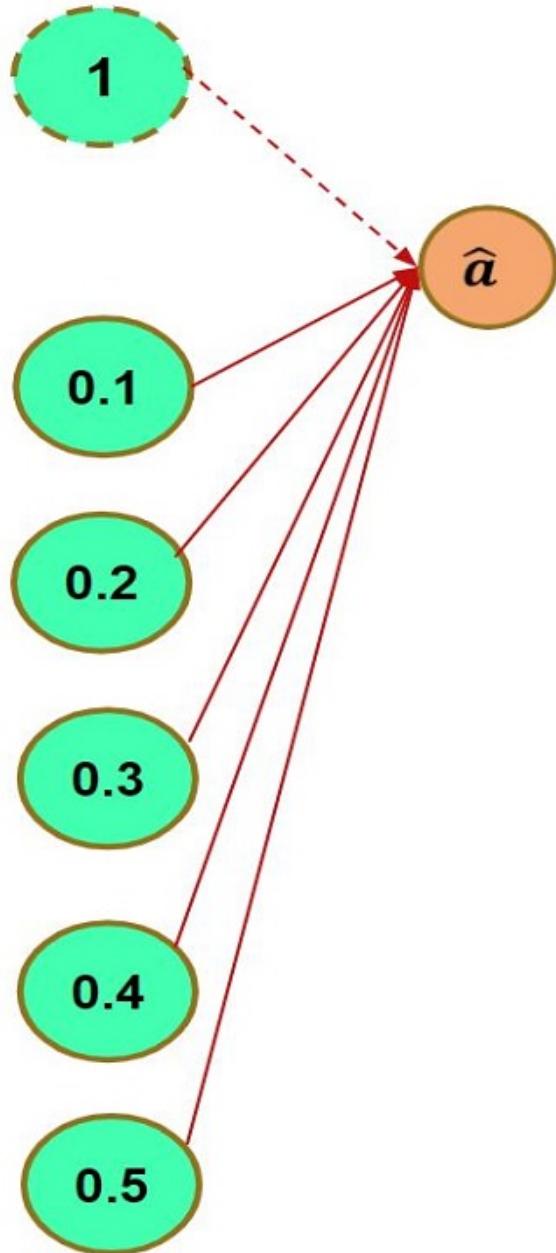
Find the output of the neuron  $\hat{a}$   
All weights are equal to 1



Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts



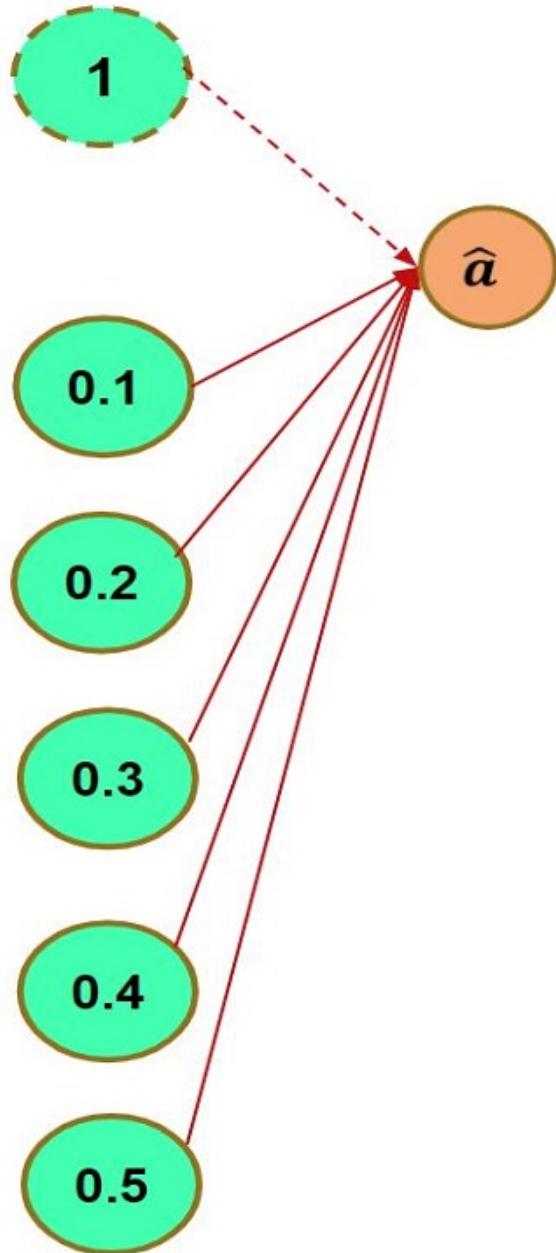
Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts

**Linear:** 
$$z = \sum w_i x_i$$





Find the output of the neuron  $\hat{a}$

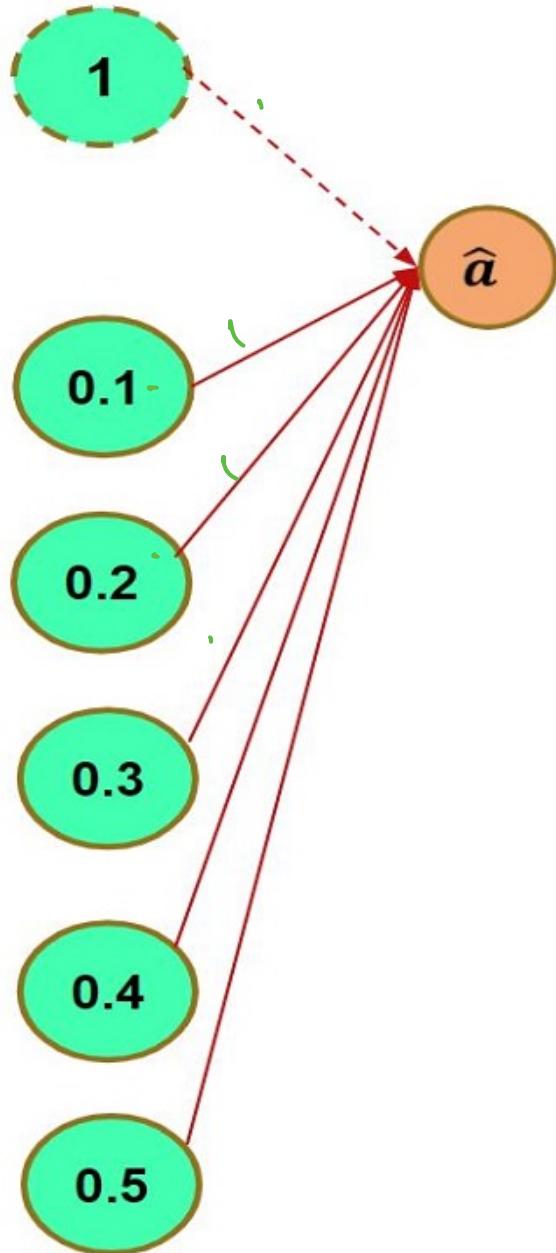
All weights are equal to 1

**Ans:** The output has two parts

**Linear:** 
$$z = \sum w_i x_i$$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$





Find the output of the neuron  $\hat{a}$

All weights are equal to 1

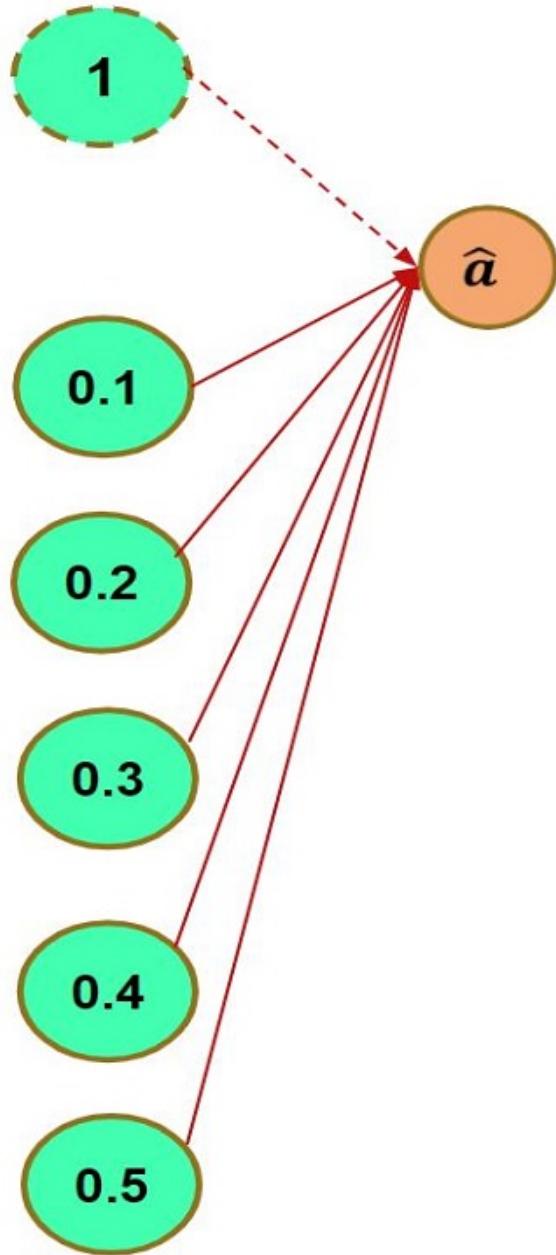
**Ans:** The output has two parts

**Linear:**  $z = \sum w_i x_i$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

ANS



Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts

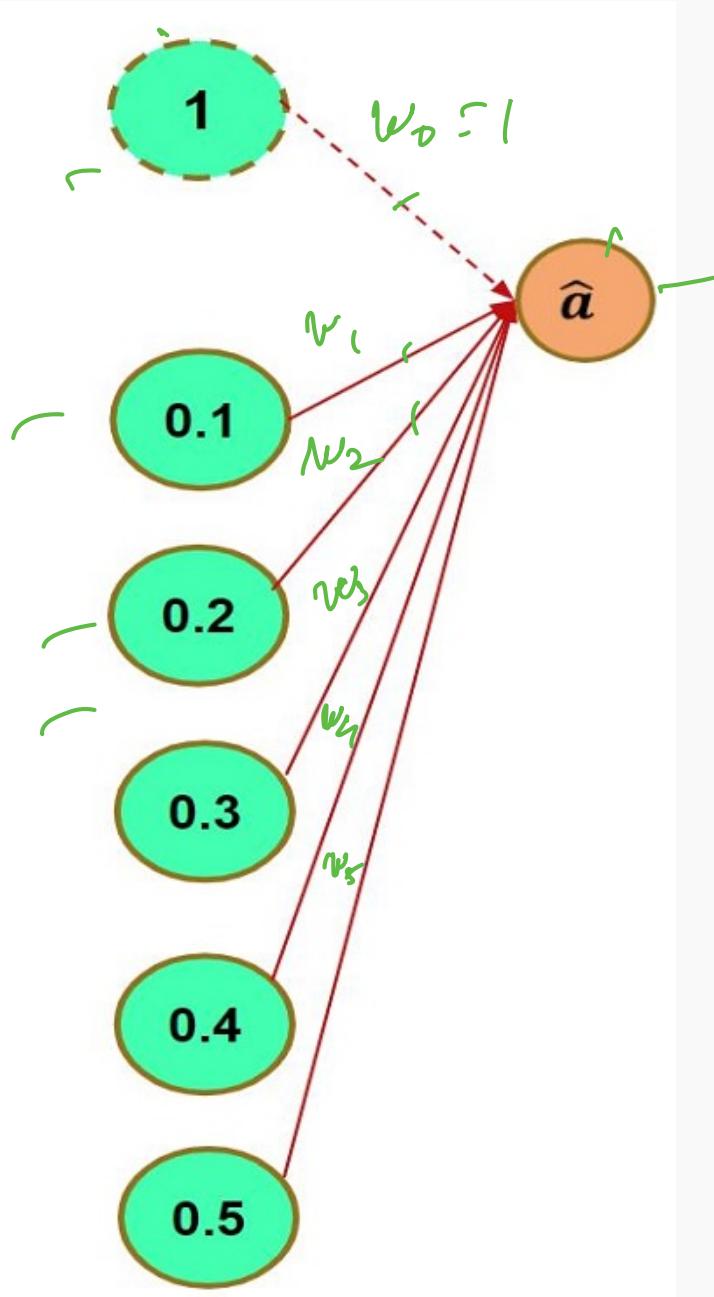
**Linear:**  $z = \sum w_i x_i$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

**Non-linear Activation:**  $\hat{a} = g(z)$





$$\sum w_i x_i + b' \quad x_0$$

$$x_0 \rightarrow 1$$

Forward Pass

Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts

**Linear:**  $z = \sum w_i x_i$

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

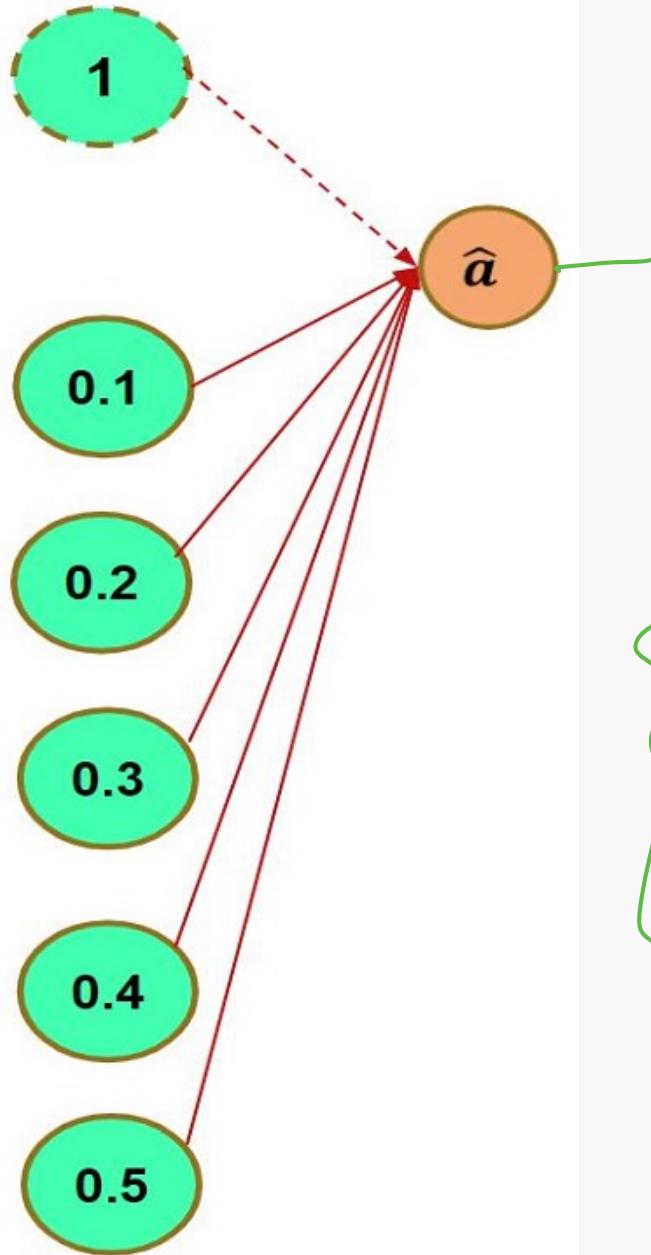
$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

**Non-linear Activation:**  $\hat{a} = g(z)$

$$\hat{a} = \frac{1}{1 + \exp(-2.5)}$$

$\equiv$

$$\sum_{i=1}^n w_i x_i + w_0 = \sum_{i=0}^n w_i x_i \quad x_0 = 1$$



Find the output of the neuron  $\hat{a}$

All weights are equal to 1

**Ans:** The output has two parts

**Linear:**  $z = \sum w_i x_i$  }  $\rightarrow$  Output

$$z = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + w_5 x_5$$

$$z = 1 + 0.1 + 0.2 + 0.3 + 0.4 + 0.5 = 2.5$$

**Non-linear Activation:**  $\hat{a} = g(z)$

$$\hat{a} = \frac{1}{1+\exp(-2.5)} \quad } \sim \text{output}$$

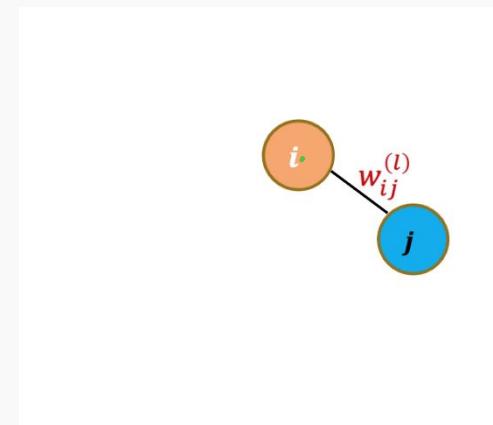
$$\hat{a} = 0.9241$$

=

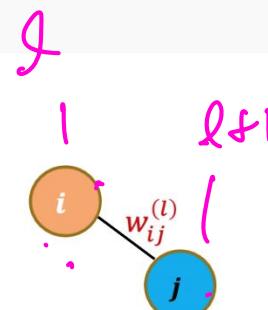
hidden layer

# Weights notation

# Weights notation

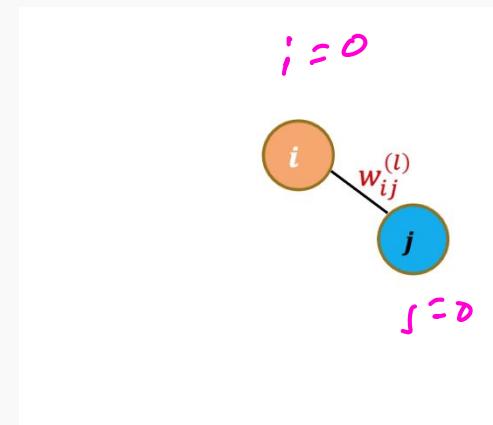


# Weights notation



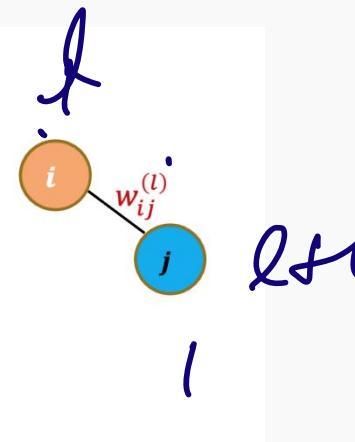
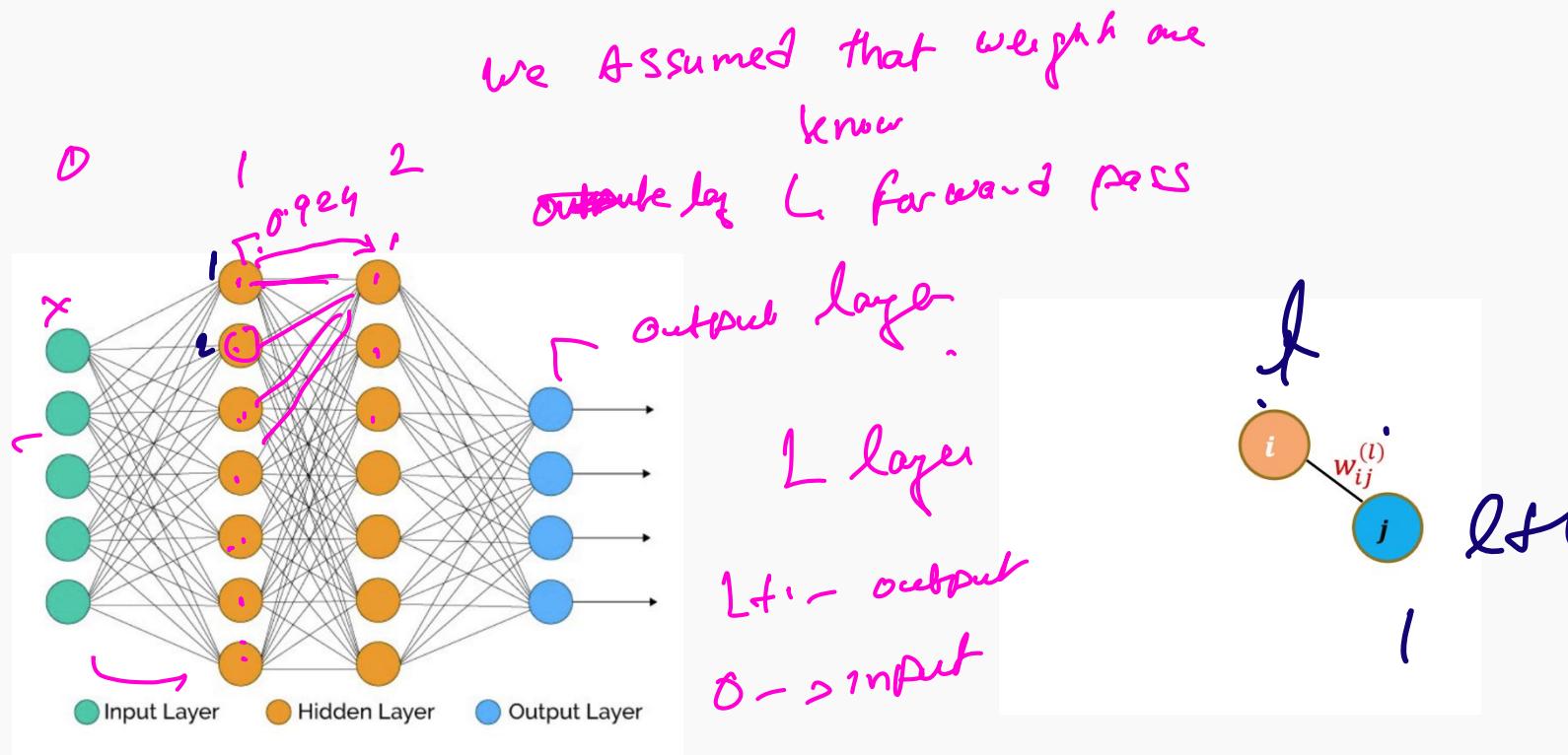
- $w_{ij}^{(l)}$  is the weight connecting Neuron  $i$  of Layer  $l$  to Neuron  $j$  of Layer  $l + 1$ .

# Weights notation



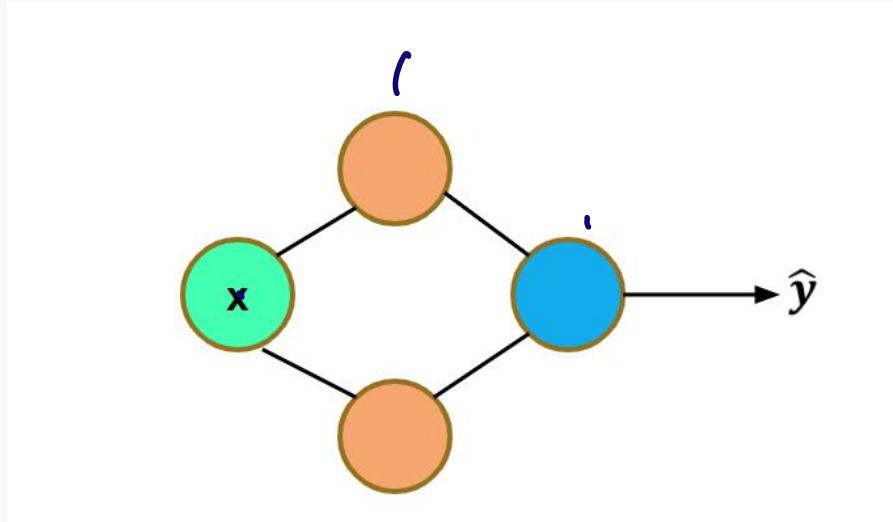
- $w_{ij}^{(l)}$  is the weight connecting Neuron  $i$  of Layer  $l$  / Neuron  $j$  of Layer  $l + 1$ .
- bias units are the  $0^{th}$  neurons.

# Weights notation



- $w_{ij}^{(l)}$  is the weight connecting Neuron  $i$  of Layer  $l$  Neuron  $j$  of Layer  $l+1$ .
- bias units are the  $0^{th}$  neurons.

A simple forward pass calculation

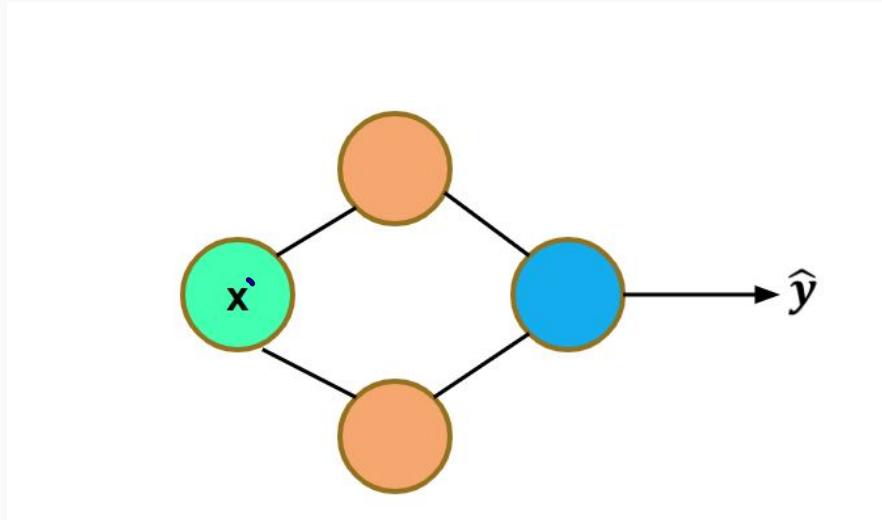


**Question:** Find the output of the given network for the following.

univariate  
input

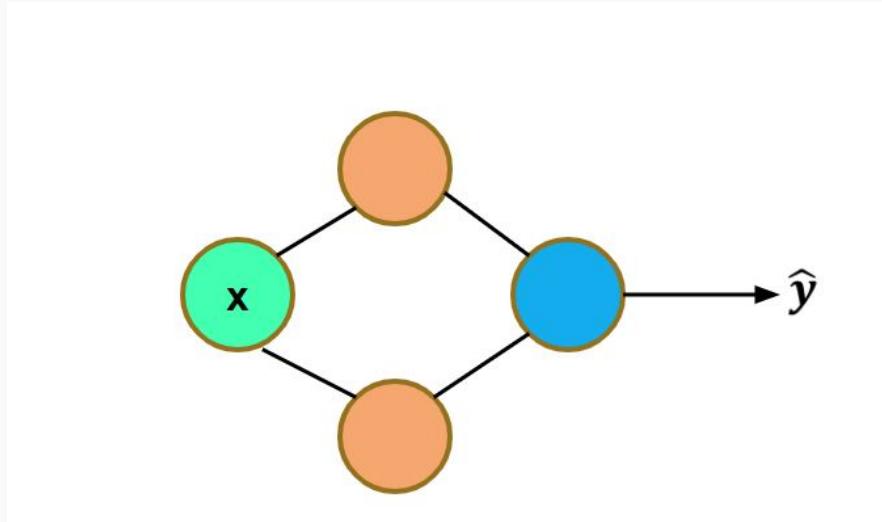
← hidden  
layer

← output ↑



**Question:** Find the output of the given network for the following.

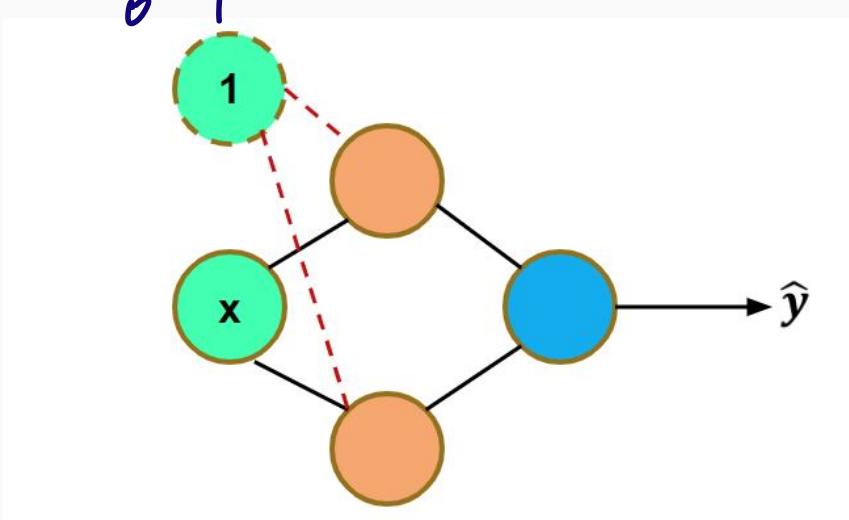
$$x = 0.5$$



**Question:** Find the output of the given network for the following.

$$x = 0.5$$

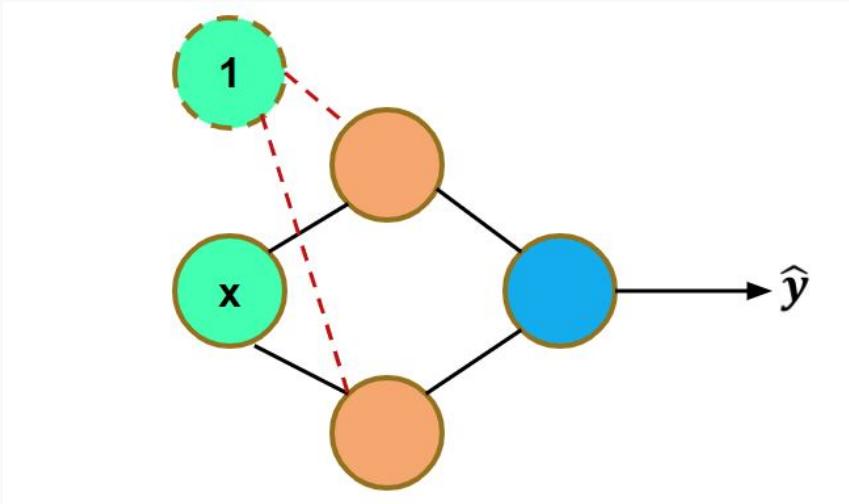
$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$



**Question:** Find the output of the given network for the following.

$$x = 0.5$$

$$\begin{matrix} & \mid \\ w_{01}^{(1)} & = 1.0 & w_{02}^{(1)} & = 0.8 \end{matrix}$$

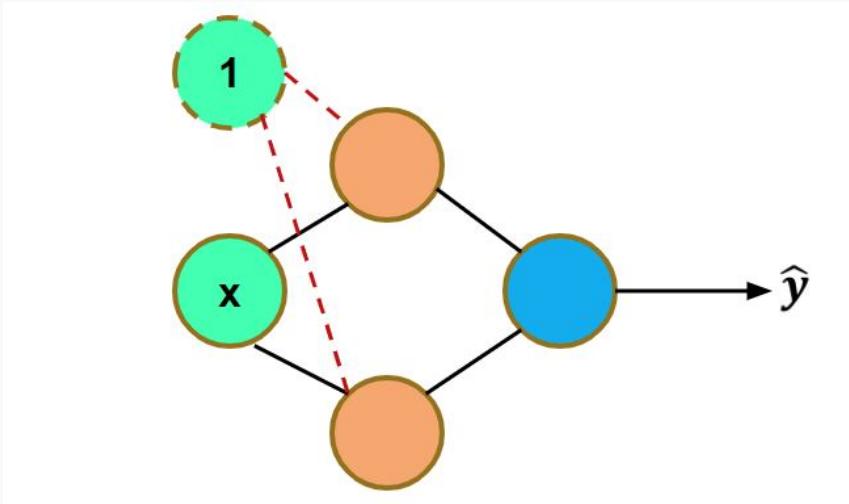


**Question:** Find the output of the given network for the following.

$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$

$$w_{11}^{(1)} = 0.7 \quad w_{12}^{(1)} = 0.9$$



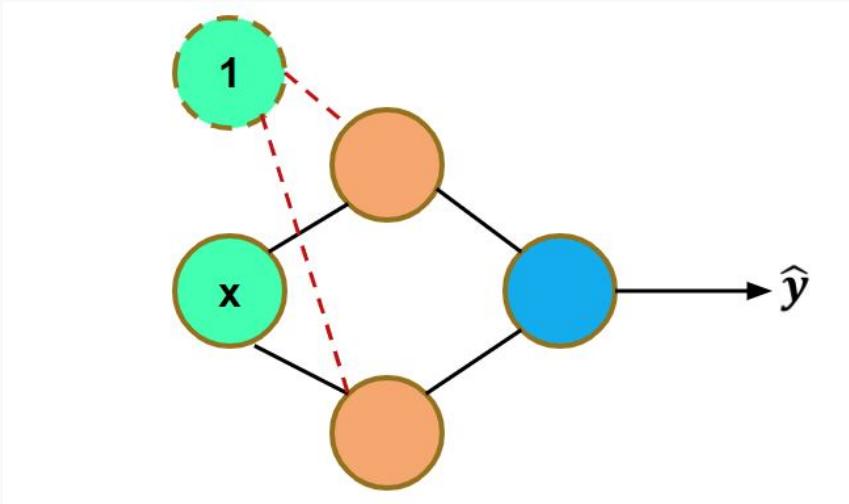
**Question:** Find the output of the given network for the following.

$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$

$$w_{11}^{(1)} = 0.7 \quad w_{12}^{(1)} = 0.9$$

$$w_{01}^{(2)} = 1$$



**Question:** Find the output of the given network for the following.

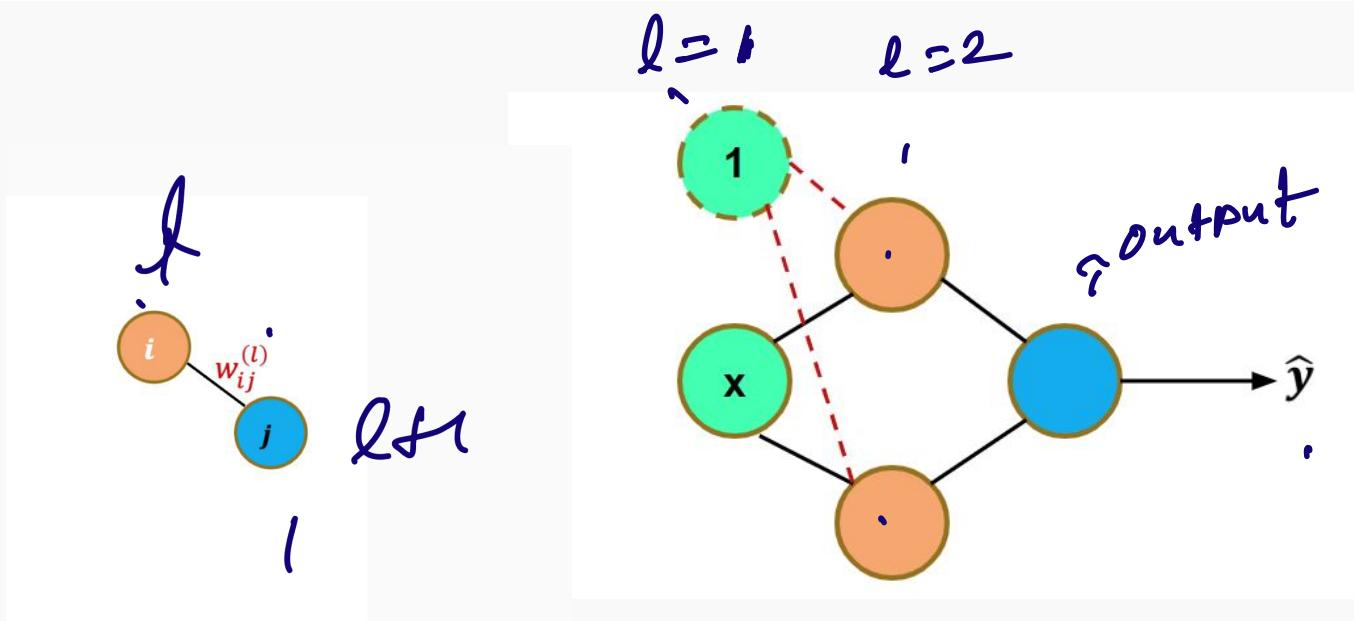
$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$

$$w_{11}^{(1)} = 0.7 \quad w_{12}^{(1)} = 0.9$$

$$w_{01}^{(2)} = 1$$

$$w_{11}^{(2)} = 1 \quad w_{21}^{(2)} = 0.9$$



**Question:** Find the output of the given network for the following.

$$x = 0.5$$

$$w_{01}^{(1)} = 1.0 \quad w_{02}^{(1)} = 0.8$$

$$w_{11}^{(1)} = 0.7 \quad w_{12}^{(1)} = 0.9$$

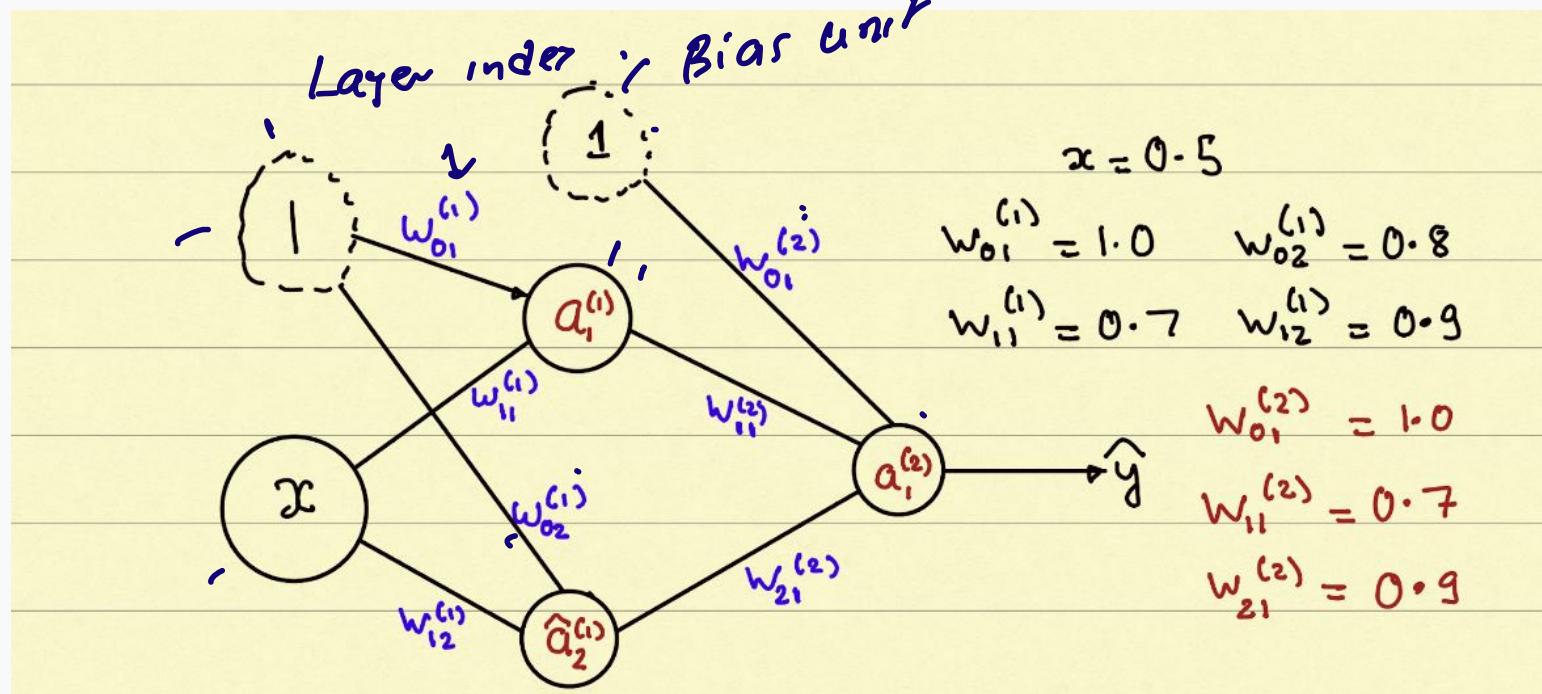
$$w_{01}^{(2)} = 1$$

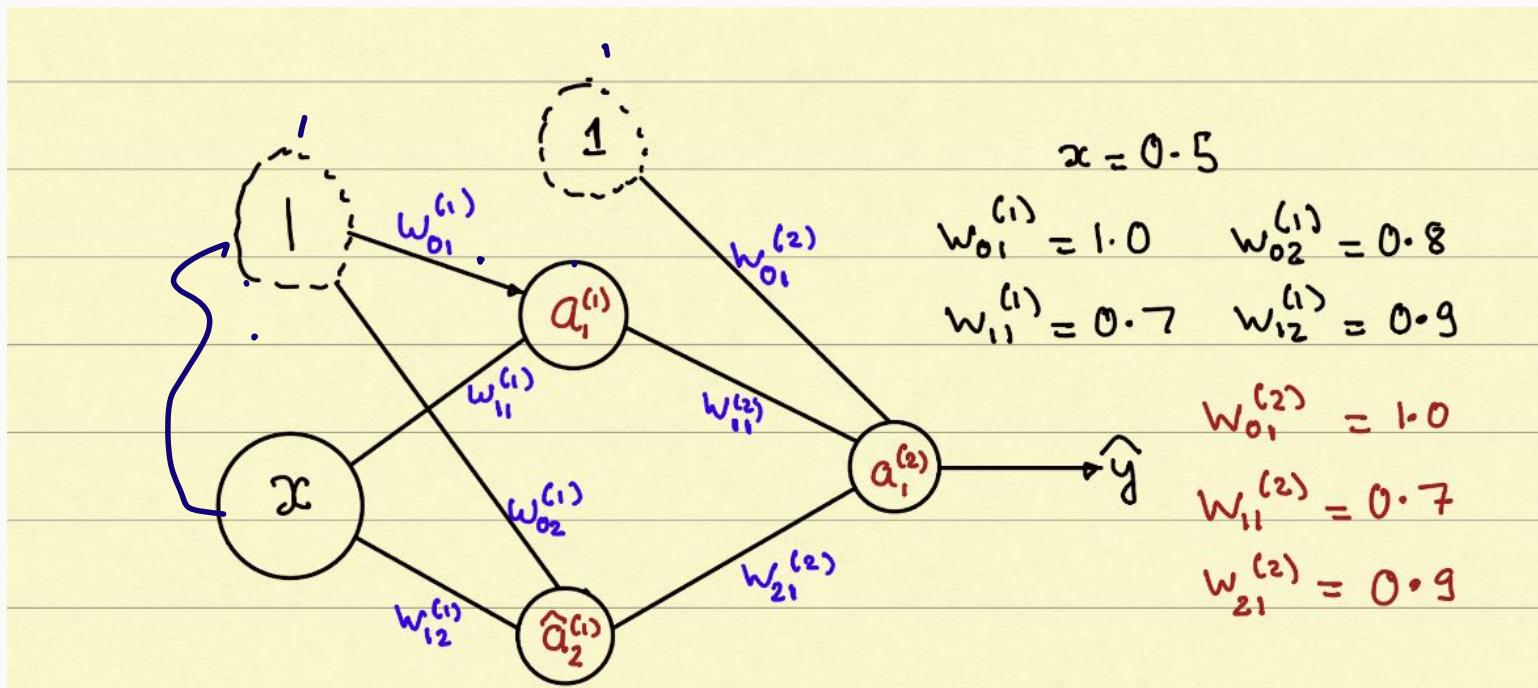
$$w_{11}^{(2)} = 1 \quad w_{21}^{(2)} = 0.9$$

Assume that the activation function in the hidden layer and the output layer is the sigmoid



SOLUTION

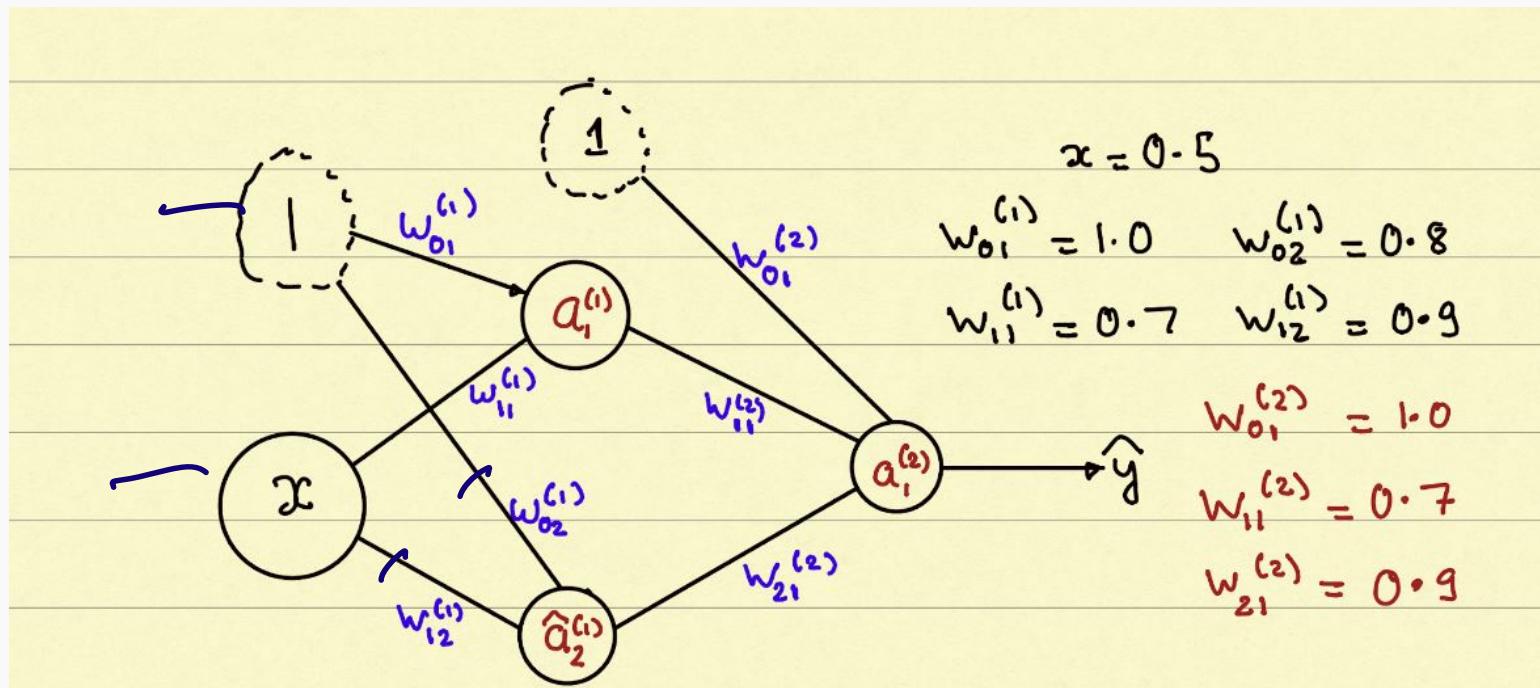




2

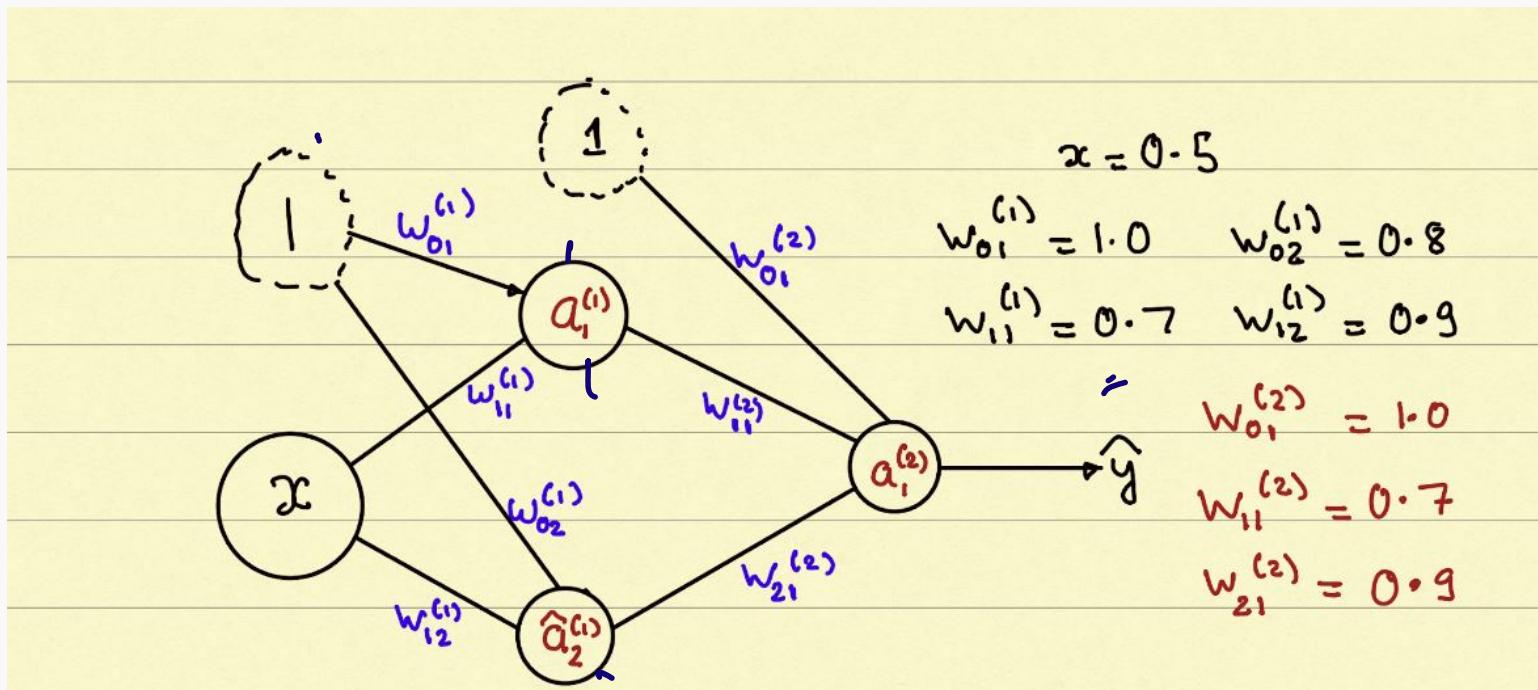
$$\hat{a}_1^{(1)} = g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x) \quad \curvearrowright$$

$g \in \sigma$



$$\hat{a}_1^{(1)} = g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x)$$

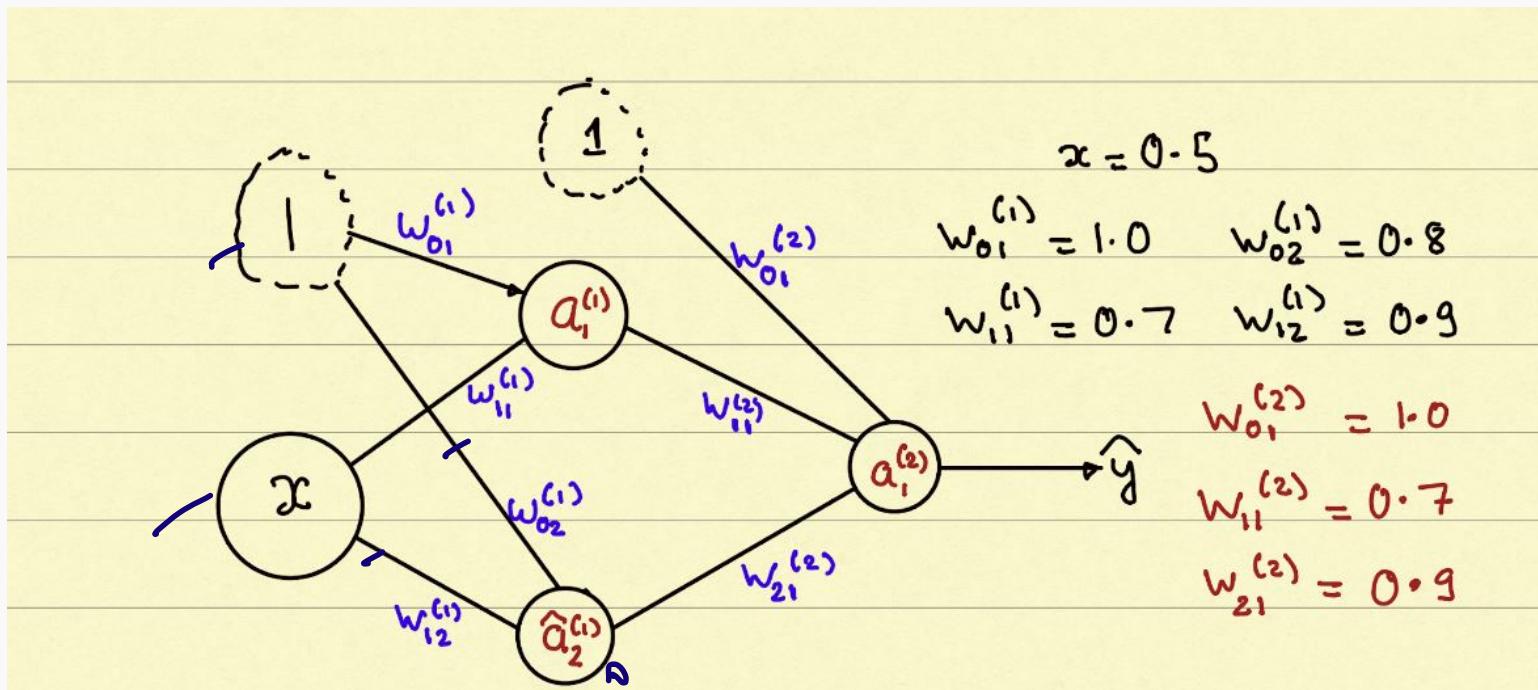
$$= g(1 \times 1 + 0.7 \times 0.5) = g(1.35) \text{ —}$$



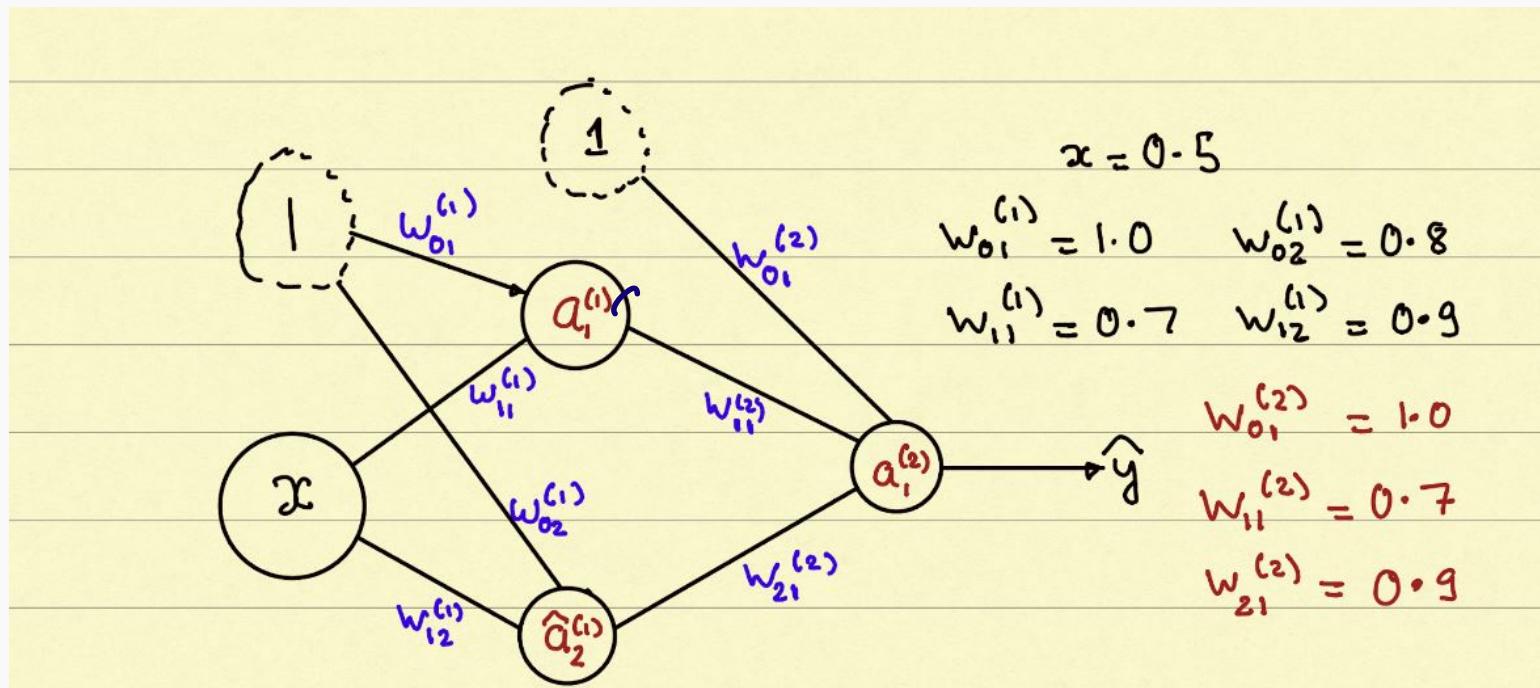
$$\hat{a}_1^{(1)} = g(z_1^{(1)}) = g(w_{01}^{(1)}x_0 + w_{11}^{(1)}x)$$

$$= g(1 \times 1 + 0.7 \times 0.5) = g(1.35)$$

$$\hat{a}_1^{(1)} = 0.7941 \leftarrow$$

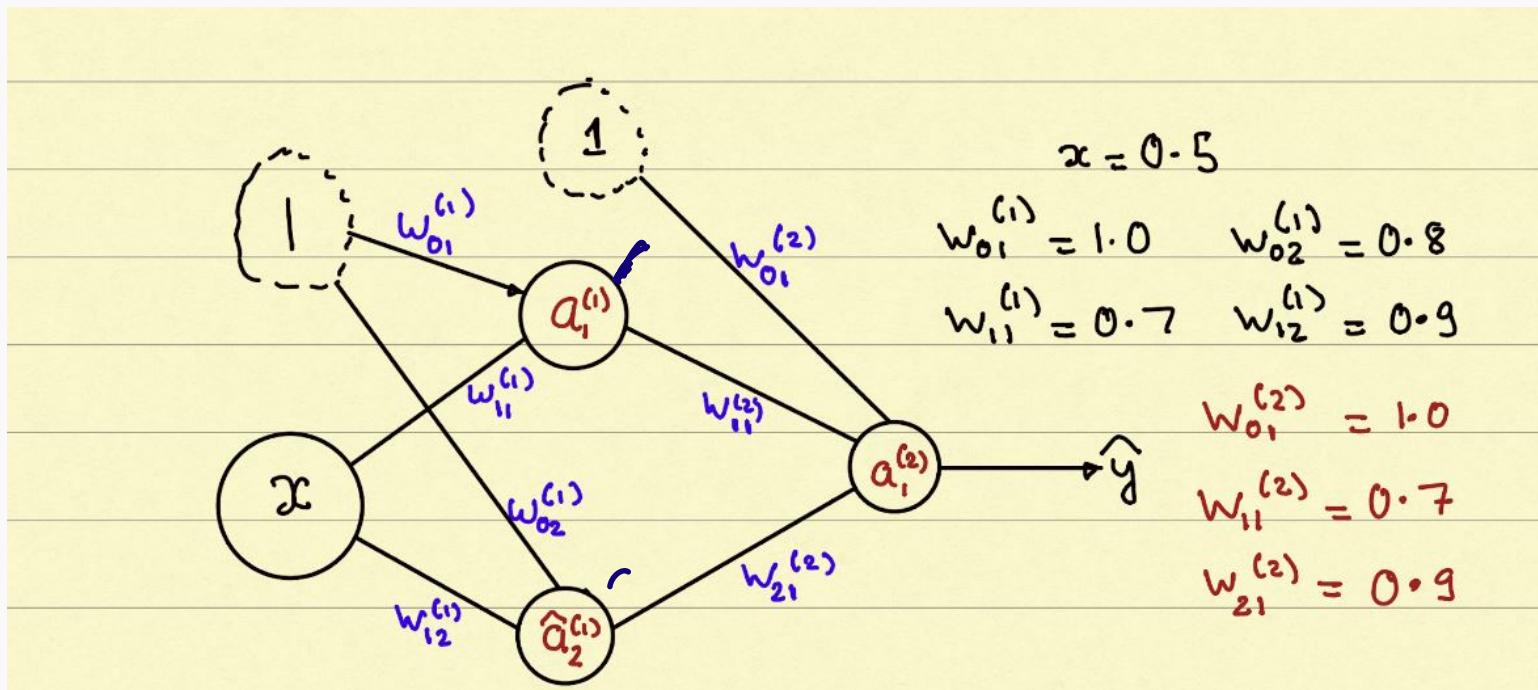


$$\hat{a}_2^{(1)} = g(z_2^{(1)}) = g(w_{02}^{(1)}x_0 + w_{12}^{(1)}x_1)$$



$$\hat{a}_2^{(1)} = g(z_2^{(1)}) = g(w_{02}^{(1)}x_0 + w_{12}^{(1)}x_1)$$

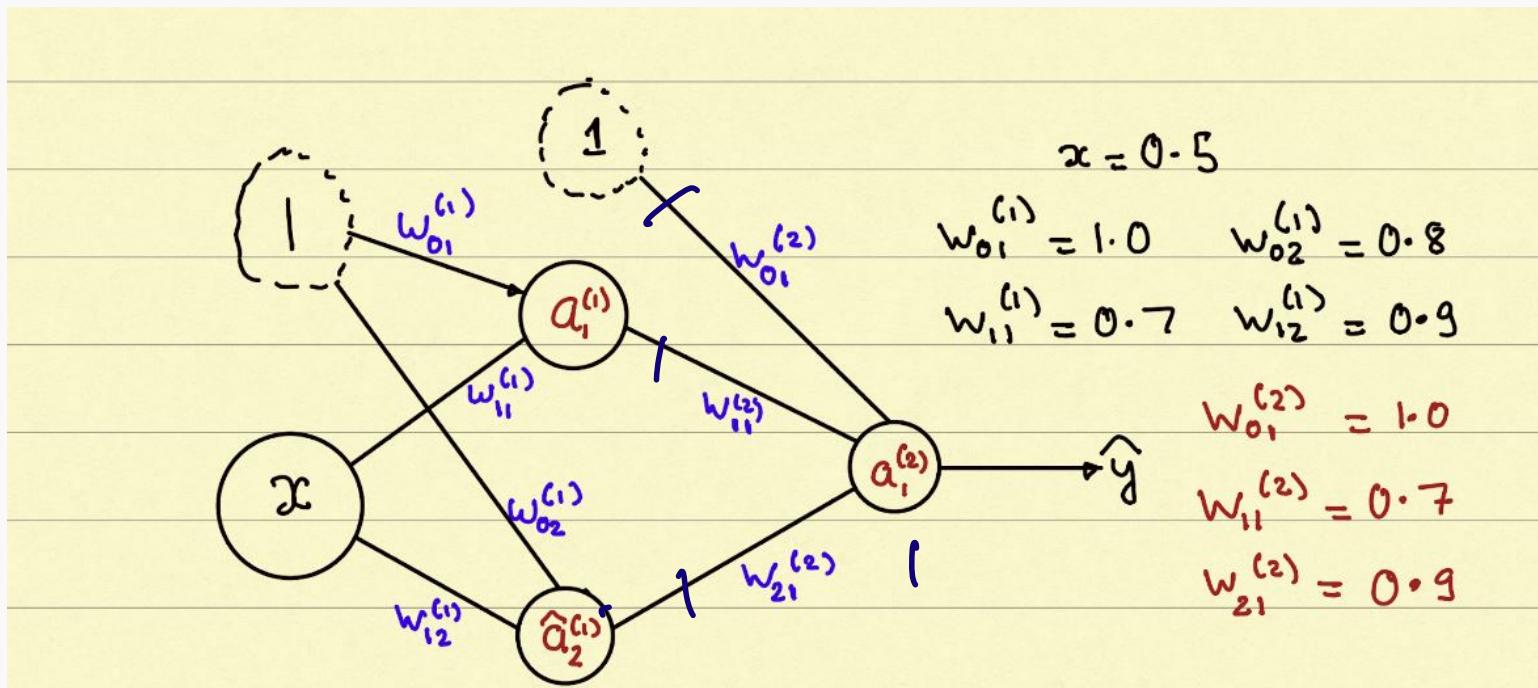
$$= g(0.8 + 0.9 \times 0.5) = g(1.25)$$



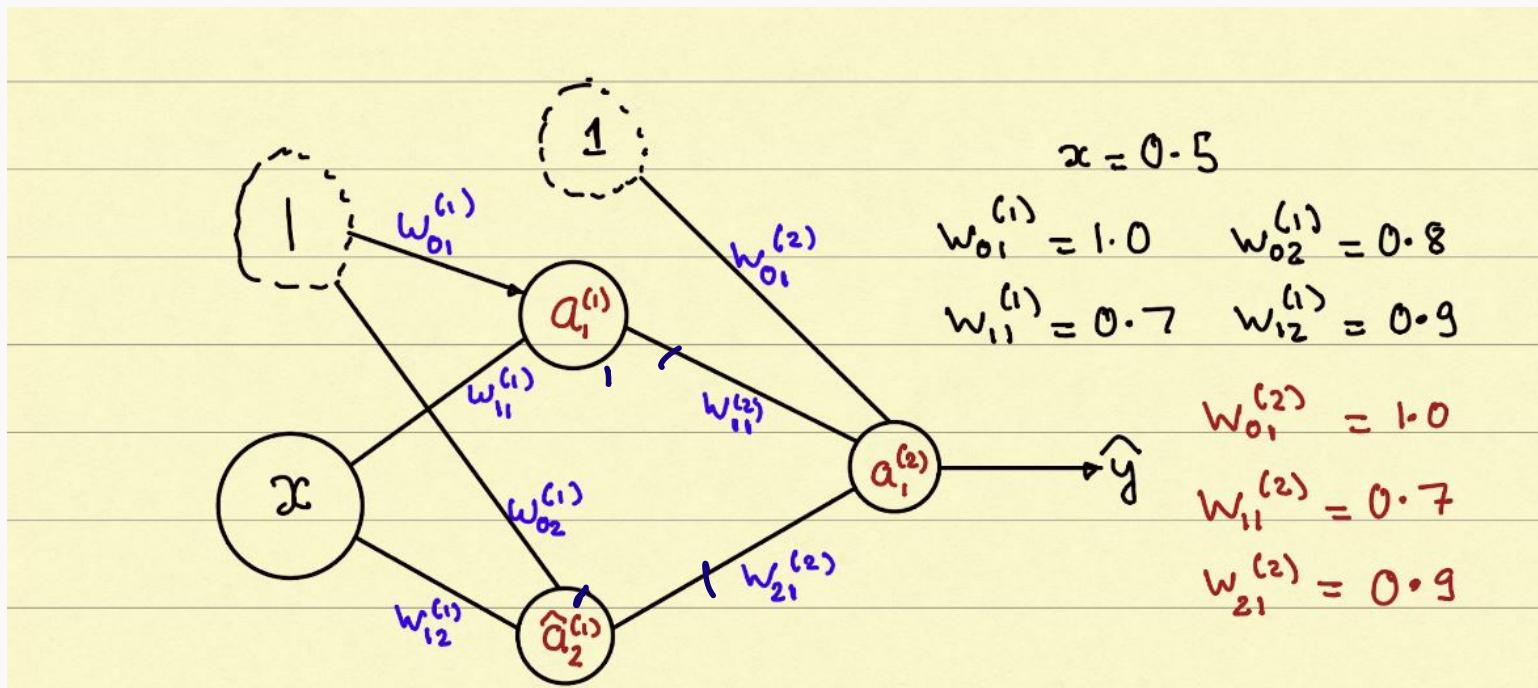
$$\hat{a}_2^{(1)} = g(z_2^{(1)}) = g(w_{02}^{(1)}x_0 + w_{12}^{(1)}x_1)$$

$$= g(0.8 + 0.9 \times 0.5) = g(1.25)$$

$$\hat{a}_2^{(1)} = 0.7773$$

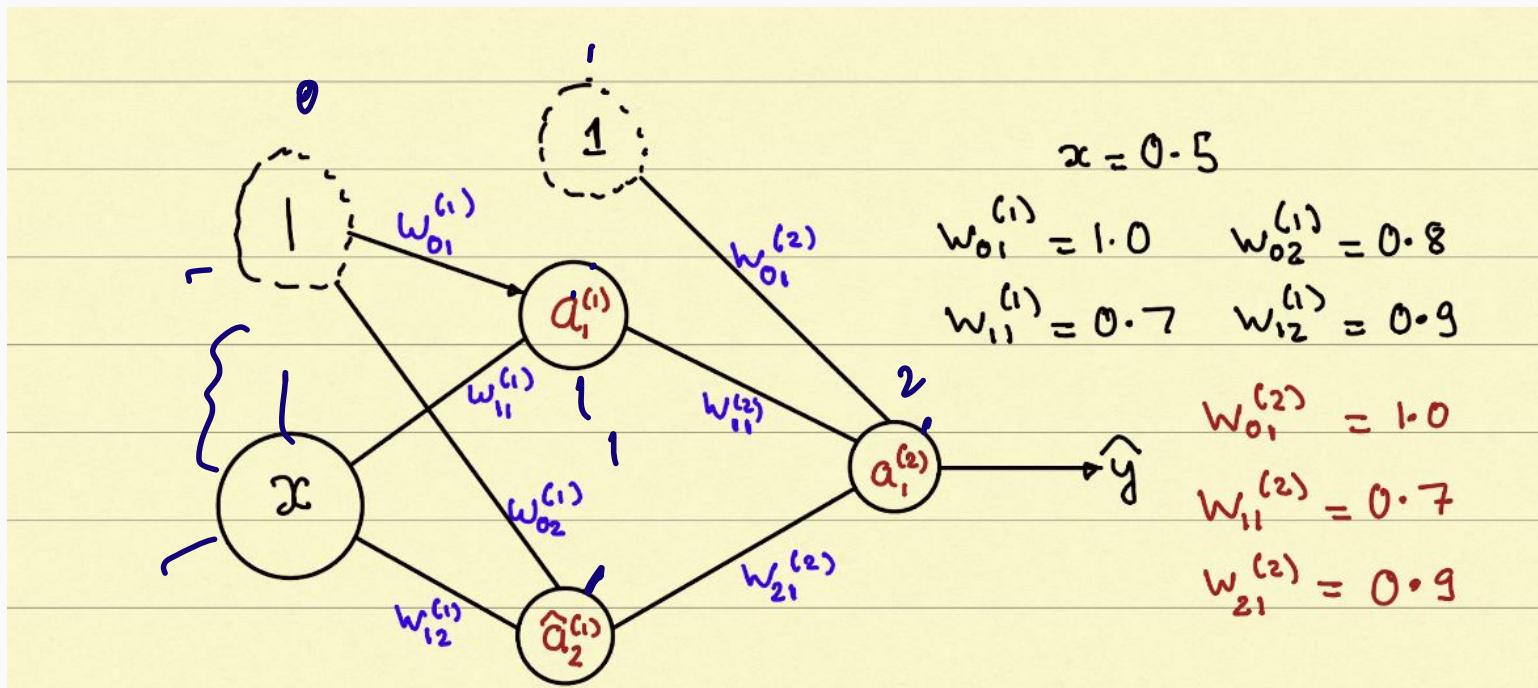


$$\hat{a}_1^{(2)} = g(z_1^{(2)}) = g(w_{01}^{(2)}x_0 + w_{11}^{(2)}a_1^{(1)} + w_{21}^{(2)}\hat{a}_2^{(1)})$$



$$\hat{a}_1^{(2)} = g(z_1^{(2)}) = g(w_{01}^{(2)}x_0 + w_{11}^{(2)}a_1^{(1)} + w_{21}^{(2)}\hat{a}_2^{(1)})$$

$$= g(1 \times 1 + 0.7 \times \underbrace{0.7941}_{\sim} + 0.9 \times 0.7773) = g(2.255)$$



$$\hat{a}_1^{(2)} = g(z_1^{(2)}) = g(w_{01}^{(2)}x_0 + w_{11}^{(2)}a_1^{(1)} + w_{21}^{(2)}\hat{a}_2^{(1)})$$

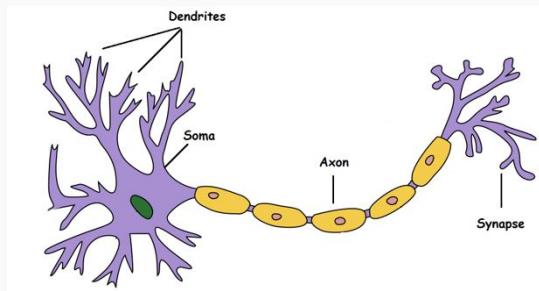
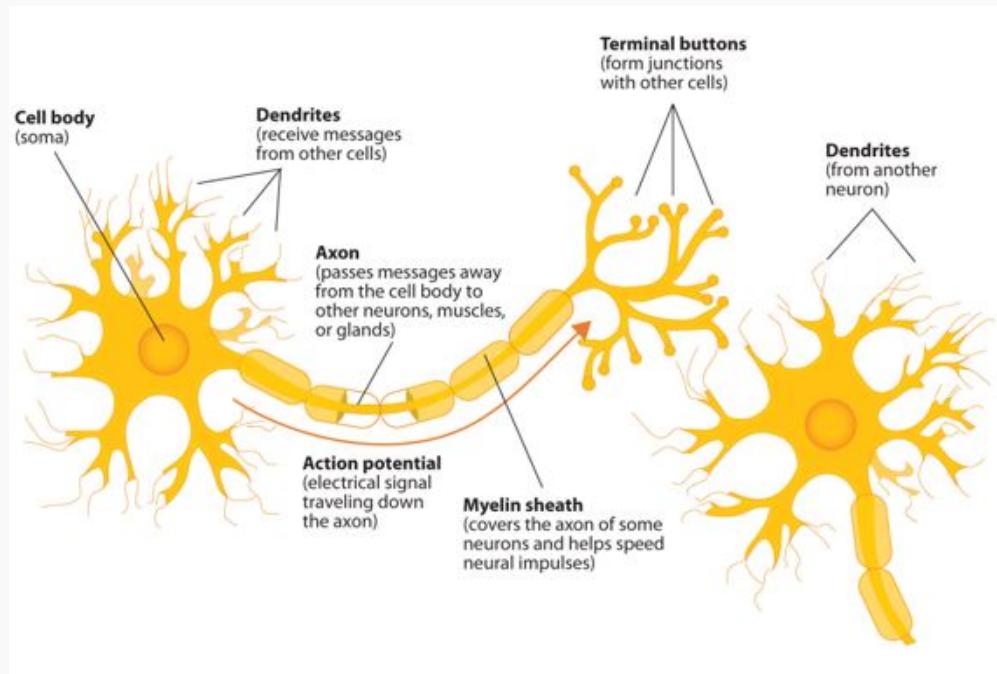
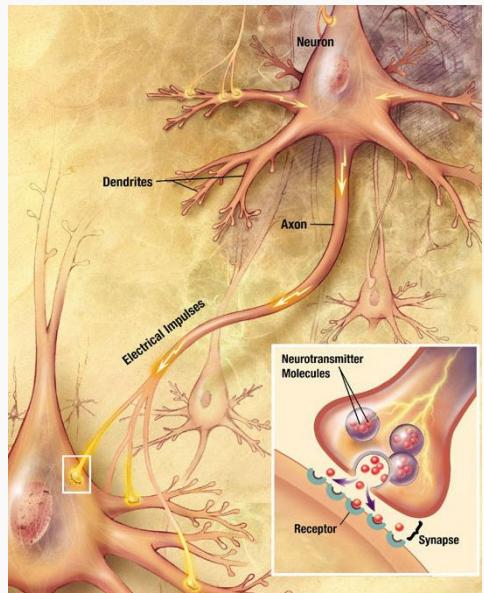
$$= g(1 \times 1 + 0.7 \times 0.7941 + 0.9 \times 0.7773) = g(2.255)$$

$=$

$$\hat{a}_1^{(2)} = 0.9051 = \hat{y}$$

$x \in W \Rightarrow \hat{y} = A'W + b$

From biology to computation



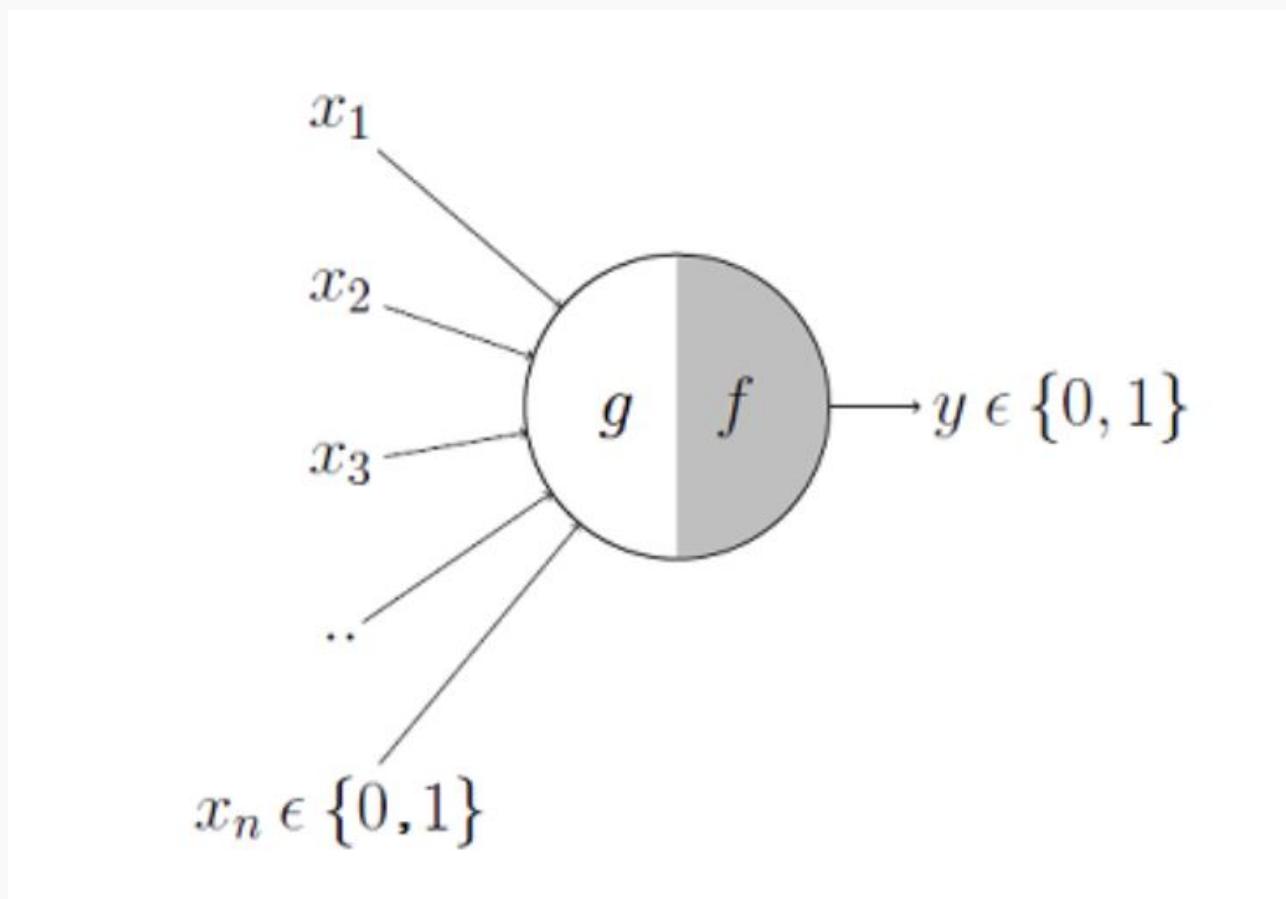
# McCulloch-Pitts Neuron

# McCulloch-Pitts Neuron

- First Mathematical model of Neuron (1943)

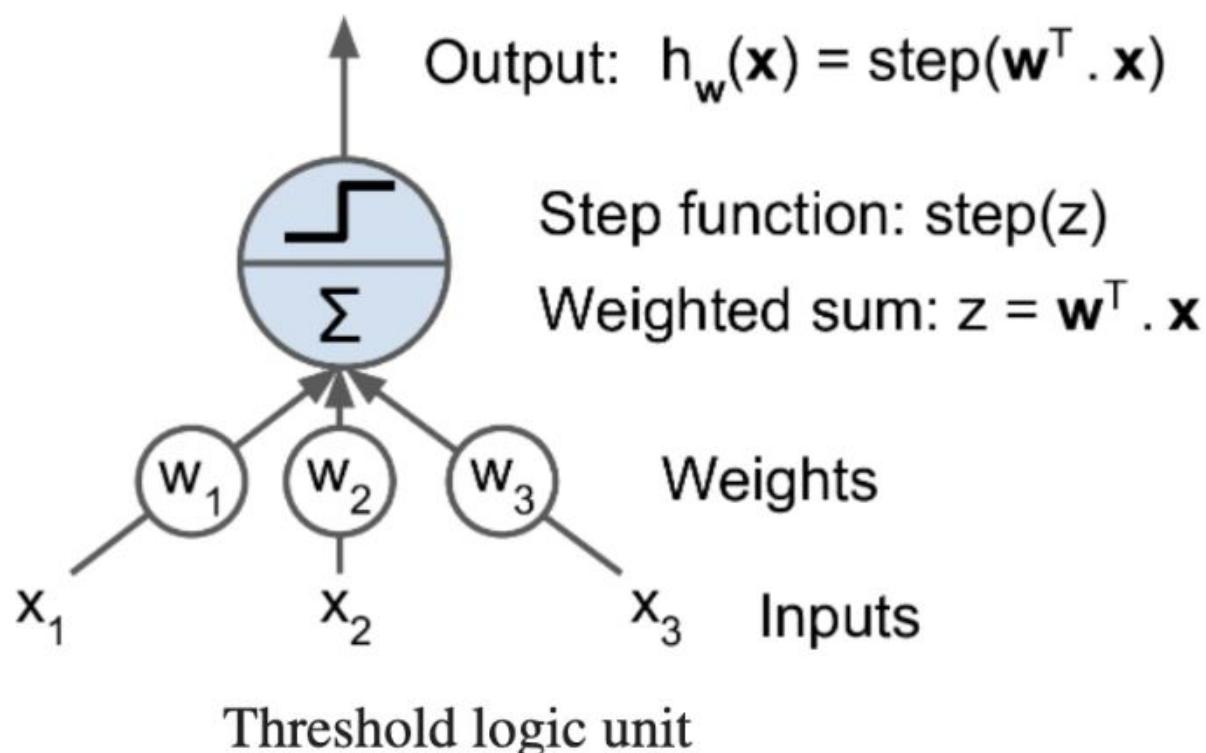
# McCulloch-Pitts Neuron

- First Mathematical model of Neuron (1943)



# Rosenblatt Perceptron (1957)

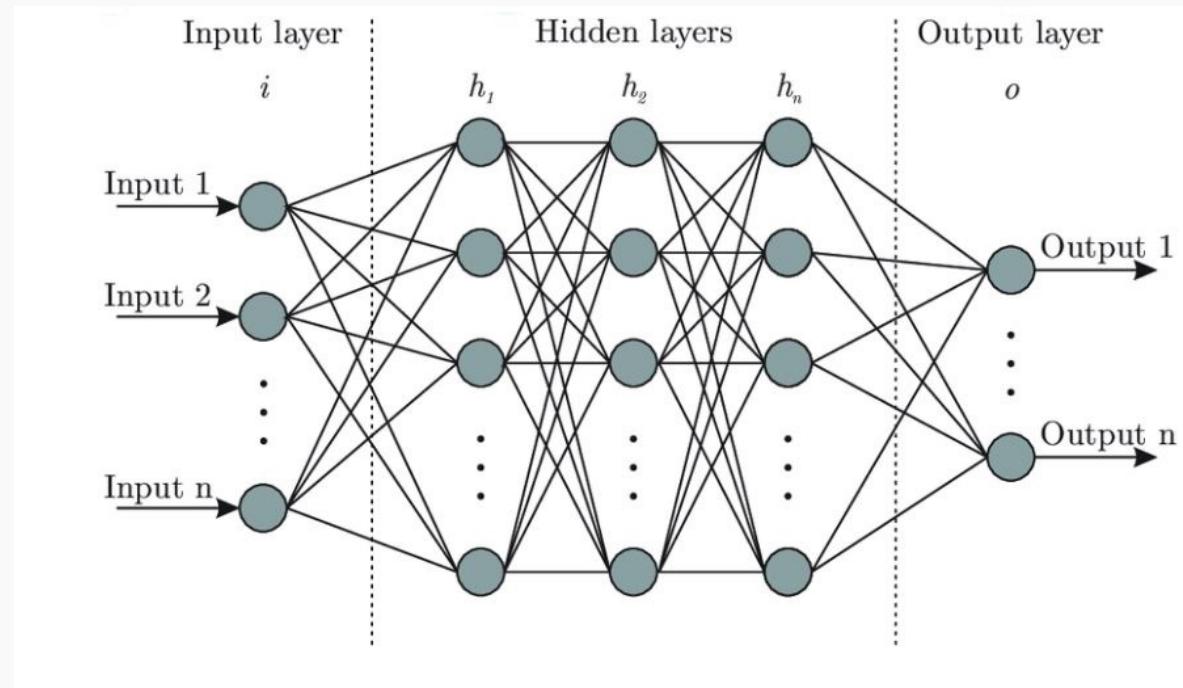
# Rosenblatt Perceptron (1957)



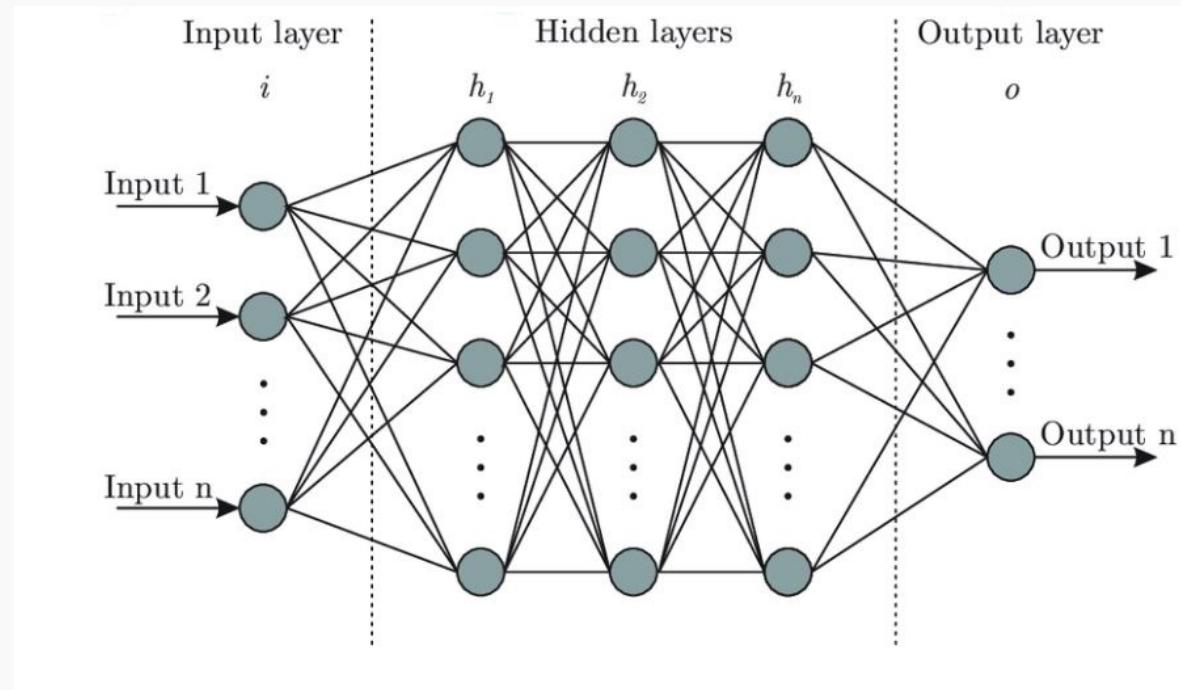
# Common non-linearities

# Multilayer Perceptrons/Deep NNs

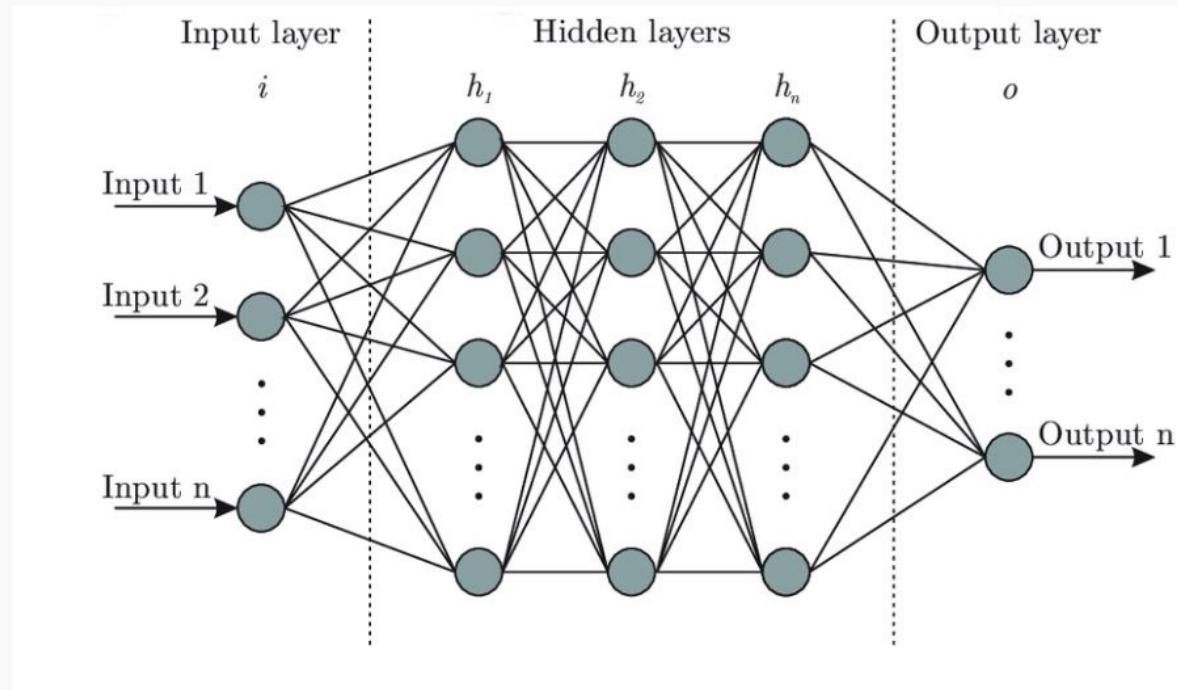
# Multilayer Perceptrons/Deep NNs



# Multilayer Perceptrons/Deep NNs

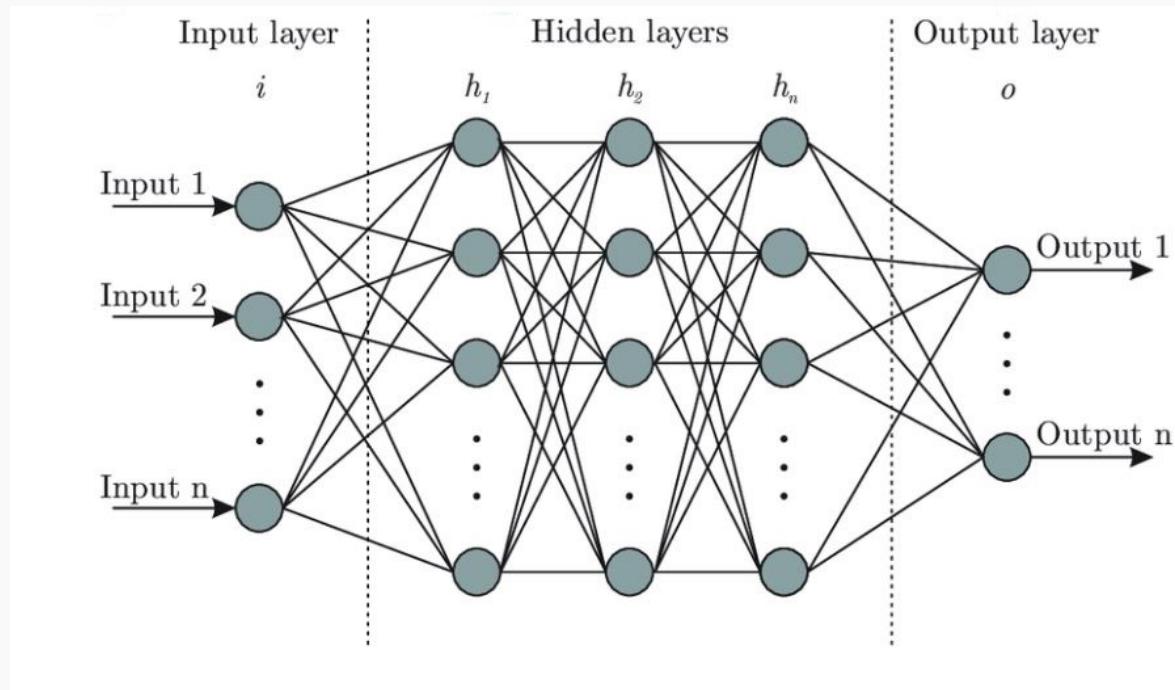


# Multilayer Perceptrons/Deep NNs



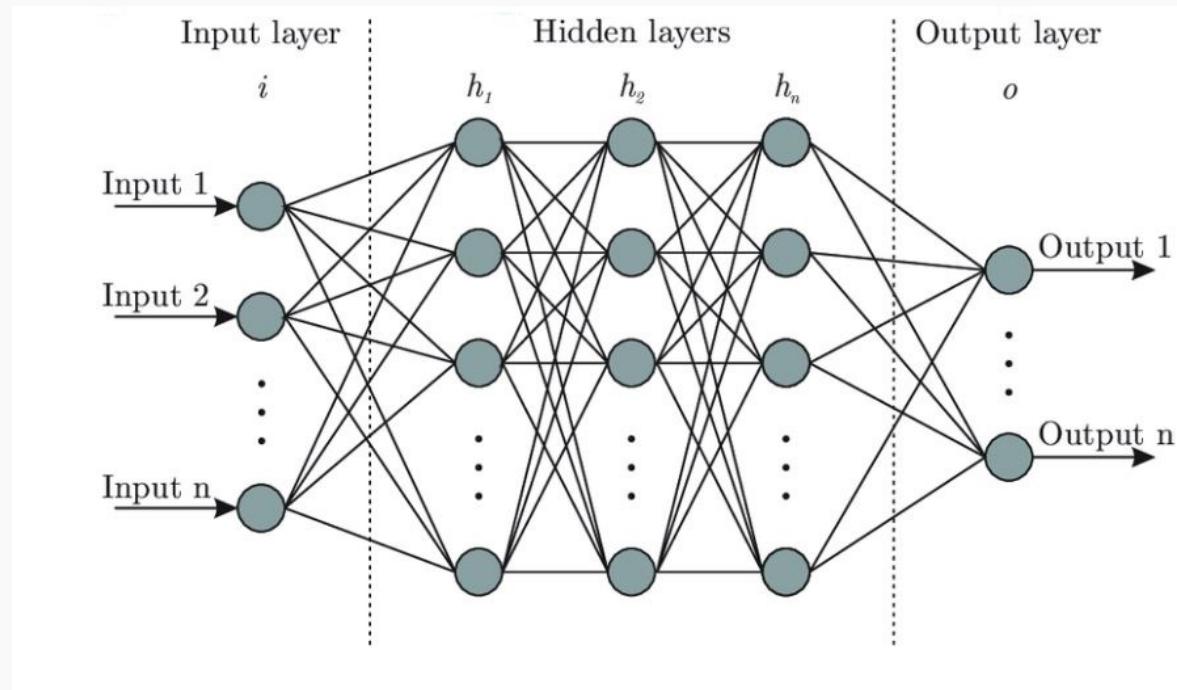
- MLPs or Fully connected Networks are very common

# Multilayer Perceptrons/Deep NNs



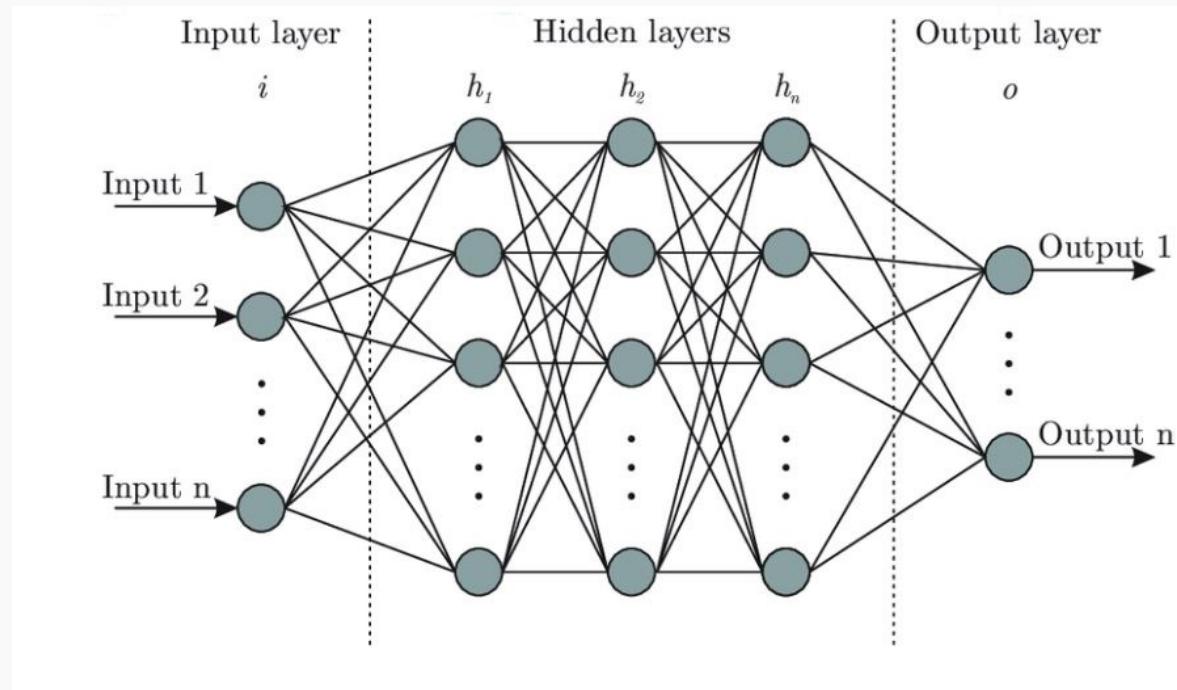
- MLPs or Fully connected Networks are very common
  - At least as parts of other, more sparse networks

# Multilayer Perceptrons/Deep NNs



- MLPs or Fully connected Networks are very common
  - At least as parts of other, more sparse networks
- Satisfy the Universal Approximation Theorem

# Multilayer Perceptrons/Deep NNs



- MLPs or Fully connected Networks are very common
  - At least as parts of other, more sparse networks
- Satisfy the Universal Approximation Theorem
  - Can approximate any function to arbitrary accuracy with just one layer.

# Types of Neural Networks

# Types of Neural Networks

- Artificial Neural Networks

# Types of Neural Networks

- Artificial Neural Networks
  - For general numerical data

# Types of Neural Networks

- Artificial Neural Networks
  - For general numerical data
- Convolutional Neural Networks (CNNs)

# Types of Neural Networks

- Artificial Neural Networks
  - For general numerical data
- Convolutional Neural Networks (CNNs)
  - For image based data

# Types of Neural Networks

- Artificial Neural Networks
  - For general numerical data
- Convolutional Neural Networks (CNNs)
  - For image based data
- Recurrent Neural Networks (RNNs)

# Types of Neural Networks

- Artificial Neural Networks
  - For general numerical data
- Convolutional Neural Networks (CNNs)
  - For image based data
- Recurrent Neural Networks (RNNs)
  - For sequential data

# From linear models to DNNs

# From linear models to DNNs