



©www.n-ix.com/

# CS395: Selected CS1 (Introduction to Machine Learning)

Associate Prof. Wessam El-Behaidy

**Fall 2021**

References:

<https://www.coursera.org/learn/machine-learning> (Andrew Ng)  
Machine learning A to Z: Kirill Eremenko ©superdatascience

# Midterm Exam

- Midterm Exam (based on announced schedule)

**Sun. 21 Nov. from 9:00 AM to 9:30 AM**

## تذكيرات هامة

- الحضور للكلية قبل موعد الامتحان بنصف ساعة على الأقل
- يتواجد معك (قلم جاف- قلم رصاص - استيكه) هام جدا جدا
- الاجابة في ورقة الاجابات الالكترونية بالقلم الرصاص وتحبر في النهاية
- كتابة البيانات على ورقة الاسئلة وورقة الاجابة
- تسليم كلا من ورقة الاسئلة وورقة الاجابة (عدم تسليم الورقتين يلغى امتحانك)
- يجب احضار الة حاسبة للامتحان

 <b>لجنة الاعداد</b>							
Name _____	Student ID _____						
Academic Year : ..... ١ ٢ ٣ ٤ ٥	- Use pencil or blue or black pen						
Department : .....	- Fill the circle completely (like this  Not like these						
Subject : .....	- Do not fill in more than one circle						
Examination Form : A B C D E	- To change your answer use the eraser						
Date : .....	- DO NOT USE THE CORRECTOR						
	- For true & False questions, (mark (A) for true & (E) for false)						
حسم الكترونول							
<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>
1 - A B C D E	11 - A B C D E	21 - A B C D E	31 - A B C D E	41 - A B C D E			
2 - A B C D E	12 - A B C D E	22 - A B C D E	32 - A B C D E	42 - A B C D E			
3 - A B C D E	13 - A B C D E	23 - A B C D E	33 - A B C D E	43 - A B C D E			
4 - A B C D E	14 - A B C D E	24 - A B C D E	34 - A B C D E	44 - A B C D E			
5 - A B C D E	15 - A B C D E	25 - A B C D E	35 - A B C D E	45 - A B C D E			
6 - A B C D E	16 - A B C D E	26 - A B C D E	36 - A B C D E	46 - A B C D E			
7 - A B C D E	17 - A B C D E	27 - A B C D E	37 - A B C D E	47 - A B C D E			
8 - A B C D E	18 - A B C D E	28 - A B C D E	38 - A B C D E	48 - A B C D E			
9 - A B C D E	19 - A B C D E	29 - A B C D E	39 - A B C D E	49 - A B C D E			
10 - A B C D E	20 - A B C D E	30 - A B C D E	40 - A B C D E	50 - A B C D E			
<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>	<b>T</b>	<b>F</b>
51 - A B C D E	61 - A B C D E	71 - A B C D E	81 - A B C D E	91 - A B C D E			
52 - A B C D E	62 - A B C D E	72 - A B C D E	82 - A B C D E	92 - A B C D E			
53 - A B C D E	63 - A B C D E	73 - A B C D E	83 - A B C D E	93 - A B C D E			
54 - A B C D E	64 - A B C D E	74 - A B C D E	84 - A B C D E	94 - A B C D E			
55 - A B C D E	65 - A B C D E	75 - A B C D E	85 - A B C D E	95 - A B C D E			
56 - A B C D E	66 - A B C D E	76 - A B C D E	86 - A B C D E	96 - A B C D E			
57 - A B C D E	67 - A B C D E	77 - A B C D E	87 - A B C D E	97 - A B C D E			
58 - A B C D E	68 - A B C D E	78 - A B C D E	88 - A B C D E	98 - A B C D E			
59 - A B C D E	69 - A B C D E	79 - A B C D E	89 - A B C D E	99 - A B C D E			
60 - A B C D E	70 - A B C D E	80 - A B C D E	90 - A B C D E	100 - A B C D E			

# الشكل العام لورقة التصحيح الإلكتروني

# البيانات الأساسية



Name : .....

Academic Year : .....       ② ③ ④ ⑤      مثال: 1st level

Department : .....

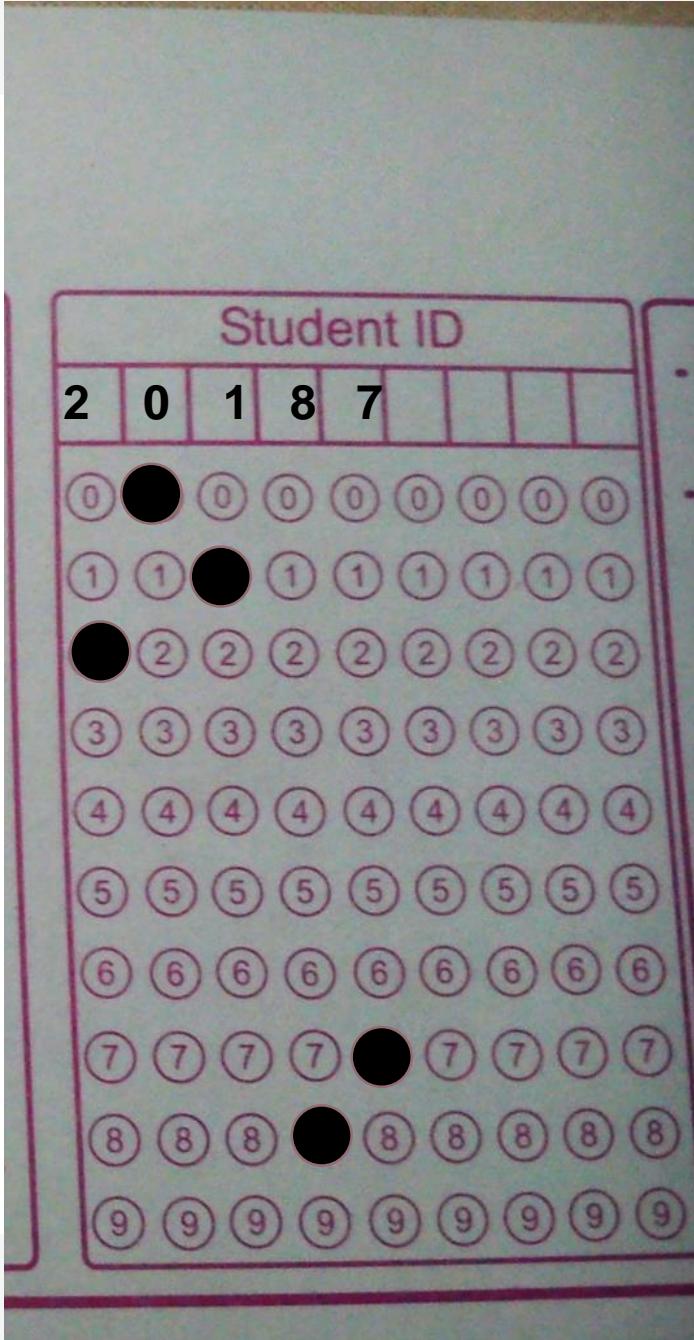
Subject: .....

Examination Form :       A  C  D  E      مثال: Model B

Date : .....

# رقم الطالب

مثال:  
رقم الجلوس : 20187



# تعليمات هامة

- Use pencil or blue or black pen
- Fill the circle completely  
( like this ● Not like these ○ ✓ ✗ )
- Do not fill in more than one circle
- To change your answer use the eraser
- DO NOT USE THE CORRECTOR
- For true & False questions,  
(mark (A) for true& (E) For false)

# مثال لإجابات مظاللة بشكل صحيح

مثال:  
اجابة  
السؤال  
الأول  
**B**

	T	F	T	F	T	F	T	F	T	F	T	F					
1 -	A	B	C	D	E	11 -	A	B	C	D	E	21 -	A	B	C	D	E
2 -		B	C	D	E	12 -	A	B	C	D	E	22 -	A	B	C	D	E
3 -	A	B	C		E	13 -	A	B	C	D	E	23 -	A	B	C	D	E
4 -	A	B	C	D	E	14 -	A	B	C	D	E	24 -	A	B	C	D	E
5 -	A	B	C	D	E	15 -	A	B	C	D	E	25 -	A	B	C	D	E
6 -	A	B	C	D	E	16 -	A	B	C	D	E	26 -	A	B	C	D	E
7 -	A	B	C	D	E	17 -	A	B	C	D	E	27 -	A	B	C	D	E
8 -	A	B	C	D	E	18 -	A	B	C	D	E	28 -	A	B	C	D	E
9 -	A	B	C	D	E	19 -	A	B	C	D	E	29 -	A	B	C	D	E
10 -		B	C	D	E	20 -	A	B	C	D	E	30 -	A	B	C	D	E
	T	F	T	F	T	F	T	F	T	F	T	F					
51 -	A	B	C	D	E	61 -	A	B	C	D	E	71 -	A	B	C	D	E
52 -	A	B	C	D	E	62 -	A	B	C	D	E	72 -	A	B	C	D	E
	T	F	T	F	T	F	T	F	T	F	T	F					
81 -	A	B	C	D	E	91 -	A	B	C	D	E						
82 -	A	B	C	D	E	92 -	A	B	C	D	E						

# مثال لإجابات مظللة بشكل خطأ

	T	F	T	F	T	F
1 -	●	Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ	11 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	21 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ
2 -	●	Ⓐ Ⓑ Ⓒ Ⓗ Ⓕ	12 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	22 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
3 -	X	Ⓐ Ⓑ Ⓒ Ⓓ Ⓔ	13 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	23 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
4 -	Ⓐ Ⓑ Ⓒ Ⓗ Ⓕ	✓	14 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	24 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
5 -	●	Ⓑ Ⓑ Ⓒ Ⓓ Ⓔ	15 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	25 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
6 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ		16 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	26 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
7 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ		17 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	27 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
8 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ		18 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	28 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ
9 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ		19 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ	29 -	Ⓐ Ⓑ Ⓒ Ⓓ Ⓕ

# Problems of Logistic Regression with non-linear classification

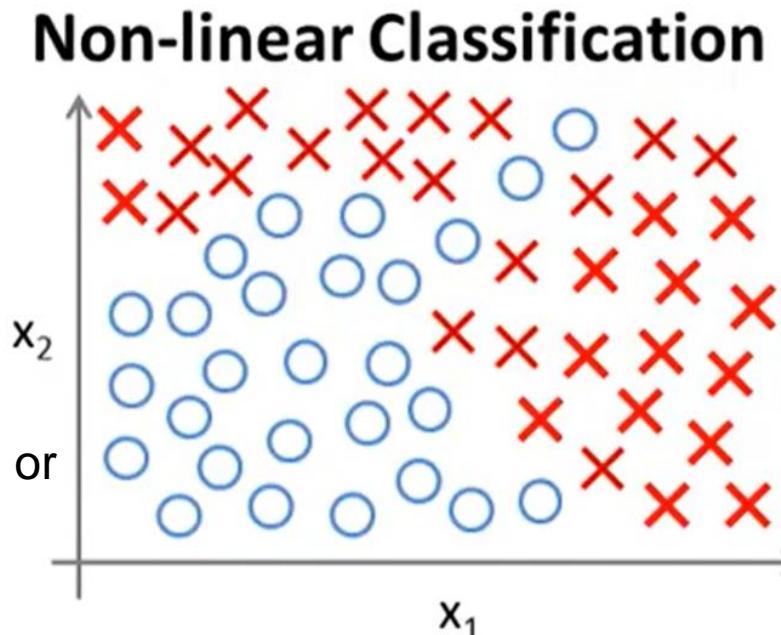
# Non-linear Classification

$$h_{\theta}(x) = g(\theta^T x) = g\left(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1^3 \dots \dots \right)$$

If n=100,  
how many **parameters** if we only use second polynomial terms or  
Quadratic terms??

What if we need to have more higher polynomial terms?  
how many parameters?

Is this a good way to build a non-linear classifier, **when n is large?**

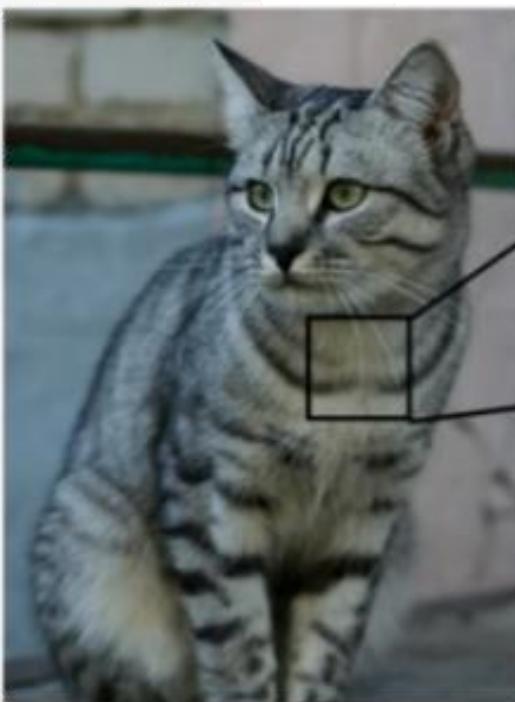


# Image Classification Problem



How many features do we need?

# Semantic Gap, What the computer sees?



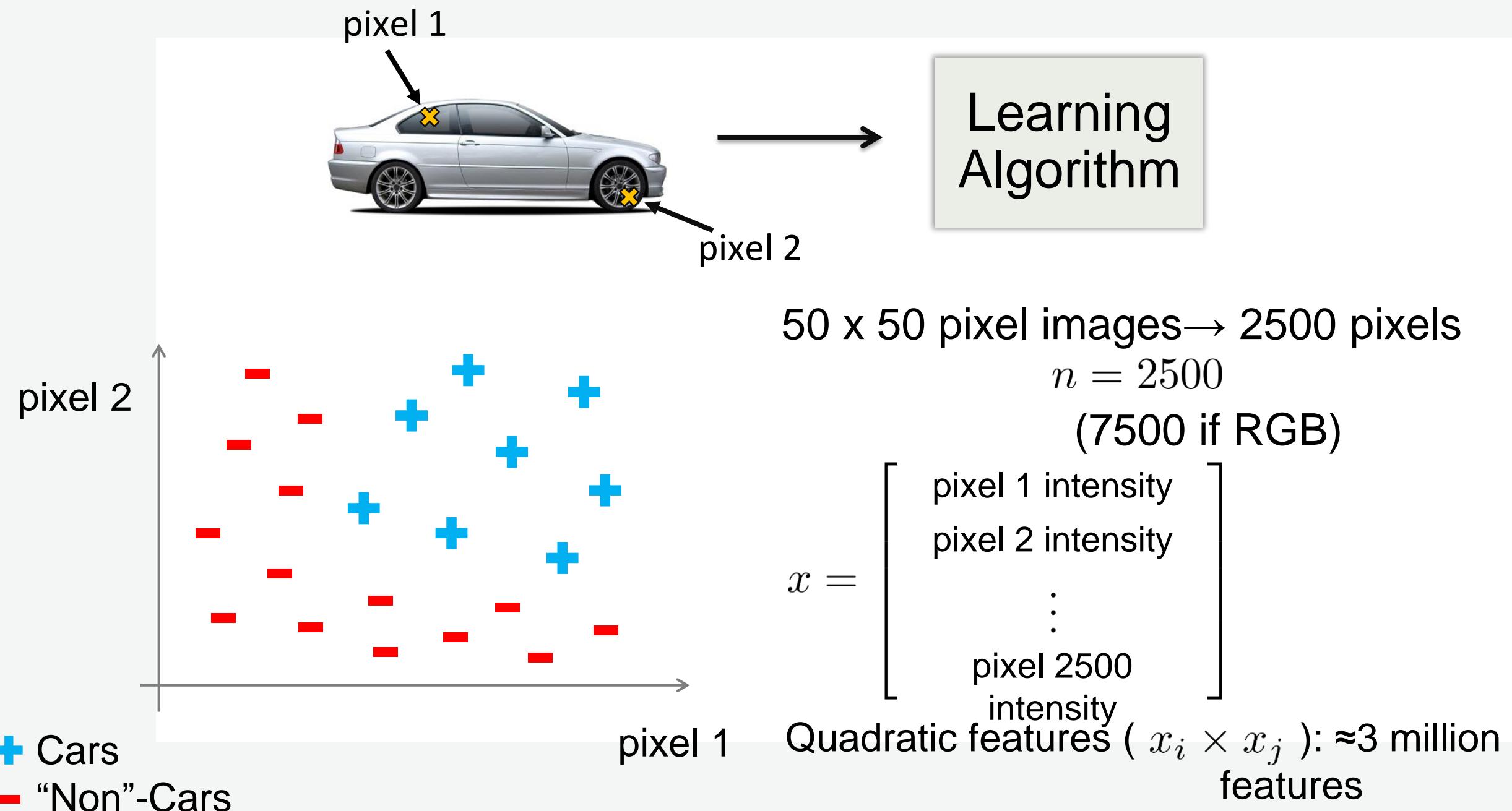
13495	332	188	111	184	99	186	99	96	183	132	119	184	97	93	187	
1	93	98	382	186	184	79	98	183	99	185	123	134	118	185	94	853
1	76	85	98	185	128	185	87	96	95	99	125	112	186	183	98	857
1	99	81	81	91	128	133	127	188	95	98	182	99	96	93	181	942
1346	91	83	64	69	91	88	85	185	187	189	98	79	84	96	953	
1314	188	89	55	55	49	64	94	64	87	132	129	98	74	84	912	
1333	337	347	183	65	81	98	85	92	94	74	84	182	93	85	823	
1328	337	244	148	189	95	86	78	62	65	63	63	68	73	86	5812	
1329	133	148	137	119	521	137	94	85	79	88	65	54	64	72	981	
1327	125	131	147	133	127	326	135	155	98	89	75	63	64	72	843	
1315	124	149	123	158	149	135	158	123	189	188	92	74	65	72	782	
1	89	93	98	97	188	147	131	118	133	114	113	189	186	95	77	882
1	43	77	86	81	77	79	182	123	117	115	125	138	115	87	119	
1	42	85	82	89	79	71	98	124	126	119	181	187	114	131	1192	
1	63	85	75	88	89	71	62	81	128	138	135	185	81	98	118	1183
1	87	85	71	87	186	85	68	45	76	138	126	187	92	84	185	1123
1218	97	82	88	117	123	106	66	41	51	95	93	89	95	182	1873	
1364	146	112	89	82	128	124	184	76	48	45	86	88	181	182	1881	
1357	178	237	128	93	86	254	132	112	97	69	55	78	82	99	942	
1338	128	134	161	139	188	189	118	325	134	114	87	65	53	69	862	
1328	112	96	117	158	144	138	115	184	187	182	93	87	81	72	782	
1329	187	96	86	83	112	153	149	122	188	284	75	88	187	132	982	
1322	521	182	89	82	86	94	117	145	148	153	182	58	78	92	1871	
1323	184	148	183	75	56	78	83	93	183	119	139	182	61	69	842	

What the computer sees

An image is just a big grid of numbers between [0, 255]:

e.g. 800 x 600 x 3

Is it easy for the machine?

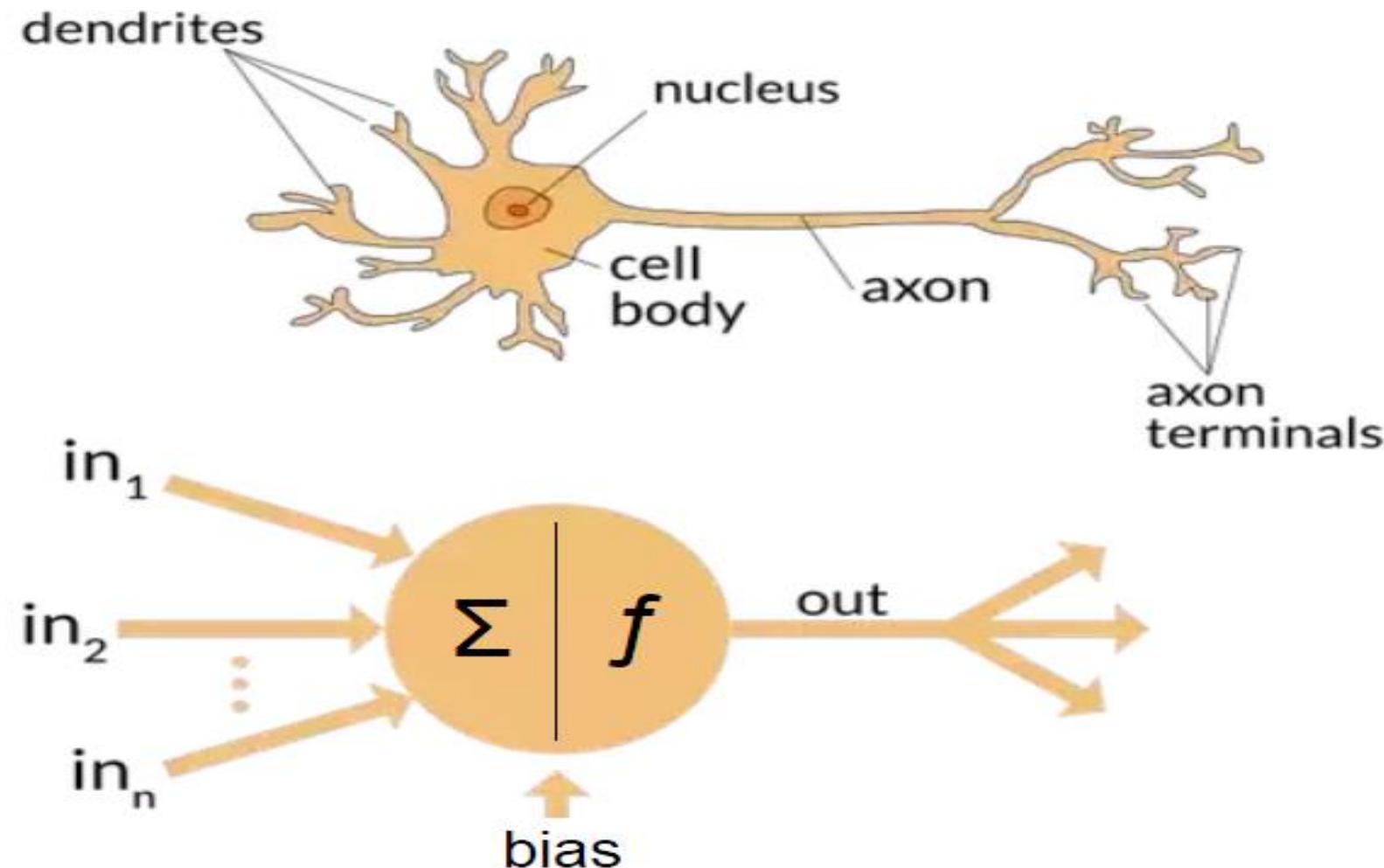


# Artificial Neural Network (ANN)

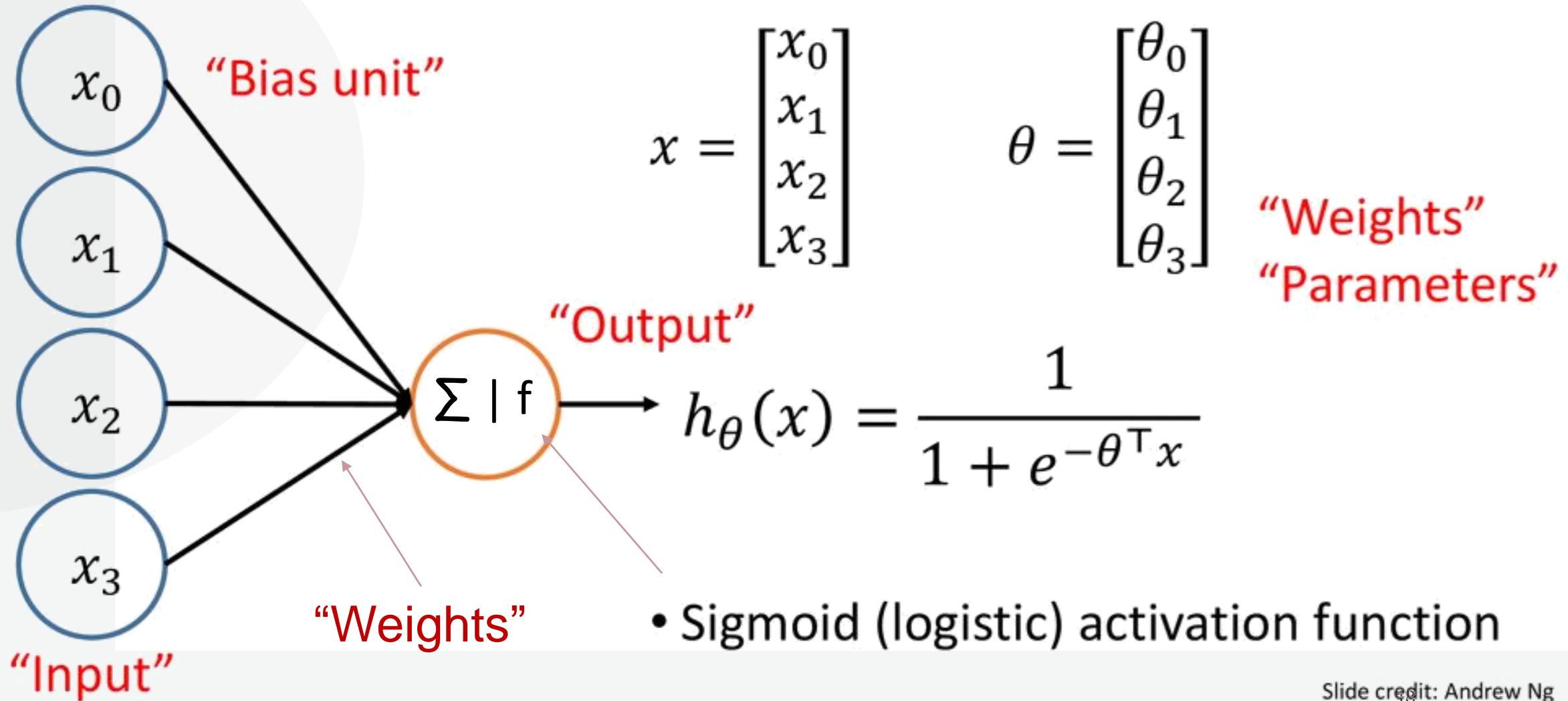
## Neural Networks, Neurons and the Brain

- Origins: Algorithms that try to mimic the brain.
- Was very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many applications

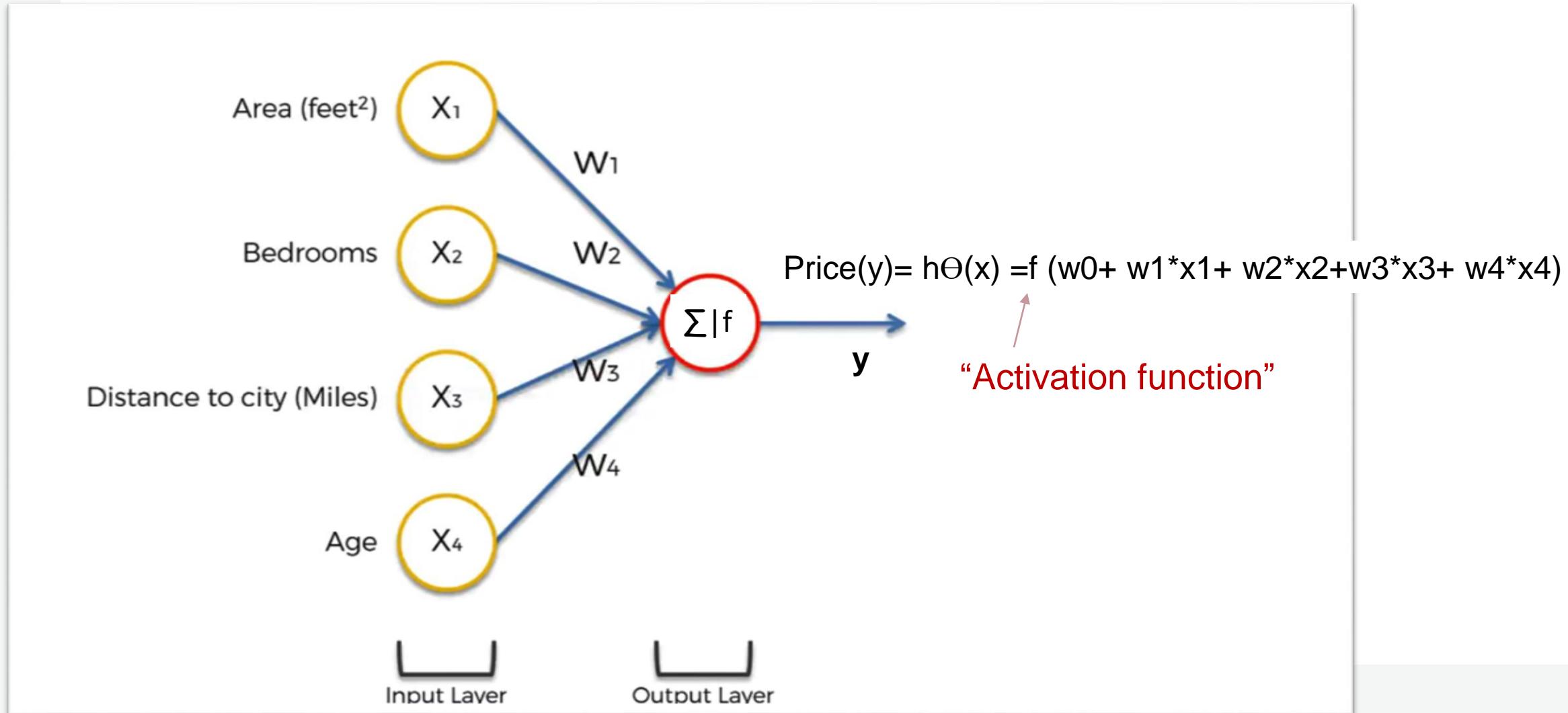
# Neural Networks, Neurons and the Brain



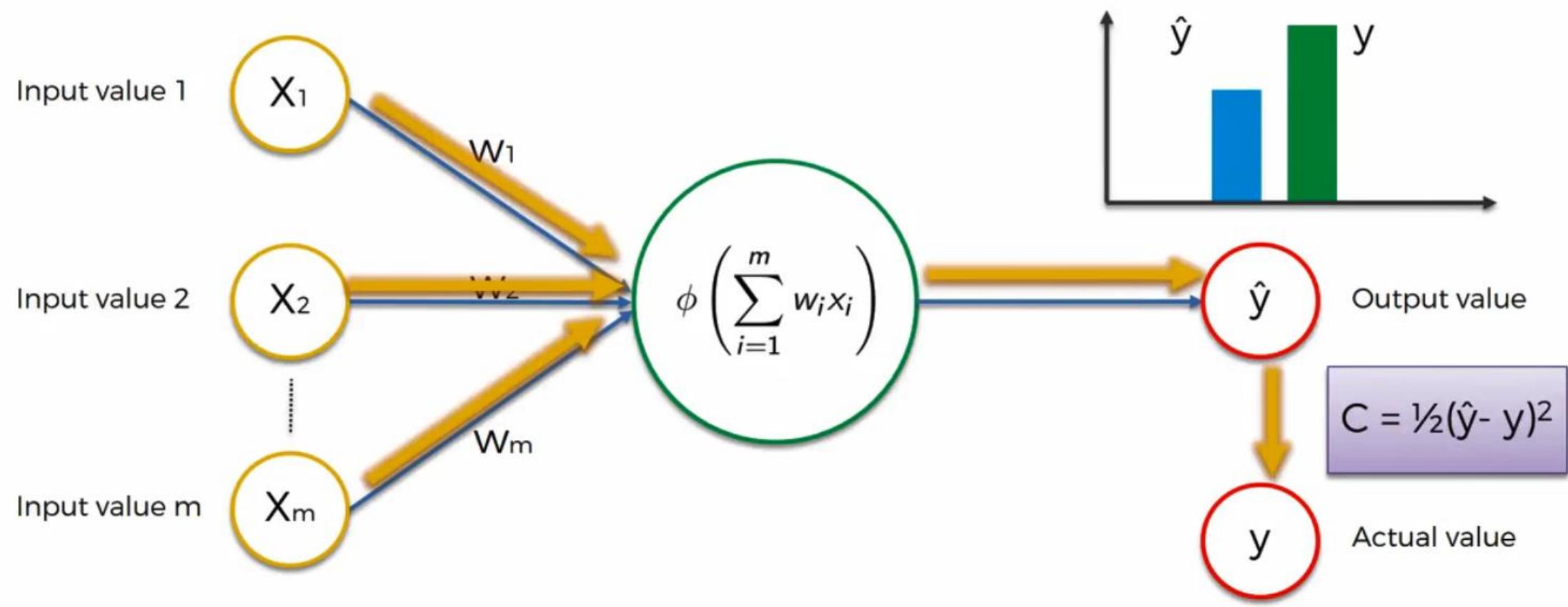
# An artificial neuron: Logistic unit



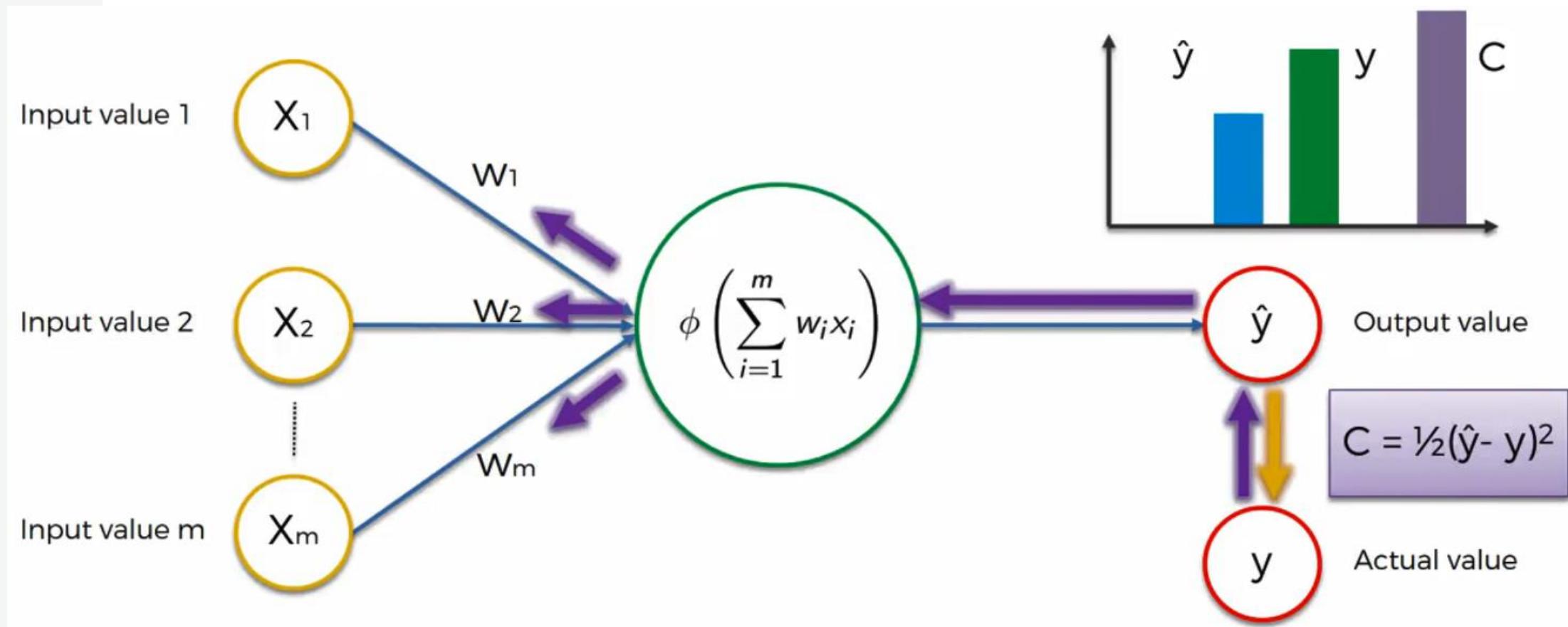
# Neural networks: Simple example



# Forward-Propagation

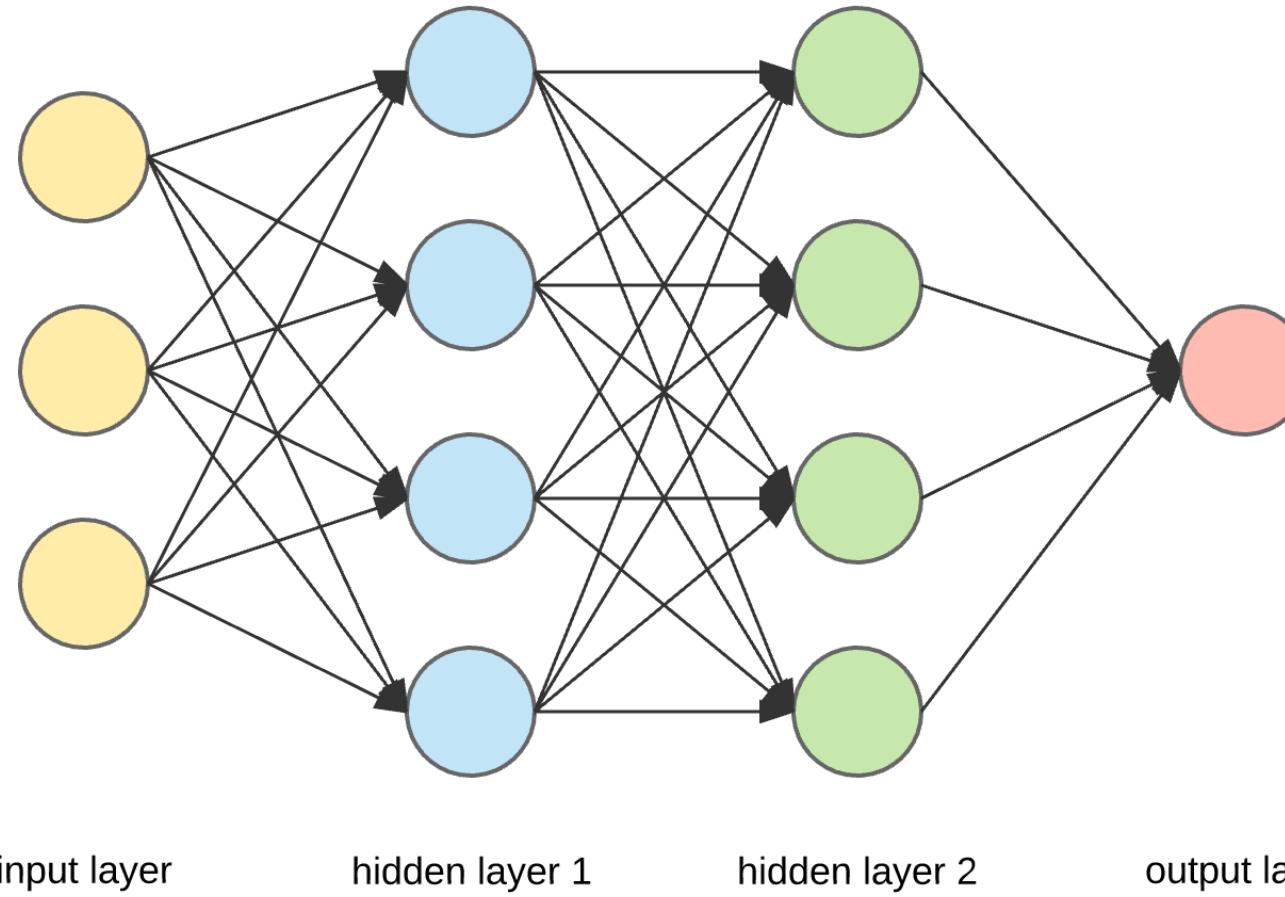


# Backward-Propagation



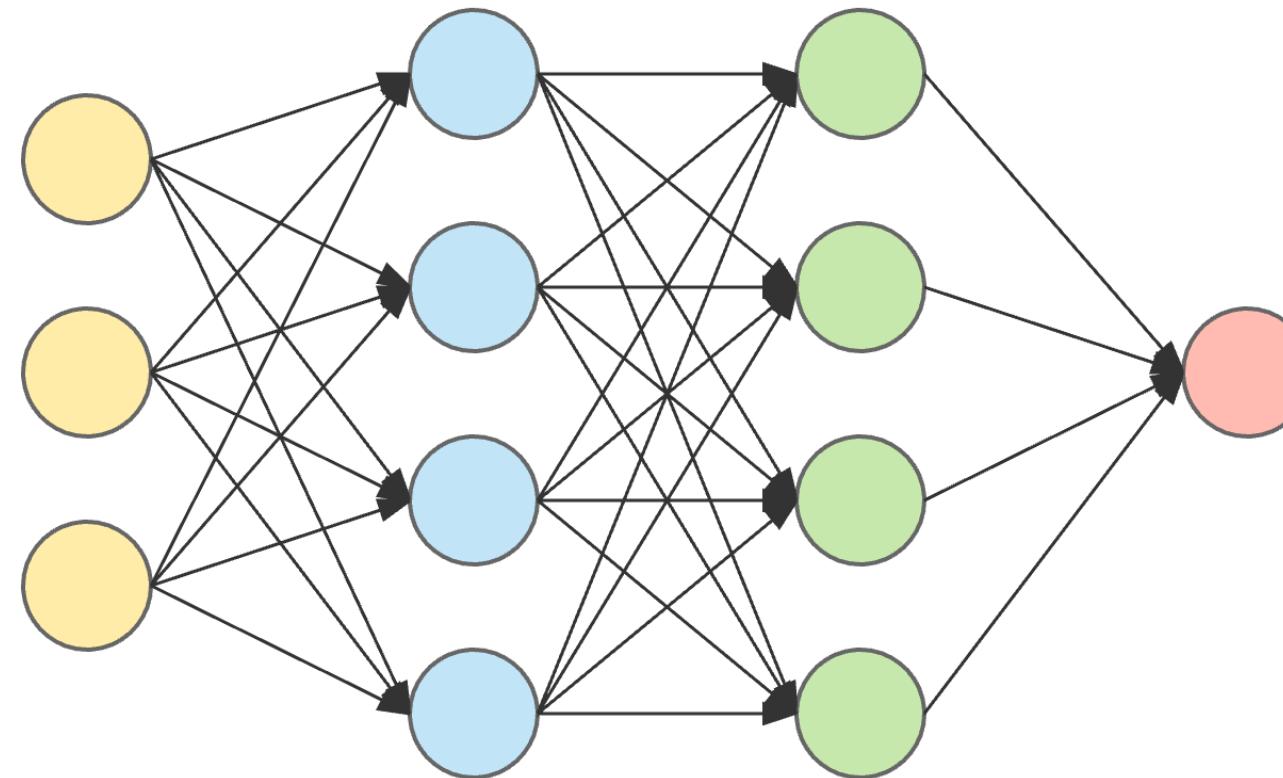
# Artificial Neural networks (ANN)

- Artificial Neural Networks (ANN) are multi-layer **fully-connected** neural nets. Every node in one layer is connected to every other node in the next layer.



# Artificial Neural networks (ANN) (Cont.)

- They consist of an **input layer**, one or multiple **hidden layers**, and an **output layer**. To make the network deeper, the number of hidden layers is increased.



input layer

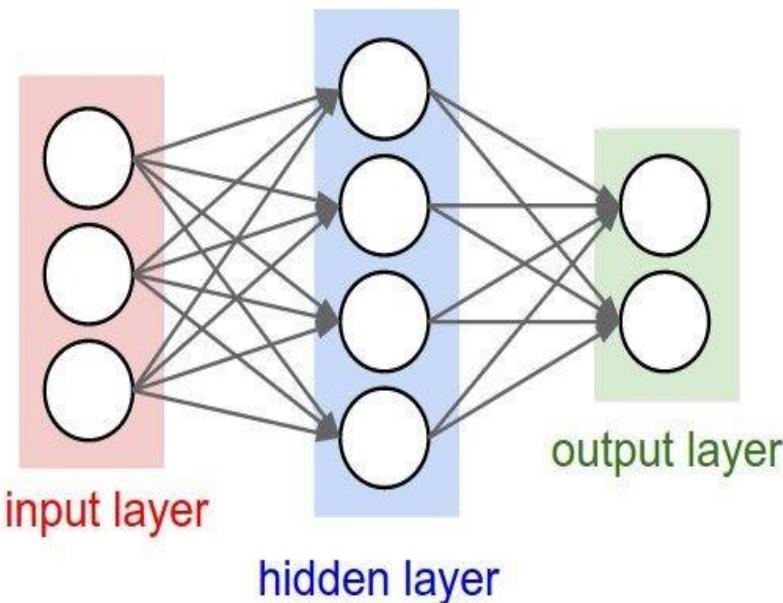
hidden layer 1

hidden layer 2

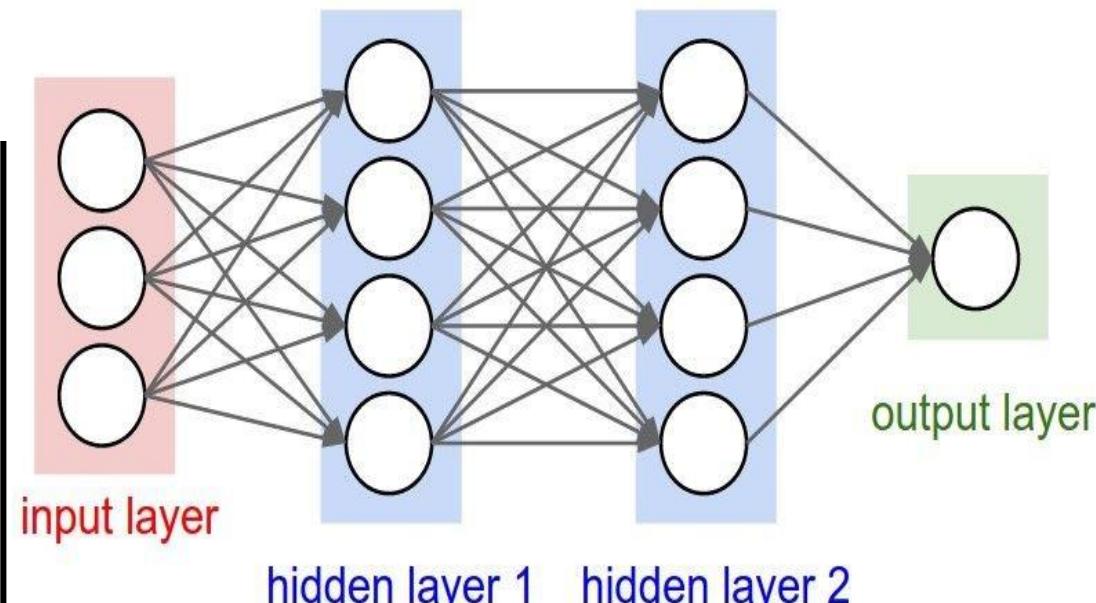
output layer

# Artificial Neural networks(ANN) Architectures

basically, anything that isn't an input layer and isn't an output layer is called a hidden layer.

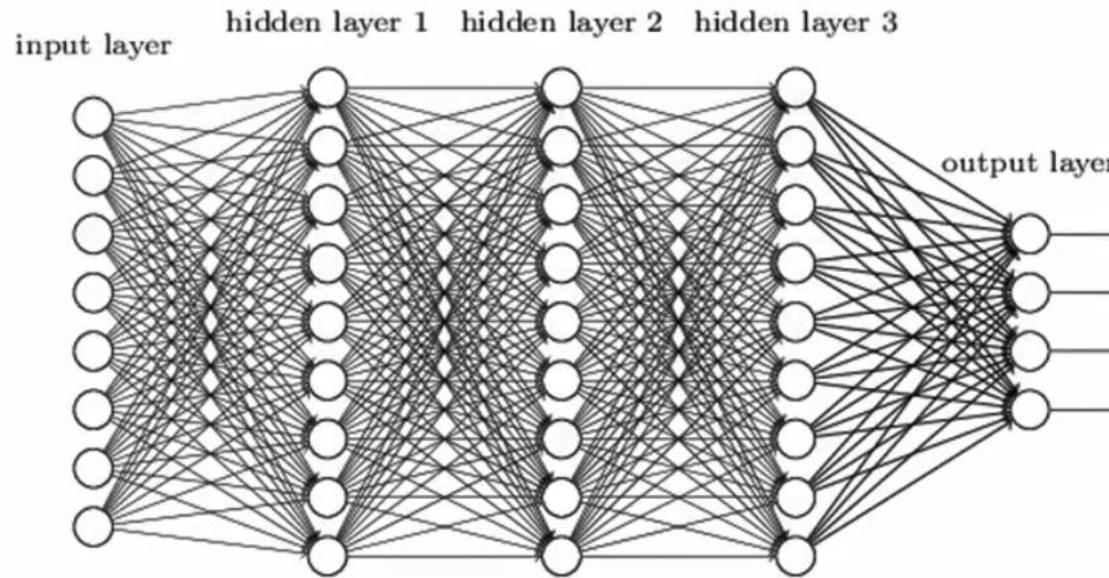


**“2-layer Neural Net”, or  
“1-hidden-layer Neural Net”**

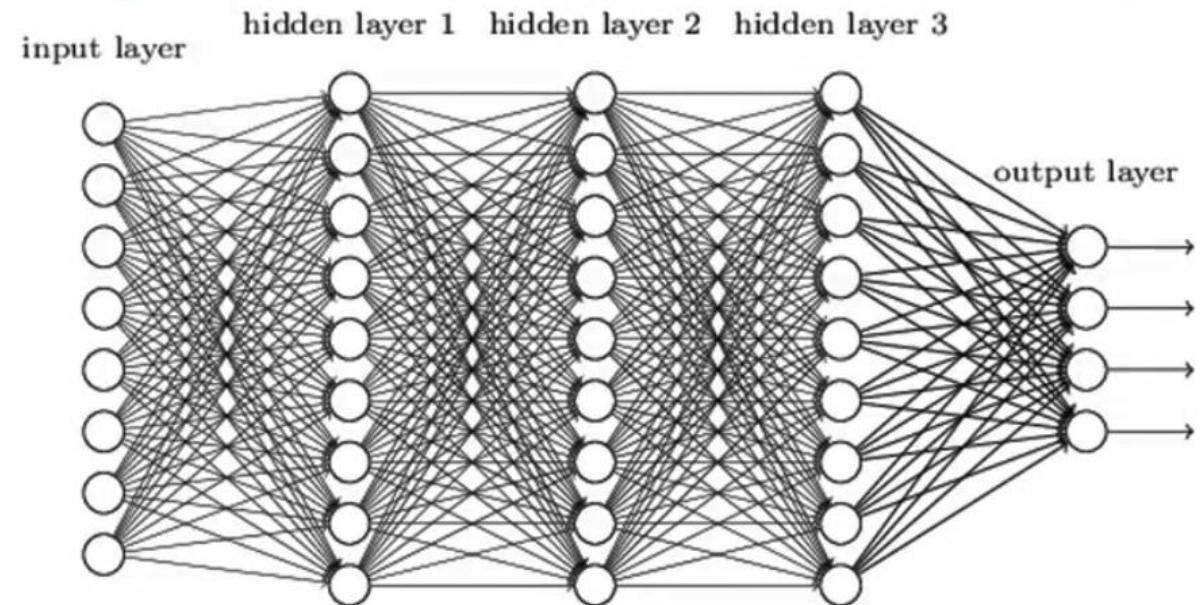


**“3-layer Neural Net”, or  
“2-hidden-layer Neural Net”**

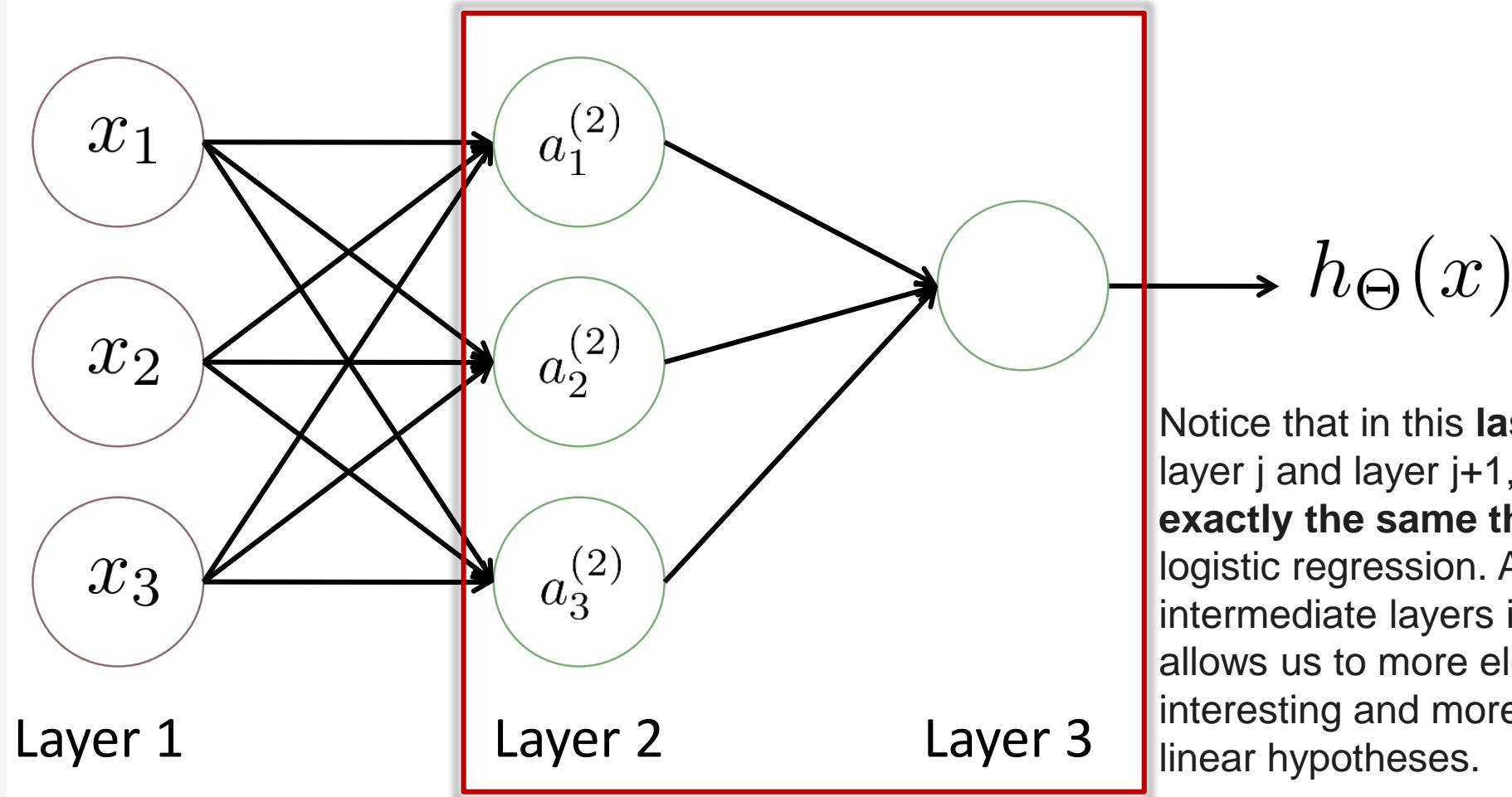
## Forward Propagation



## Backpropagation



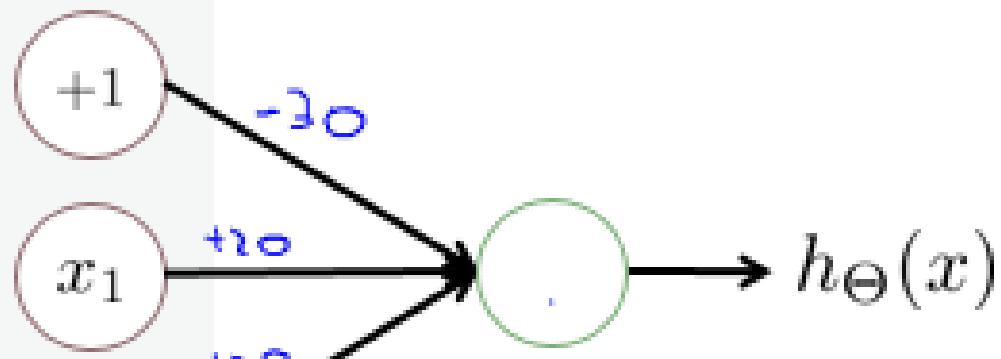
# Neural Network learning its own features



## Simple example: AND

$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2$$

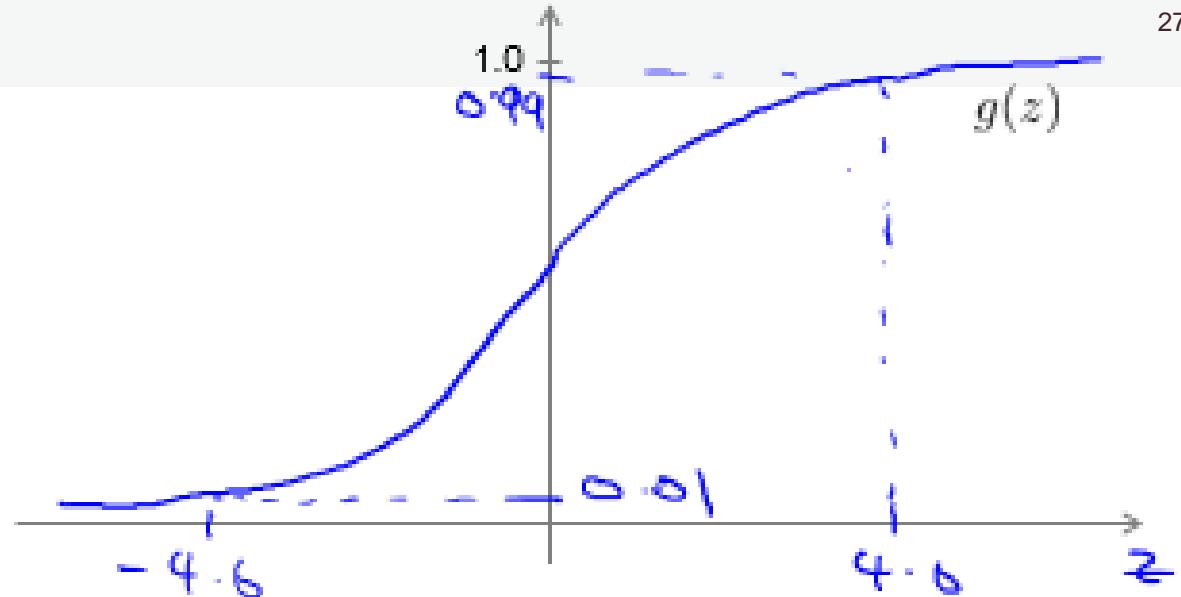


$$h_{\Theta}(x) = g\left(\frac{-30}{\cancel{\pi}} + \frac{10x_1}{\cancel{\pi}} + \frac{10x_2}{\cancel{\pi}}\right)$$

$$\Theta_{1,0}^{(1)}$$

$$\Theta_{1,1}^{(1)}$$

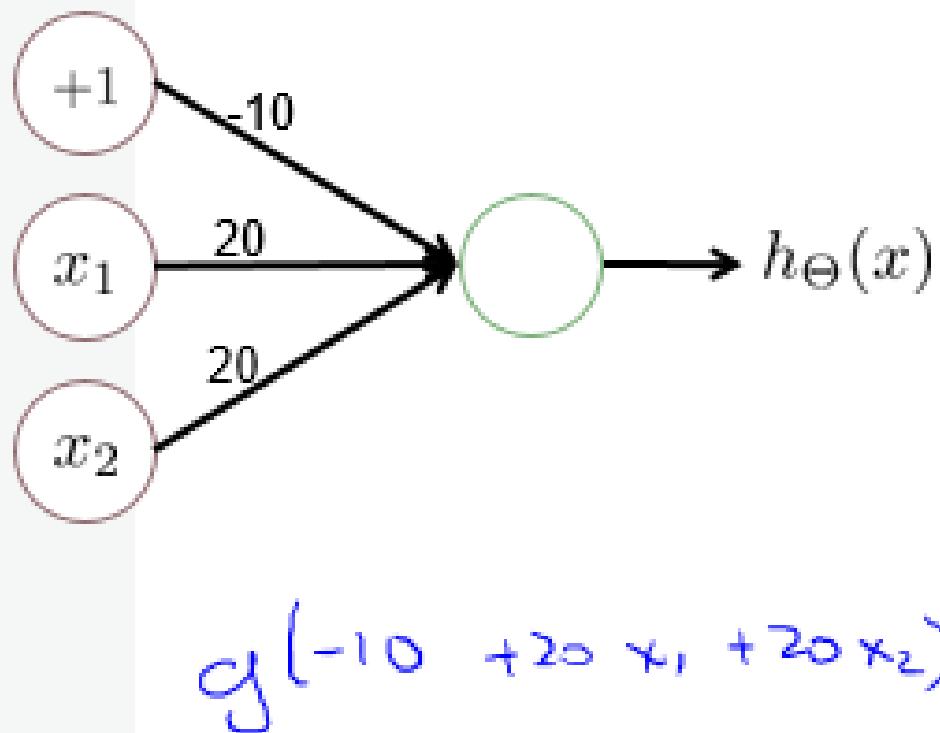
$$\Theta_{1,2}^{(1)}$$



$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$

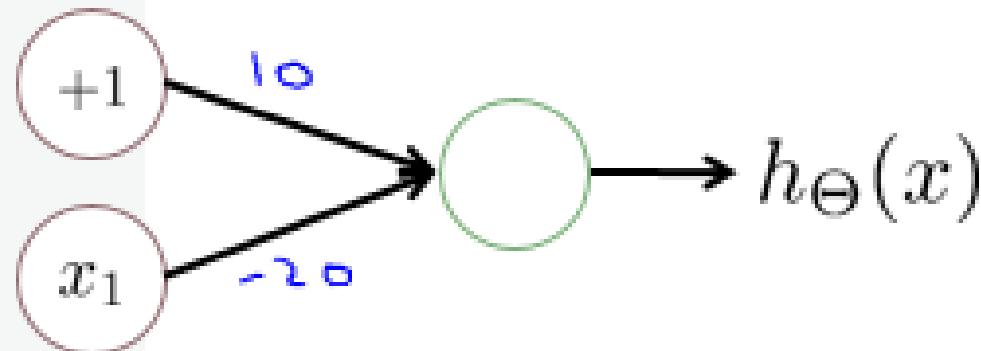
## Example: OR function



$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$\approx 1$
1	1	$\approx 1$

## Negation:

NOT  $x_1$



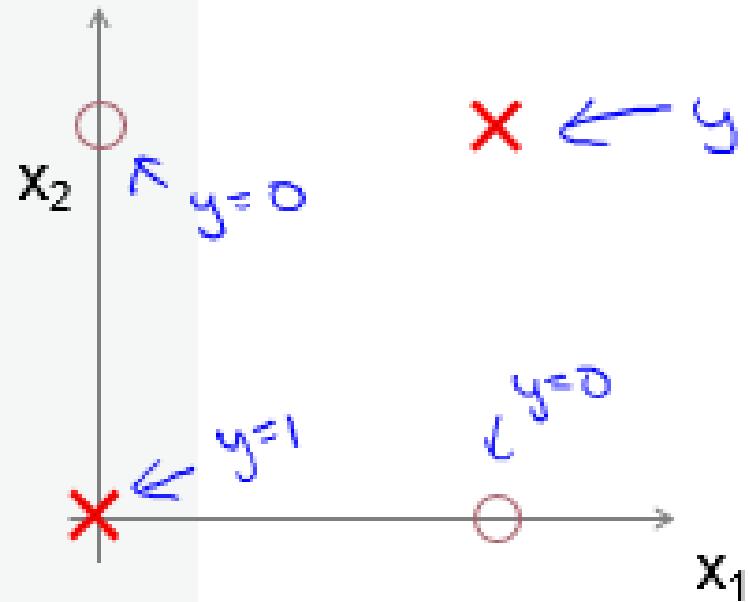
$$h_\Theta(x) = g(10 - 20x_1)$$

$x_1$	$h_\Theta(x)$
0	$g(10) \approx 1$
1	$g(-20) \approx 0$

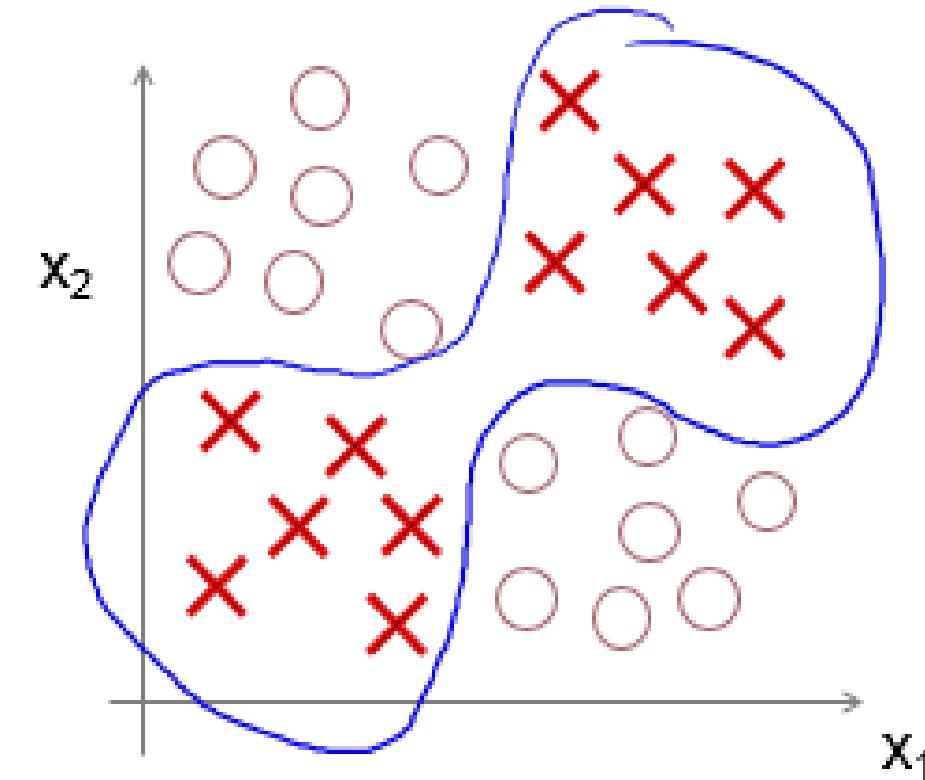
(NOT  $x_1$ ) AND (NOT  $x_2$ ) ?

# Non-linear classification example: XOR/XNOR

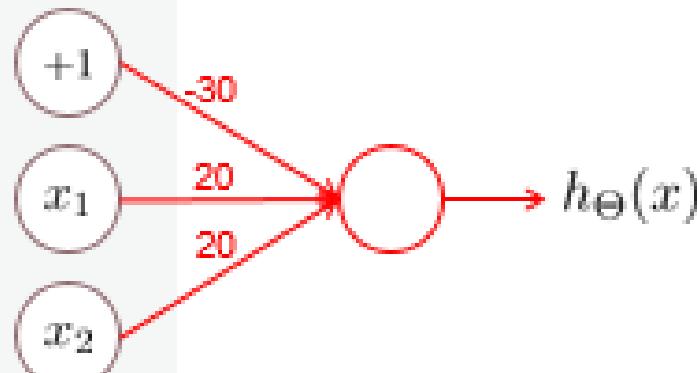
$x_1, x_2$ , are binary (0 or 1).



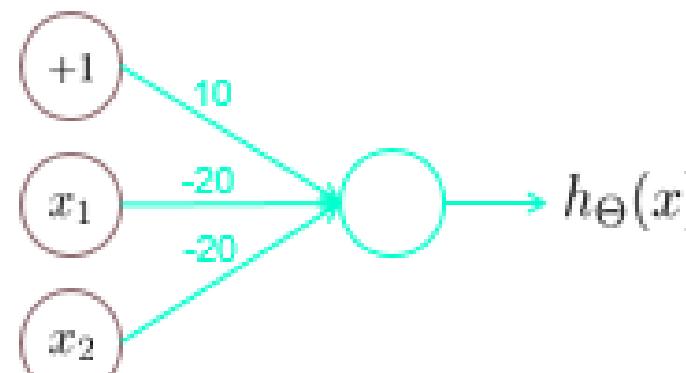
$$y = x_1 \text{ XOR } x_2$$



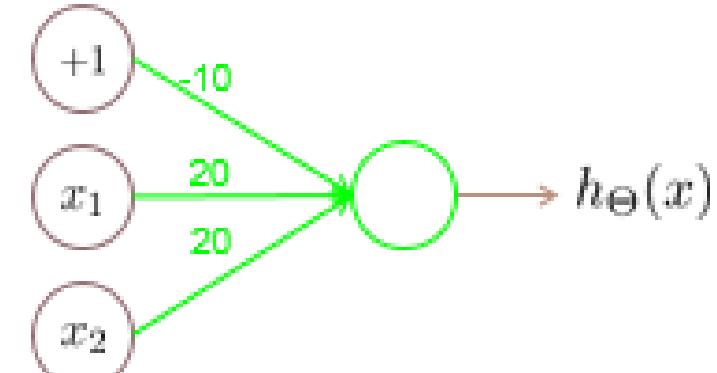
# Putting it together: $x_1 \text{ XNOR } x_2$



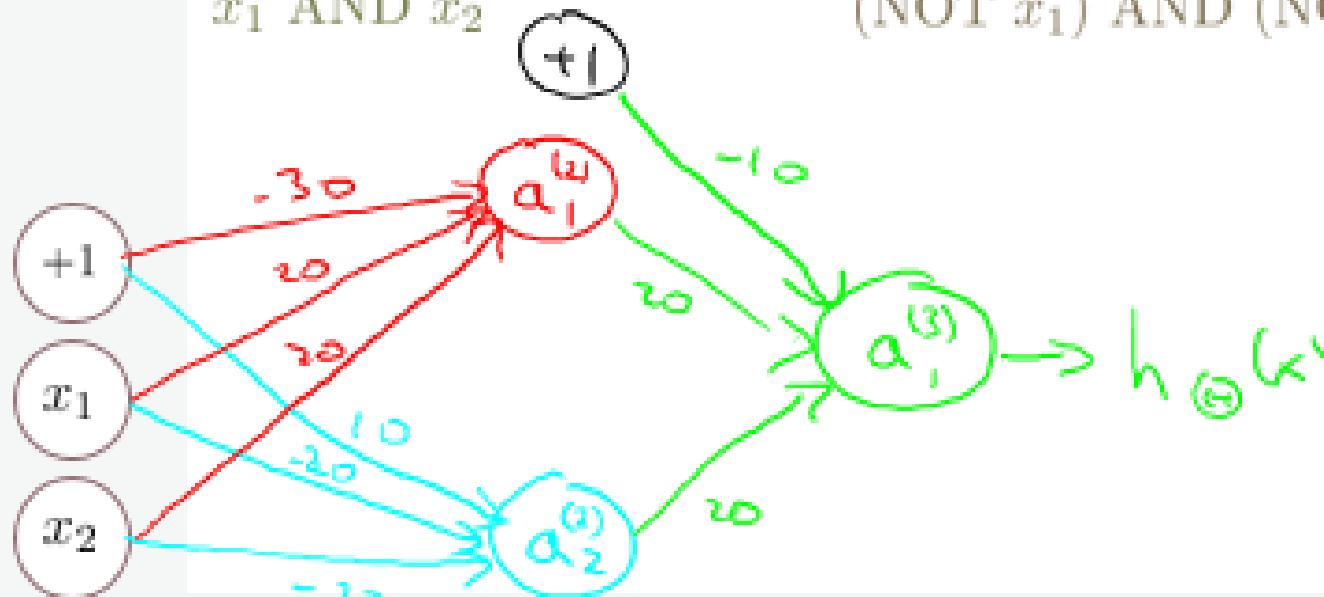
$x_1 \text{ AND } x_2$



(NOT  $x_1$ ) AND (NOT  $x_2$ )



$x_1 \text{ OR } x_2$

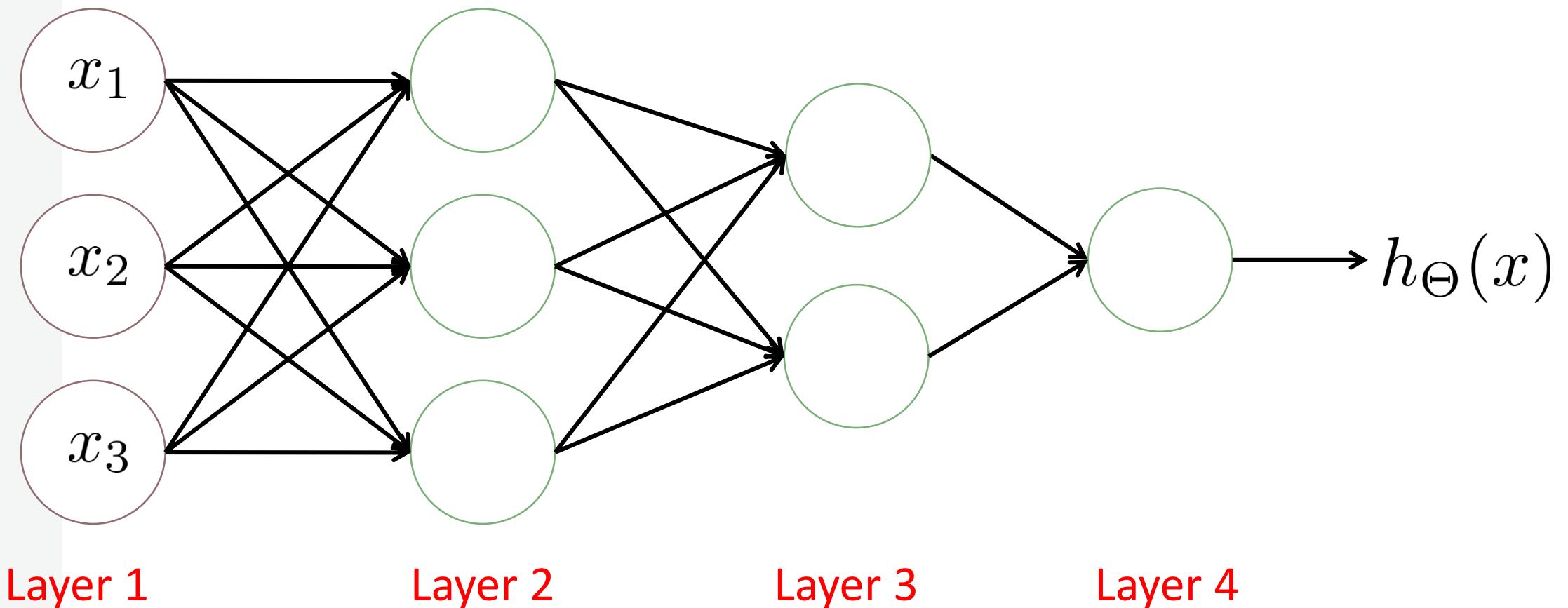


$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

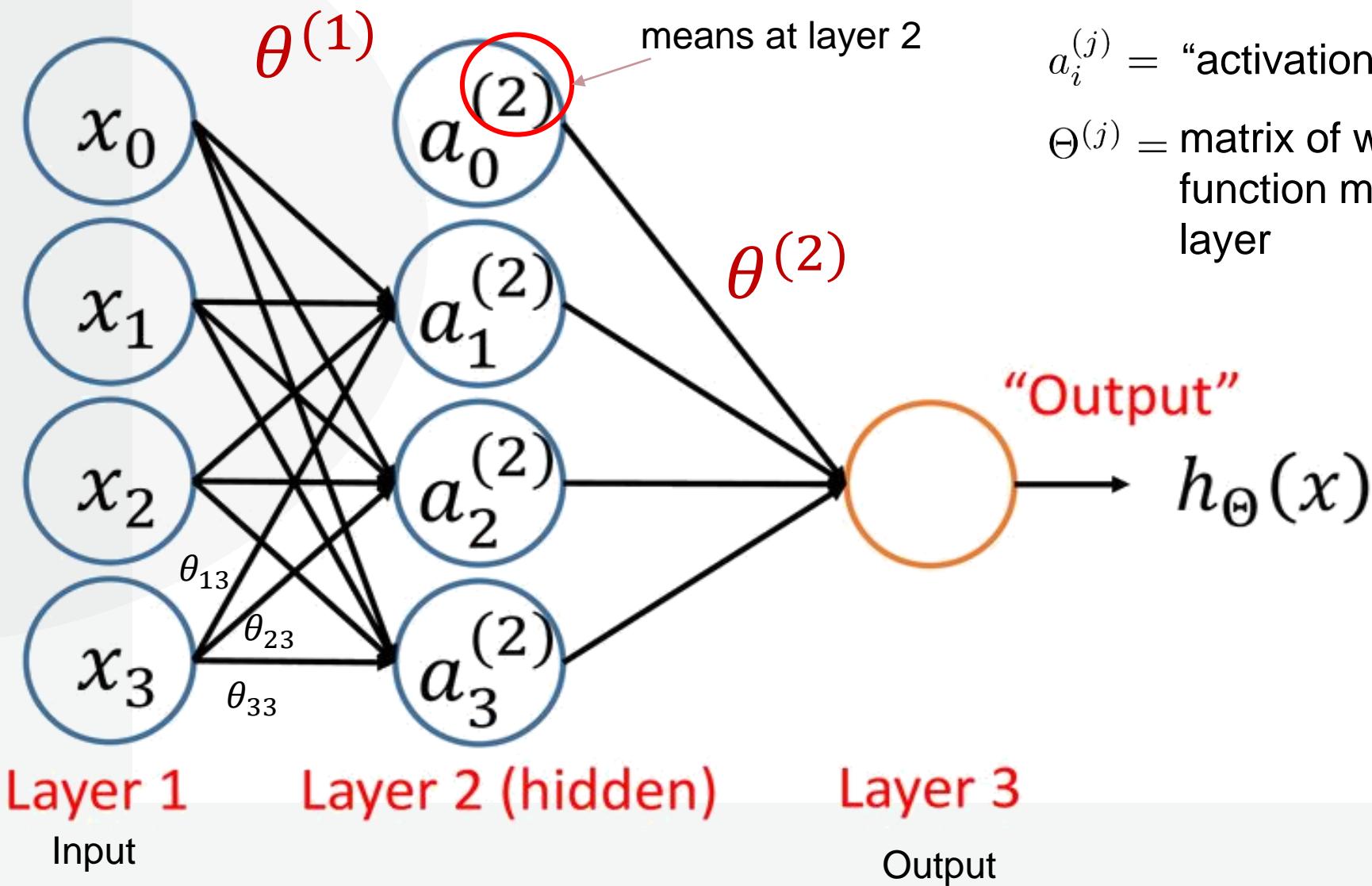
# ANN Notation

Forward and Backward propagation

# Layers Notation



# Neural network – Multilayer



- $a_i^{(j)}$  = “activation” of unit in layer
- $\Theta^{(j)}$  = matrix of weights controlling function mapping from layer to layer

“Output”  
→  $h_{\Theta}(x)$

## Layer 1

## Layer 2 (hidden)

## Layer 3

Slide credit: Andrew Ng

What is the dimension of  $\Theta^{(1)}$ ?

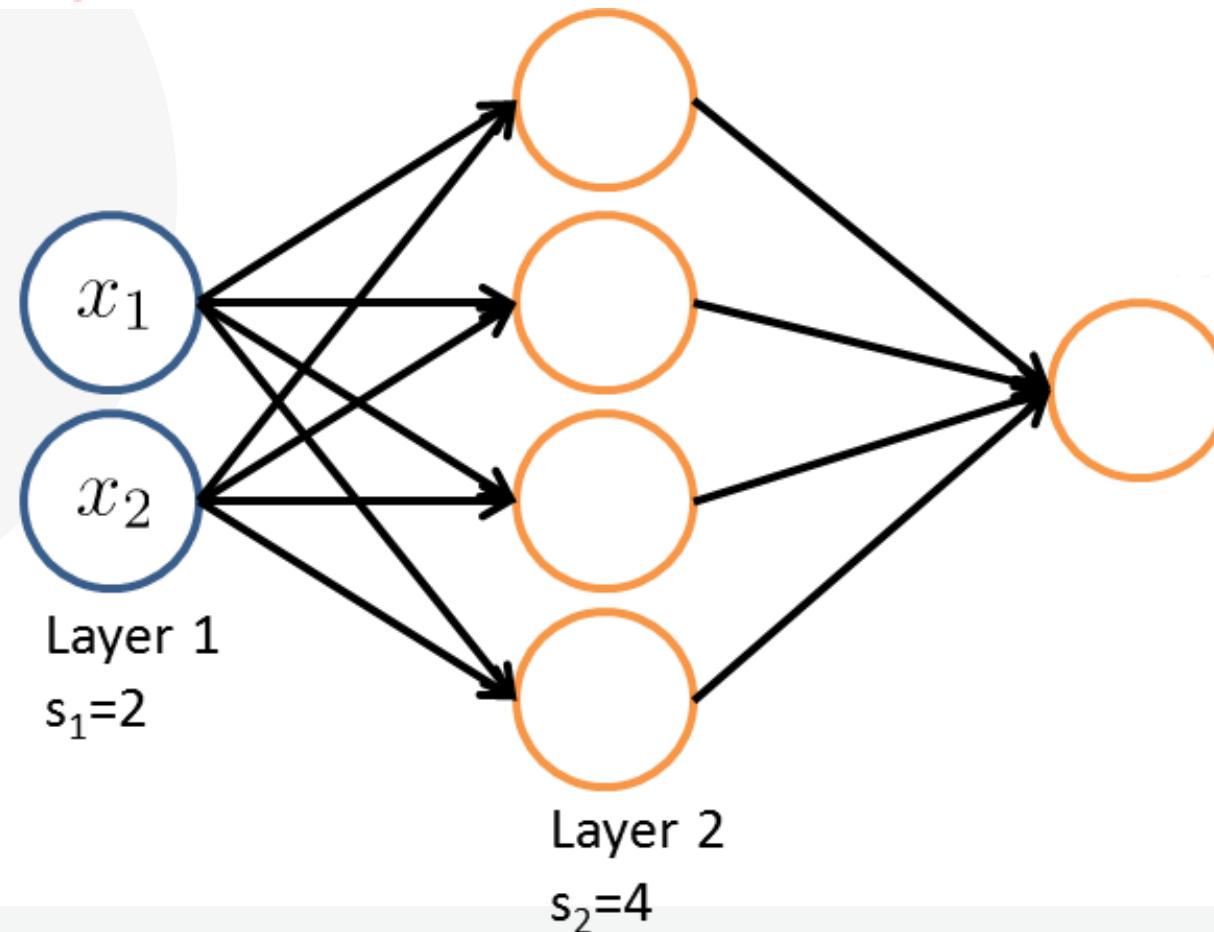
$$s_{j+1} \times (s_j + 1) = 4 \times 3$$

$s_j$  unit in layer  $j$

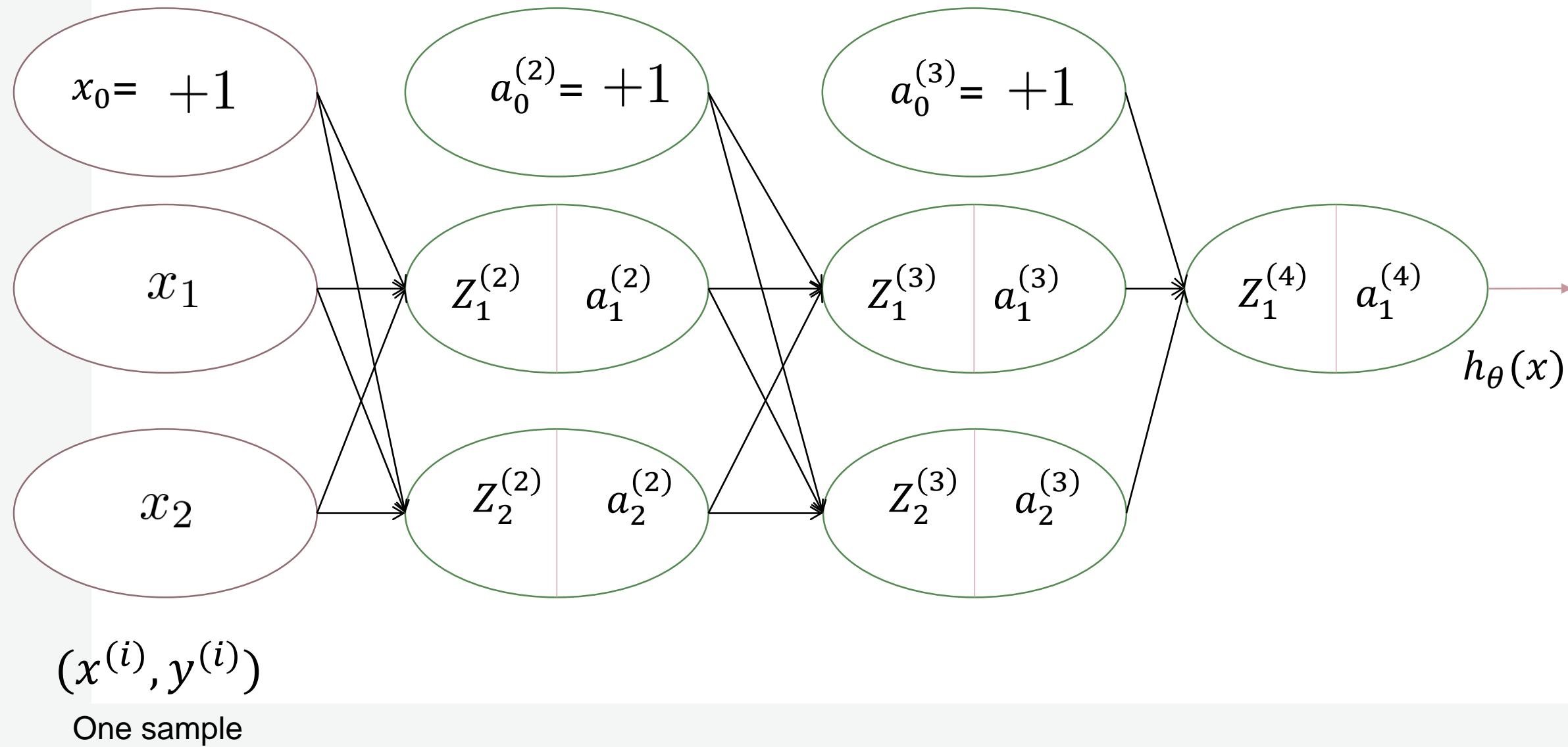
$s_{j+1}$  units in layer  $j + 1$

Size of  $\Theta^{(j)}$ ?

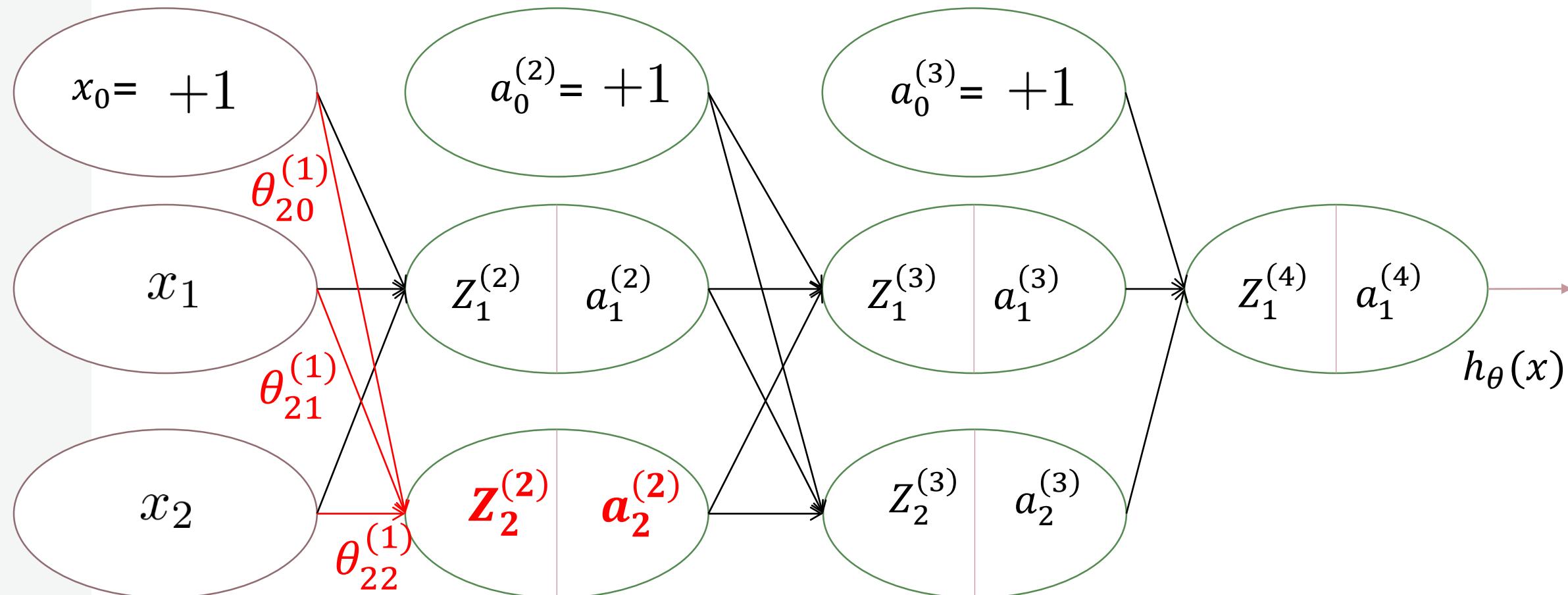
$$s_{j+1} \times (s_j + 1)$$



# Forward Propagation



# Forward Propagation

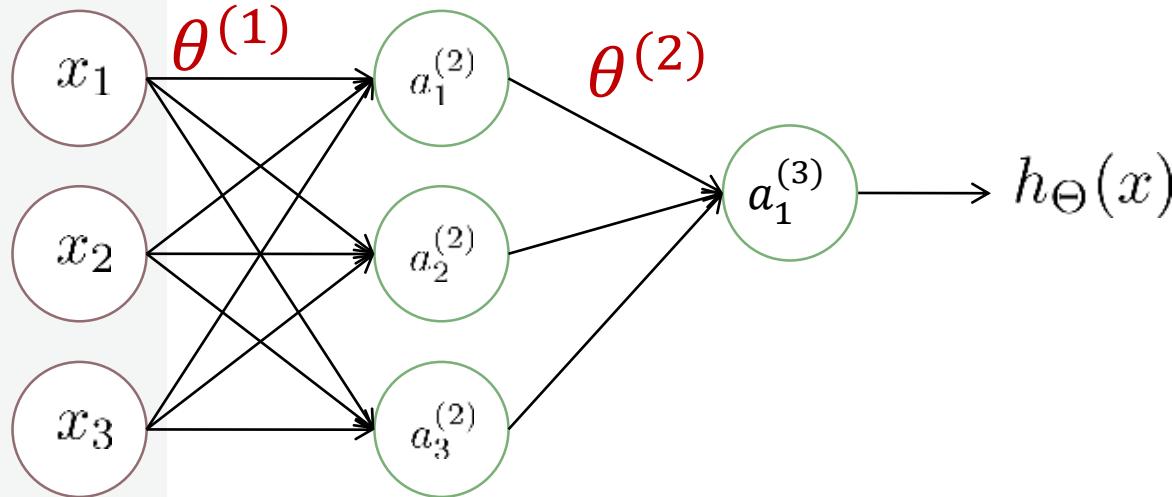


$(x^{(i)}, y^{(i)})$   
One sample

$$Z_2^{(2)} = \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2$$

$$a_2^{(2)} = g(Z_2^{(2)}) = \frac{1}{1 + e^{-Z_2^{(2)}}}$$

# Forward propagation: Vectorized implementation



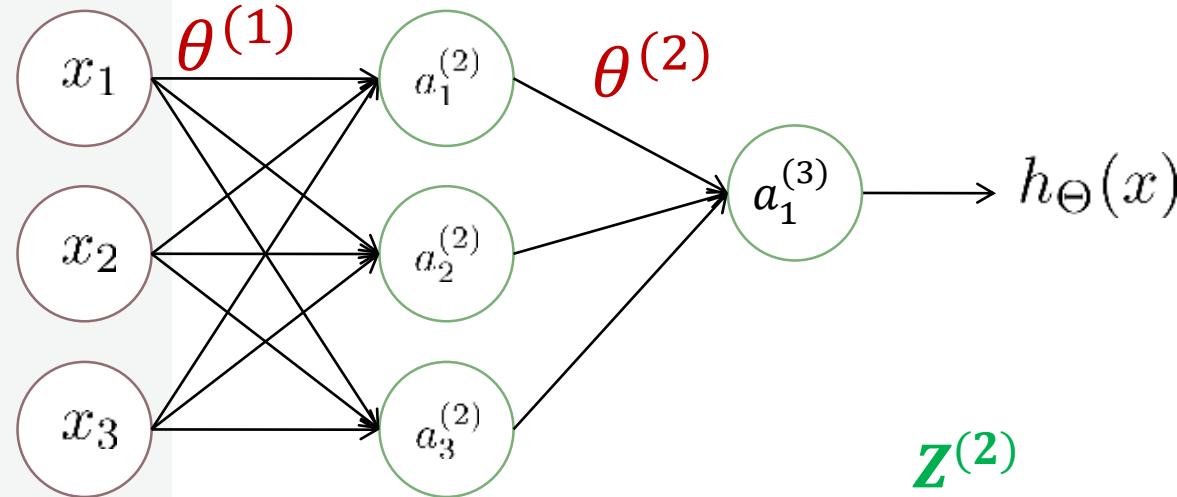
$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$h_{\Theta}(x) = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

# Forward propagation: Vectorized implementation



$$a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$h_{\Theta}(x) = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

$Z^{(3)}$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

## Forward propagation

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

$$z^{(2)} = \Theta^{(1)} x$$

$$a^{(2)} = g(z^{(2)})$$

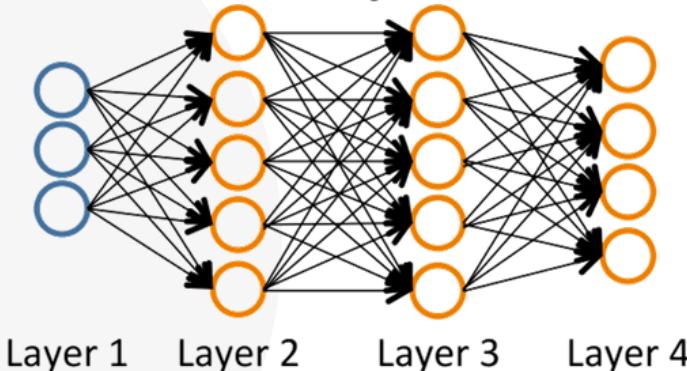
$$\text{Add } a_0^{(2)} = 1$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$$

# Neural Networks: Learning

## Neural Network (Classification)



### Binary classification

$y = 0$  or  $1$

1 output unit

$$K = 2$$

$$s_L = 1 \text{ (output layer)}$$

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

$L$  = total no. of layers in network

$s_l$  = no. of units (not counting bias unit) in layer  $l$

$$s_1 = 3, s_2 = 5, s_3 = 5, s_4 = 4$$

### Multi-class classification (K classes)

$$y \in \mathbb{R}^K \quad \text{E.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

pedestrian   car   motorcycle   truck

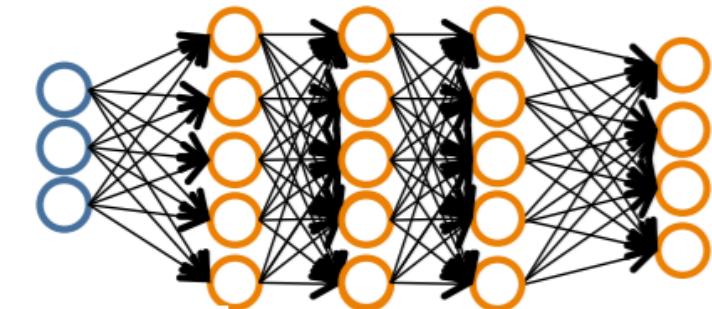
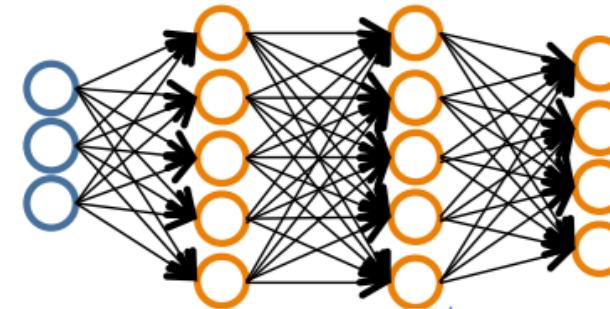
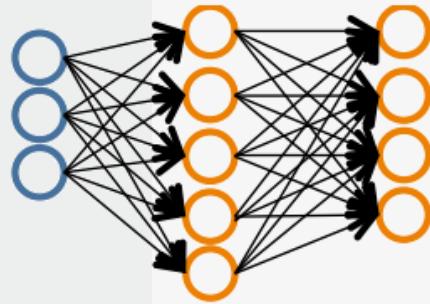
K output units

$$s_L = K$$

$$K \geq 3$$

## Training a neural network

Pick a network architecture (connectivity pattern between neurons)



No. of input units: Dimension of features  $x^{(i)}$

No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

# Gradient computation

$a_i^{(j)}$  = “activation” of unit  $i$  in layer  $j$   
 $\Theta^{(j)}$  = matrix of weights controlling  
 function mapping from layer  $j$  to layer  $j + 1$

Given one training example  $(x, y)$ :

Forward propagation:

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

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

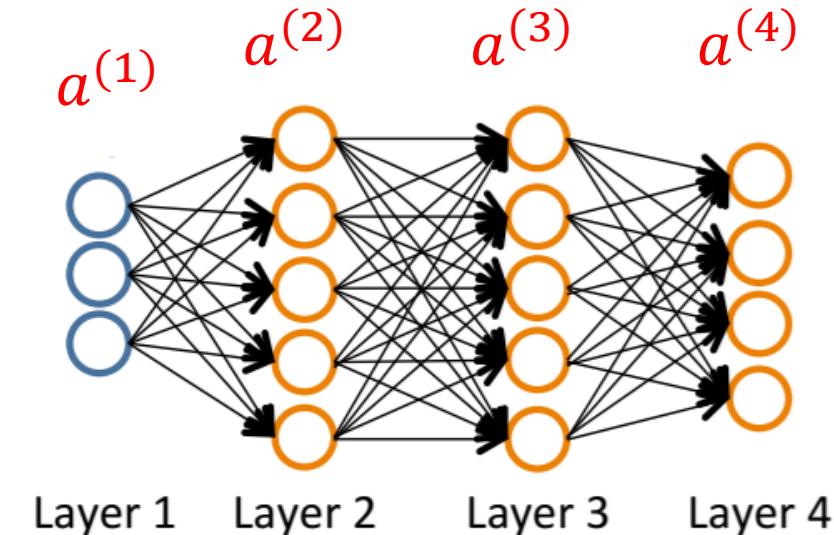
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\Theta}(x) = g(z^{(4)})$$



# Cost function

Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Neural network:

$$h_\Theta(x) \in \mathbb{R}^K \quad (h_\Theta(x))_i = i^{th} \text{ output}$$

$$\begin{aligned} J(\Theta) &= -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k) \right] \\ &\quad + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2 \end{aligned}$$

# Gradient computation

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

$$\min_{\Theta} J(\Theta)$$

Need code to compute:

- $J(\Theta)$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$

# Gradient computation: Backpropagation algorithm

Intuition:  $\delta_j^{(l)}$  = “error” of node  $j$  in layer  $l$ .

Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$  (for  $j \geq 0$ )

**for each output unit** (layer  $L = 4$ )

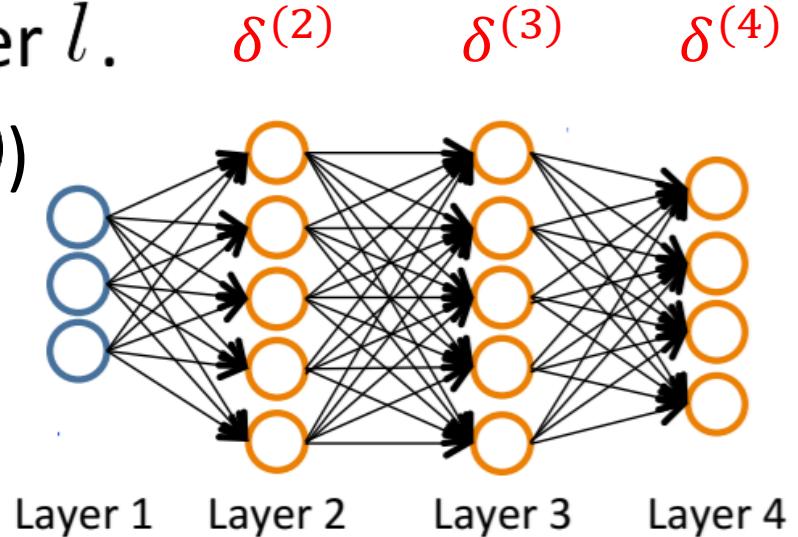
$$\delta_j^{(4)} = a_j^{(4)} - y_j$$

$h_{\theta}(x)_j$

$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} . * g'(z^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} . * g'(z^{(2)})$$

*element wise multiplication*



$$g'(z^{(l)}) = a^{(l)} . * (1 - a^{(l)})$$

**Note:**  
We don't calculate the error of cost for **bias units** and **input units**

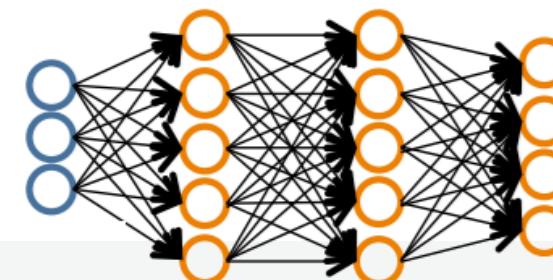
## Training a neural network

1. Randomly initialize weights (**Don't initialize with zero**)
2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
3. Implement code to compute cost function  $J(\Theta)$
4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

**for**  $i = 1:m$

    Perform forward propagation and backpropagation using  
    example  $(x^{(i)}, y^{(i)})$

    (Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ ).



## Training a neural network

5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$  computed using backpropagation vs. using numerical estimate of gradient of  $J(\Theta)$ .  
Then disable gradient checking code.
6. Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\Theta)$  as a function of parameters  $\Theta$

## In next Example

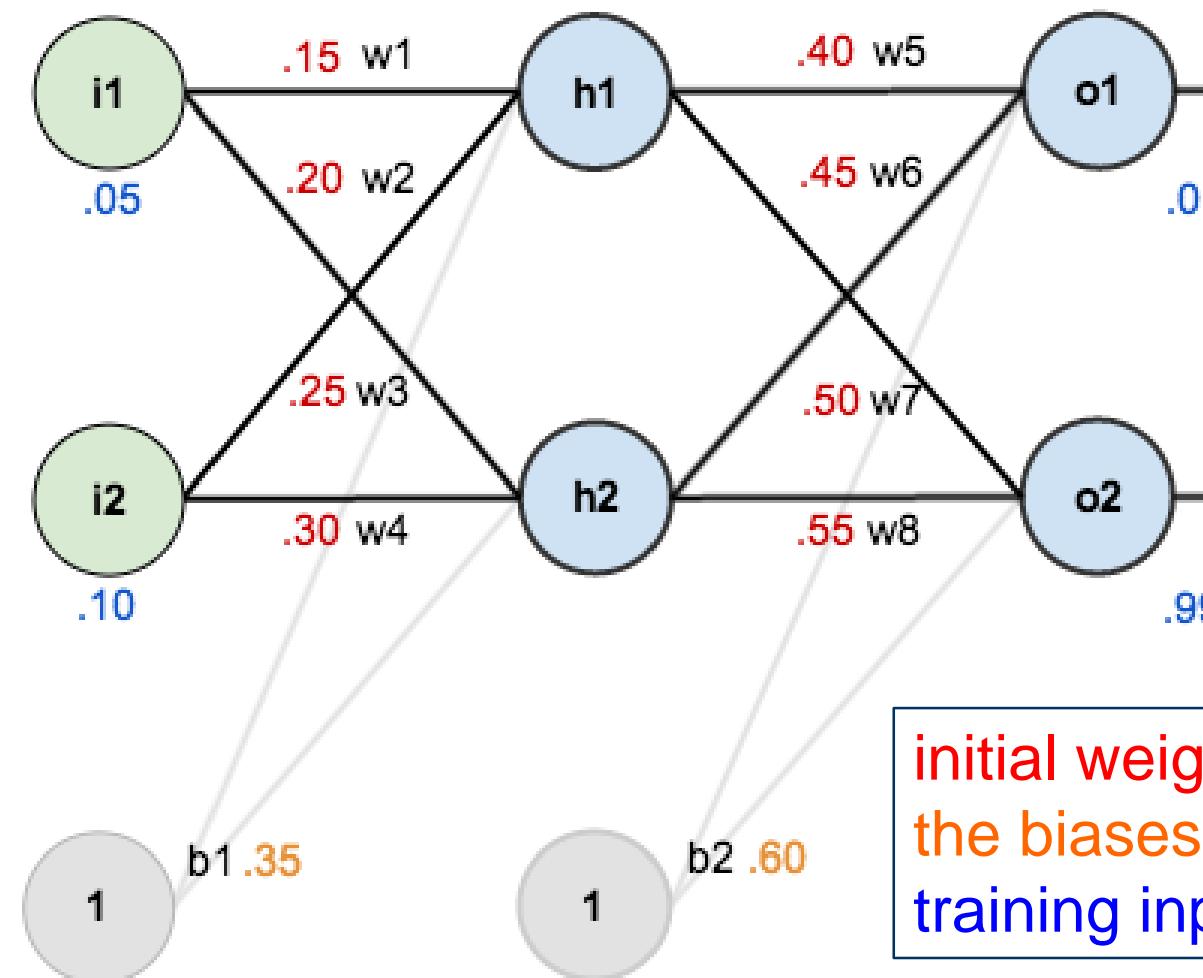
$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(h_\Theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - (h_\Theta(x^{(i)}))) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

Focusing on a **single example**  $x^{(i)}$ ,  $y^{(i)}$ , and ignoring regularization ( $\lambda = 0$ ),

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$$

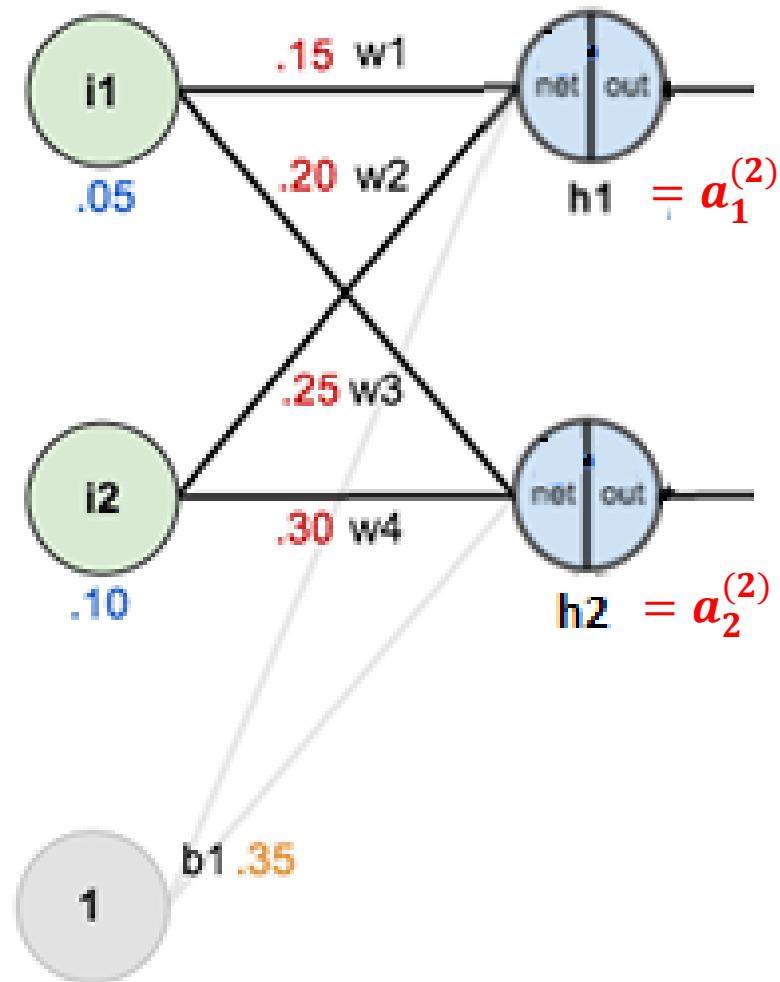
(Think of  $\text{cost}(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2$ )

# Numerical Example



- <https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>

# Numerical Example: Forward Pass-<sup>hidden layer</sup>



- Calculate total net input for  $h_1$  ( $\text{net}_{h_1}$ )

$$\text{net}_{h_1} = w_1 * i_1 + w_2 * i_2 + b_1 * 1 \quad Z_1^{(2)}$$

$$\text{net}_{h_1} = 0.15 * 0.05 + 0.2 * 0.1 + 0.35 * 1 = 0.3775$$

- Apply the sigmoid function to get the output of  $h_1$  ( $\text{out}_{h_1}$ )

$$a_1^{(2)} = g(Z_1^{(2)})$$

$$\text{out}_{h_1} = \frac{1}{1+e^{-\text{net}_{h_1}}} = \frac{1}{1+e^{-0.3775}} = 0.593269992$$

- Carrying out the same process for  $h_2$ , we get:

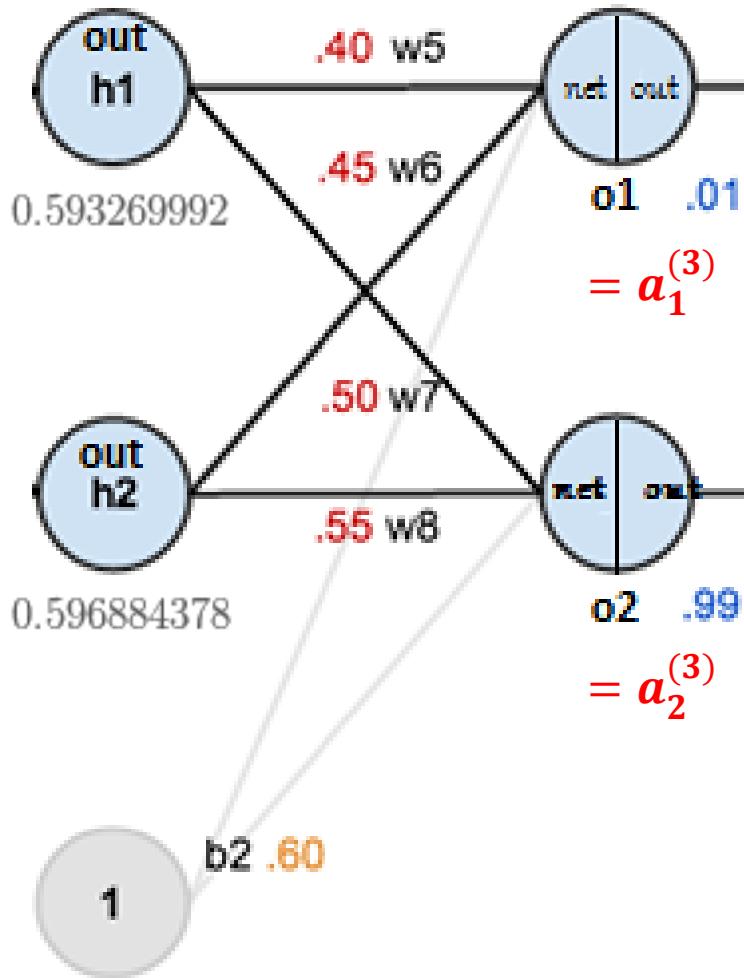
$$\text{out}_{h_2} = 0.596884378$$

$$a_2^{(2)}$$

# Numerical Example: Forward Pass-

output layer

53



Repeat same process for output layer:  
Calculate total net input for  $o_1$  ( $net_{o1}$ )

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\begin{aligned} net_{o1} &= 0.4 * 0.593269992 + 0.45 * 0.596884378 + 0.6 * 1 \\ &= 1.105905967 \end{aligned}$$

Apply the sigmoid function to get the output of  $o_1$  ( $out_{o1}$ )

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}} = \frac{1}{1+e^{-1.105905967}} = 0.75136507 \quad a_1^{(3)} = g(z_1^{(3)})$$

Carrying out the same process for  $o_2$ , we get:

$$out_{o2} = 0.772928465 \quad a_2^{(3)}$$

# Numerical Example: Calculate Total Error<sup>54</sup>

- Calculate the error for each output neuron using the **squared error function** and sum them to get the total error:

$$E_{total} = \sum \frac{1}{2}(target - output)^2$$

Used for an approximate cost for one sample

- For example, the target output for o1 is 0.01 but the neural network output 0.75136507, therefore its error is:

$$E_{o1} = \frac{1}{2}(target_{o1} - output_{o1})^2 = \frac{1}{2}(0.01 - 0.75136507)^2 = 0.274811083$$

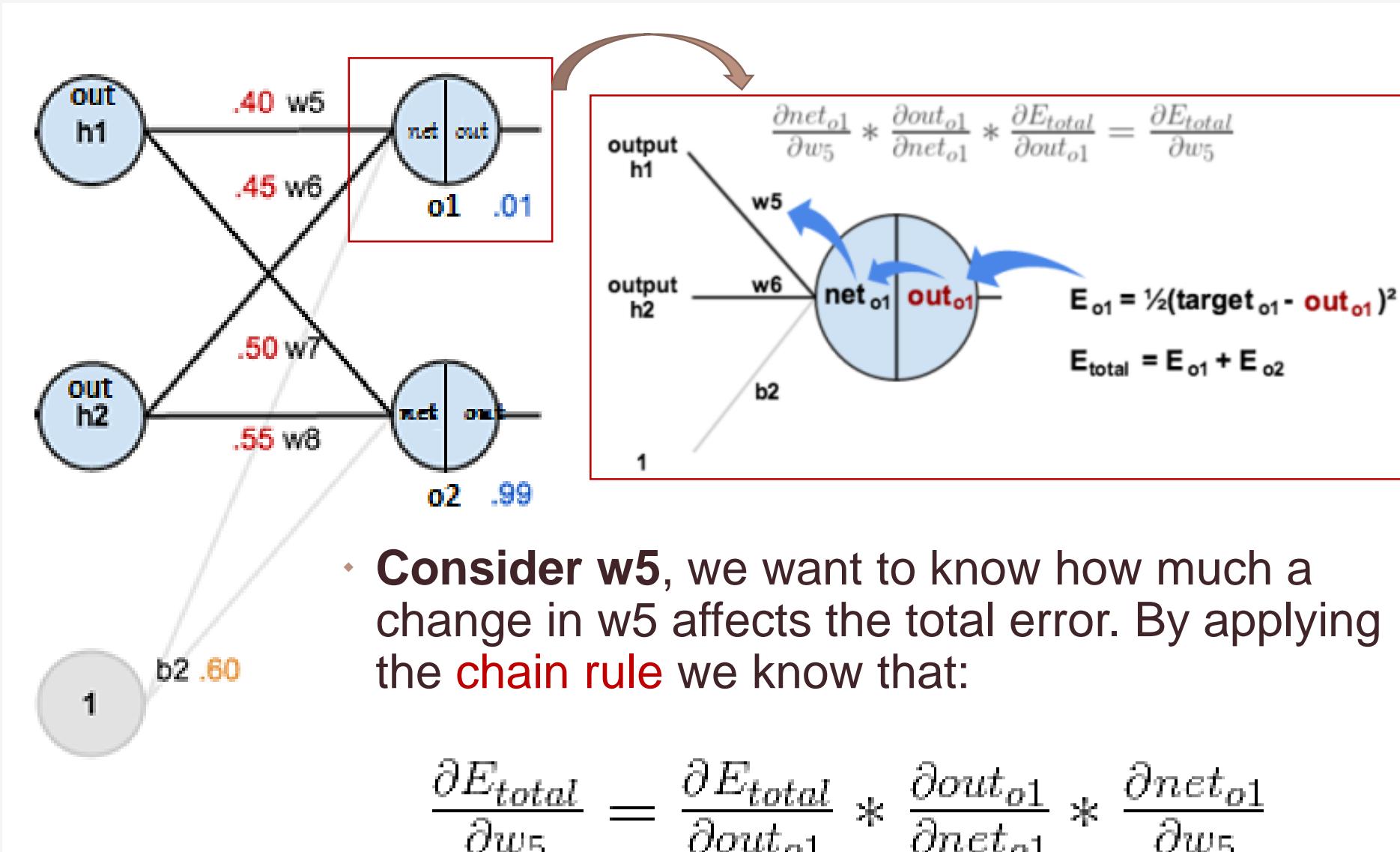
- Repeating this process for o2 (remembering that the target is 0.99) we get:

$$E_{o2} = 0.023560026$$

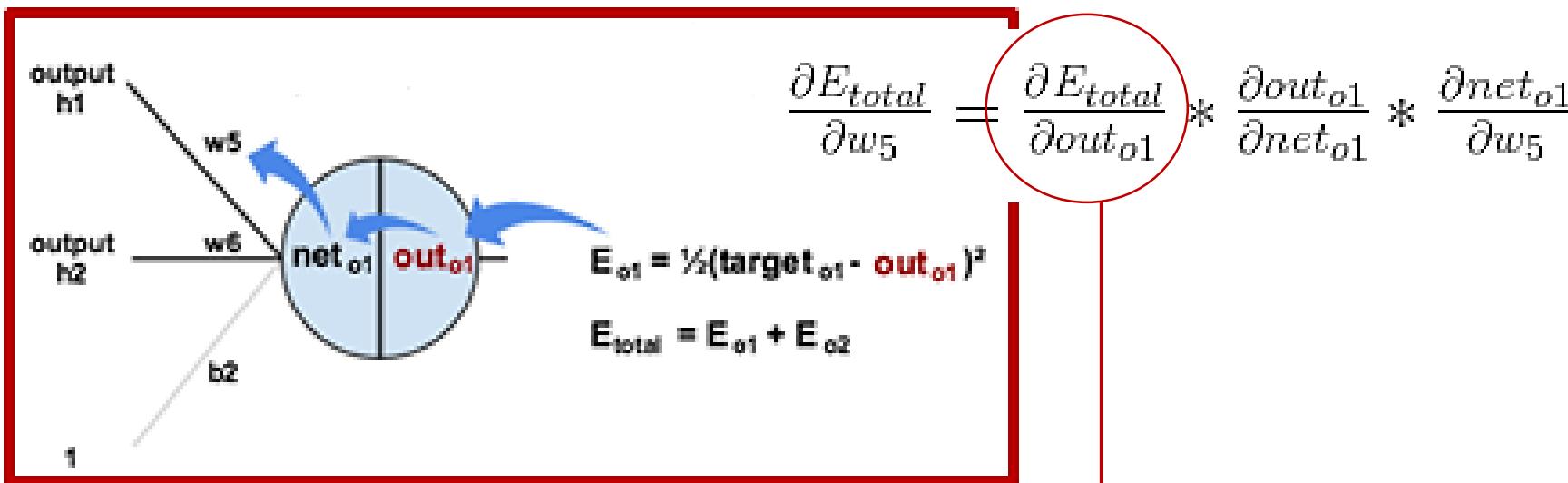
- The total error for the neural network is the sum of these errors:

$$E_{total} = E_{o1} + E_{o2} = 0.274811083 + 0.023560026 = 0.298371109$$

# Numerical Example: Backward Pass-<sub>output layer</sub><sup>55</sup>



# Numerical Example: Backward Pass-<sub>output layer</sub><sup>56</sup>



- 1- The total error change with respect to the output  $o_1$

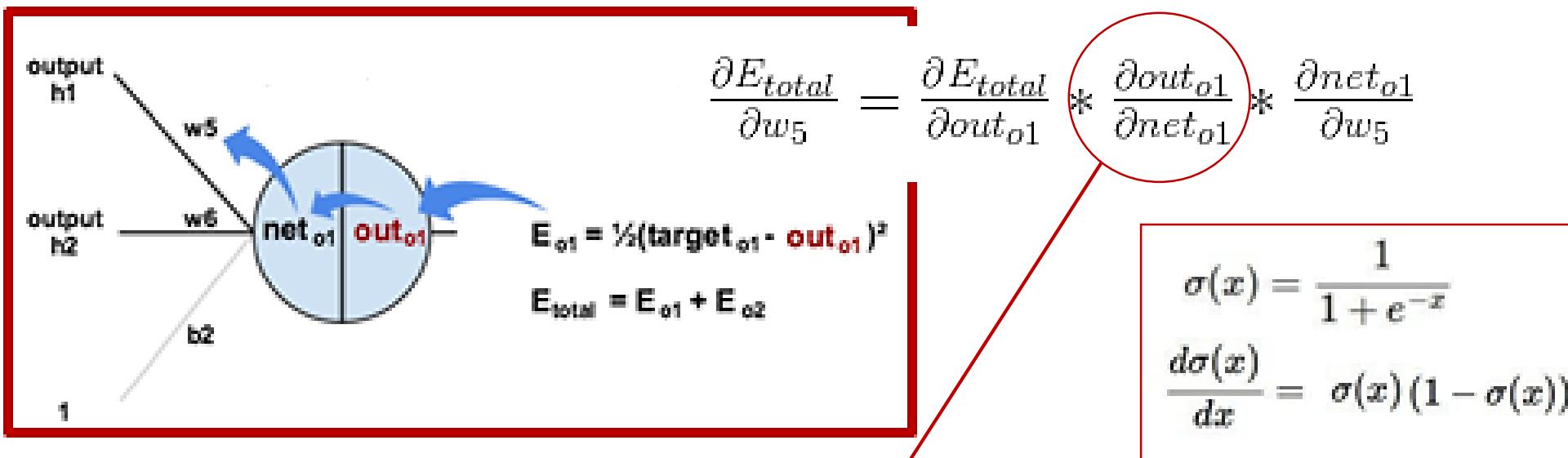
$$E_{total} = \frac{1}{2}(target_{o1} - out_{o1})^2 + \frac{1}{2}(target_{o2} - out_{o2})^2$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = 2 * \frac{1}{2}(target_{o1} - out_{o1})^{2-1} * -1 + 0$$

$$\frac{\partial E_{total}}{\partial out_{o1}} = -(target_{o1} - out_{o1}) = -(0.01 - 0.75136507) = 0.74136507$$

$-(target - out)$  is sometimes expressed as  $out - target$

# Numerical Example: Backward Pass-<sub>output layer</sub><sup>57</sup>

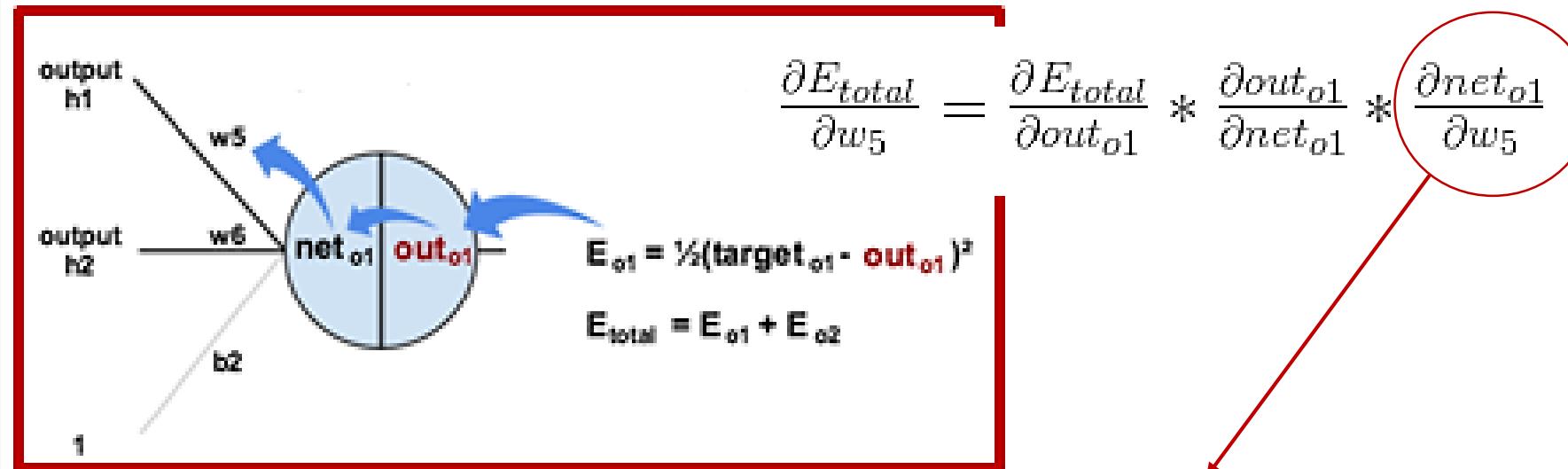


- 2- The output  $o1$  change with respect to its total net input

$$out_{o1} = \frac{1}{1+e^{-net_{o1}}}$$

$$\frac{\partial out_{o1}}{\partial net_{o1}} = out_{o1}(1 - out_{o1}) = 0.75136507(1 - 0.75136507) = 0.186815602$$

# Numerical Example: Backward Pass-<sub>output layer</sub><sup>58</sup>

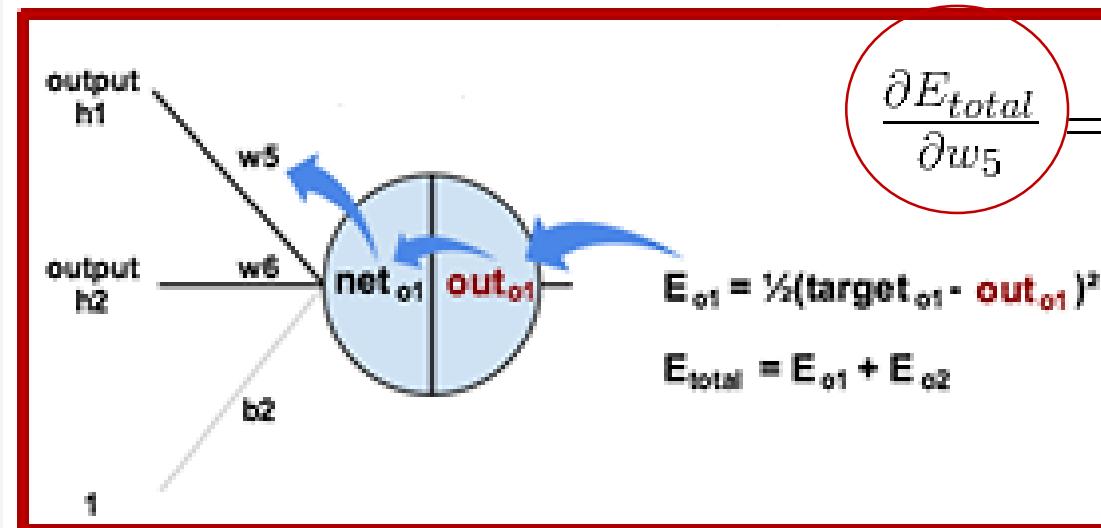


- 3- The total net input of o1 change with respect to w5

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial w_5} = 1 * out_{h1} * w_5^{(1-1)} + 0 + 0 = out_{h1} = 0.593269992$$

# Numerical Example: Backward Pass-<sub>output layer</sub><sup>59</sup>



$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

- 4- Putting all together

$$\frac{\partial E_{total}}{\partial w_5} = \frac{\partial E_{total}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial w_5}$$

$$\frac{\partial E_{total}}{\partial w_5} = 0.74136507 * 0.186815602 * 0.593269992 = 0.082167041$$

# Numerical Example: Backward Pass-<sub>output layer</sub>

- 5- To decrease the error, we then subtract this value from the current weight (optionally multiplied by some **learning rate**, eta, which we'll set to **0.5**):

$$w_5^+ = w_5 - \eta * \frac{\partial E_{total}}{\partial w_5} = 0.4 - 0.5 * 0.082167041 = 0.35891648$$

- We can repeat this process to all weights of output layer to get the new weights w6, w7 and w8:

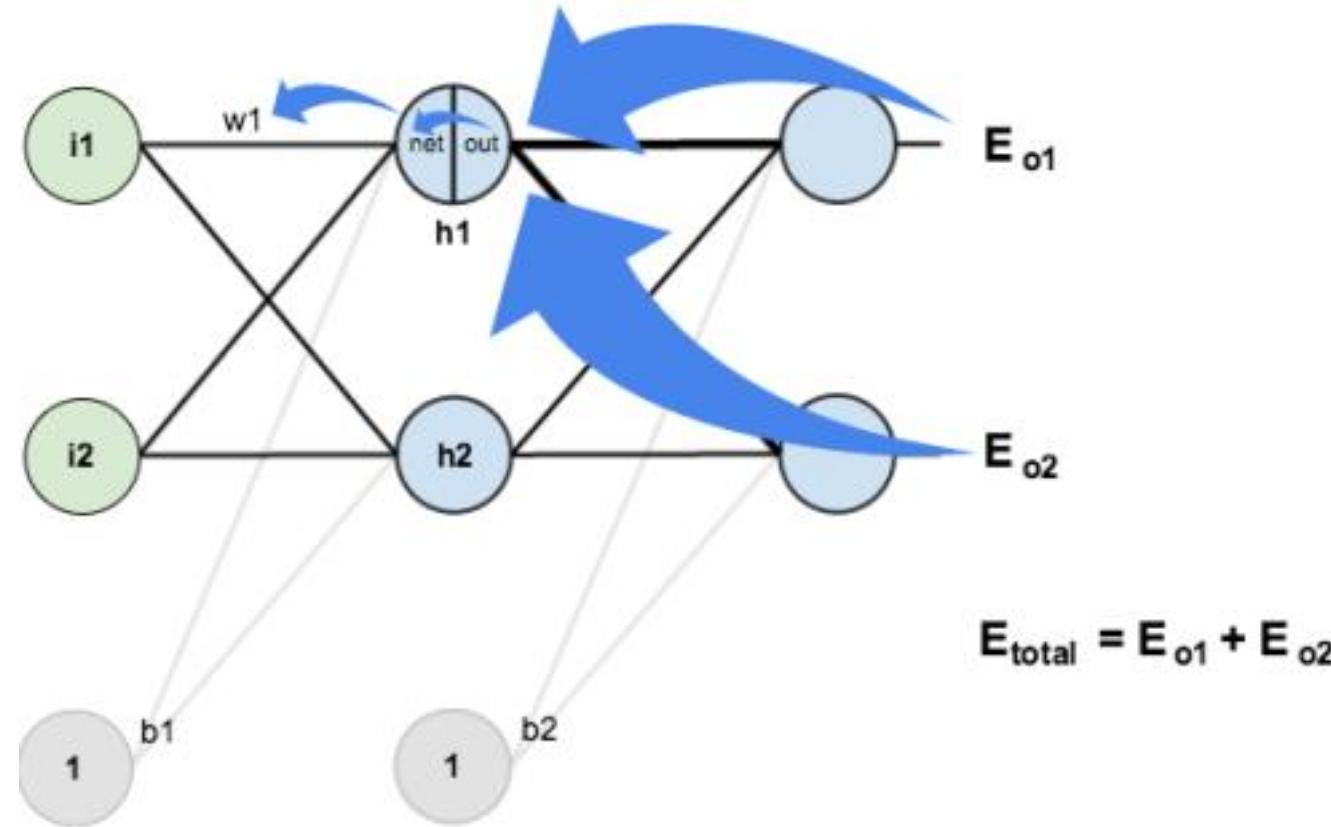
$$w_6^+ = 0.408666186$$

$$w_7^+ = 0.511301270$$

$$w_8^+ = 0.561370121$$

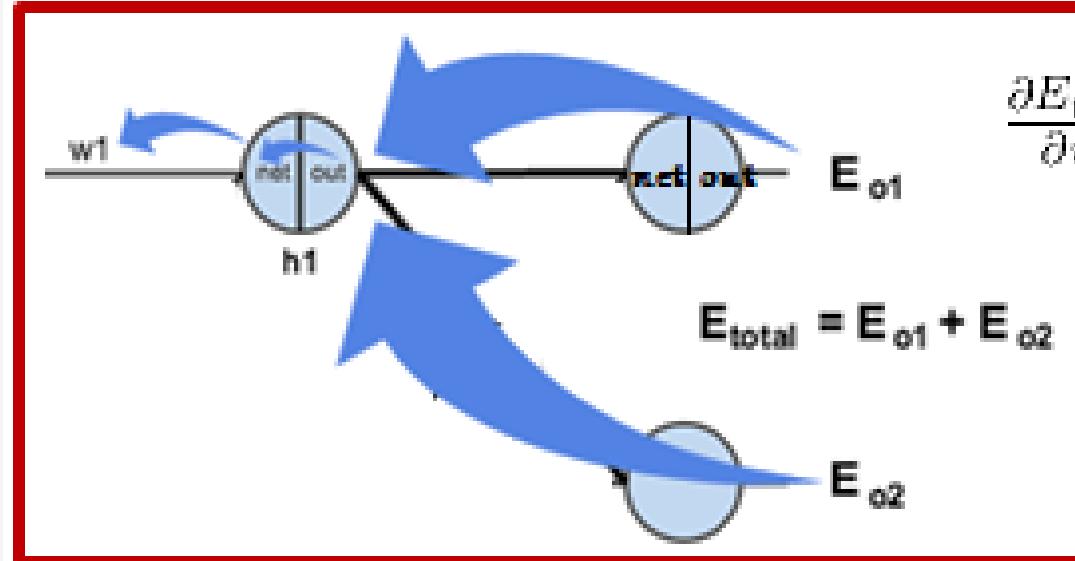
# Numerical Example: Backward Pass-<sub>hidden layer</sub><sup>61</sup>

- we'll continue the backwards pass by calculating new values for w1, w2, w3 and w4



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1} \text{ where, } \frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

# Numerical Example: Backward Pass - hidden layer<sup>62</sup>



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

Where,

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

- 1- Starting with  $\frac{\partial E_{o1}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{o1}} * \frac{\partial out_{o1}}{\partial net_{o1}} * \frac{\partial net_{o1}}{\partial out_{h1}}$

$$= 0.74136507 * 0.186815602$$

using values calculated earlier

$$= 0.138498562$$

$$net_{o1} = w_5 * out_{h1} + w_6 * out_{h2} + b_2 * 1$$

$$\frac{\partial net_{o1}}{\partial out_{h1}} = w_5 = 0.40$$

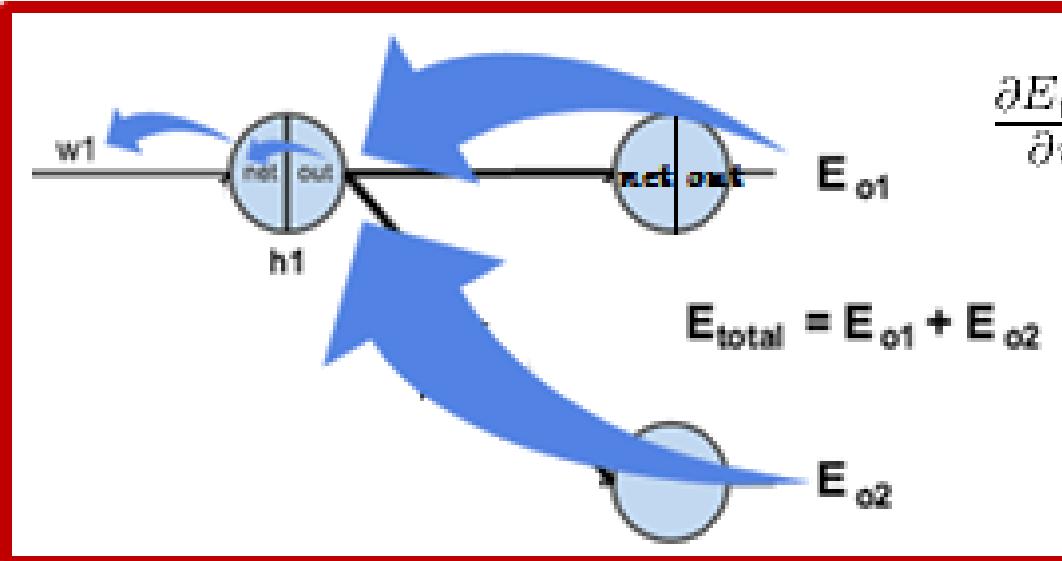
Not updated weight

$$\frac{\partial E_{o1}}{\partial out_{h1}} = 0.138498562 * 0.40 = 0.055399425$$

# Numerical Example: Backward Pass-

63

hidden layer



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

Where,

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

Following the same process for  $\frac{\partial E_{o2}}{\partial out_{h1}}$ , we get:

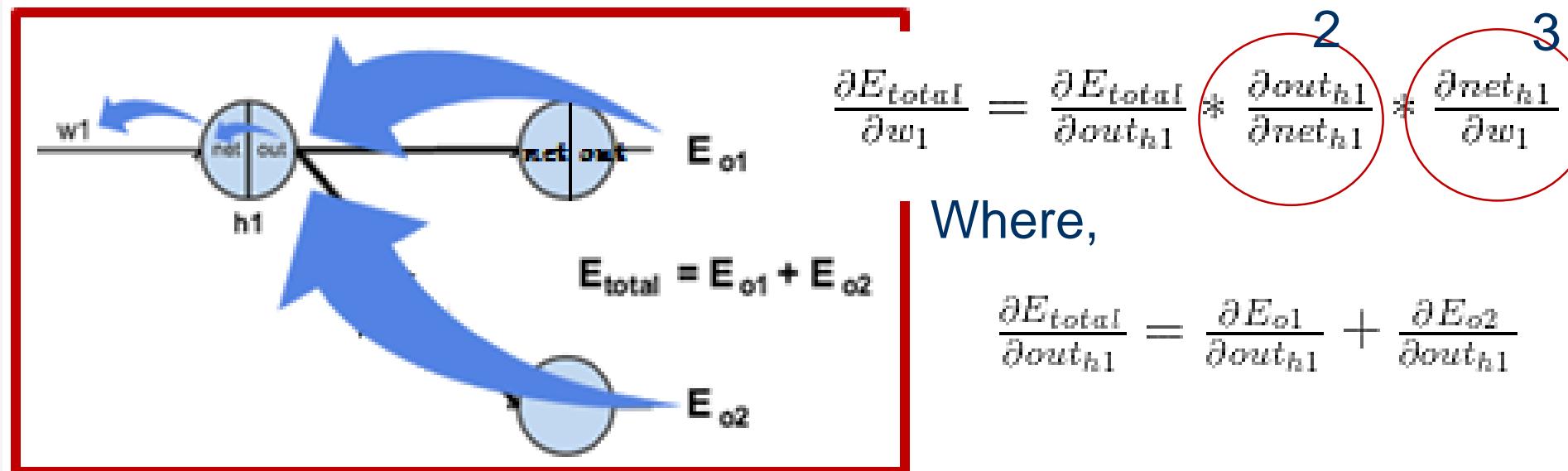
$$\frac{\partial E_{o2}}{\partial out_{h1}} = -0.019049119$$

Therefore:

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}} = 0.055399425 + -0.019049119 = 0.036350306$$

# Numerical Example: Backward Pass-

64



$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

Where,

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

- 2-  $out_{h1} = \frac{1}{1+e^{-net_{h1}}}$

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$

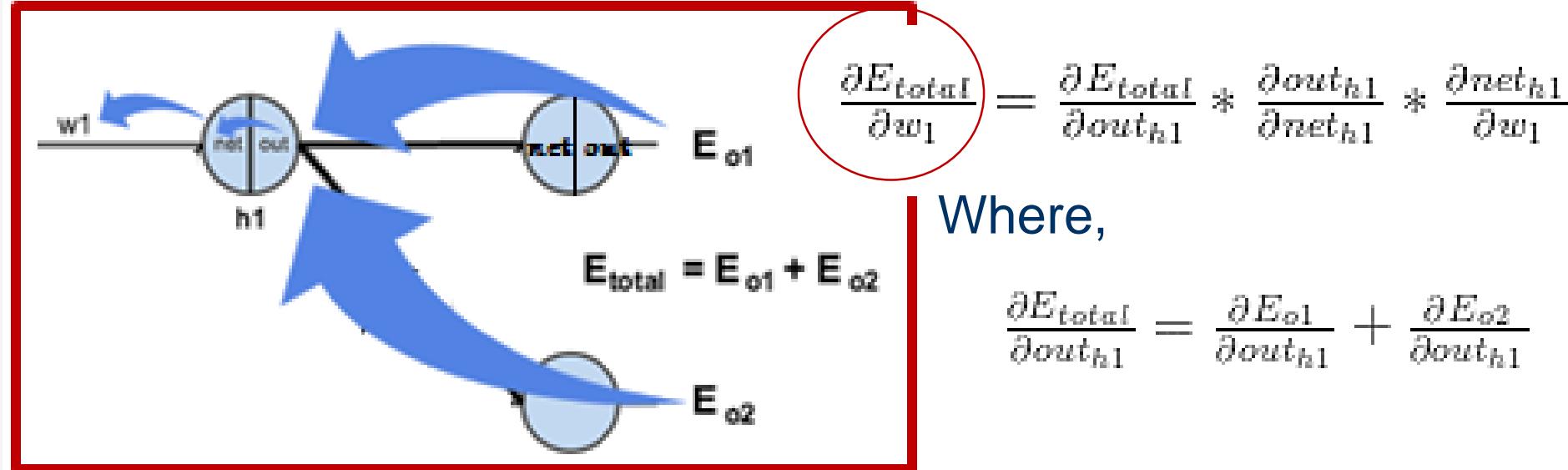
$$\frac{\partial out_{h1}}{\partial net_{h1}} = out_{h1}(1 - out_{h1}) = 0.59326999(1 - 0.59326999) = 0.241300709$$

- 3-  $net_{h1} = w_1 * i_1 + w_3 * i_2 + b_1 * 1$

$$\frac{\partial net_{h1}}{\partial w_1} = i_1 = 0.05$$

# Numerical Example: Backward Pass-

hidden layer



Where,

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$

- 4- Put all together

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

$$\frac{\partial E_{total}}{\partial w_1} = 0.036350306 * 0.241300709 * 0.05 = 0.000438568$$

# Numerical Example: Backward Pass-<sup>hidden layer</sup>

- 5- We can update w1:

$$w_1^+ = w_1 - \eta * \frac{\partial E_{total}}{\partial w_1} = 0.15 - 0.5 * 0.000438568 = 0.149780716$$

- We can repeat this process to all weights of hidden layer to get the new weights w2, w3 and w4:

$$w_2^+ = 0.19956143$$

$$w_3^+ = 0.24975114$$

$$w_4^+ = 0.29950229$$

# Finally,

- When we fed forward the 0.05 and 0.1 inputs originally,
  - Total error was 0.298371109.
- After this first round of backpropagation,
  - Total error is now 0.291027924.
- But after repeating this process 10,000 times, for example,
  - Total error drops to 0.0000351085.
- At this point, when we feed forward 0.05 and 0.1, the two outputs neurons generate
  - 0.015912196 (vs 0.01 target) and
  - 0.984065734 (vs 0.99 target).

# Derivative Rules

# Derivative rules

- **Constant Rule:**  $f(x) = c$  then  $f'(x) = 0$
- **Constant Multiple Rule:**  $g(x) = c \cdot f(x)$  then  $g'(x) = c \cdot f'(x)$
- **Power Rule:**  $f(x) = x^n$  then  $f'(x) = nx^{n-1}$
- **Sum and Difference Rule:**  $h(x) = f(x) \pm g(x)$  then  $h'(x) = f'(x) \pm g'(x)$
- **Product Rule:**  $h(x) = f(x)g(x)$  then  $h'(x) = f'(x)g(x) + f(x)g'(x)$
- **Quotient Rule:**  $h(x) = \frac{f(x)}{g(x)}$  then  $h'(x) = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}$
- **Chain Rule:**  $h(x) = f(g(x))$  then  $h'(x) = f'(g(x))g'(x) \frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$

# Derivative rules

- **Trig Derivatives:**

- $f(x) = \sin(x)$  then  $f'(x) = \cos(x)$
- $f(x) = \cos(x)$  then  $f'(x) = -\sin(x)$
- $f(x) = \tan(x)$  then  $f'(x) = \sec^2(x)$
- $f(x) = \sec(x)$  then  $f'(x) = \sec(x) \tan(x)$
- $f(x) = \cot(x)$  then  $f'(x) = -\csc^2(x)$
- $f(x) = \csc(x)$  then  $f'(x) = -\csc(x) \cot(x)$

- **Exponential Derivatives**

- $f(x) = a^x$  then  $f'(x) = \ln(a)a^x$
- $f(x) = e^x$  then  $f'(x) = e^x$
- $f(x) = a^{g(x)}$  then  $f'(x) = \ln(a)a^{g(x)}g'(x)$
- $f(x) = e^{g(x)}$  then  $f'(x) = e^{g(x)}g'(x)$

# Derivative rules

- Logarithm Derivatives

- $f(x) = \log_a(x)$  then  $f'(x) = \frac{1}{\ln(a)x}$

- $f(x) = \ln(x)$  then  $f'(x) = \frac{1}{x}$

- $f(x) = \log_a(g(x))$  then  $f'(x) = \frac{g'(x)}{\ln(a)g(x)}$

- $f(x) = \ln(g(x))$  then  $f'(x) = \frac{g'(x)}{g(x)}$

Thanks