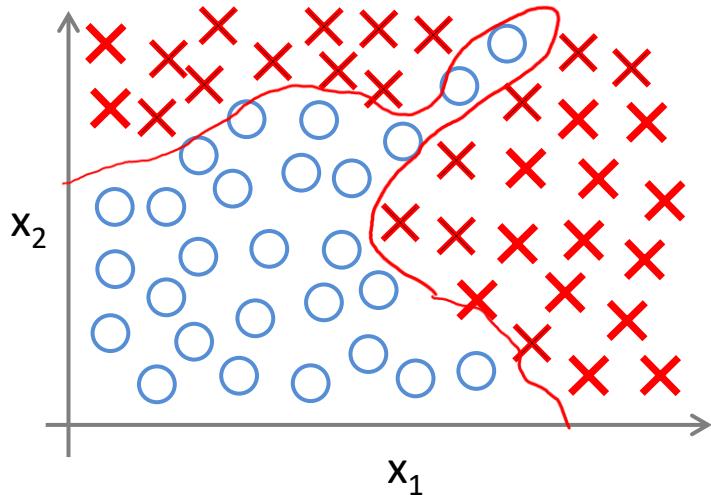


Machine Learning

Neural Networks: Representation

Non-linear hypotheses

Non-linear Classification



x_1 = size

x_2 = # bedrooms

x_3 = # floors

x_4 = age

...

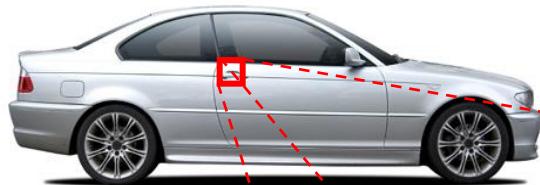
x_{100}

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

$$\begin{aligned} & \theta_0 + \theta_1 x_1 + \dots + \theta_{100} x_{100} \\ & + \theta_{101} x_1^2 + \theta_{102} x_1 x_2 + \theta_{103} x_1 x_3 \\ & + \theta_{104} x_2^2 + \theta_{105} x_2 x_3 + \dots \\ & + x_{100}^2 \end{aligned}$$

What is this?

You see this:



But the camera sees this:

194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	87	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	68	99	63	45	60	82	58	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	68	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	78	76	73	59	75	69	50



Computer Vision: Car detection



Cars

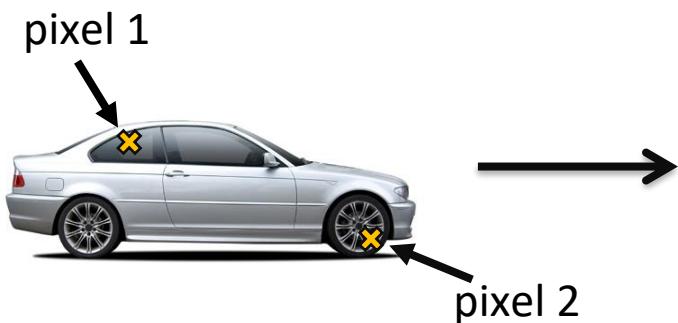


Not a car

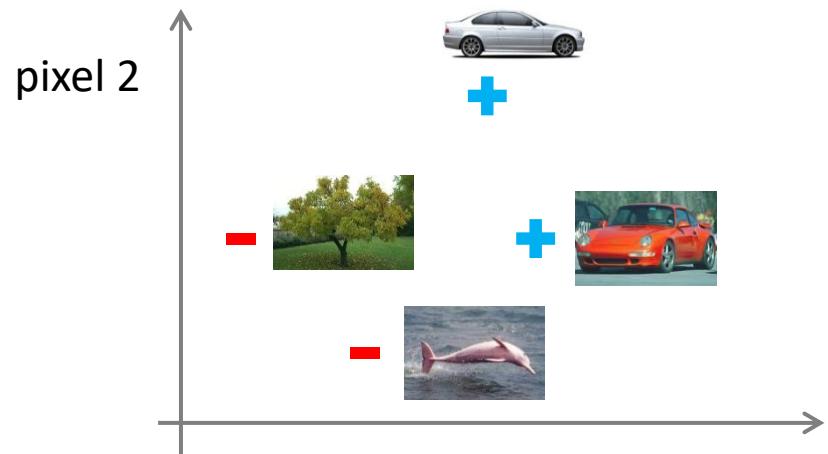
Testing:



What is this?



Learning
Algorithm



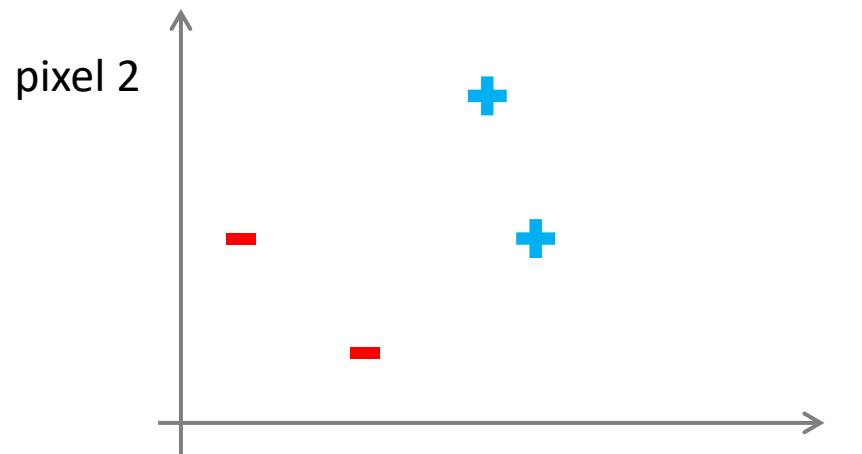
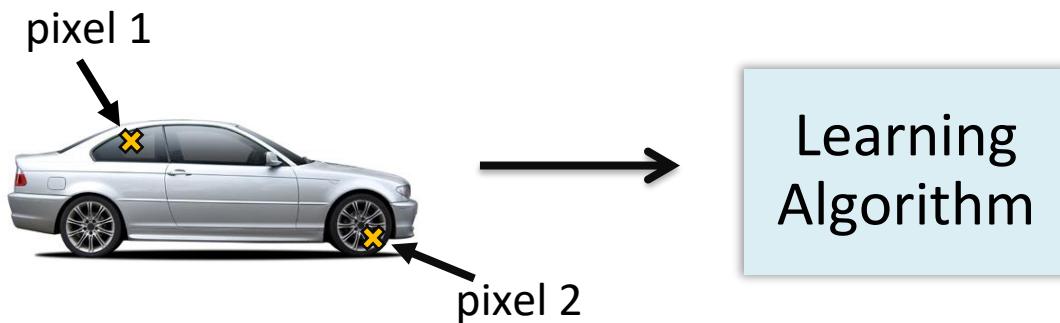
+

Cars

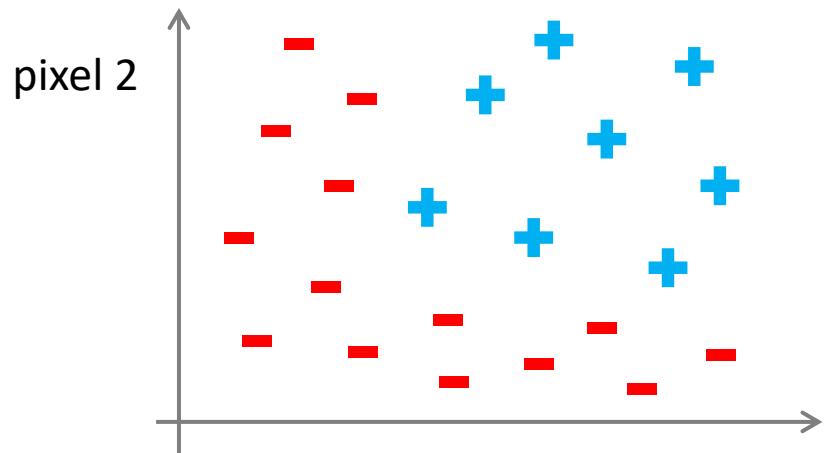
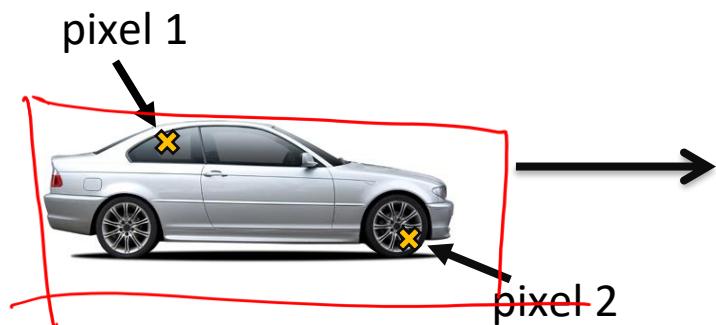
-

"Non"-Cars

pixel 1



Cars
 "Non"-Cars

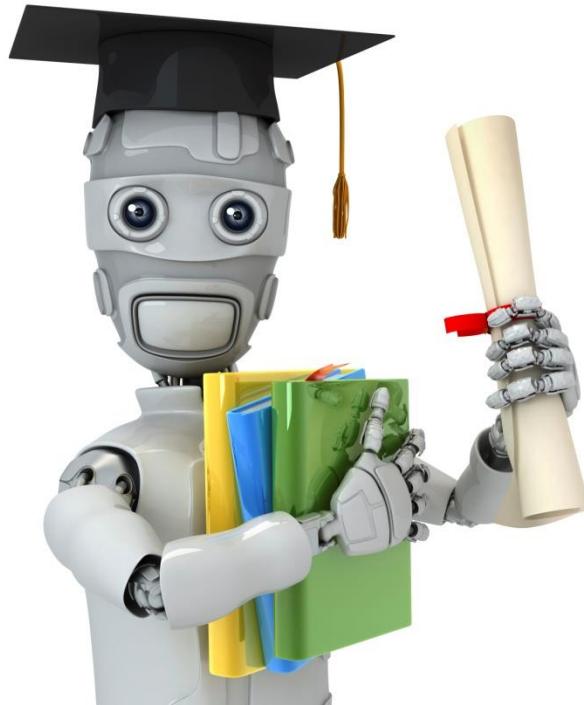


+ Cars
- "Non"-Cars

50 x 50 pixel images \rightarrow 2500 pixels
 $n = 2500$ (7500 if RGB)

$$x = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix}$$

Quadratic features ($x_i \times x_j$): ≈ 3 million features



Machine Learning

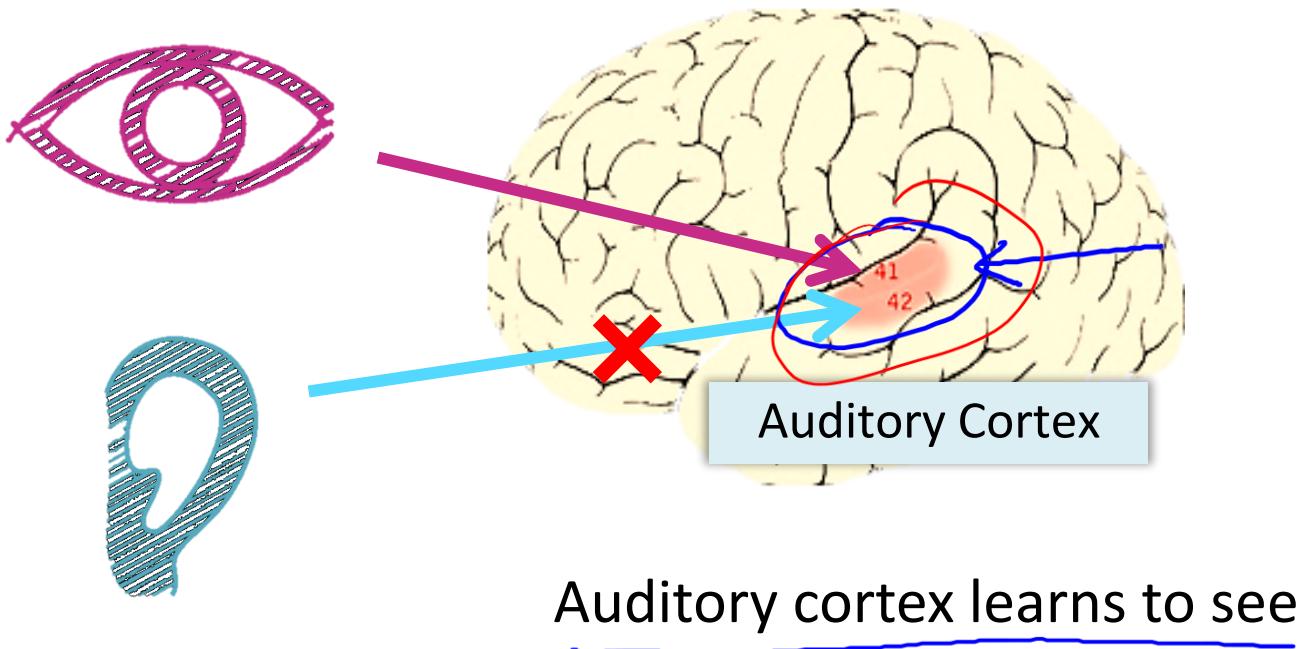
Neural Networks: Representation

Neurons and the brain

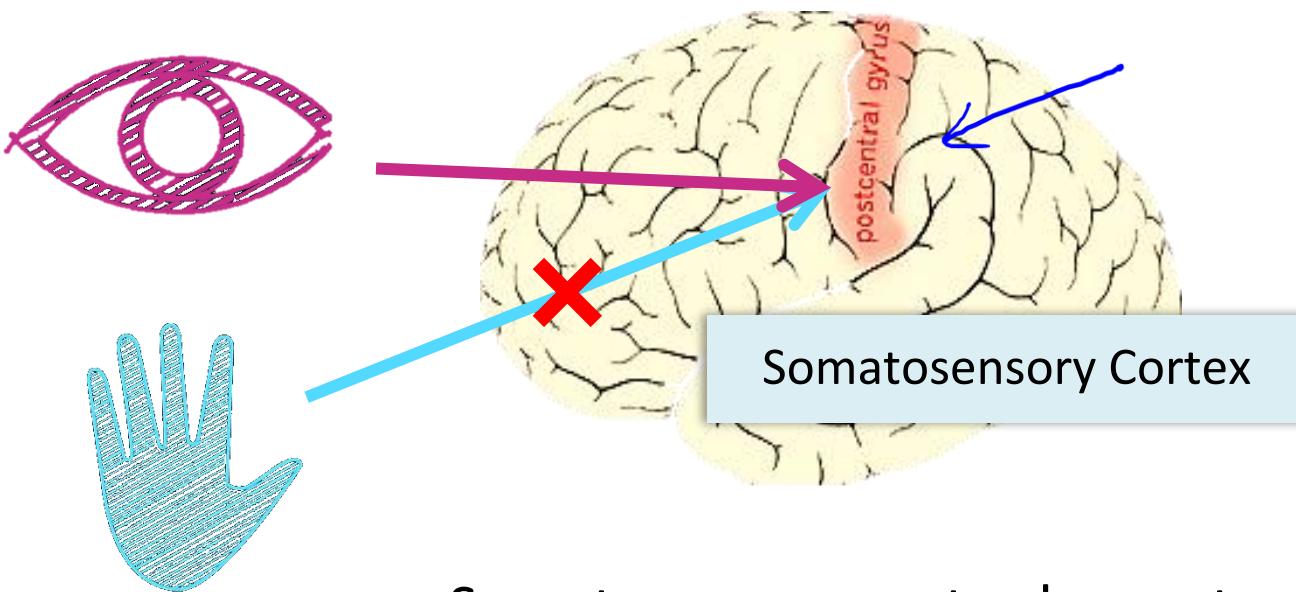
Neural Networks

- 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

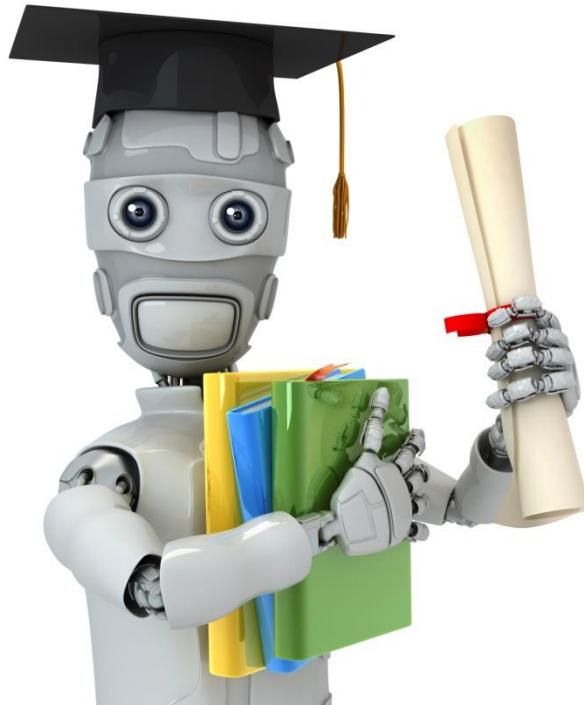
The “one learning algorithm” hypothesis



The “one learning algorithm” hypothesis



Somatosensory cortex learns to see

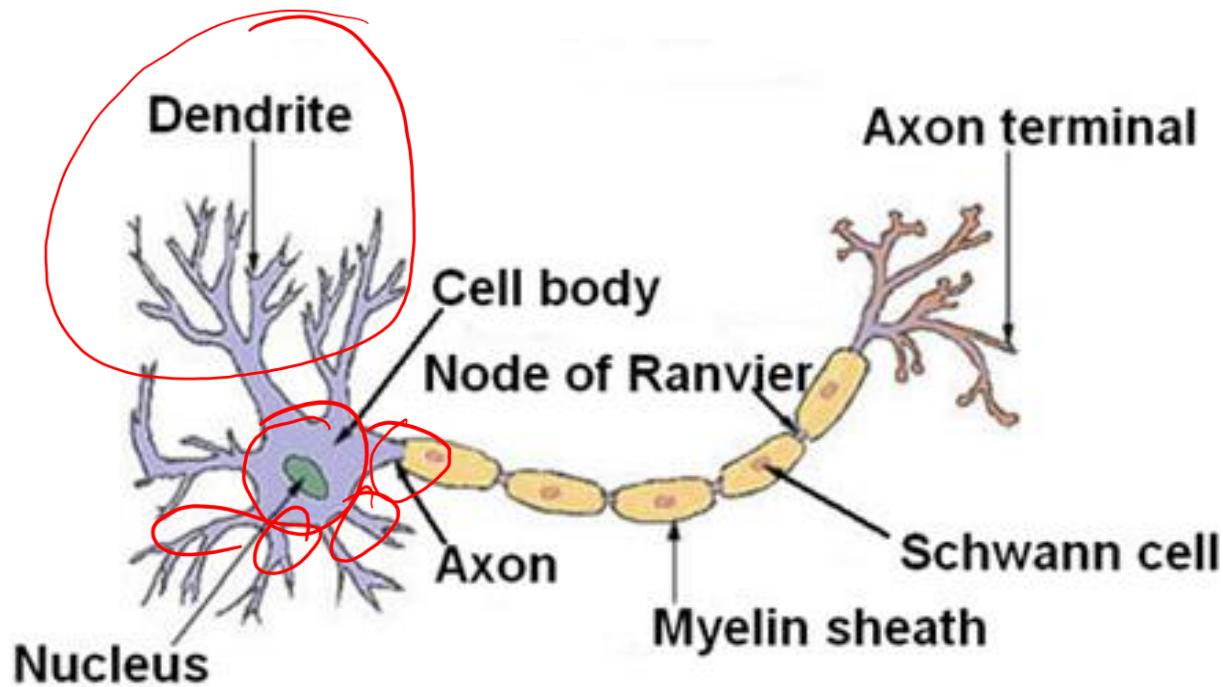


Machine Learning

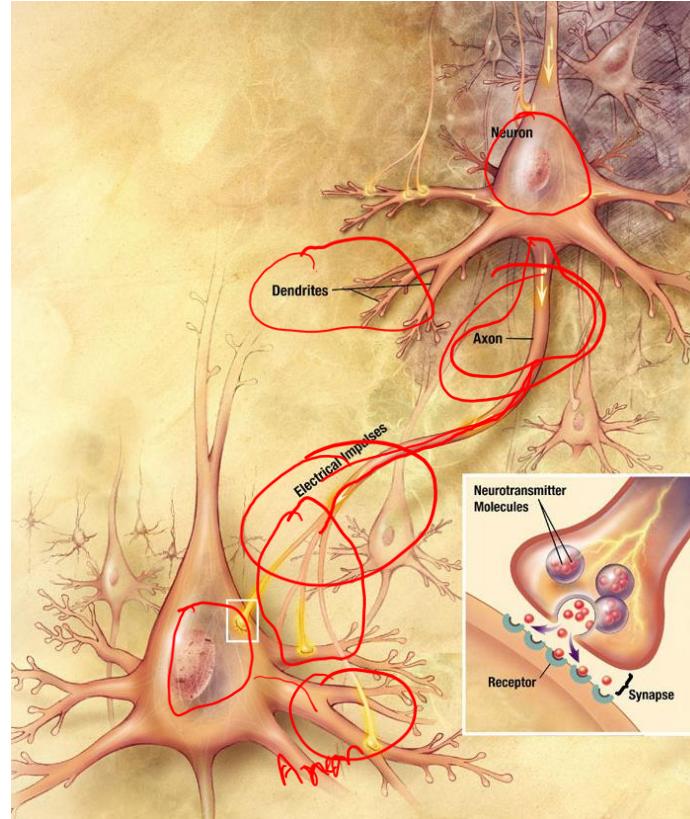
Neural Networks: Representation

Model representation I

Neuron in the brain

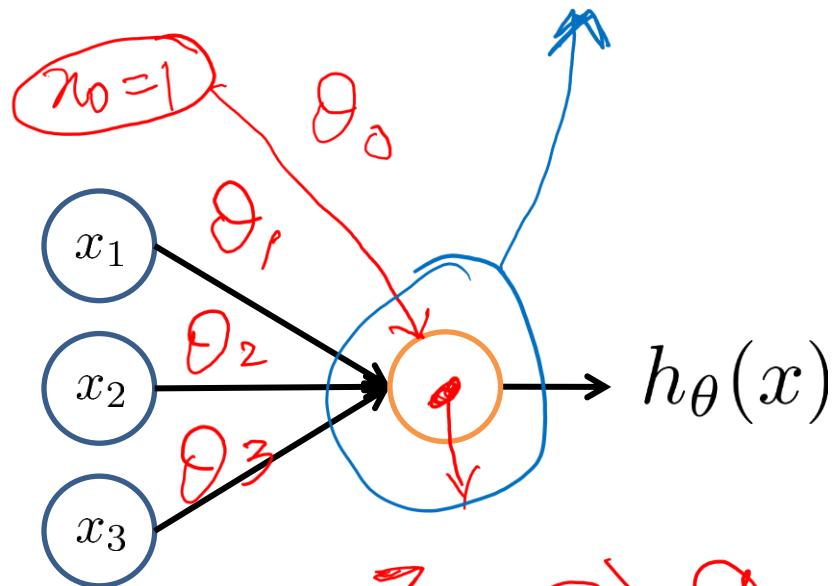


Neurons in the brain



[Credit: US National Institutes of Health, National Institute on Aging]

Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$\frac{1}{1 + e^{-\theta^T x}}$$

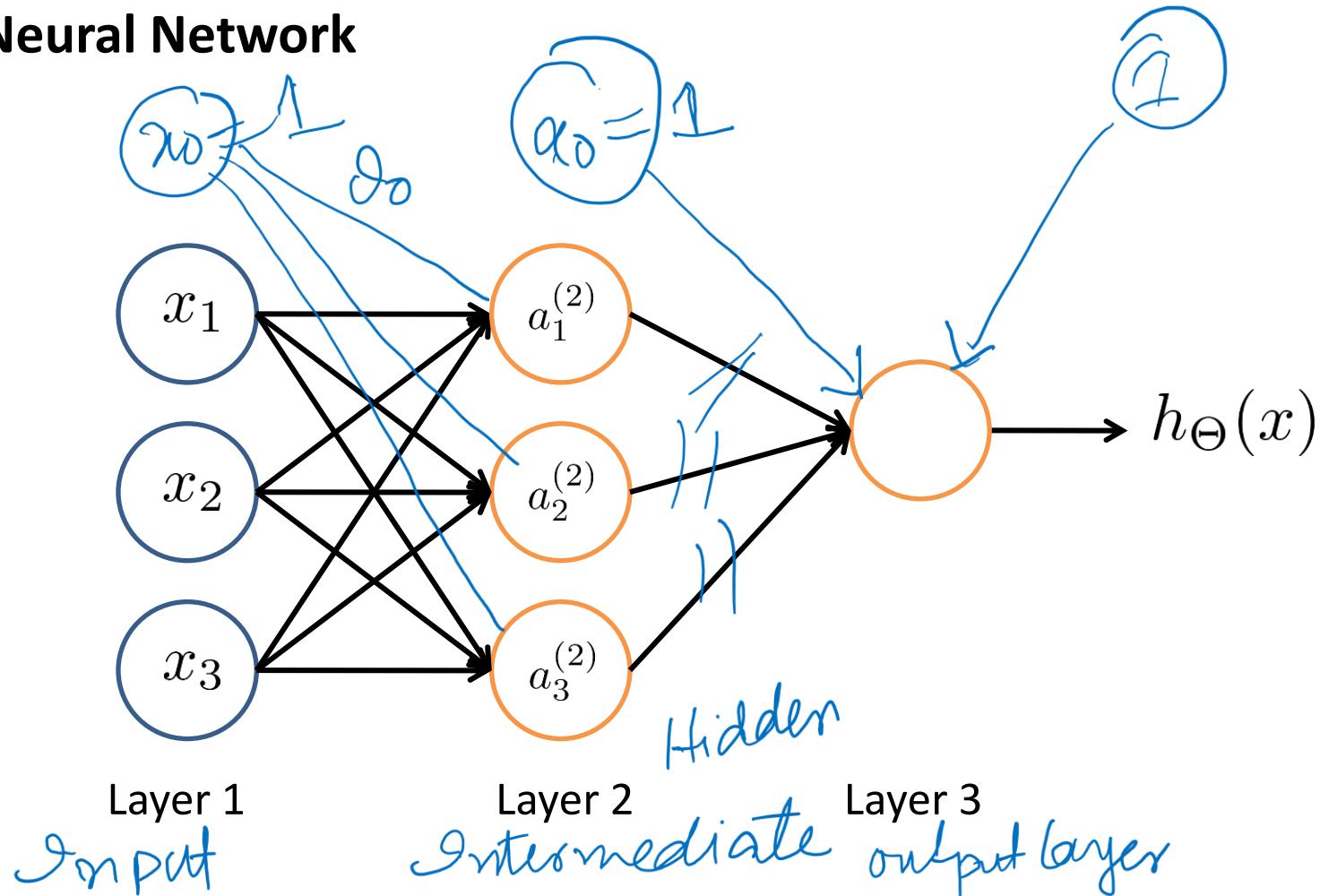
Sigmoid (logistic) activation function.
Layer 1 Layer 2

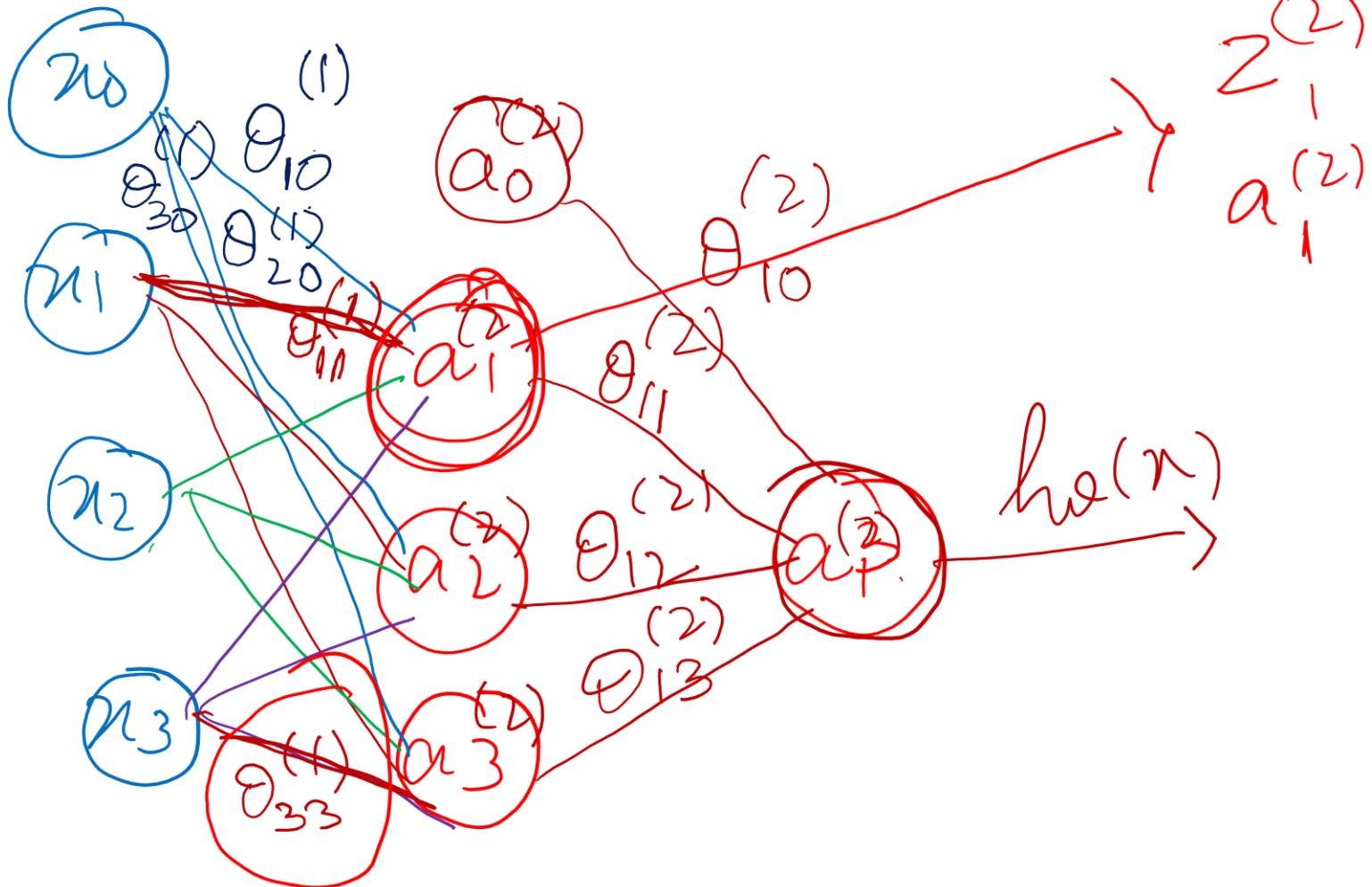
$$z = \theta^T x$$

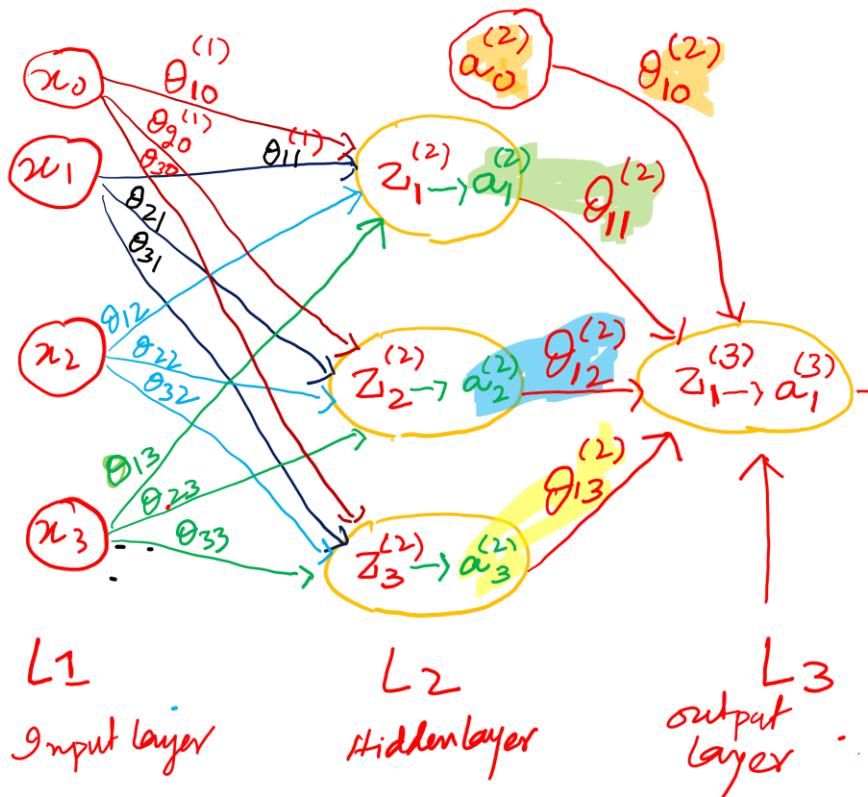
$$a = \sigma(z)$$

Perceptron

Neural Network







$$z_1^{(2)} = \theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3$$

$$z_2^{(2)} = \theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3$$

$$z_3^{(2)} = \theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3$$

g : Sigmoid. σ

ReLU

$$a_1^{(2)} = g(z_1^{(2)})$$

$$a_2^{(2)} = g(z_2^{(2)})$$

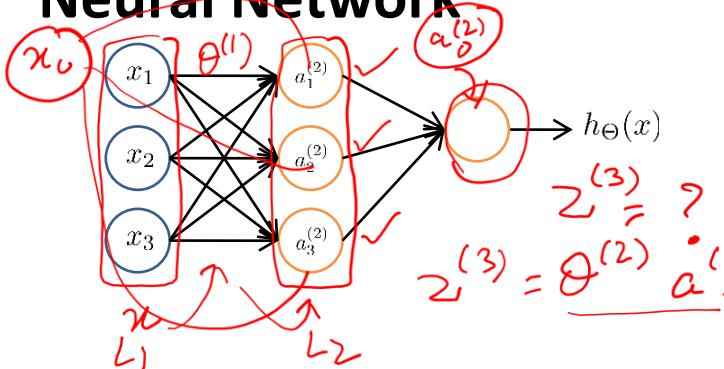
$$a_3^{(2)} = g(z_3^{(2)})$$

$$z_1^{(3)} = \theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}$$

$$a_1^{(3)} = g(z_1^{(3)})$$

(i) Sigmoid.
(ii) Softmax

Neural Network



$a_i^{(j)}$ = “activation” of unit i in layer j

$\underline{\Theta}^{(j)}$ = matrix of weights controlling
function mapping from layer j to

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

increase dimension by 1

$$\underline{\Theta}^{(1)} = \begin{bmatrix} \Theta_{10} & \Theta_{11} & \Theta_{12} & \Theta_{13} \\ \Theta_{20} & \Theta_{21} & \Theta_{22} & \Theta_{23} \\ \Theta_{30} & \Theta_{31} & \Theta_{32} & \Theta_{33} \end{bmatrix}$$

$\underline{x} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad 3 \times 4$

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

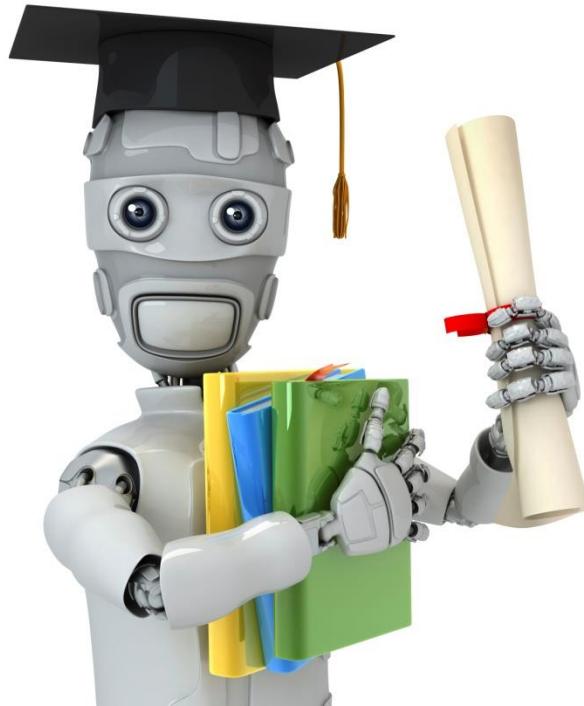
$$a_2^{(2)} = g(\underline{\Theta}_{20}^{(1)} x_0 + \underline{\Theta}_{21}^{(1)} x_1 + \underline{\Theta}_{22}^{(1)} x_2 + \underline{\Theta}_{23}^{(1)} x_3)$$

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

$$h_{\Theta}(x) = a_1^{(3)} = g(\underline{\Theta}_{10}^{(2)} a_0^{(2)} + \underline{\Theta}_{11}^{(2)} a_1^{(2)} + \underline{\Theta}_{12}^{(2)} a_2^{(2)} + \underline{\Theta}_{13}^{(2)} a_3^{(2)})$$

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

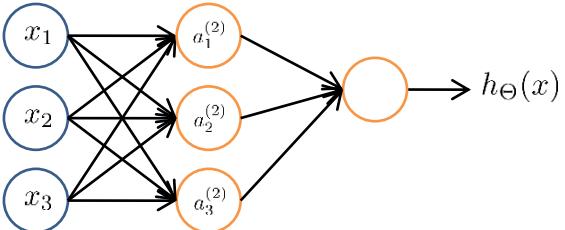
If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\underline{\Theta}^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.



Machine Learning

Neural Networks: Representation Model representation II

Forward propagation: Vectorized implementation



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

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

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

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

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

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

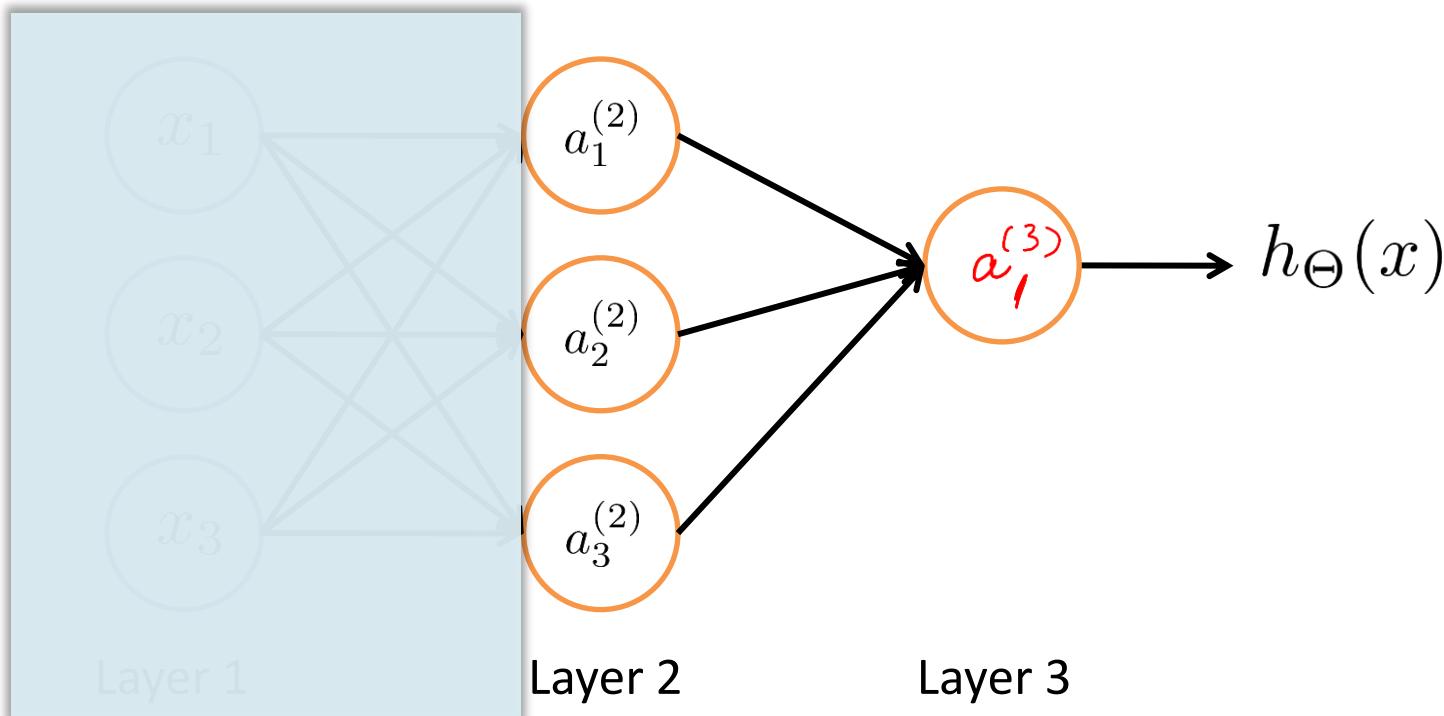
$$\underline{a^{(2)} = g(z^{(2)})}$$

Add $a_0^{(2)} = 1$.

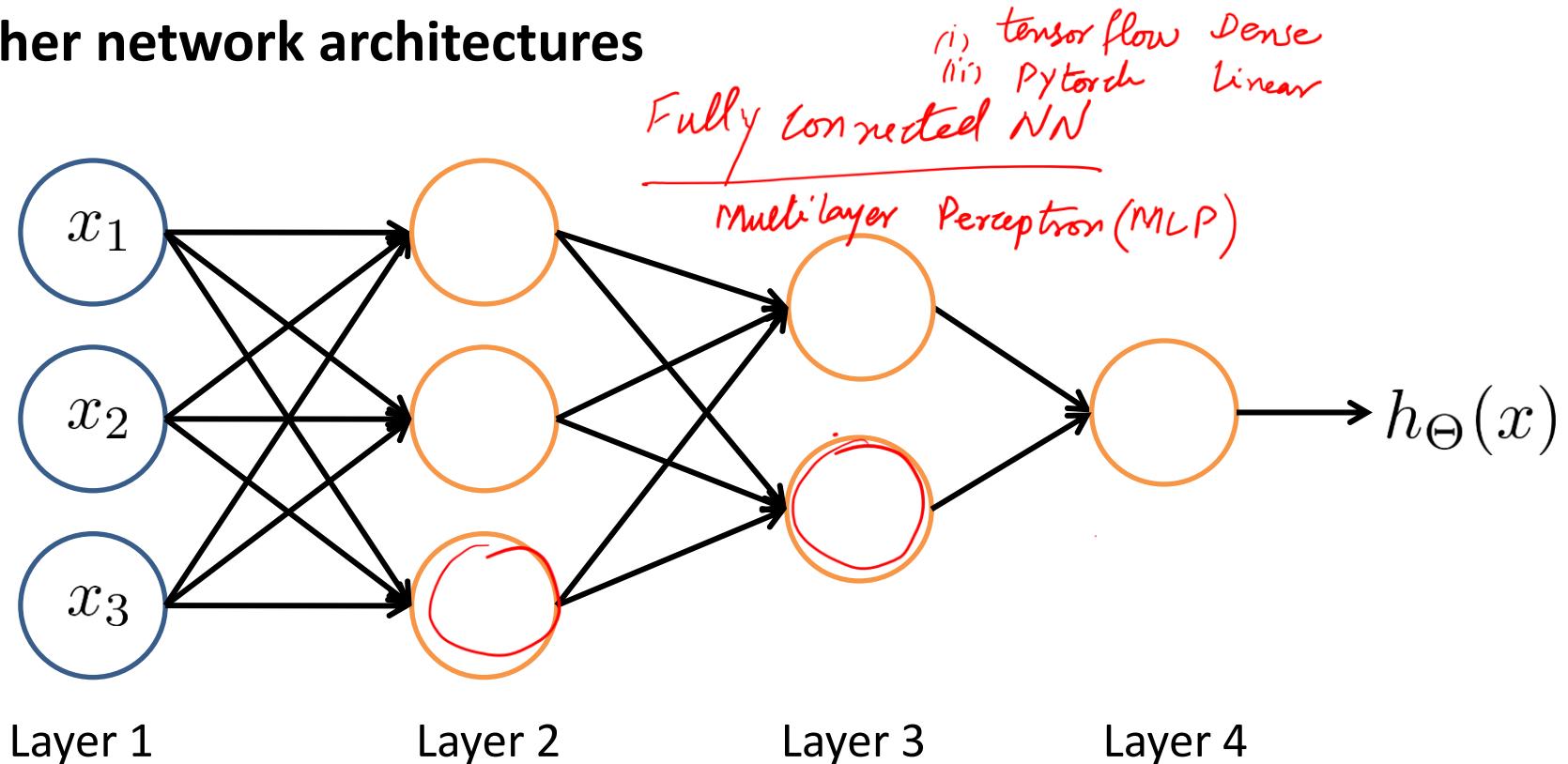
$$\underline{\underline{z^{(3)} = \Theta^{(2)} a^{(2)}}}$$

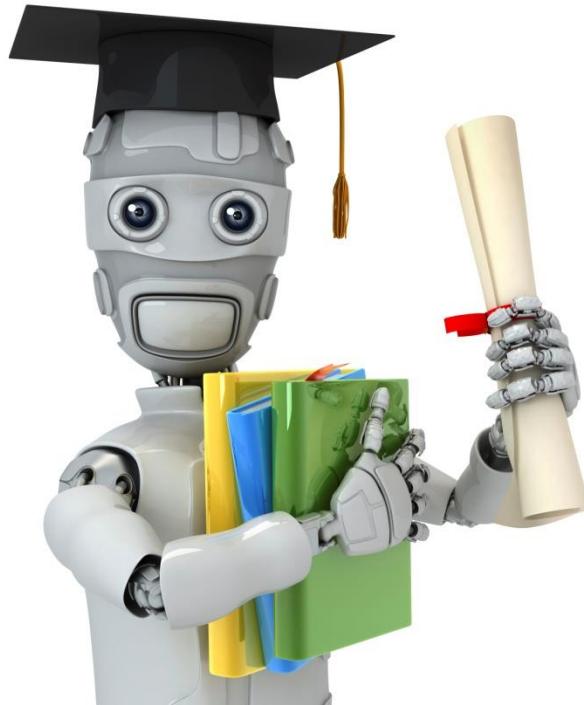
$$\underline{\underline{h_\Theta(x) = a^{(3)} = g(z^{(3)})}}$$

Neural Network learning its own features



Other network architectures





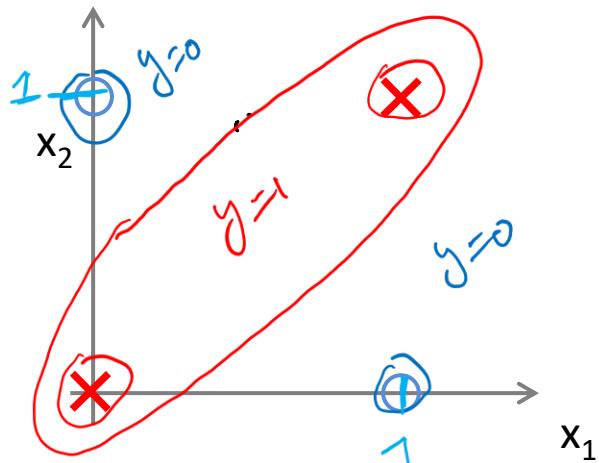
Machine Learning

Neural Networks: Representation

Examples and intuitions I

Non-linear classification example: XOR/XNOR

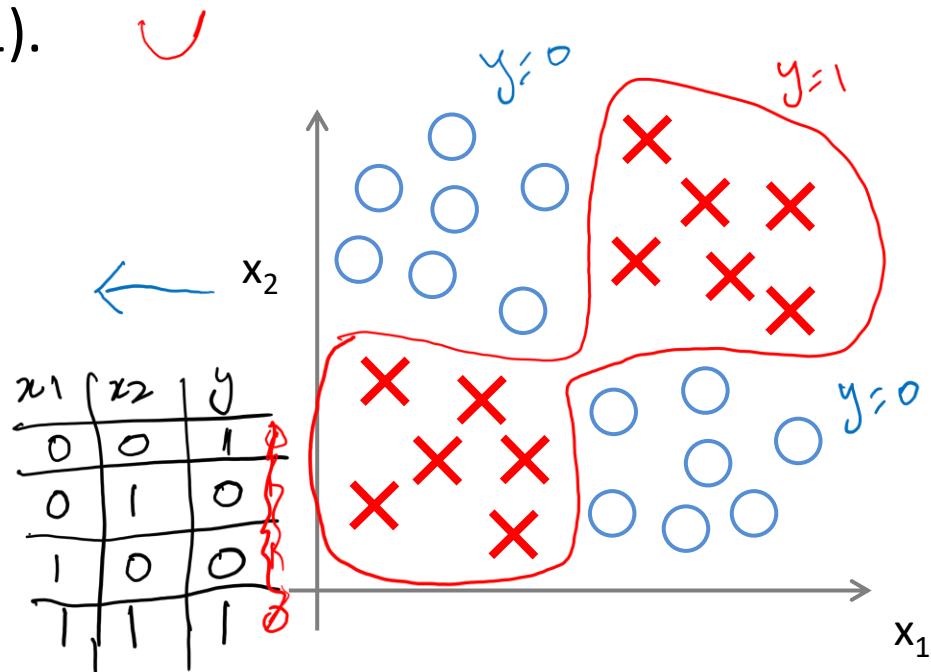
x_1, x_2 are binary (0 or 1).



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

$$x_1 \text{ XNOR } x_2$$

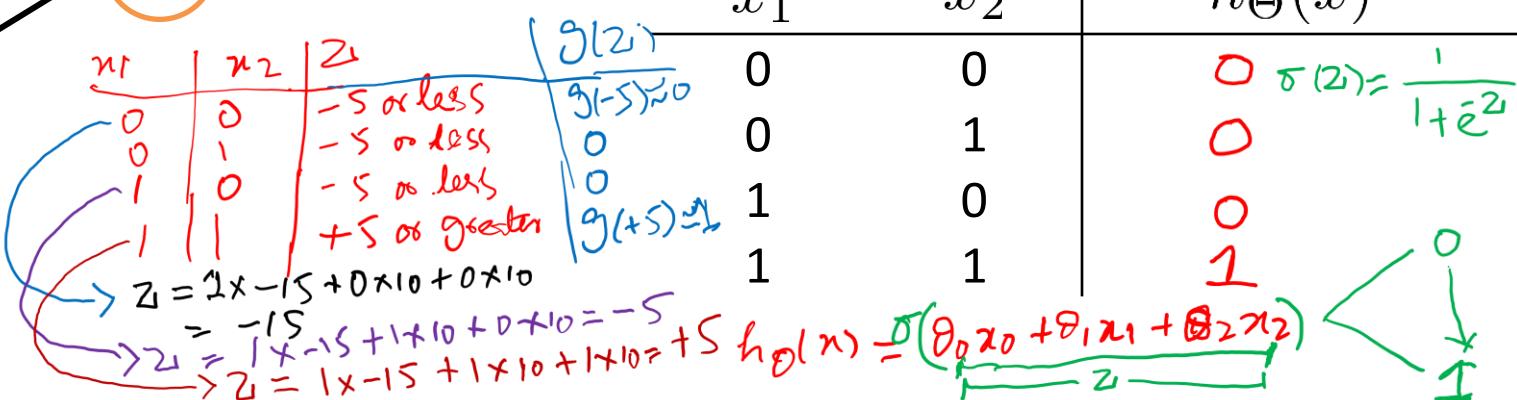
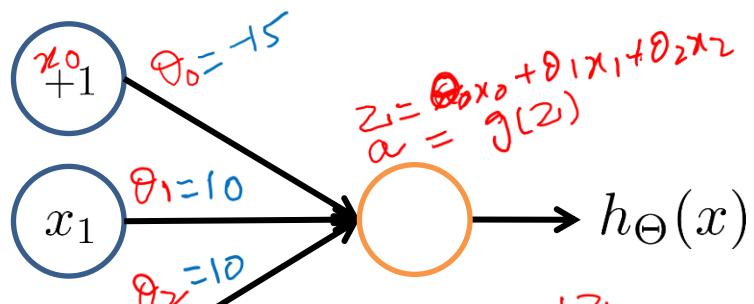
$$\text{NOT } (x_1 \text{ XOR } x_2)$$



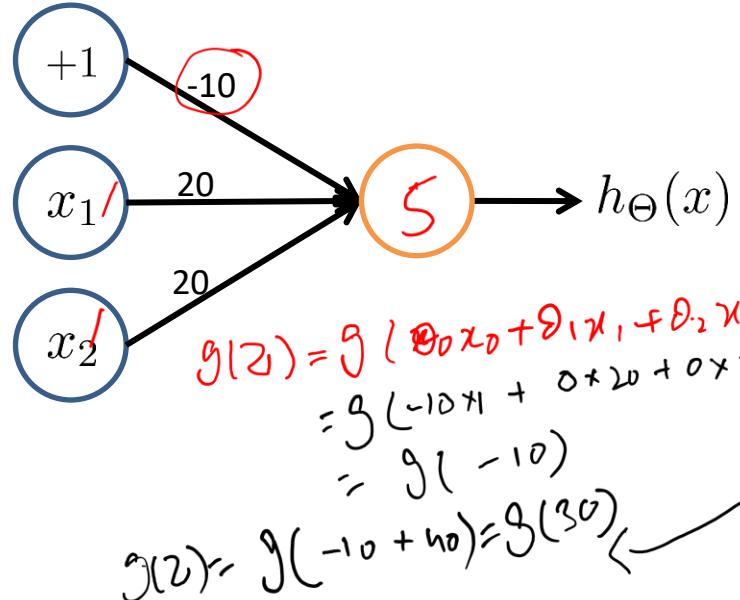
Simple example: AND

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

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



Example: OR function

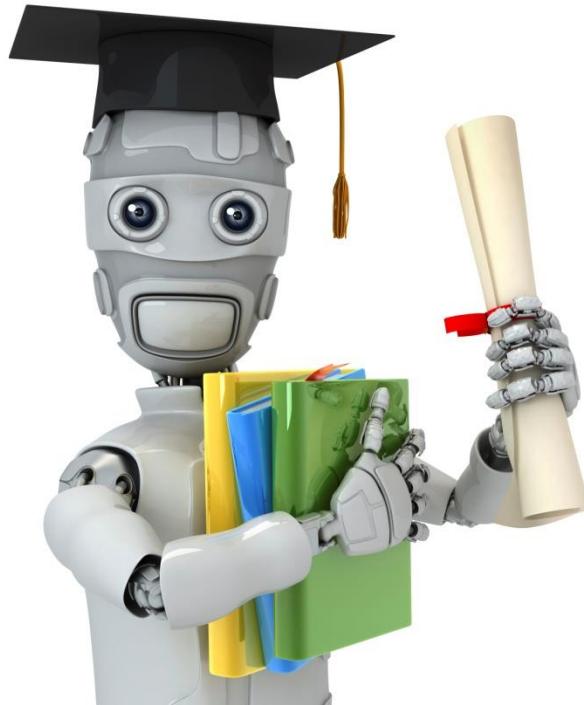


$$z = 4.6 \rightarrow g(z) = 0.98 \approx 1$$

$$z = -4.6 \rightarrow g(z) = 0.02 \approx 0$$

		x_1	x_2	z	$a = g(z)$
				$h_{\Theta}(x)$	
		0	0	-5 or less	0
		0	1	+5	1
		1	0	+5	1
		1	1	+5	1

Annotations: 'less' with an arrow pointing to the first row; 'greater' with an arrow pointing to the second, third, and fourth rows.



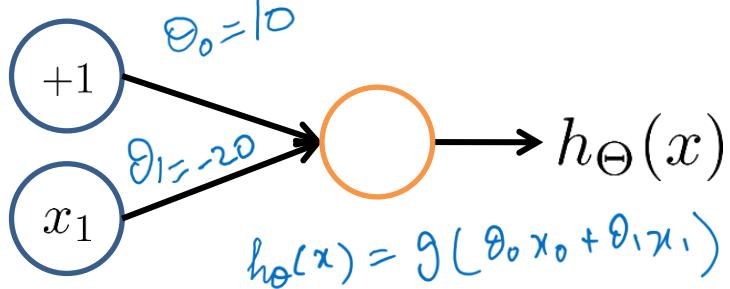
Machine Learning

Neural Networks: Representation

Examples and intuitions II

x_1 AND x_2

Negation:



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

x_1 OR x_2

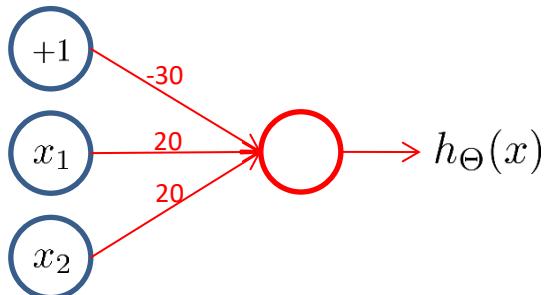
		x_0	x_1	and	nand	and'	/
0	0	0	0	0	1	1	1
0	1	0	0	0	1	0	0
1	0	1	0	0	1	0	0
1	1	1	1	1	0	0	0

x_1	z	$h_{\Theta}(x)$
0	+5 or greater	1 ✓
1	-5 or less	0 ✓

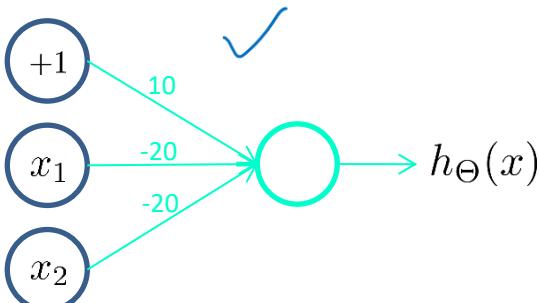
(NOT x_1) AND (NOT x_2)



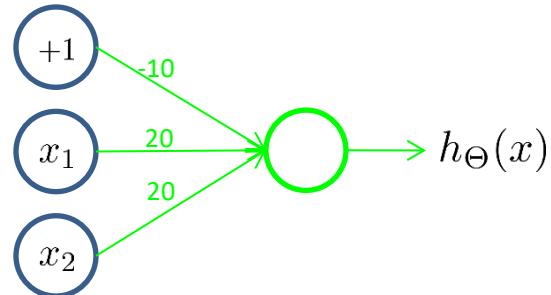
Putting it together: x_1 XNOR x_2



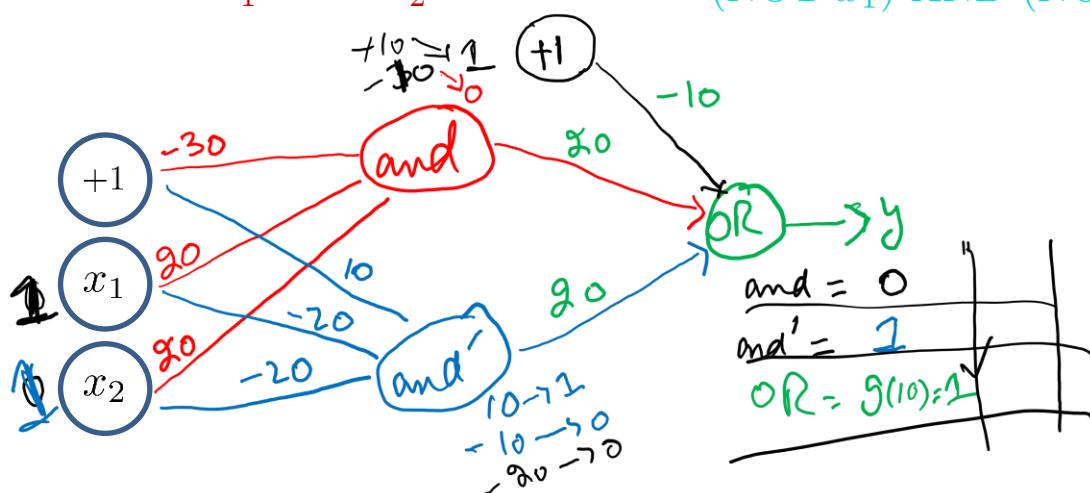
x_1 AND x_2



(NOT x_1) AND (NOT x_2)

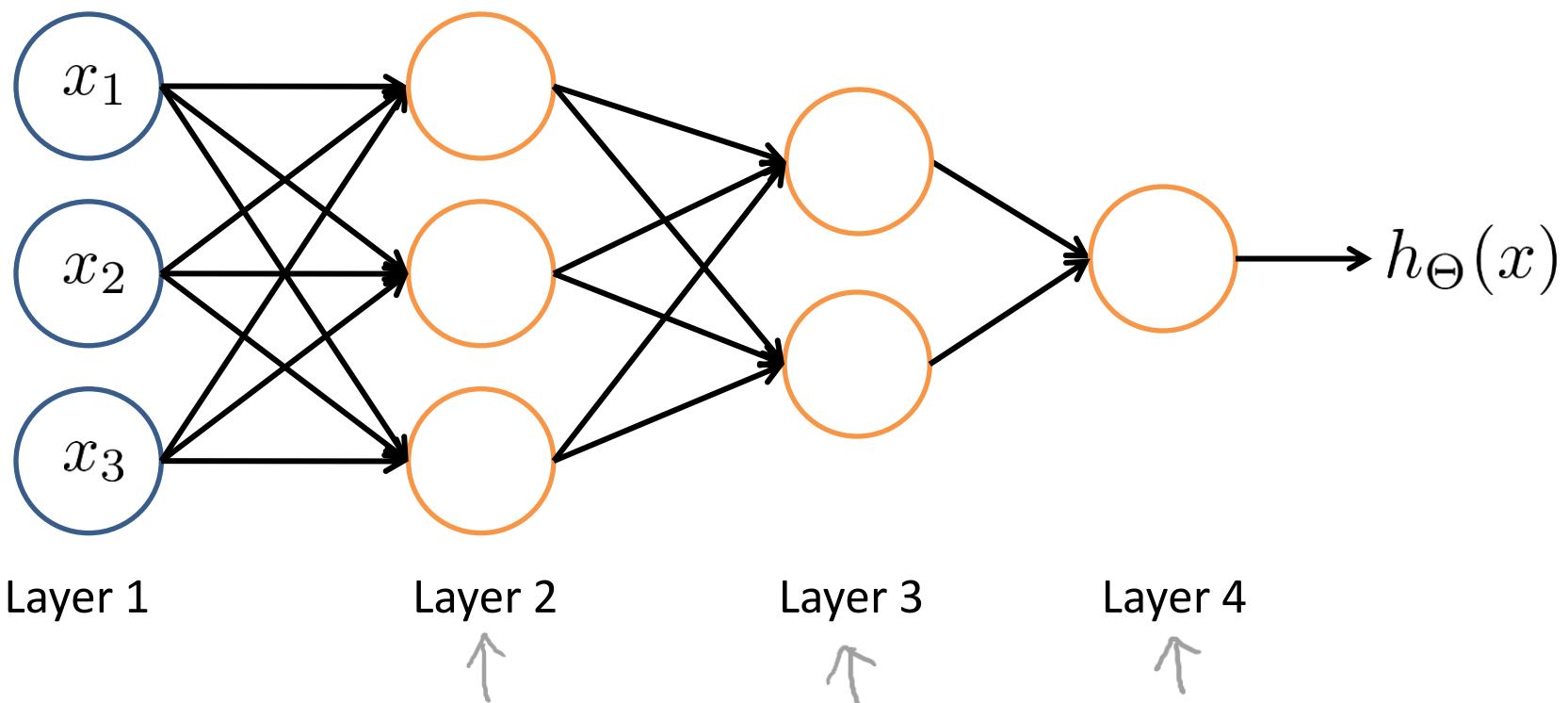


x_1 OR x_2

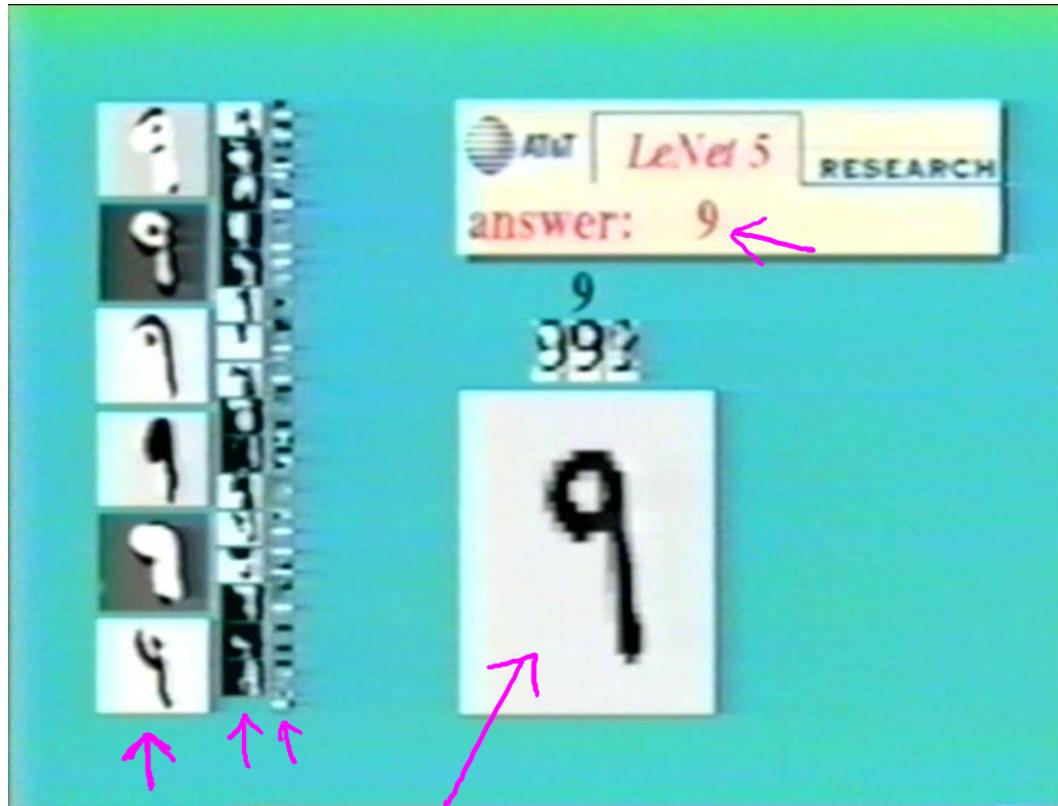


x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1 ✓
0	1	0	0	0 ✓
1	0	0	0	0 ✓
1	1	1	0	1 ✓

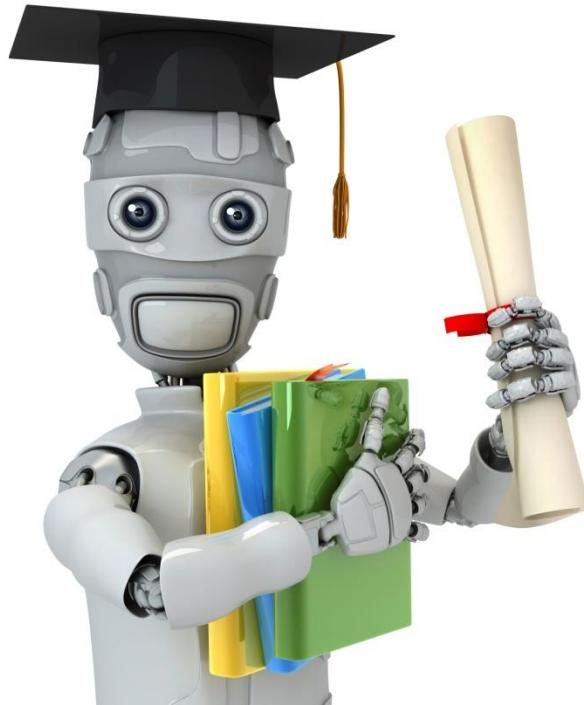
Neural Network intuition



Handwritten digit classification



[Courtesy of Yann LeCun]



Machine Learning

Neural Networks: Representation

Multi-class classification

Multiple output units: One-vs-all.



Jeep



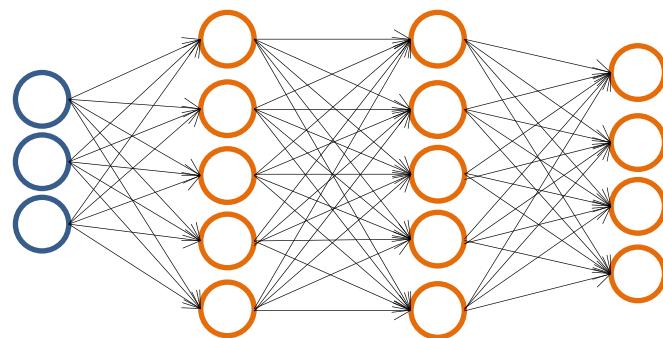
Truck



Light



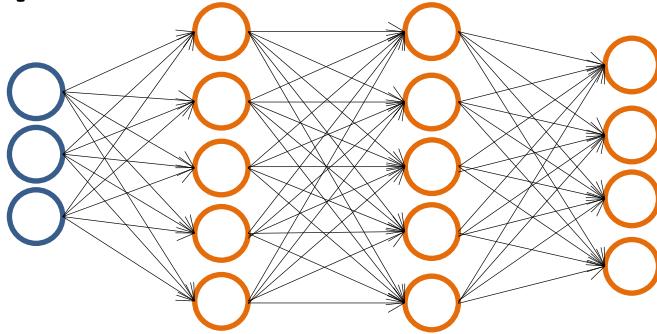
Pedestrian



$$h_{\Theta}(x) \in \mathbb{R}^4$$

Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
when jeep when truck when light

Multiple output units: One-vs-all.

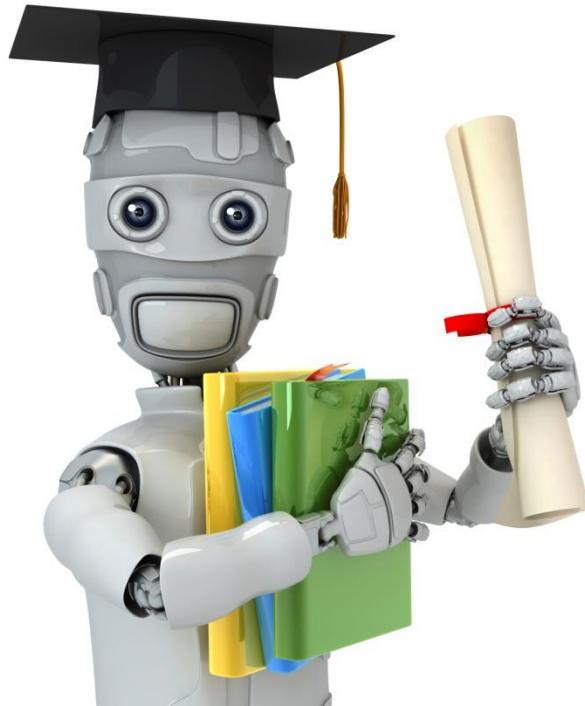


$$h_{\Theta}(x) \in \mathbb{R}^4$$

Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
when Jeep when truck when light

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

$y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
Jeep truck light pedestrian

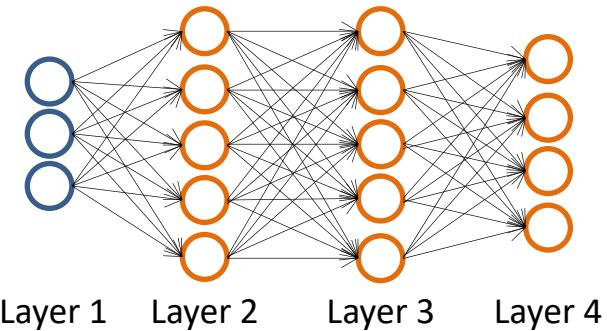


Machine Learning

Neural Networks: Learning

Cost function

Neural Network (Classification)



Binary classification

$y = 0 \text{ or } 1$

1 output unit

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

L = total no. of layers in network

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

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

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:

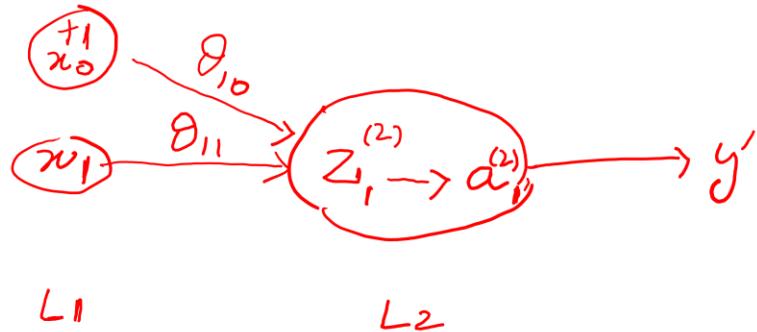
$$h_\Theta(x) \in \mathbb{R}^K \quad (h_\Theta(x))_i = i^{th} \text{ output}$$

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

Neural Network:

- 1- Forward Propagation (Predict ~~as~~ a label)
- 2- Cost computation (Finds the distance between actual and Predicted)
- 3- Backward Propagation (optimizes the parameters and finds the ~~correct~~ predicted label as close as ~~correct~~ label)
 - Learning / training
 - by finding the optimal values of the parameters.
we use an optimizer.
 - Gradient Descent Algorithm

1. Forward Propagation:



$$z_1^{(2)} = \theta_{1,0} x_0 + \theta_{1,1} x_1$$

$$a_1^{(2)} = g(z_1^{(2)})$$

$$y' = a_1^{(2)}$$

g is an activation function

(2) Cost Computation

(a) Binary Cross Entropy

$$J(\theta) = -[y \log(a_1^{(2)}) + (1-y) \log(1-a_1^{(2)})]$$

(b) Squared Error function

$$J(\theta) = (a_1^{(2)} - y)^2$$

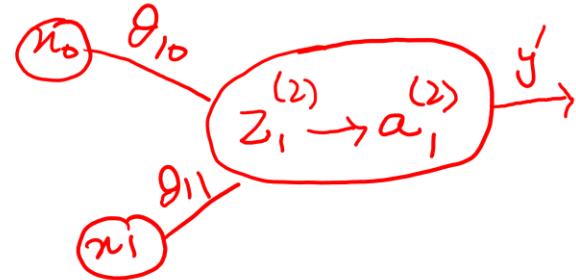
(3) Backward / Back Propagation :

- (i) Find lost by using either
 - (a) Binary Cross Entropy
 - or (b) Squared Error
- (ii) use an optimizer : (GD)

$$\theta_{10}, \theta_{11}$$

$$\theta_{10} := \theta_{10} - \alpha \frac{\partial J(\theta)}{\partial \theta_{10}}$$

$$\theta_{11} := \theta_{11} - \alpha \frac{\partial J(\theta)}{\partial \theta_{11}}$$



Chain Rule :
if $y = f(u)$

Chain Rule

if $y = f(u)$ and
 $u = g(x)$

then $y = f(g(x))$

and $\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}$

- (i) $J(\theta) = f(a_1^{(2)}) \rightarrow f:$ Binary crossEntropy
Squared Error
- (ii) $a_1^{(2)} = g(z_1^{(2)})$
- (iii) $z_1^{(2)} = \theta_{10}x_0 + \theta_{11}x_1$

Therefore,

$$\frac{\partial J(\theta)}{\partial \theta_{10}} = \frac{\partial J(\theta)}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial \theta_{10}}$$

A
 B
 C

A; In case of binary loss Entropy function;

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \alpha_1^{(2)}} &= \frac{\partial}{\partial \alpha_1^{(2)}} [-y \log(\alpha_1^{(2)}) - (1-y) \log(1-\alpha_1^{(2)})] \\&= -y \frac{\partial}{\partial \alpha_1^{(2)}} \log(\alpha_1^{(2)}) - (1-y) \frac{\partial}{\partial \alpha_1^{(2)}} \log(1-\alpha_1^{(2)}) \\&= \frac{-y}{\alpha_1^{(2)}} \times 1 - \frac{(1-y)}{1-\alpha_1^{(2)}} \frac{\partial (1-\alpha_1^{(2)})}{\partial \alpha_1^{(2)}} \\&= \frac{-y}{\alpha_1^{(2)}} - \frac{(1-y)}{1-\alpha_1^{(2)}} * -1 \\&= (\alpha_1^{(2)} - y)\end{aligned}$$

(A) : In case of Squared Error function:

$$\begin{aligned}\frac{\partial J(\theta)}{\partial a_1^{(2)}} &= \frac{\partial}{\partial a_1^{(2)}} (a_1^{(2)} - y)^2 \\ &= 2(a_1^{(2)} - y) \frac{\partial}{\partial a_1^{(2)}} (a_1^{(2)} - y) \\ &= 2(a_1^{(2)} - y)\end{aligned}$$

B)

$$\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = \frac{\partial}{\partial z_1^{(2)}} g(z_1^{(2)})$$

$$\boxed{a_1^{(2)} = g(z_1^{(2)})}$$
$$a_1^{(2)} = \sigma(z_1^{(2)})$$

$$\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = \frac{\partial}{\partial z_1^{(2)}} \sigma(z_1^{(2)})$$

$$\boxed{\sigma' = \sigma(1-\sigma)}$$

$$= \sigma(z_1^{(2)}) (1 - \sigma(z_1^{(2)}))$$

$$\boxed{\frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = a_1^{(2)} (1 - a_1^{(2)})} .$$

② :

$$\frac{\partial z_1^{(2)}}{\partial \theta_{10}} = \frac{\partial (\theta_{10}x_0 + \theta_{11}x_1)}{\partial \theta_{10}}$$

$$\boxed{\frac{\partial z_1^{(2)}}{\partial \theta_{10}} = x_0 = 1}$$

$$z_1^{(2)} = \theta_{10}x_0 + \theta_{11}x_1$$

..

So, we combine ①, ②, ③

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_{10}} &= \left[\frac{-y}{a_1^{(2)}} + \frac{(1-y)}{1-a_1^{(2)}} \right] \times \left[a_1^{(2)} (1 - a_1^{(2)}) \right] \times [1] \\ &= -y(1 - a_1^{(2)}) + (1-y) a_1^{(2)} \\ &= \boxed{\quad}\end{aligned}$$

$$\frac{\partial \bar{J}(\theta)}{\partial \theta_{11}} = ?$$

$$\frac{\partial \bar{J}(\theta)}{\partial \theta_{11}} = (a_1^{(2)} - y) x_1$$

$$\frac{\partial \bar{J}(\theta)}{\partial \theta_{10}} = (a_1^{(2)} - y)$$

$$\checkmark \quad \theta_{10} := \theta_{10} - \lambda (a_1^{(2)} - y)$$

$$\checkmark \quad \theta_{11} := \theta_{11} - \lambda (a_1^{(2)} - y) x_1$$

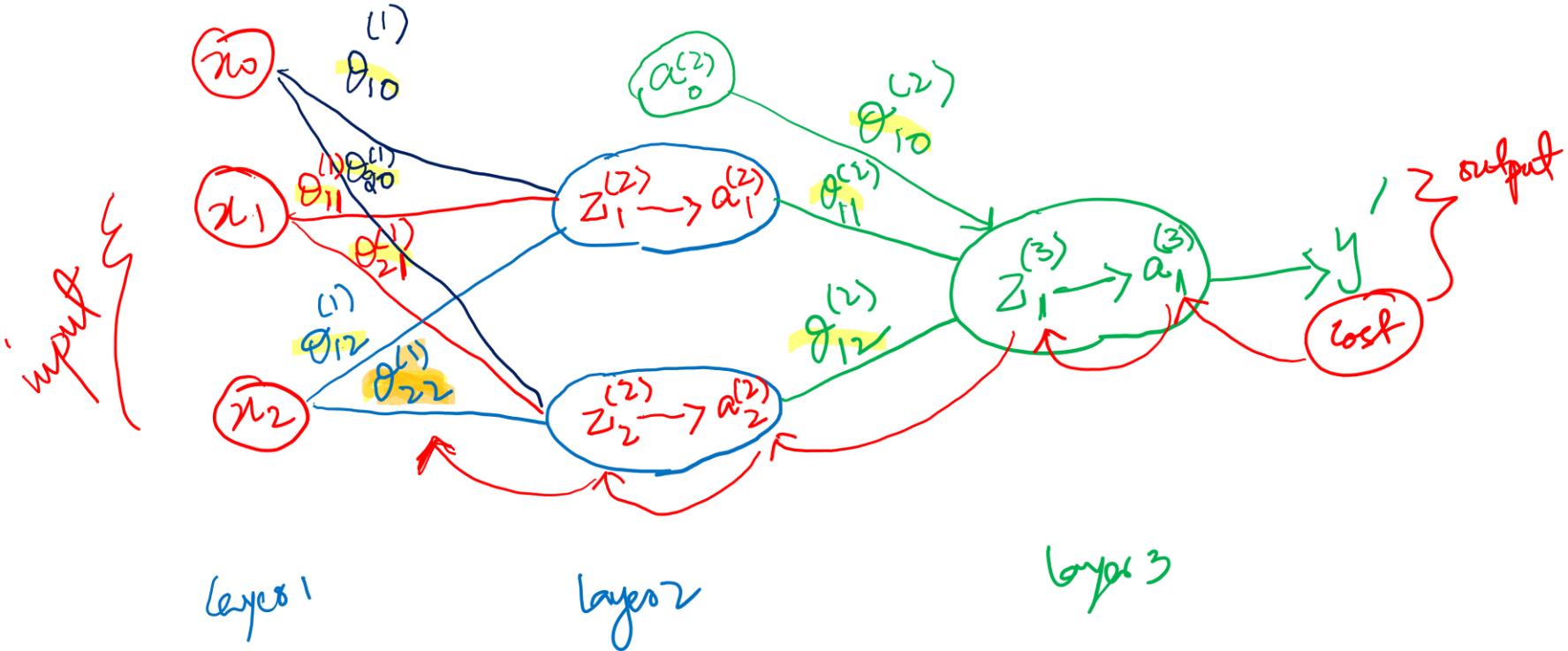
Forward Propagation: $a_1^{(2)} = \sigma (\theta_{10} x_0 + \theta_{11} x_1)$

epoch:

if you have $(x^{(1)})$,
 $(x^{(1)}, y^{(1)})$, ..., $(x^{(m)}, y^{(m)})$
 m examples.

~~After~~ Gradient Computation:

- After m examples: Batch Gradient Descent.
- After $K < m$ examples: mini-batch
- After each example: Stochastic gradient descent



let's compute gradients for $\theta_{22}^{(1)}$:

chain rule :

$$\frac{\partial J(\theta)}{\partial \theta_{22}} = \frac{\partial J(\theta)}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \cdot \frac{\partial z_1^{(2)}}{\partial \theta_{22}^{(1)}}$$

① $\frac{\partial z_2^{(2)}}{\partial \theta_{22}^{(1)}} = \frac{\partial}{\partial \theta_{22}^{(1)}} (\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2)$

$$= x_2$$

$$\textcircled{2} \quad \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} = \frac{\partial}{\partial z_1^{(2)}} \left(\frac{1}{1 + \bar{e}^{z_1^{(2)}}} \right) \\ = a_1^{(2)} (1 - a_1^{(2)})$$

$$\textcircled{3} \quad \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} = \frac{\partial}{\partial a_1^{(2)}} (\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)}) \\ = \theta_{11}^{(2)}$$

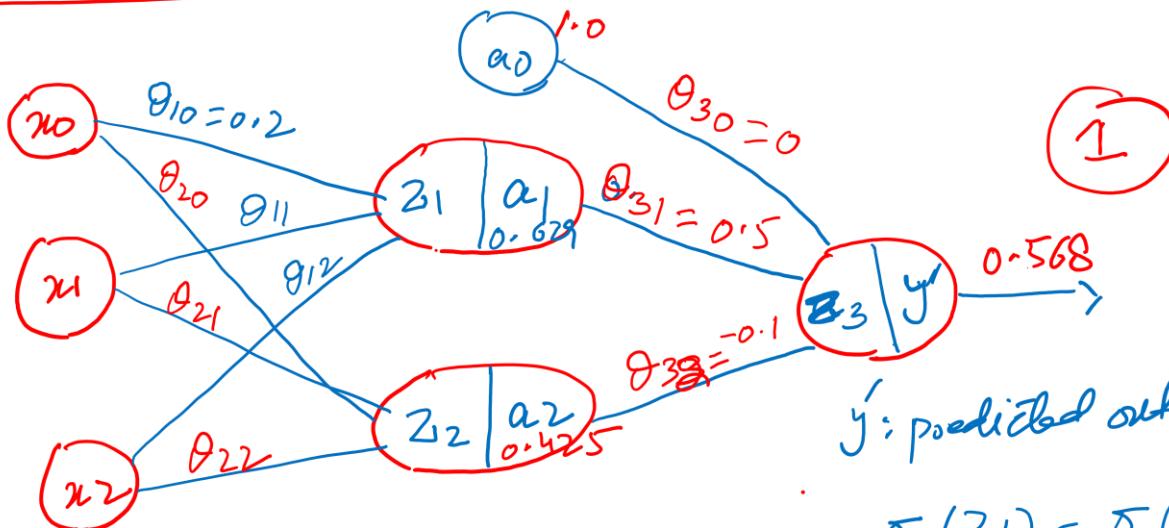
$$\textcircled{4} \quad \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} = a_1^{(3)} (1 - a_1^{(3)}) \quad \cdot \quad \cdot \quad .$$

⑤ $\frac{\partial J(\theta)}{\partial a_1^{(3)}}$: ① Binary cross ✓
 Entropy ?
 ② Squared Error Function X ?

$$\begin{aligned}
 (i) \frac{\partial J(\theta)}{\partial a_1^{(3)}} &= \frac{\partial}{\partial a_1^{(3)}} \left[-y \log(a_1^{(3)}) - (1-y) \log(1-a_1^{(3)}) \right] \\
 &= -y \frac{\partial}{\partial a_1^{(3)}} (\log(a_1^{(3)})) - (1-y) \frac{\partial}{\partial a_1^{(3)}} (\log(1-a_1^{(3)})) \\
 &= \frac{-y}{a_1^{(3)}} + \frac{(1-y)}{1-a_1^{(3)}} \quad ; \quad \dots \quad .
 \end{aligned}$$

$$\boxed{\frac{\partial J(\theta)}{\partial \theta_{22}^{(1)}} = \left[\frac{-y}{a_1^{(3)}} + \frac{(1-y)}{1-a_1^{(3)}} \right] \times a_1^{(3)}(1-a_1^{(3)}) \times \theta_{11}^{(2)} \times a_1^{(2)}(1-a_1^{(2)}) \times x_2}$$

Neural Network Example :



$$z_1 = \theta_{10} \cdot 1 + \theta_{11} \cdot x_1 + \theta_{12} \cdot x_2$$

$$z_1 = 0.2 \cdot 1 + 0.4 \cdot 0.6 + 0.1 \cdot 0.9$$

$$z_1 = 0.2 + 0.24 + 0.09 = 0.53$$

$$a_1 = \sigma(z_1) = \sigma(0.53) = 0.629$$

\therefore

$$z_2 = \theta_{20} \cdot 1 + \theta_{21} \cdot x_1 - \theta_{22} \cdot x_2$$

$$= -0.3 + 0.18 + 0.18$$

$$= -0.3$$

$$a_2 = \sigma(-0.3) \approx 0.425$$

$$z_3 = 0.272$$

$$y' = \sigma(0.272) \approx 0.568$$

② $J(\theta) = E = \text{Cost computation: } \frac{1}{2} (y' - t)^2$

$$J(\theta) = E = \frac{1}{2} * (0.568 - 1)^2$$

$$J(\theta) = E = 0.093$$

③ Optimization / Backpropagation

$$\theta_{ji} = \theta_{ji} - \alpha \frac{\partial J(\theta)}{\partial \theta_{ji}}$$

$$\theta_{31} = \theta_{31} - \alpha \frac{\partial J(\theta)}{\partial \theta_{31}}$$

Chain Rule

$$\frac{\partial J(\theta)}{\partial \theta_{31}} = \frac{\partial J(\theta)}{\partial y'} \frac{\partial y'}{\partial z_3} \frac{\partial z_3}{\partial \theta_{31}}$$

$$\begin{aligned}\frac{\partial J(\theta)}{\partial y'} &= \frac{1}{2} (y' - t)^2 \\ &= \frac{1}{2} (y' - t) \times \cancel{\frac{\partial (y' - t)}{\partial y'}} \\ &= y' - t \\ &= 0.568 - 1 \\ &= -0.432\end{aligned}$$

$$\frac{\partial y'}{\partial z_3} = y'(1-y')$$

$$= 0.568(1-0.568)$$

$$= \boxed{s_1}$$

$$\begin{aligned}\frac{\partial z_3}{\partial \theta_{31}} &= \theta_{30} \cancel{x_0} + \boxed{\theta_{31} \cdot a_1} + \theta_{32} \cancel{x_2} \\ \frac{\partial z_3}{\partial \theta_{31}} &= a_1 = 0.629\end{aligned}$$

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_{31}} &= (-0.432) \times 0.568 \times (-432) \times 0.629 \\ \frac{\partial J(\theta)}{\partial \theta_{31}} &= \boxed{\delta_{31}}\end{aligned}$$

$$\frac{\partial J(\theta)}{\partial \theta_{32}}, \frac{\partial J(\theta)}{\partial \theta_{30}}, \partial \theta_{10}, \theta_{11}, \theta_{12}, \theta_{20}, \theta_{21}, \theta_{22}$$

$\alpha = 0.01$

$$\theta_{31} = \theta_{31} - \alpha (\delta_{31})$$

$$\boxed{\theta_{31} = 0.5 - (0.01) \times \delta_{31}}$$

θ_{30}
 θ_{31}
 θ_{11}
 θ_{12}
 θ_{20}
 θ_{21}
 θ_{22}

2nd iteration of
 - forward propagation
 - cost.
 - backward prop.