

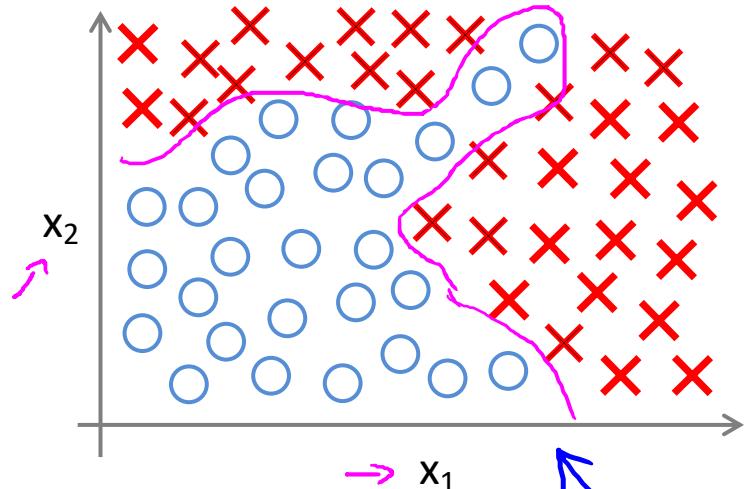
Machine Learning

# Neural Networks: Representation

---

## Non-linear hypotheses

# Non-linear Classification



$\rightarrow \underline{x_1}$  = size  
 $\underline{x_2}$  = # bedrooms  
 $\underline{x_3}$  = # floors  
 $x_4$  = age  
 $\dots$   
 $x_{100}$  -

$\left. \right\} h=100$

$$\begin{aligned}
 & \downarrow \\
 & g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 \\
 & + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 \\
 & + \theta_5 x_1^3 x_2 + \underline{\theta_6 x_1 x_2^2} + \dots)
 \end{aligned}$$

$$\begin{aligned}
 & \rightarrow \underline{x_1^2}, \underline{x_1 x_2}, \underline{x_1 x_3}, \underline{x_1 x_4} \dots \underline{x_1 x_{100}} \\
 & \underline{x_2^2}, \underline{x_2 x_3} \dots
 \end{aligned}$$

$\approx 5000$  feature

$$\begin{aligned}
 & O(n^2) \\
 & \frac{n^2}{2} \\
 & 10
 \end{aligned}$$

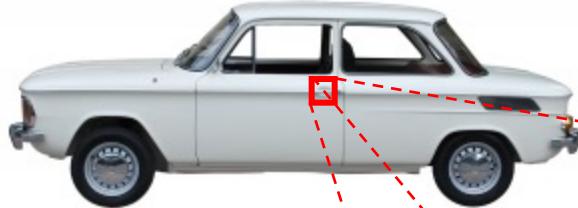
$$\rightarrow \underline{x_1^2}, \underline{x_2^2}, \underline{x_3^2}, \dots, \underline{x_{100}^2}$$

$$\begin{aligned}
 & \rightarrow \underline{x_1 x_2 x_3}, \underline{x_1^2 x_2}, \underline{x_{10} x_{11} x_{12}}, \dots \\
 & O(n^3)
 \end{aligned}$$

170,000

# 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

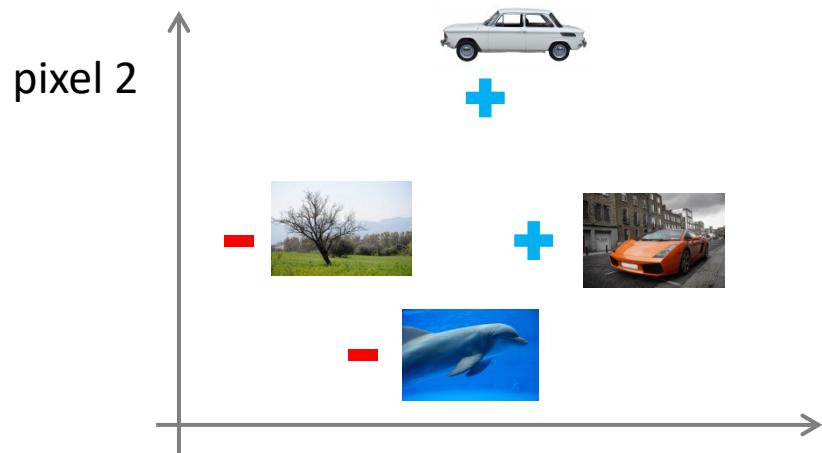
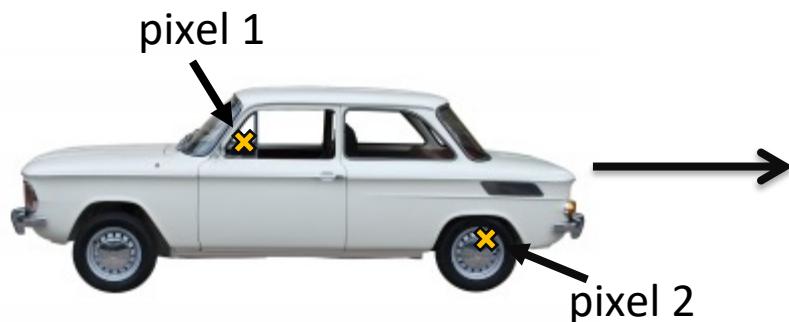


Not a car

Testing:

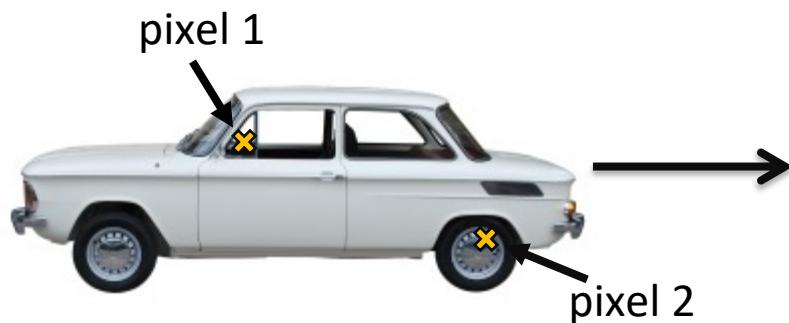


What is this?

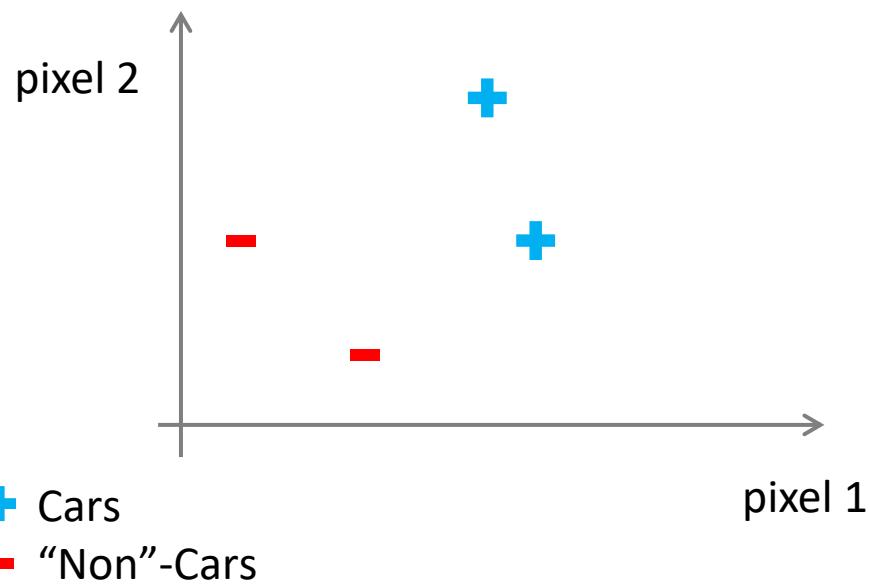


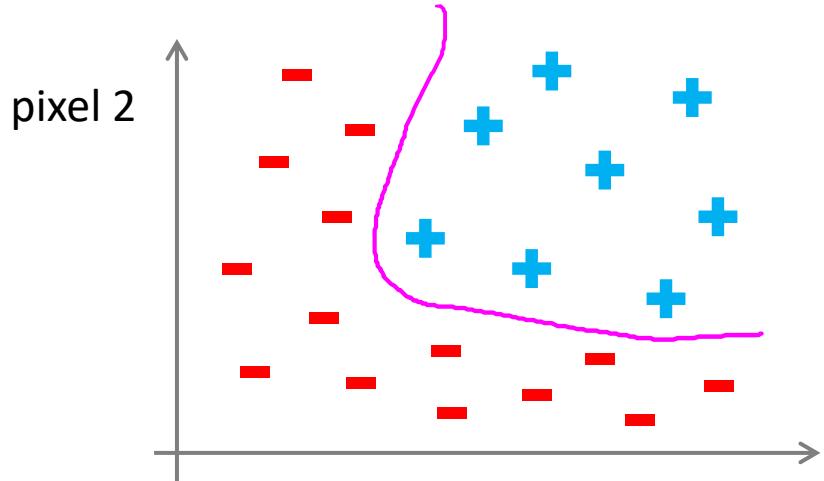
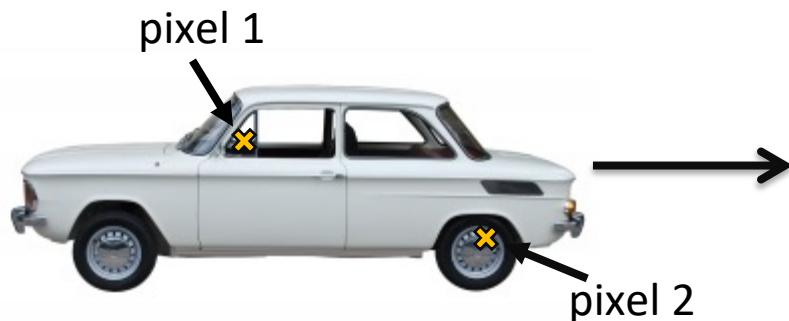
**+** Cars

**-** "Non"-Cars



Learning  
Algorithm





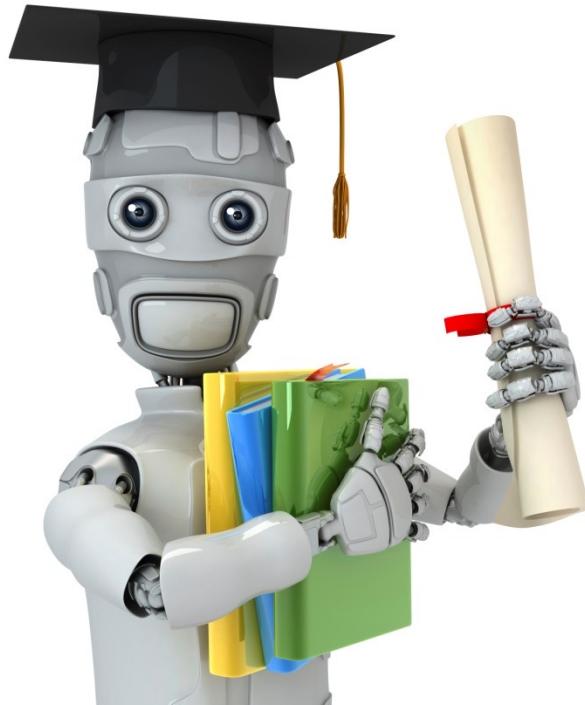
$50 \times 50$  pixel images  $\rightarrow$  2500 pixels  
 $n = 2500$  (7500 if RGB)

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

0 - 255

+ Cars  
- "Non"-Cars

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



Machine Learning

# Neural Networks: Representation

---

## Neurons and the brain

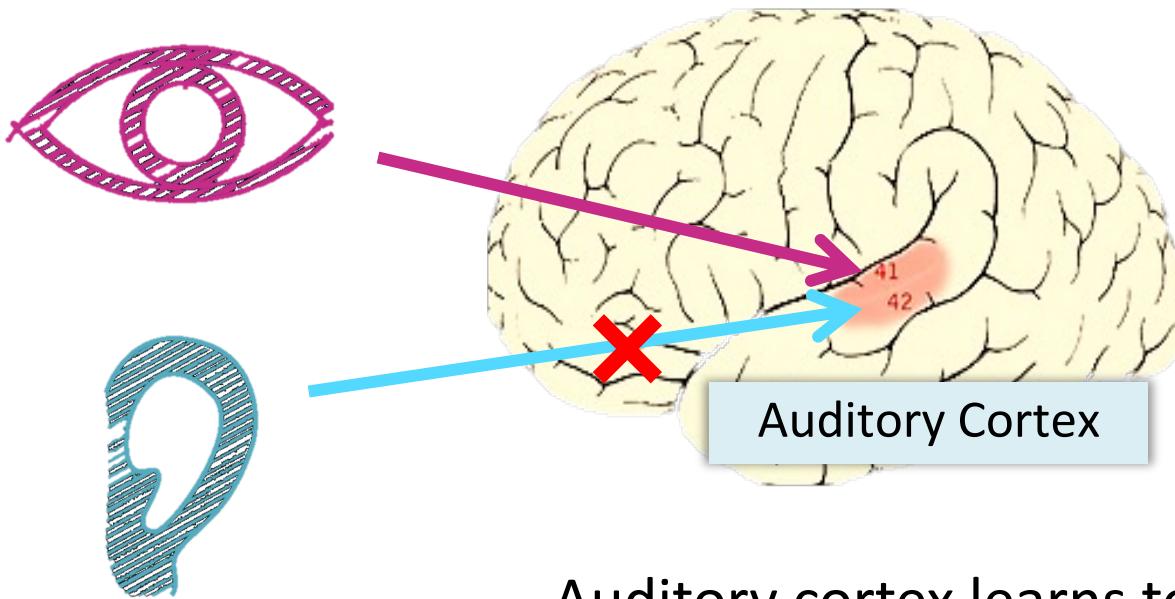
# Neural Networks

Origins: Algorithms that try to mimic the brain.

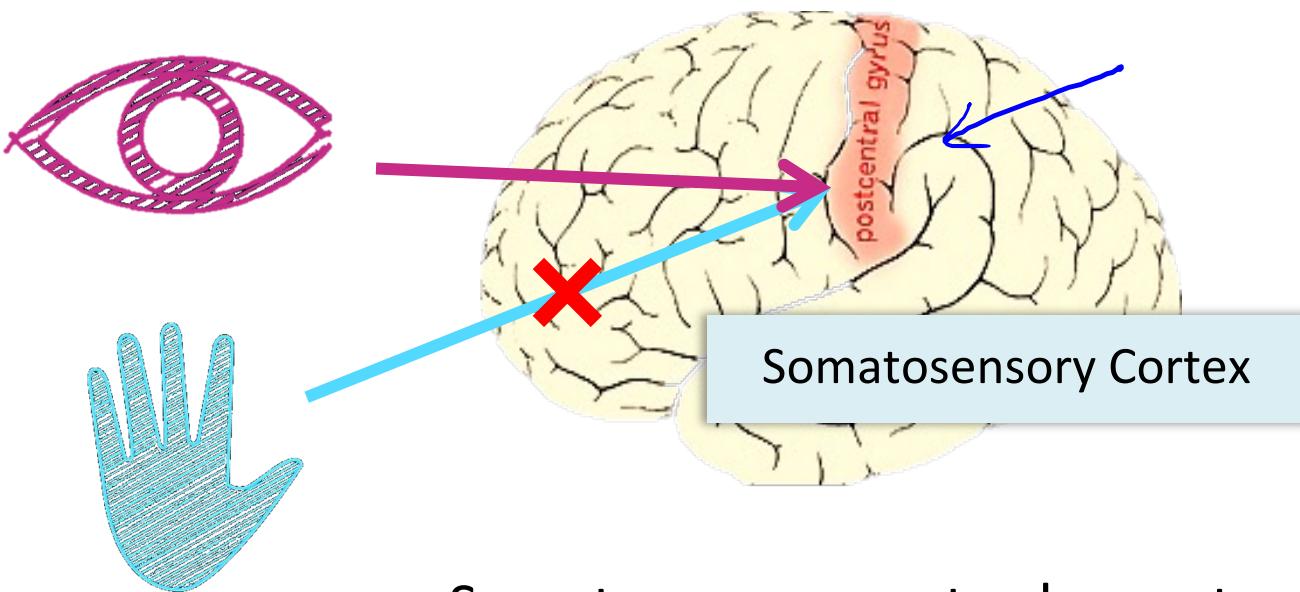
Was very widely used in 80s and early 90s; popularity diminished in late 90s.

Recent resurgence: State-of-the-art technique for many applications

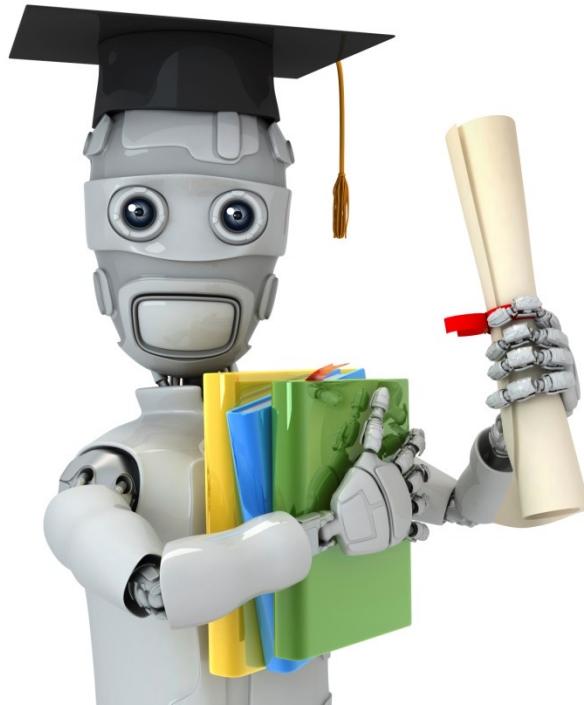
# The “one learning algorithm” hypothesis



# The “one learning algorithm” hypothesis



Somatosensory cortex learns to see



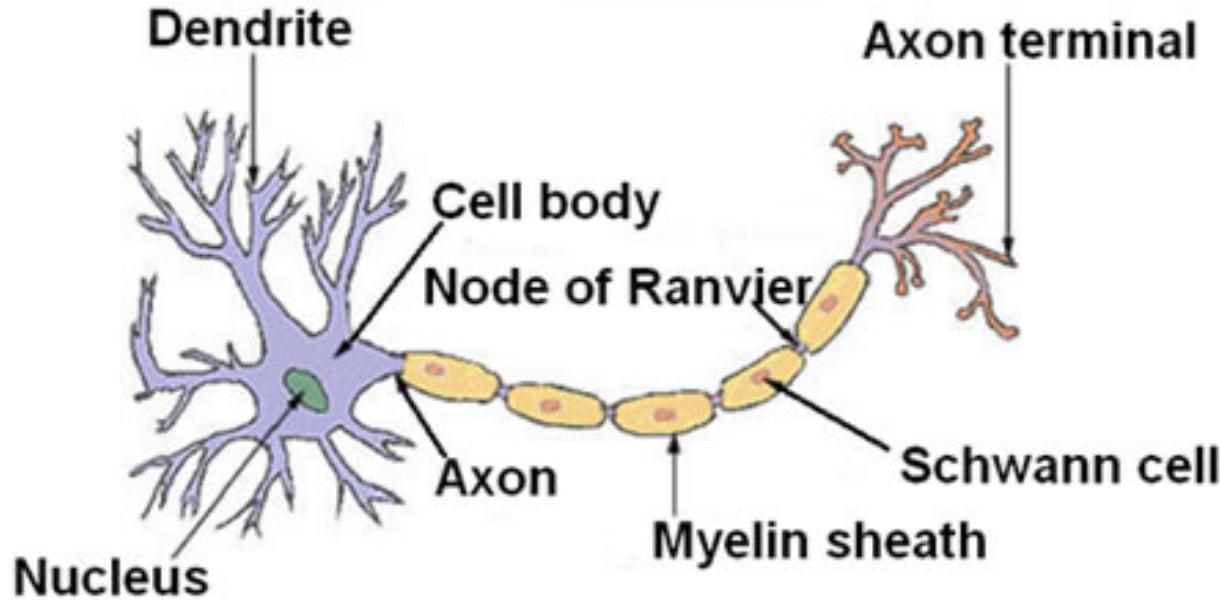
Machine Learning

# Neural Networks: Representation

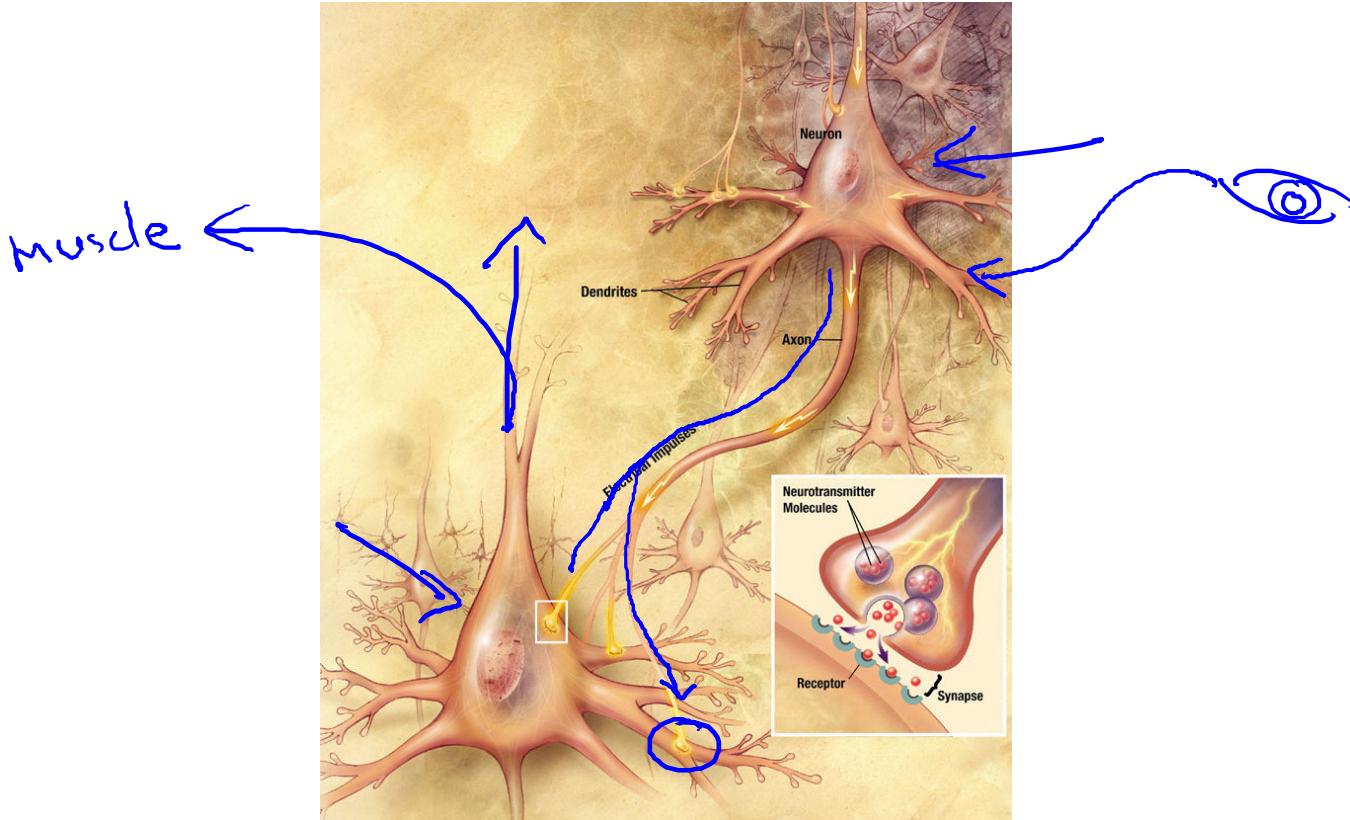
---

## Model representation I

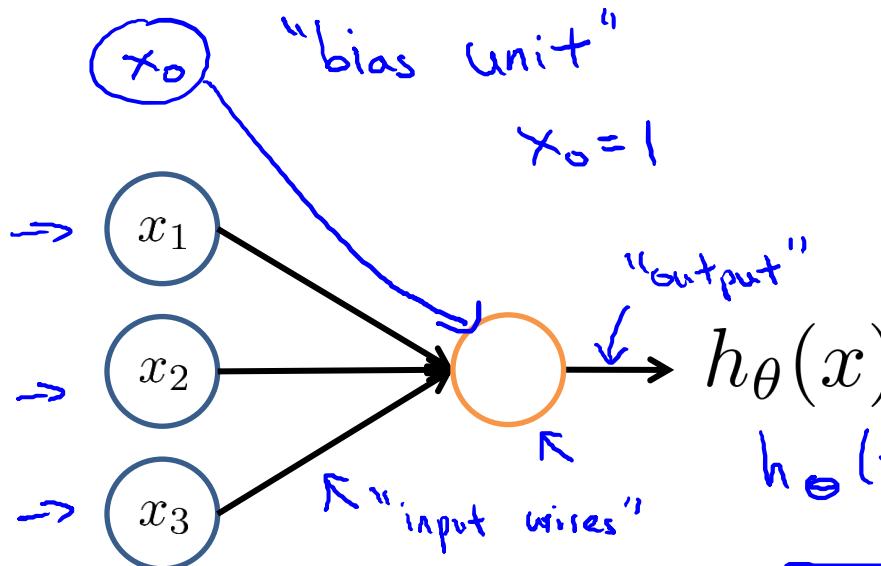
# Neuron in the brain



# Neurons in the brain



# Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

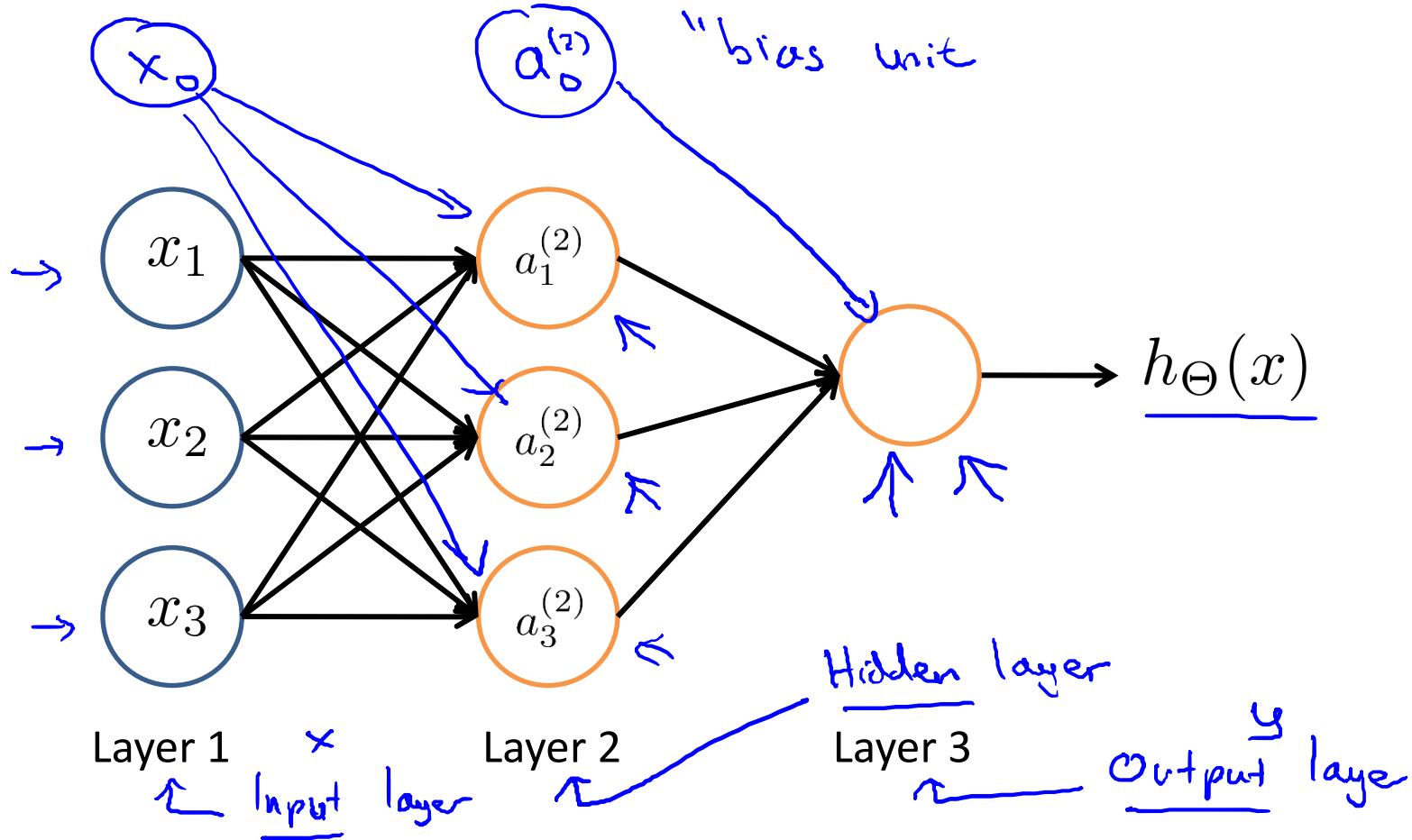
↑  
"weights" ←  
(parameters ←

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

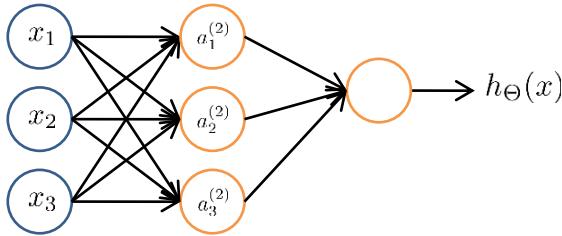
Sigmoid (logistic) activation function.

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

# Neural Network



# Neural Network



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

$\Theta^{(j)}$  = matrix of weights controlling  
function mapping from layer  $j$  to  
layer  $j + 1$

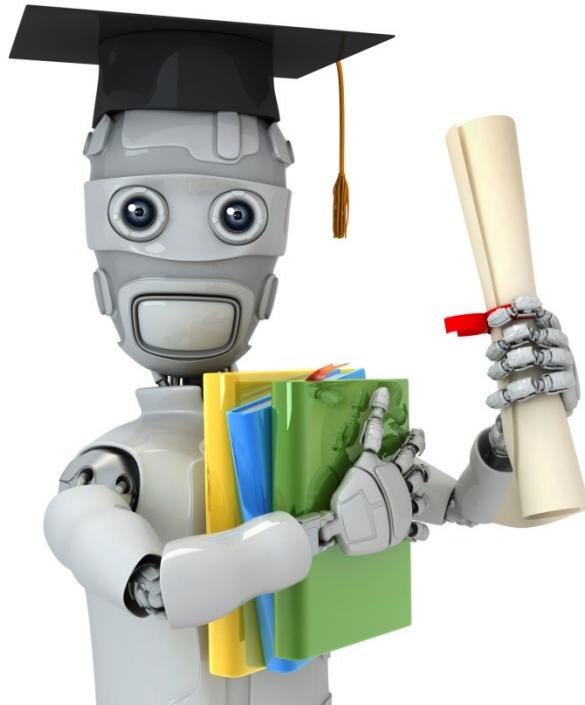
$$a_1^{(2)} = g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3)$$

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

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

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

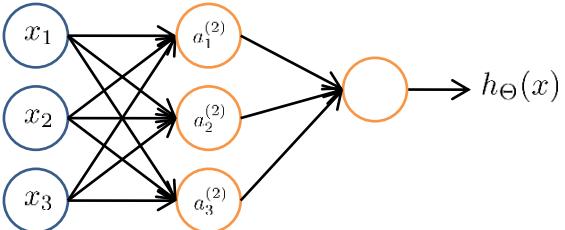
If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j + 1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .



Machine Learning

# Neural Networks: Representation --- Model representation II

# Forward propagation: Vectorized implementation



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

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

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

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

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

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

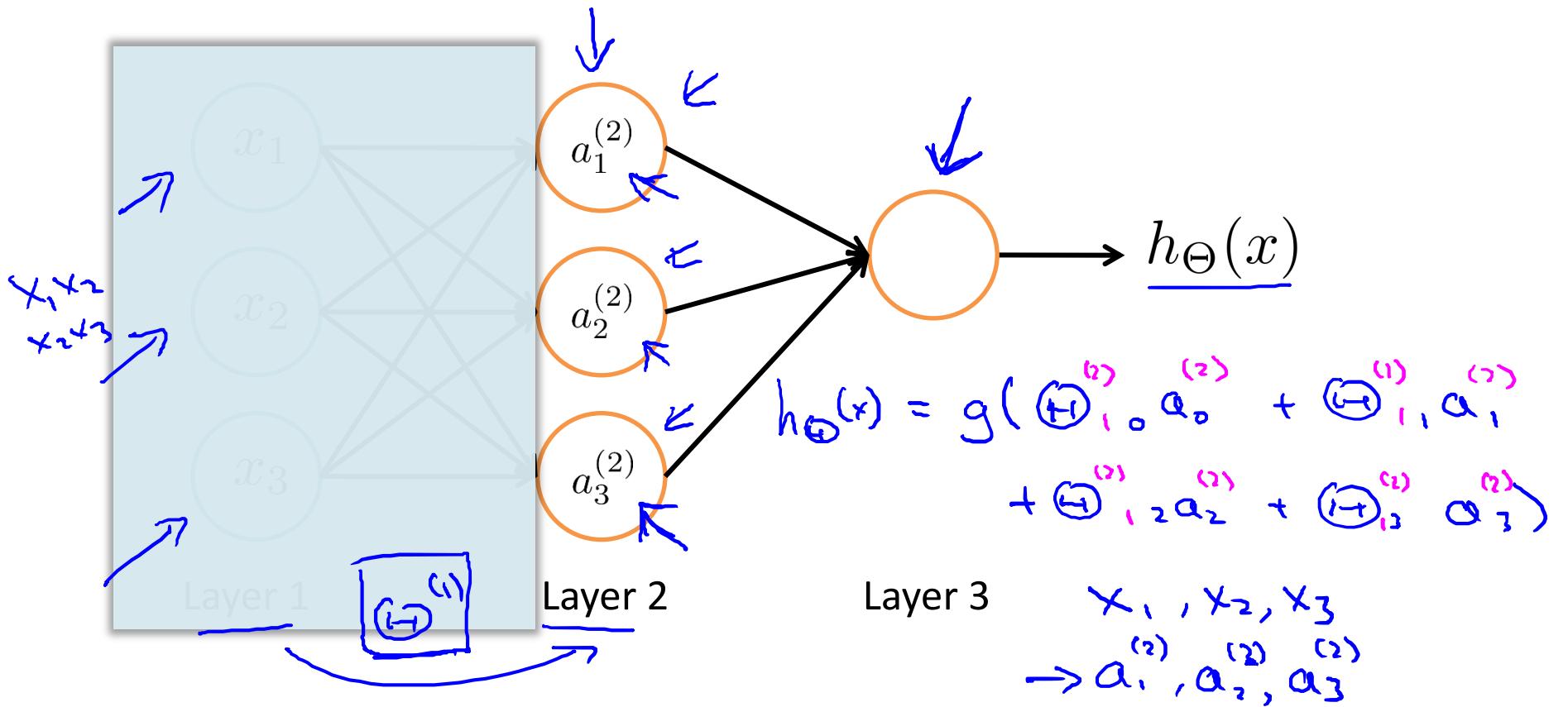
$$a^{(2)} = g(z^{(2)})$$

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

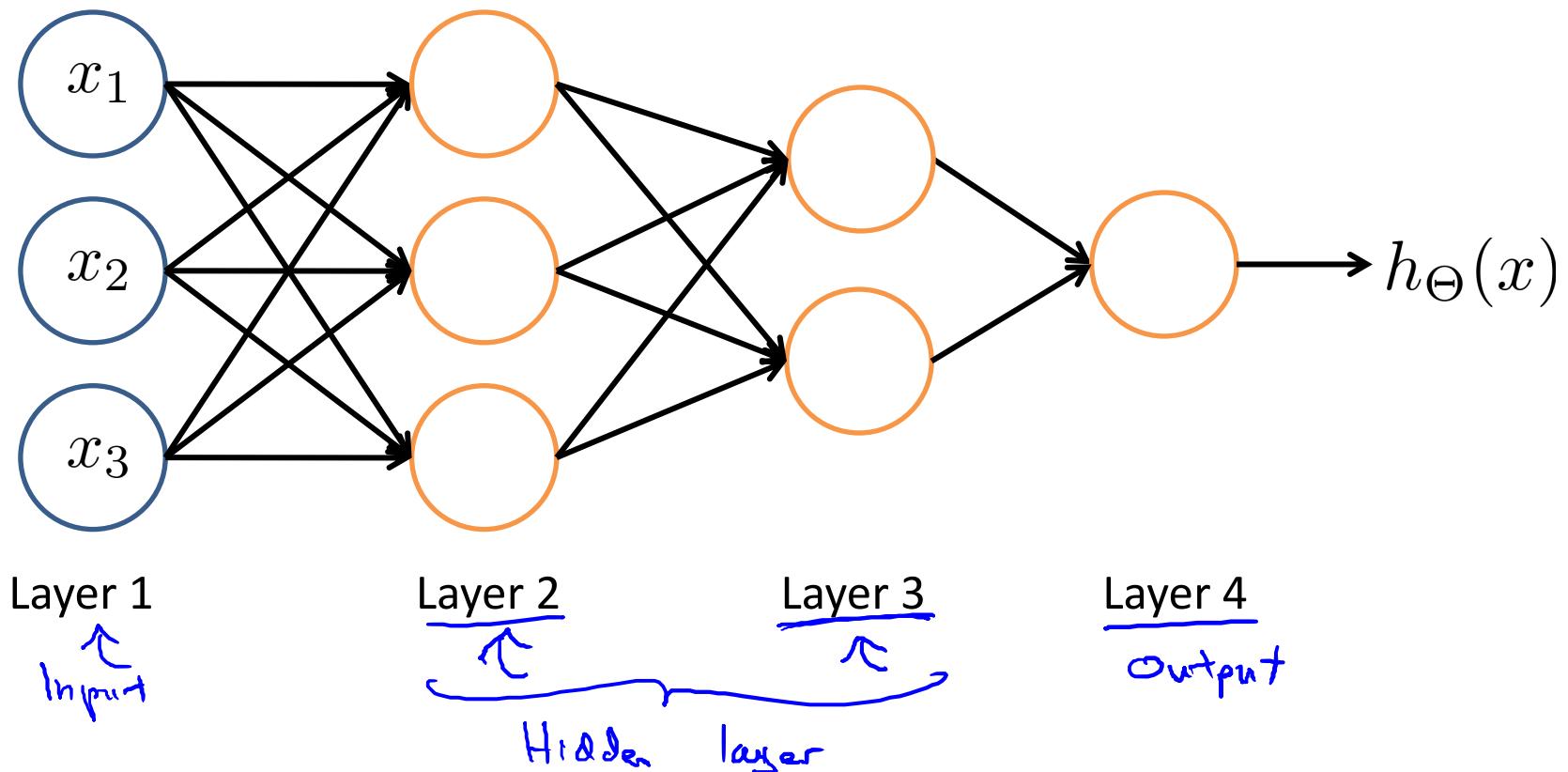
$$z^{(3)} = \Theta^{(2)} a^{(2)}$$

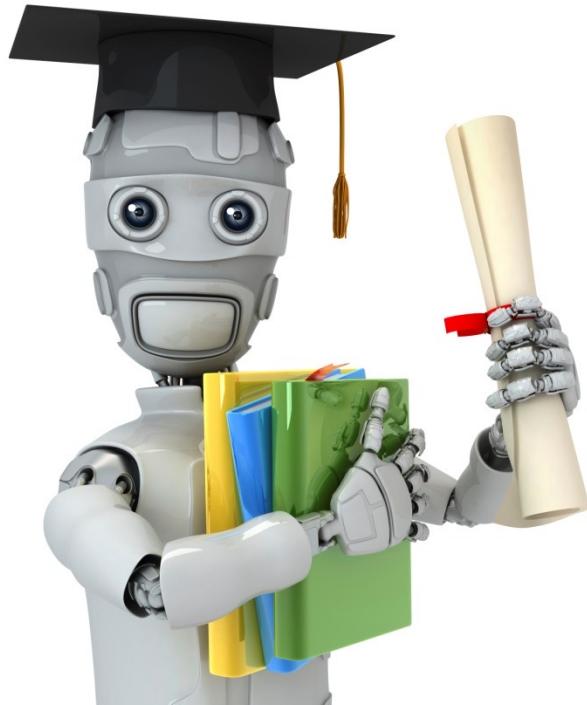
$$h_\Theta(x) = a^{(3)} = g(z^{(3)})$$

# Neural Network learning its own features



## Other network architectures





Machine Learning

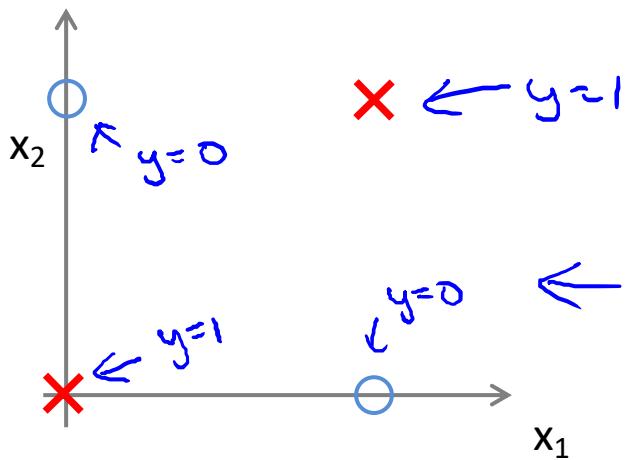
# Neural Networks: Representation

---

## Examples and intuitions I

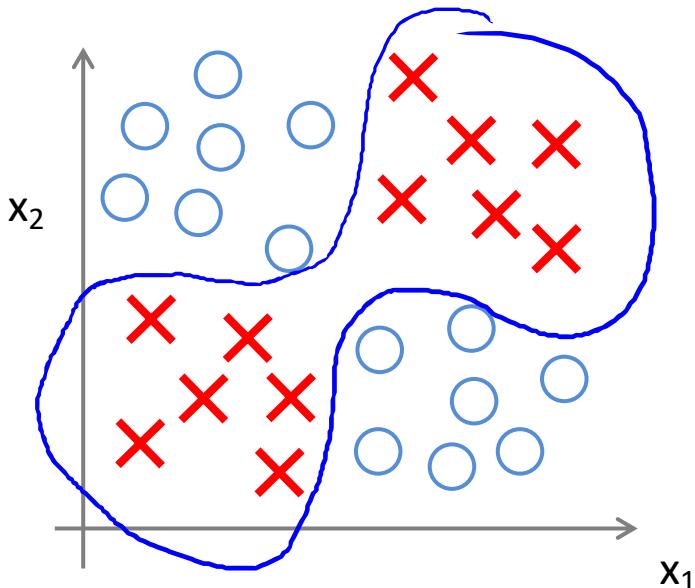
## Non-linear classification example: XOR/XNOR

→  $x_1, x_2$  are binary (0 or 1).



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

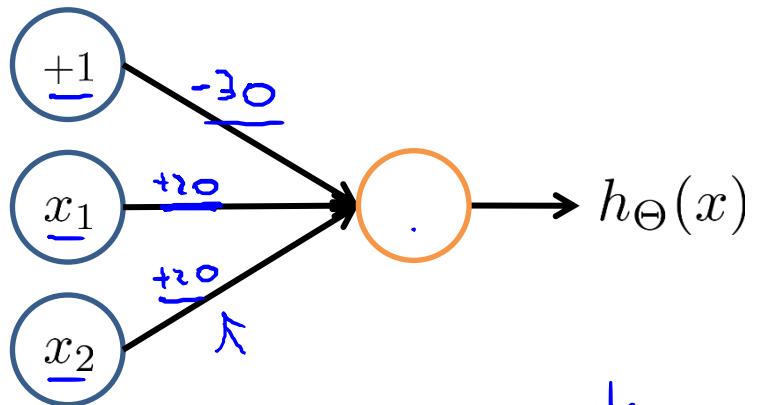
→  $\underline{x_1 \text{ XNOR } x_2}$  ←  
→ NOT (x<sub>1</sub> XOR x<sub>2</sub>)



## Simple example: AND

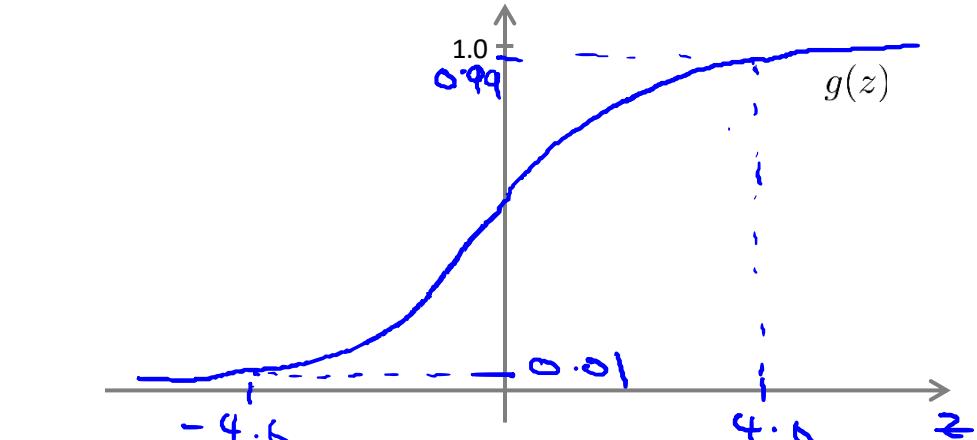
$$\rightarrow x_1, x_2 \in \{0, 1\}$$

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



$$\rightarrow h_{\Theta}(x) = g\left(\frac{-30}{\pi} + \frac{20}{\pi}x_1 + \frac{20}{\pi}x_2\right)$$

$\Theta_{1,0}$        $\Theta_{1,1}$        $\Theta_{1,2}$

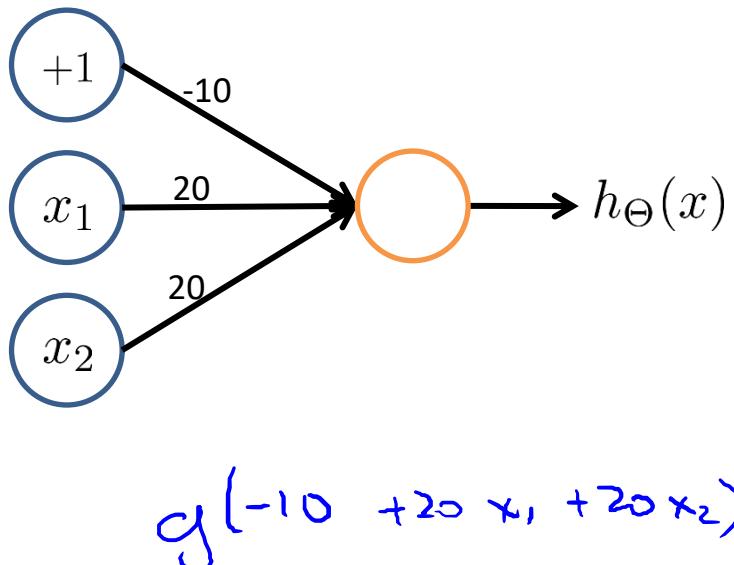


$\leftarrow$

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

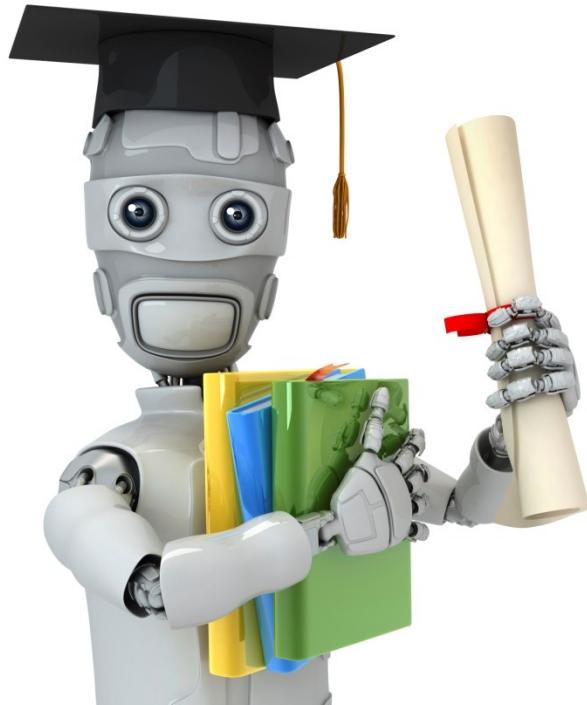
$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$

## 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	$\approx 1$

The table shows the output of the function  $h_{\Theta}(x)$  for all combinations of  $x_1$  and  $x_2$ . The values are approximated as follows:  $g(-10) \approx 0$ ,  $g(10) \approx 1$ , and both  $x_1=1$  and  $x_2=1$  result in  $\approx 1$ . A pink hand-drawn bracket highlights the row where  $x_1=1$  and  $x_2=1$ .



Machine Learning

# Neural Networks: Representation

---

## Examples and intuitions II

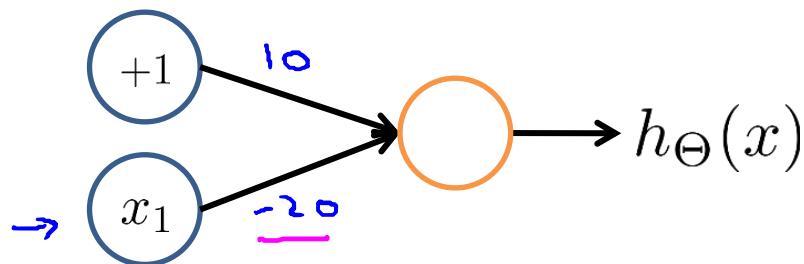
$\rightarrow x_1 \text{ AND } x_2$

$\rightarrow x_1 \text{ OR } x_2$

$\{0, 1\}$ .

## Negation:

NOT  $x_1$



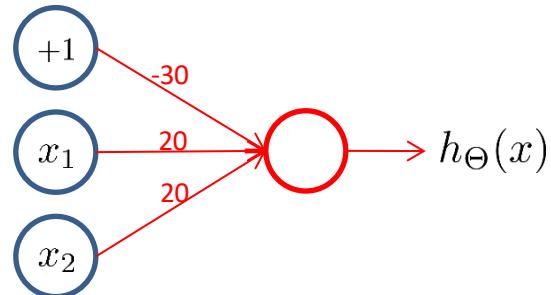
$x_1$	$h_{\Theta}(x)$
0	$g(10) \approx 1$
1	$g(-20) \approx 0$

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

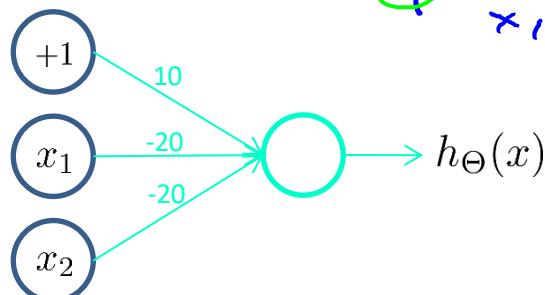
$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

$\begin{cases} = 1 & \text{if and only if} \\ & \rightarrow x_1 = x_2 = 0 \end{cases}$

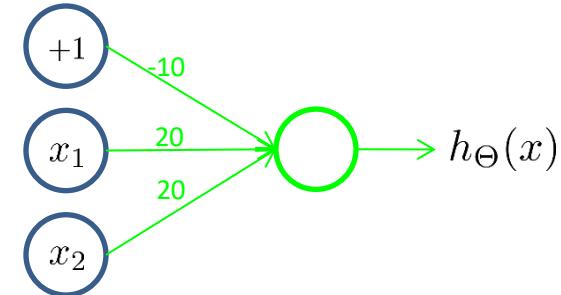
## Putting it together: $x_1$ XNOR $x_2$



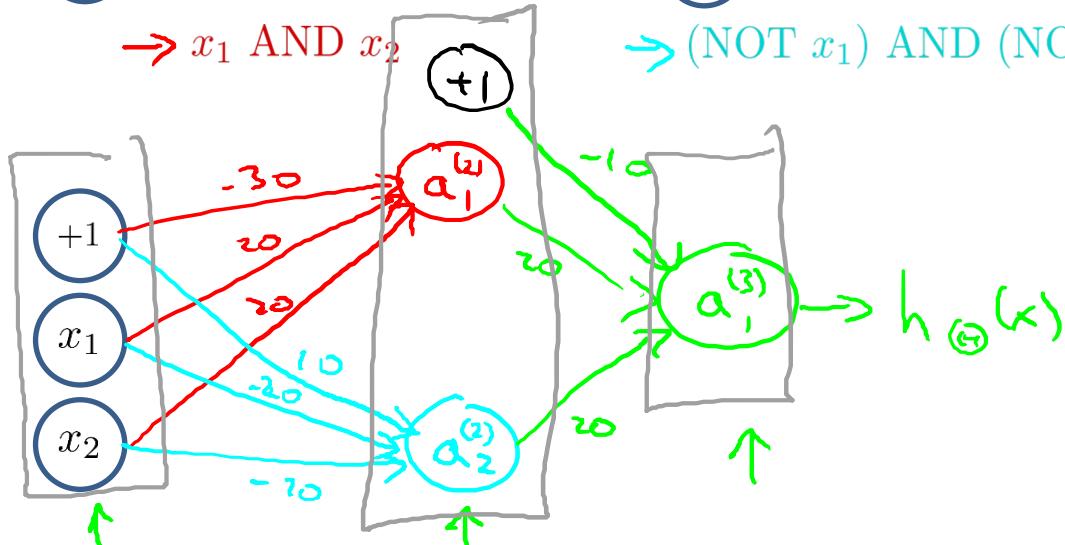
$\rightarrow x_1$  AND  $x_2$



$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

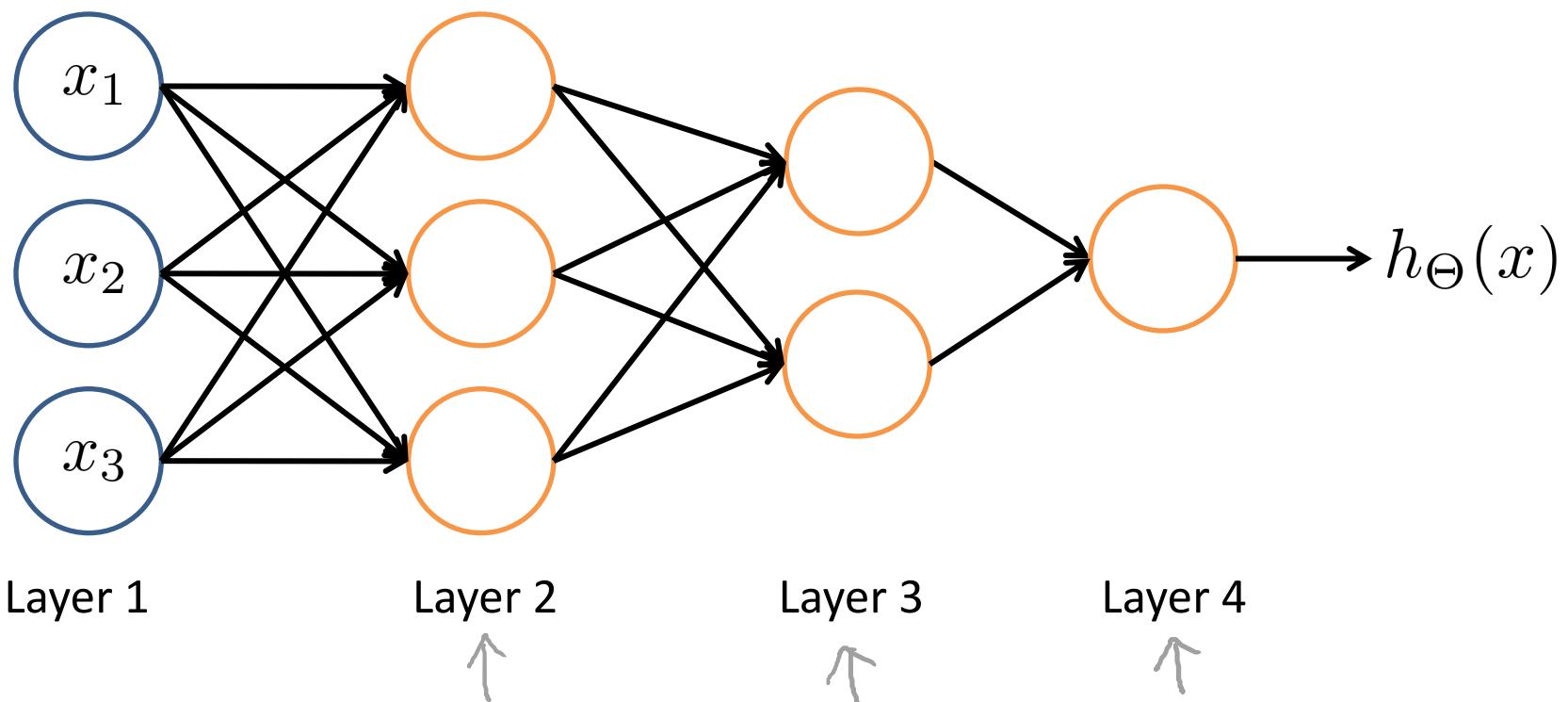


$\rightarrow x_1$  OR  $x_2$

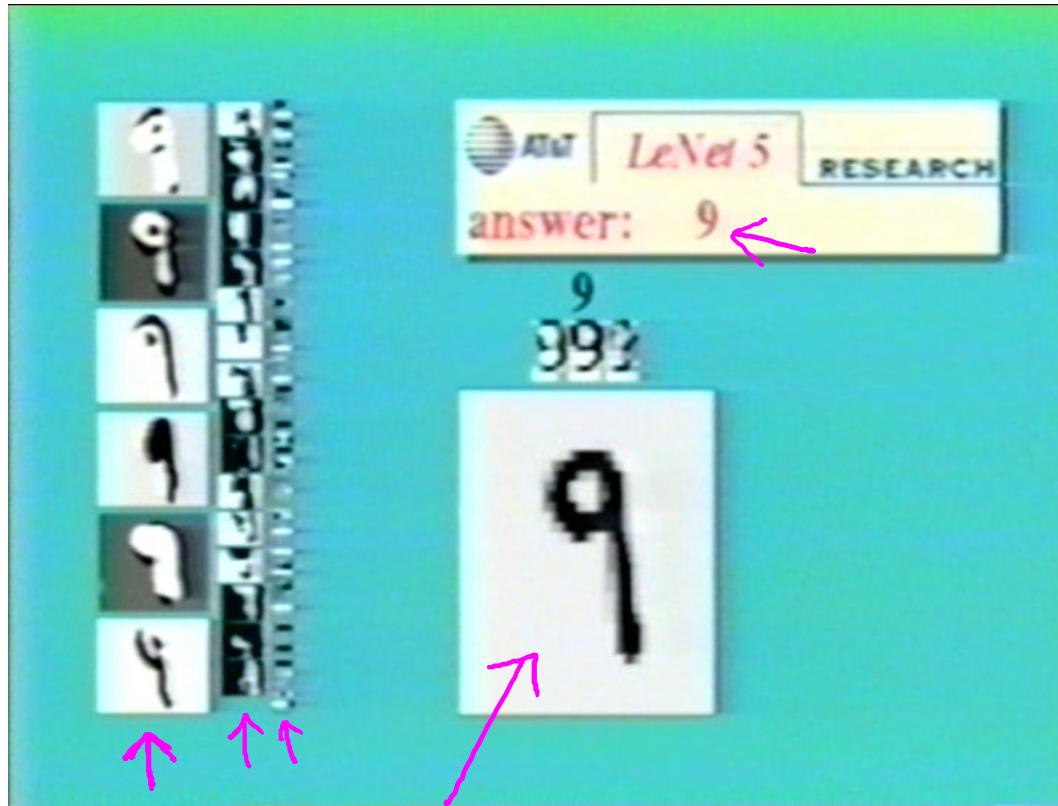


$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1 ←
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1 ←

# Neural Network intuition



# Handwritten digit classification

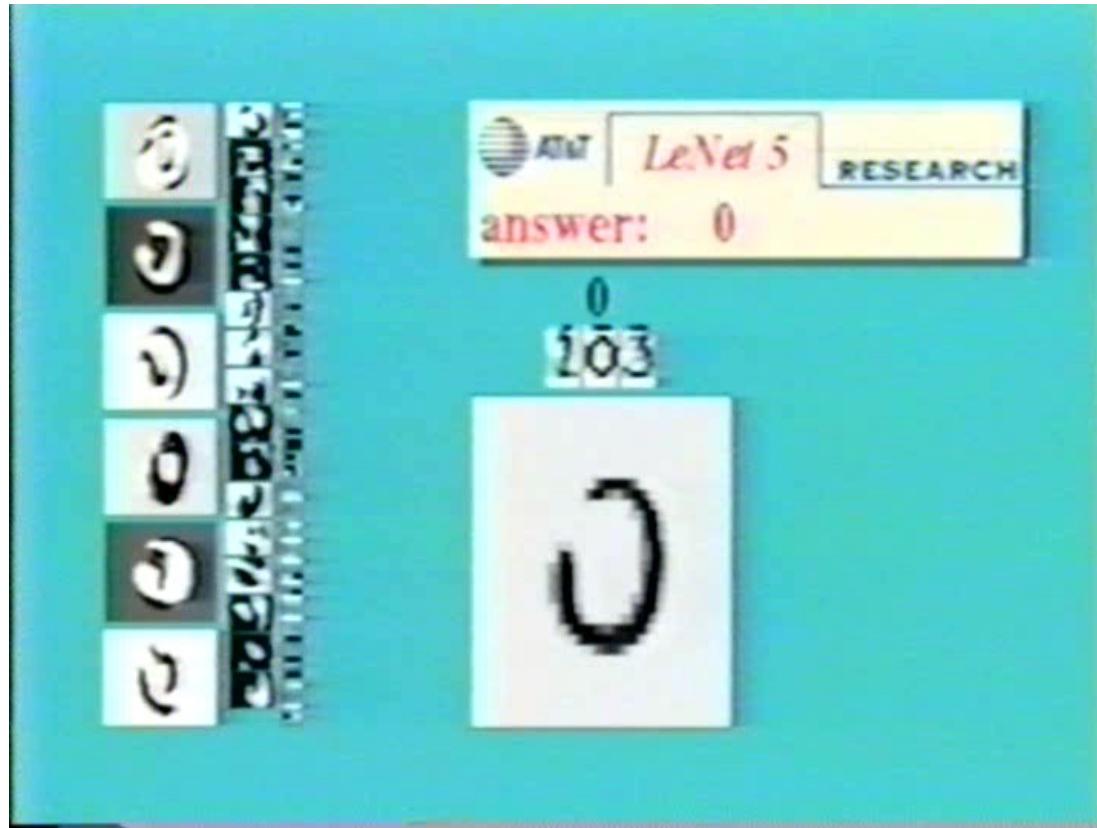


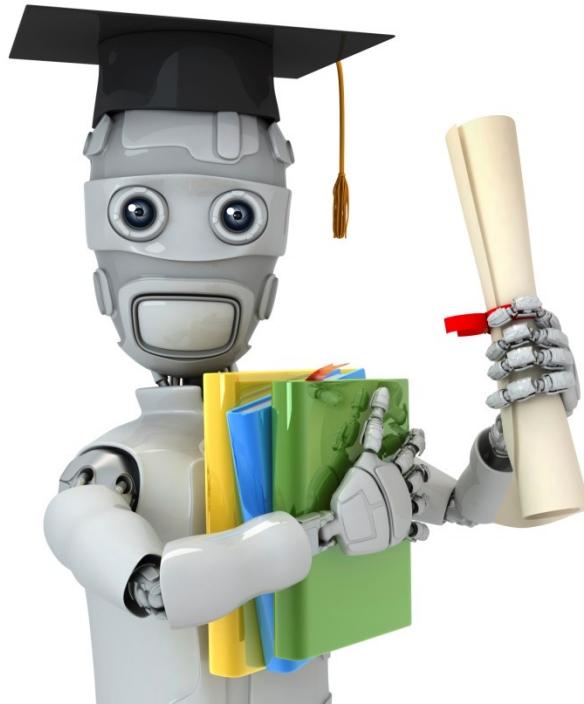
[Courtesy of Yann LeCun]



Andrew Ng

# Handwritten digit classification





Machine Learning

# Neural Networks: Representation

---

## Multi-class classification

# Multiple output units: One-vs-all.



Pedestrian



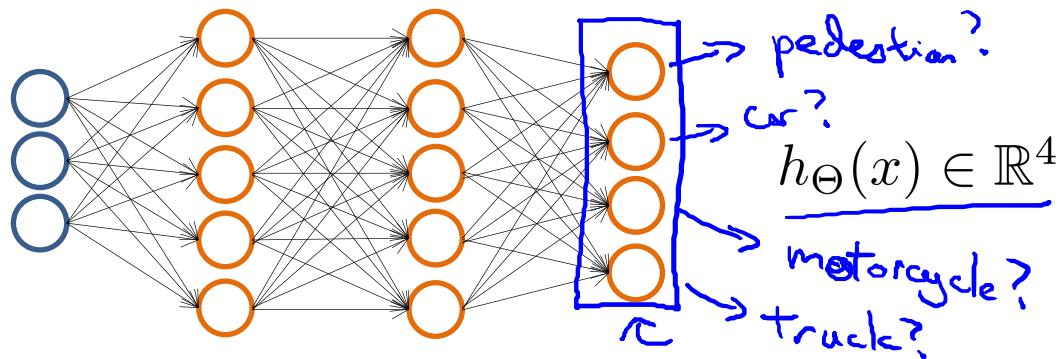
Car



Motorcycle



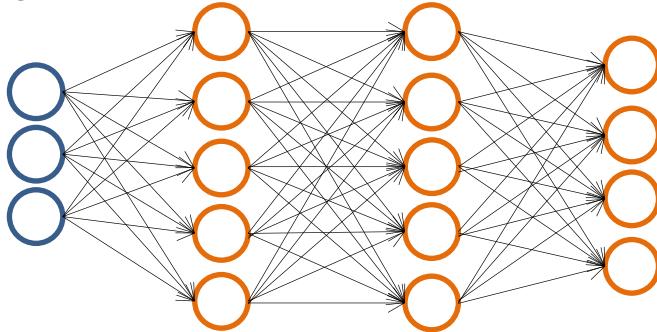
Truck



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 pedestrian                          when car                          when motorcycle

## Multiple output units: One-vs-all.



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

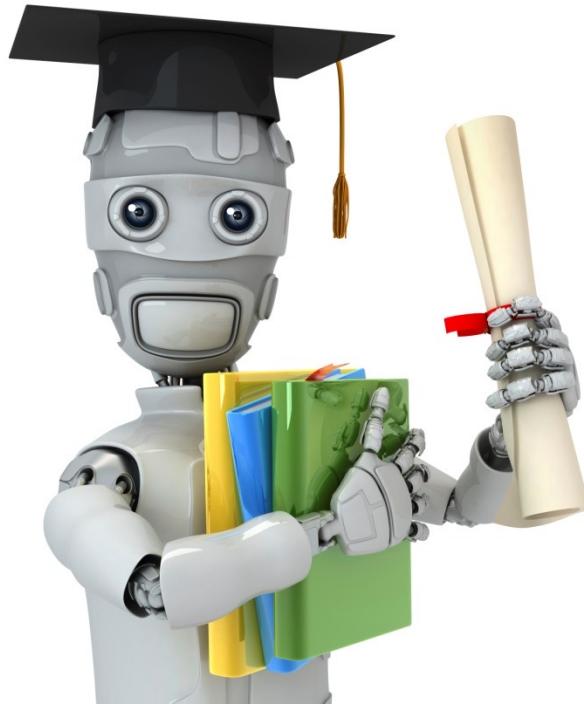
Want  $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ , etc.  
when pedestrian      when car      when motorcycle

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

→  $y^{(i)}$  one of  
pedestrian      car      motorcycle      truck

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

~~Previously~~  
 $y \in \{1, 2, 3, 4\}$   
 $\underline{h_{\Theta}(x^{(i)}) \approx y^{(i)}} = \mathbb{R}^4$



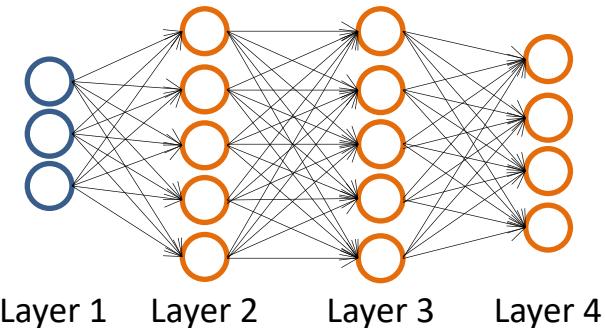
Machine Learning

# 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

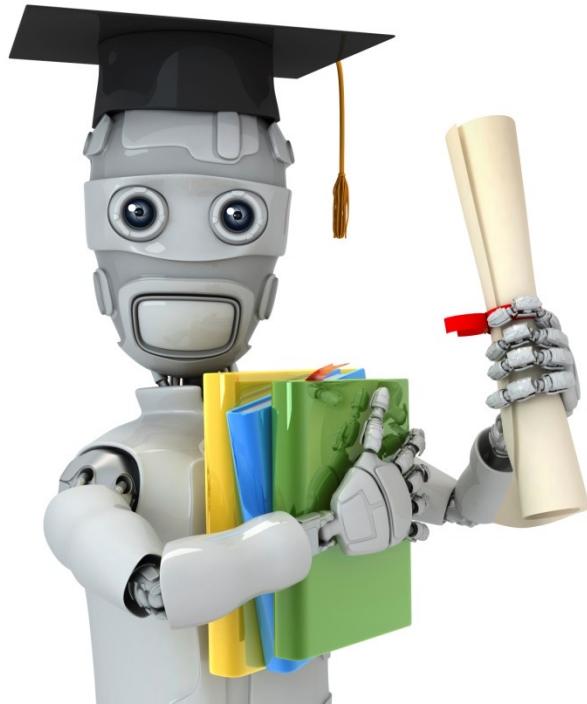
Logistic regression:

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

Neural network:

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

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



Machine Learning

# Neural Networks: Learning

---

## Backpropagation algorithm

# Gradient computation

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

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

Need code to compute:

$$\rightarrow -\underline{J(\Theta)}$$
  
$$\rightarrow -\underline{\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)} \quad \leftarrow$$

$$H_{ij}^{(l)} \in \mathbb{R}$$

# Gradient computation

Given one training example  $(\underline{x}, \underline{y})$ :

Forward propagation:

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

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

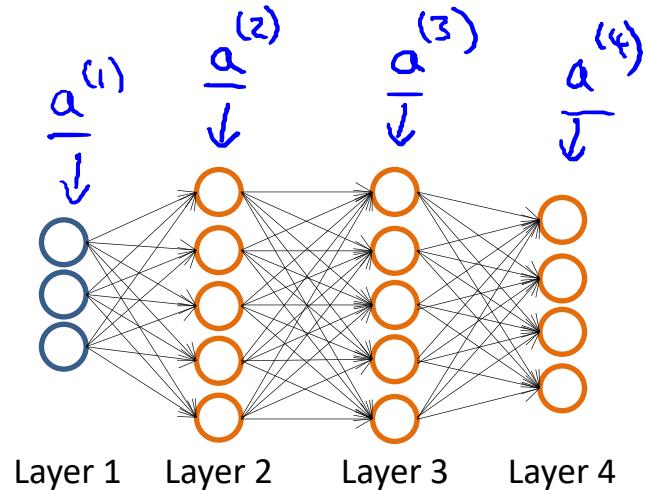
$$\rightarrow \underline{a}^{(2)} = g(\underline{z}^{(2)}) \quad (\text{add } \underline{a}_0^{(2)})$$

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

$$\rightarrow \underline{a}^{(3)} = g(\underline{z}^{(3)}) \quad (\text{add } \underline{a}_0^{(3)})$$

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

$$\rightarrow \underline{a}^{(4)} = \underline{h}_{\Theta}(\underline{x}) = g(\underline{z}^{(4)})$$

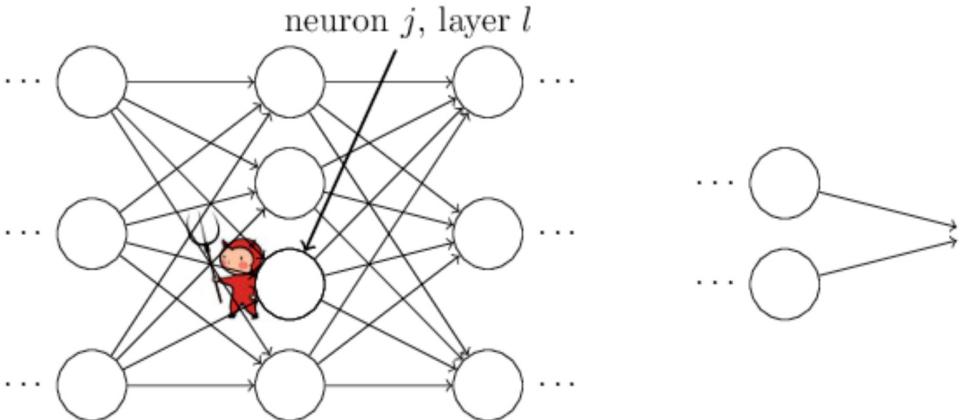


## Gradient computation: Error

- Adds a little change  $\Delta z_i^l$
- Instead of outputting  $g(z_i^l)$ , neuron output  $g(z_i^l + \Delta z_i^l)$
- This change causing the overall cost to change by:

$$\frac{\partial J}{\partial z_i^l} \Delta z_i^l$$

- Define error  $\delta_j^l = \frac{\partial J}{\partial z_i^l} = \frac{\partial J}{\partial a_i^l} g'(z_i^l)$

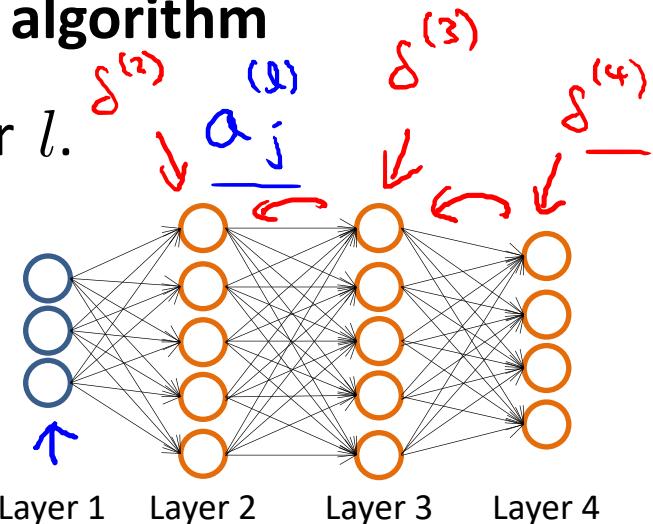


# Gradient computation: Backpropagation algorithm

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

For each output unit (layer  $L = 4$ )

$$\underline{\delta_j^{(4)}} = \underline{a_j^{(4)}} - \underline{y_j} \quad (\underline{h_{\Theta}(x)})_j \quad \underline{\delta^{(4)}} = \underline{a} - \underline{y}$$



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

$$\frac{a^{(3)}}{a^{(3)}} * \frac{(1-a^{(3)})}{(1-a^{(3)})}$$

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

$$\frac{a^{(2)}}{a^{(2)}} * \frac{(1-a^{(2)})}{(1-a^{(2)})}$$

(No  $\delta^{(1)}$ )

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(ignoring  $\lambda$ ; if  
 $\lambda = 0$ ) ↵

# Backpropagation algorithm

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ). (use to compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ )

For  $i = 1$  to  $m$  ←  $(\underline{x}^{(i)}, \underline{y}^{(i)})$ .

Set  $\underline{a}^{(1)} = \underline{x}^{(i)}$

→ Perform forward propagation to compute  $\underline{a}^{(l)}$  for  $l = 2, 3, \dots, L$

→ Using  $\underline{y}^{(i)}$ , compute  $\delta^{(L)} = \underline{a}^{(L)} - \underline{y}^{(i)}$

→ Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

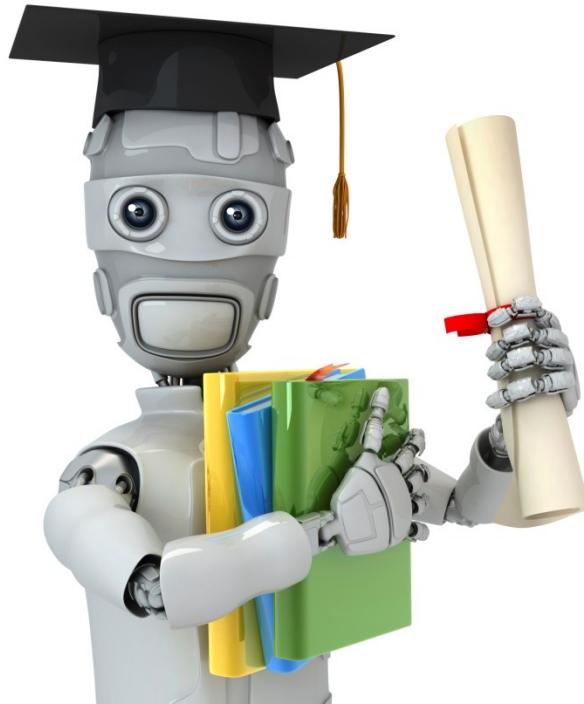
$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \delta^{(l+1)} (\underline{a}^{(l)})^T$ .

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$  if  $j \neq 0$

$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$  if  $j = 0$

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

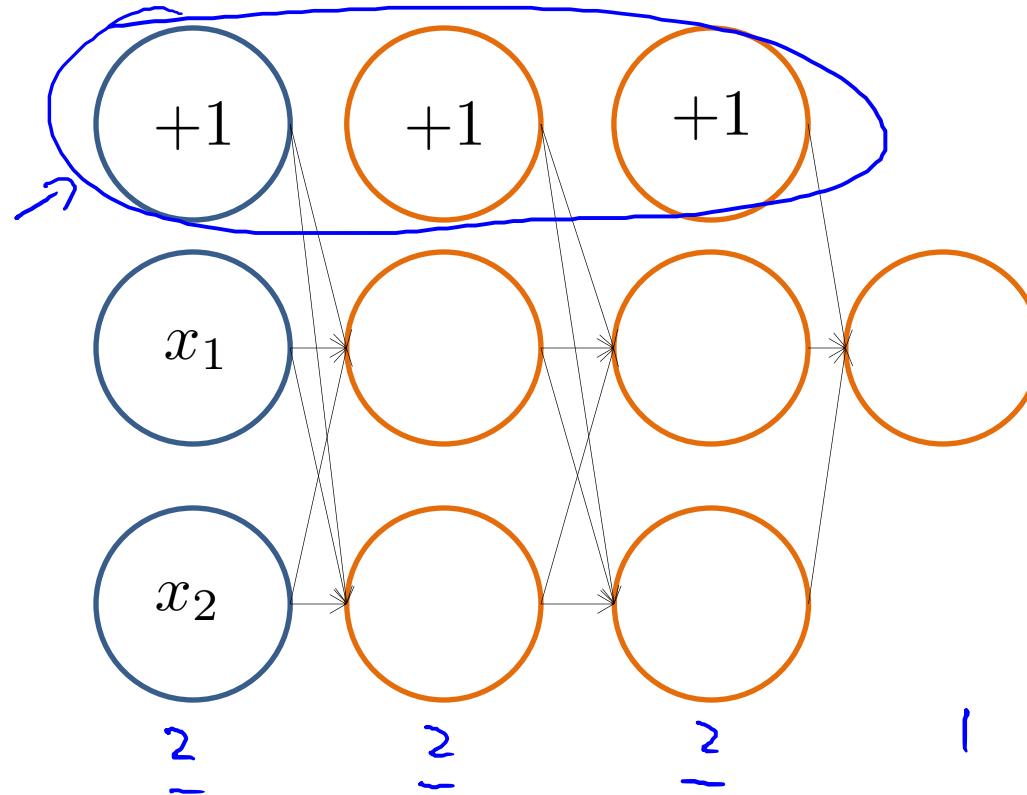


Machine Learning

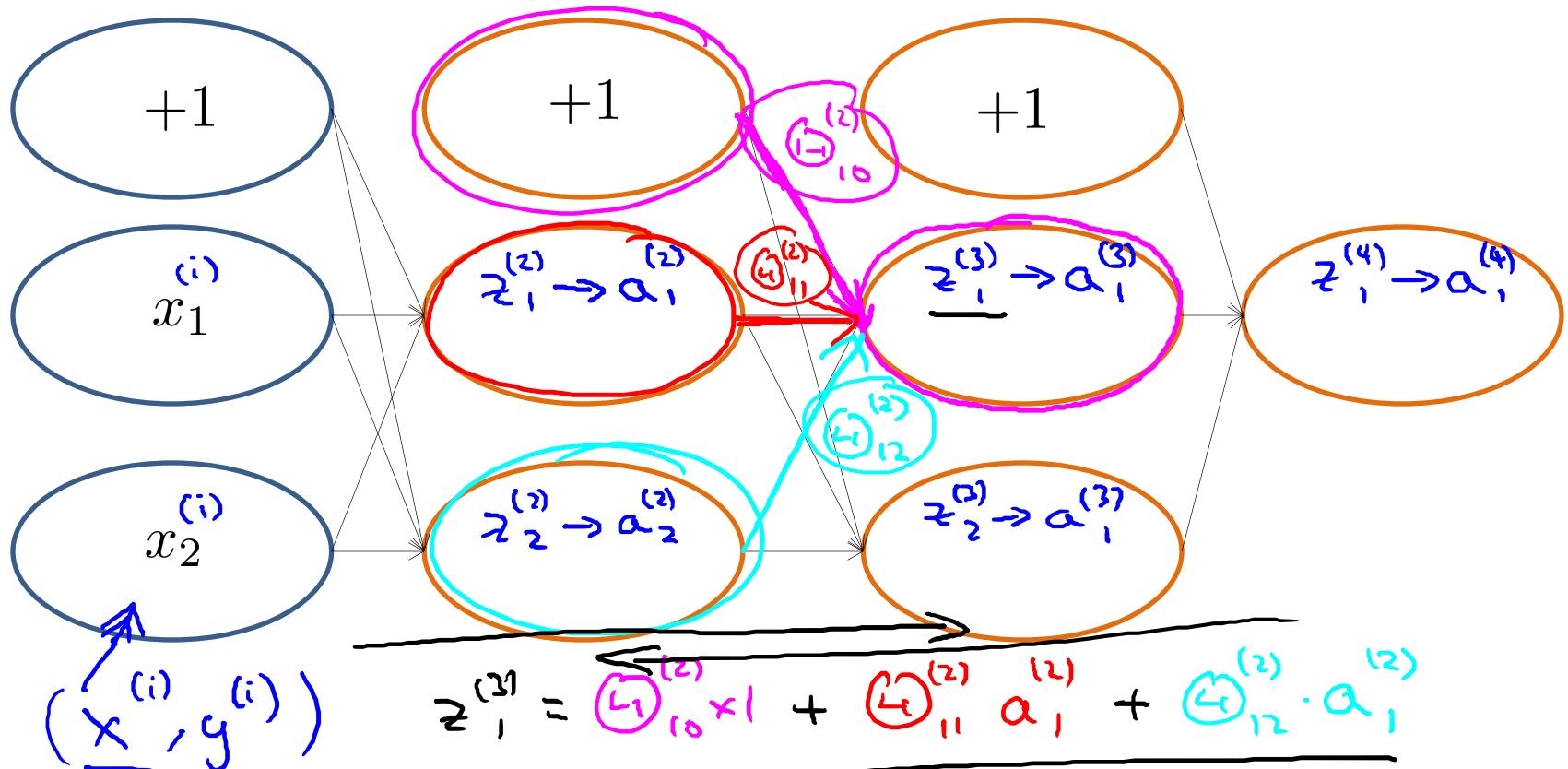
# Neural Networks: Learning

## Backpropagation intuition

## Forward Propagation



# Forward Propagation



# What is backpropagation doing?

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

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

Focusing on a single example  $x^{(i)}$ ,  $y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda = 0$ ),

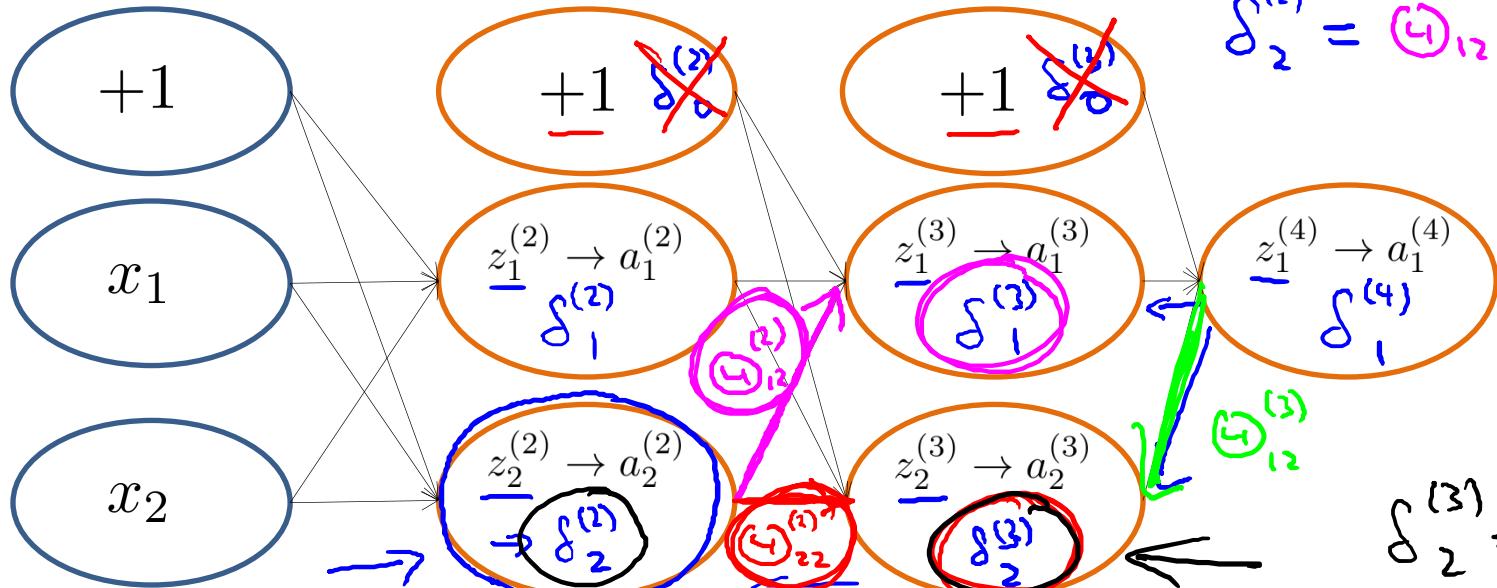
$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log h_\Theta(x^{(i)})$$

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

i.e. how well is the network doing on example i?

$y^{(i)}$

## Forward Propagation



$\rightarrow \delta_j^{(l)}$  = “error” of cost for  $a_j^{(l)}$  (unit  $j$  in layer  $l$ ).

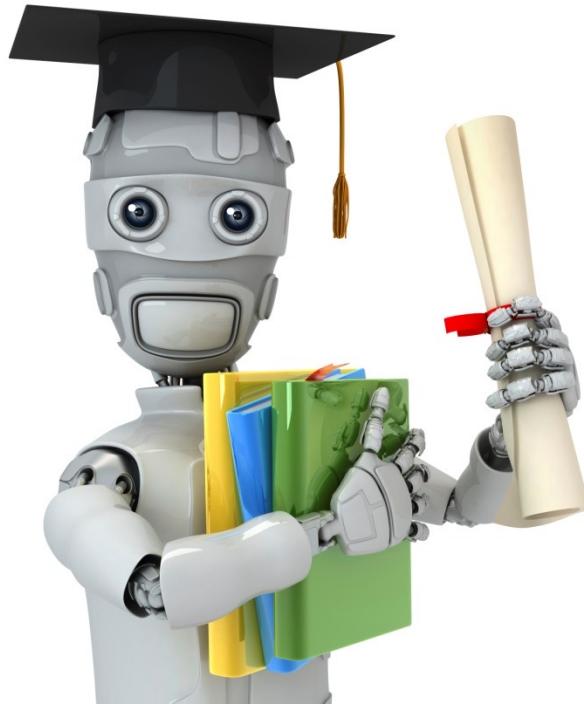
Formally,  $\delta_j^{(l)} = \frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$  (for  $j \geq 0$ ), where

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

$$\delta_1^{(4)} = \underline{y^{(i)}} - \underline{a_1^{(4)}}$$

$$\delta_2^{(2)} = (\underline{4})_{12} \delta_1^{(3)} + (\underline{4})_{22} \cdot \delta_2^{(3)}$$

$$\delta_2^{(3)} = (\underline{4})_{12} \cdot \delta_1^{(4)}.$$



Machine Learning

# Neural Networks: Learning

---

## Random initialization

## Initial value of $\Theta$

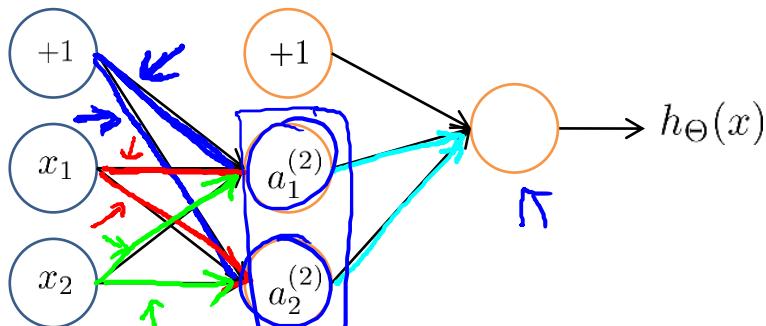
For gradient descent and advanced optimization method, need initial value for  $\Theta$ .

```
optTheta = fminunc(@costFunction,  
                    initialTheta, options)
```

Consider gradient descent

Set initialTheta = zeros(n,1) ?

## Zero initialization



$$\Rightarrow \Theta_{ij}^{(l)} = 0 \text{ for all } i, j, l.$$

Also  $\delta_1^{(2)} = \delta_2^{(2)}$ .

$$\frac{\partial}{\partial \Theta_{01}^{(2)}} J(\Theta) = \frac{\partial}{\partial \Theta_{02}^{(2)}} J(\Theta)$$

$$\underline{\Theta_{01}^{(1)}} = \underline{\Theta_{02}^{(1)}}$$

After each update, parameters corresponding to inputs going into each of two hidden units are identical.

$$\underline{a_1^{(2)}} = \underline{a_2^{(2)}}$$

## Random initialization: Symmetry breaking

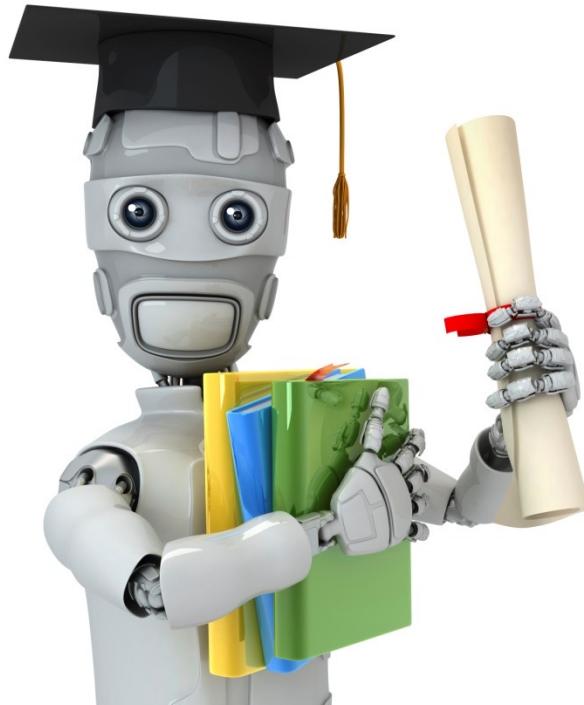
- Initialize each  $\Theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$   
(i.e.  $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$ )

E.g.

Random  $10 \times 11$  matrix (betw. 0 and 1)

$$\Theta_1 = \boxed{\text{rand}(10, 11) * (2 * \text{INIT\_EPSILON})} - \underline{\text{INIT\_EPSILON}}; \quad [-\epsilon, \epsilon]$$

$$\Theta_2 = \boxed{\text{rand}(1, 11) * (2 * \text{INIT\_EPSILON})} - \underline{\text{INIT\_EPSILON}};$$



Machine Learning

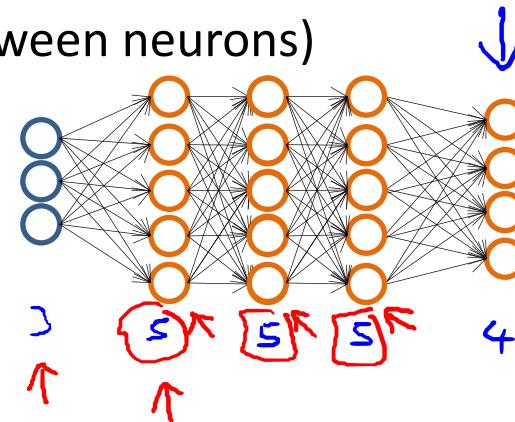
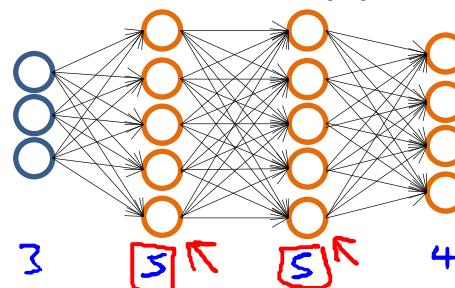
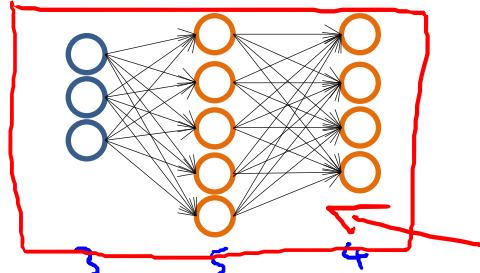
# Neural Networks: Learning

---

# Putting it together

# 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

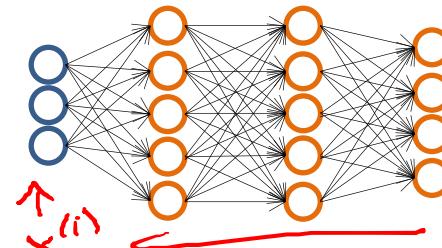
[Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)]

$$y \in \{1, 2, 3, \dots, 10\}$$

~~$y = 5$~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

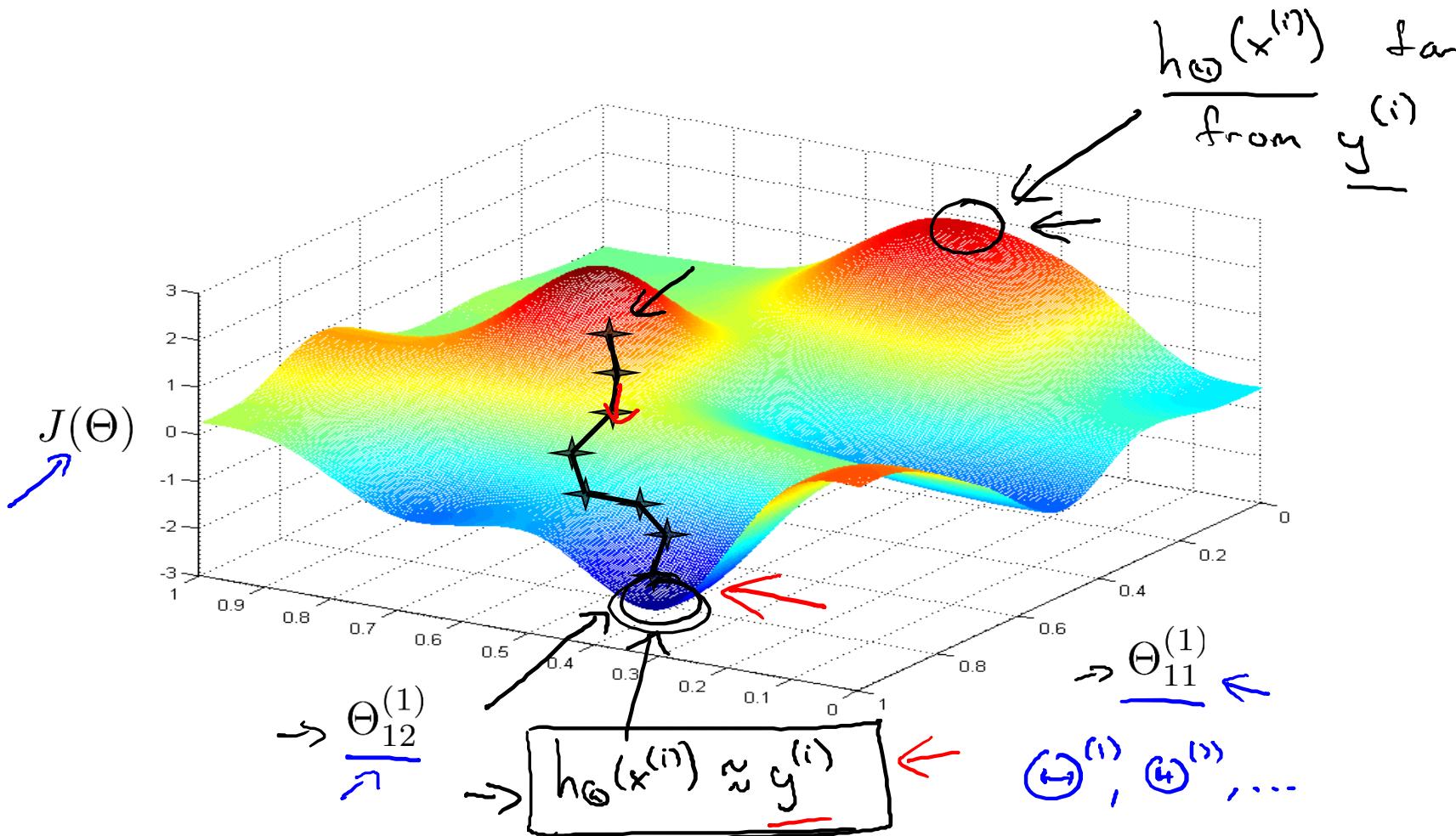
# Training a neural network

- 1. Randomly initialize weights
  - 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  $\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta)$
- for  $i = 1:m$  {  $(x^{(1)}, y^{(1)})$      $(x^{(2)}, y^{(2)})$  , ... ,  $(x^{(m)}, y^{(m)})$  }
- Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$
- (Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ ).
- $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (\alpha^{(l)})^T$
- ...  
}
- compute  $\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta)$ .
- 

## 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  $\Theta$

$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta) \quad \text{non-convex.}$$



# Neural Network Libraries

# Scikit-learn

- [https://scikit-learn.org/stable/modules/generated/sklearn.  
neural network.MLPClassifier.html](https://scikit-learn.org/stable/modules/generated/sklearn.neural_network.MLPClassifier.html)

# Keras

- [https://keras.io/guides/sequential\\_model/](https://keras.io/guides/sequential_model/)
- [https://keras.io/api/models/model\\_training\\_and\\_evaluation APIs/](https://keras.io/api/models/model_training_and_evaluation APIs/)

# Pytorch

- [https://pytorch.org/tutorials/beginner/basics/  
buildmodel\\_tutorial.html](https://pytorch.org/tutorials/beginner/basics/buildmodel_tutorial.html)

# Colab

- <https://colab.research.google.com/>

# Bài tập – scikit-learn

Trong dữ liệu ex7data.mat chưa dữ liệu lưu dưới dạng dict gồm:

X: 5000x400 là 5000 ảnh nhị phân chữ số viết tay có kích thước 20x20

y: 5000x1 là nhãn của các ảnh tương ứng

Các bạn làm các công việc sau:

- Đọc dữ liệu, reshape một số ảnh về kích thước 20x20 rồi show ra màn hình (plt.imshow)
- Chia dữ liệu thành 70% train, 30% test (train\_test\_split) đảm bảo tính ngẫu nhiên và đồng đều về nhãn.
- Train một mạng neural network với các thông số sau:
  - + 2 lớp ẩn: 100, 50
  - + Hàm kích hoạt: sigmoid
  - + Learning rate khởi tạo: 0.1
  - + Trộn dữ liệu mỗi vòng lặp
  - + Phương pháp tối ưu: stochastic gradient descent
  - + Hệ số regularization: 0.1
- Show đường cong loss trong quá trình học
- Show độ chính xác trên tập test