

with some manipulation our update rule can also be represented as:

$$\theta_j = \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

The first term in the above equation,  $\left(1 - \alpha \frac{\lambda}{m}\right)$  will be less than 1, intuitively you can see it as reducing the value of  $\theta_j$  by some amount on every update. Notice that the second term is now exactly the same as it was before.

### A Normal Equation

Now let's approach regularization using the alternate method of non-iterative normal equation.

To add in regularization, the equation is the same as our original, except that we add another term inside the parentheses:

$$\theta = (X^T X + \lambda L)^{-1} X^T y$$

where  $L = \begin{bmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix}$

Dimension of matrix  $(n+1) \times (n+1)$

$L$  is a matrix with 0's at the top left

$$(x_0, x_1, x_2, \dots, x_{n-1}, 1)$$

and 1's down the diagonal, with 0's everywhere else. It should have dimension  $(n+1) \times (n+1)$ , intuitively, this is the identity matrix (though we are not including  $x_0$ ), multiple with a ~~for~~ single real number  $\lambda$

Example of L dimensions:  $(\lambda \cdot I)$

Let  $n=2$

then  $L = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$  dimension of matrix  $(n+1) \times (n+1)$

recall that if  $m < n$ , then  $X^T X$  is non-invertible. However, when we add the term  $\lambda \cdot I$ , then  $(X^T X + \lambda \cdot I)$  becomes invertible.

## \* Regularized Logistic Regression

We can regularize logistic regression in a similar way that we regularize linear regression.

$$\begin{array}{c|c} \text{non-regularized} & \text{regularized function} \\ \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} & \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \end{array}$$

- non-regularized function; over-fit.  $\rightarrow$  Just Right
- $g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \dots)$

Cost function: ~~gradient descent~~

(start) gradient descent =梯度下降, 优化目标函数

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

+  $\left[ \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \right]$  (正则化项)

Regularize this equation by adding.

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

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j$$

The second sum,  $\sum_{j=1}^n \theta_j^2$  means to explicitly exclude the bias term,  $\theta_0$ .

→ Gradient descent

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

(i) 梯度下降  $\theta_0$  的更新公式

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad (j=1, 2, 3, \dots, n)$$

}

algorithm is identical but definition is different  
 (partial derivative of Regularized ~~not~~ linear regression cost function and partial derivative of Regularized logistic regression cost function is totally different)

for regularized cost function  
on ("fminunc(C)")

\*enlivio

Date : \_\_\_\_\_ / \_\_\_\_\_ / \_\_\_\_\_

Page: \_\_\_\_\_

## \* Advanced optimization. (Octave) (cont)

function [Jval, gradient] = costfunction(theta)

→ Jval = [Code to compute  $J(\theta)$ ];

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

$$\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

→ gradient(1) = [Code to Compute  $\frac{\partial J(\theta)}{\partial \theta_0}$ ];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

→ gradient(j) = [Code to Compute  $\frac{\partial J(\theta)}{\partial \theta_j}$ ];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$j \in \{2, 3, \dots, n+1\}$

In above function gradient(1) will compute

$\frac{\partial J(\theta)}{\partial \theta_0}$ , and therefore in gradient(j)

$$\frac{\partial J(\theta)}{\partial \theta_j}, x^{(i)} - (h_\theta(x^{(i)}))$$

$i \in \{1, 2, 3, \dots, m\}$  and  $j \in \{2, 3, 4, \dots, n+1\}$

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

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

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

Now, above Cost function can be parameterize  
in "fminunc(C)".

# \* Neural Networks :

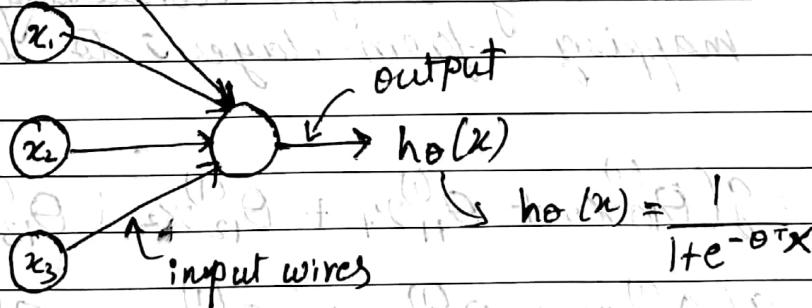
## → Neural Networks

- 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.

## \* Model representation

- neuron model: Logistic unit.

bias unit  $\rightarrow x_0$        $x_0 = 1$  (always)

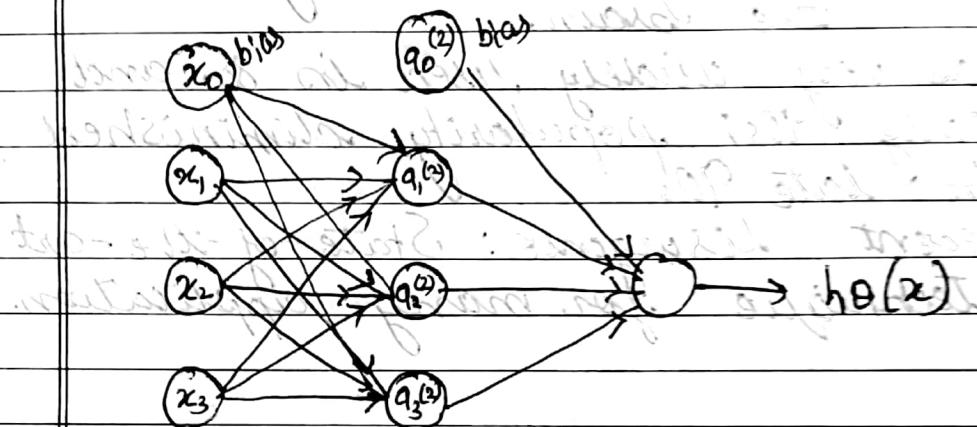


$$\begin{aligned} x &= \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} & \theta &= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \\ && \text{weights} & \\ && (\text{parameters}) & \end{aligned}$$

# Sigmoid (logistic) activation function.

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

→ Neural Network



layer 1      layer 2      layer 3.  
 (Input layer)    (hidden layer)    (output layer)

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

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

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

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

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

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

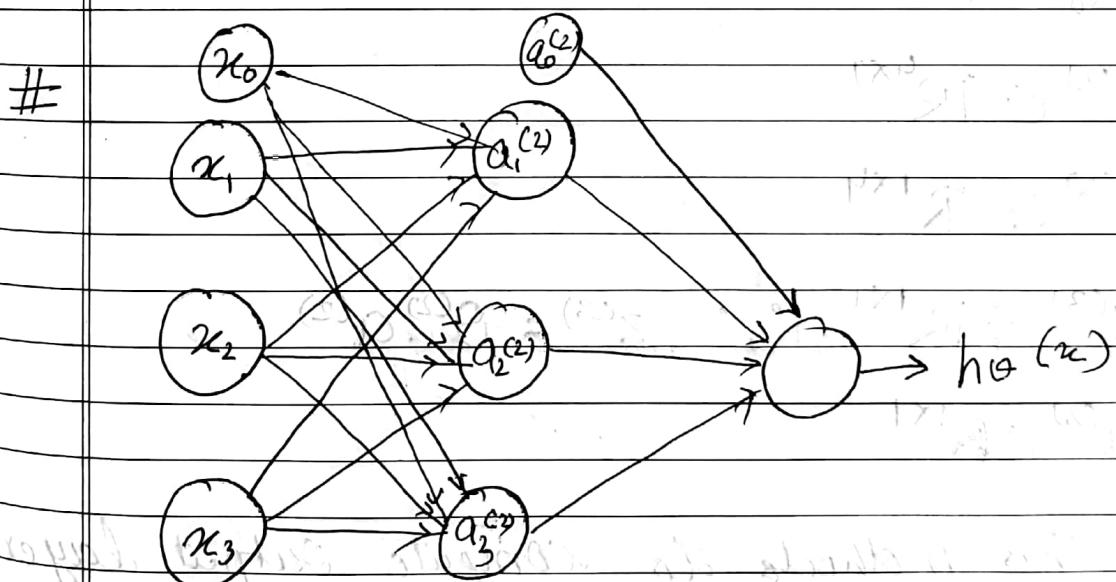
→ if network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j+1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

→ it means that dimension of  $\Theta^{(1)}$  will be "No. of unit in next layer (excluding bias unit)"  $\times$  "No. of unit in layer 1 (including bias unit)"

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

$$\Theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

Example : if layer 1 has 2 input and layer 2 has 4 activation nodes. Dimension of  $\Theta^{(1)}$  is going to be  $4 \times 3$  where  $s_j = 2$  and  $s_{j+1} = 4$ , so,  $s_{j+1} \times (s_j + 1) = 4 \times 3$ .



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

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

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

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

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

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

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

$$z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$$x \in \mathbb{R}^{4 \times 1}$$

$$\theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

$$\therefore z^{(2)} \in \mathbb{R}^{3 \times 1}$$

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

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

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

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

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

$$\therefore a^{(3)} \in \mathbb{R}^{3 \times 1}$$

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

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

$$\therefore a^{(2)} \in \mathbb{R}^{4 \times 1}$$

$$\theta^{(2)} \in \mathbb{R}^{1 \times 4}$$

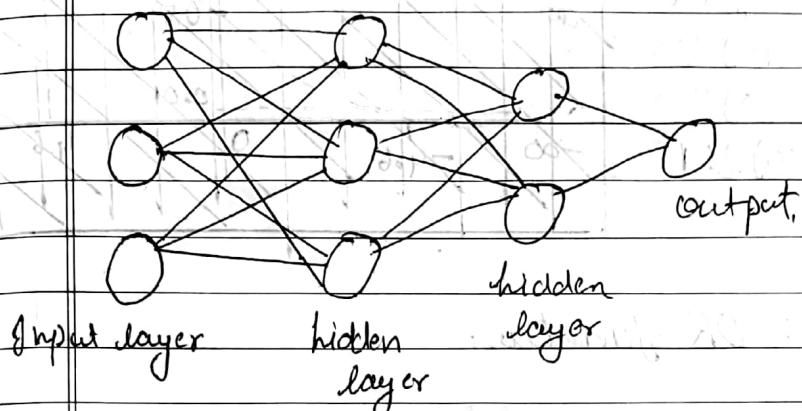
$$\therefore z^{(3)} \in \mathbb{R}^{1 \times 1}$$

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

$$\therefore a^{(3)} \in \mathbb{R}^{1 \times 1}$$

This method to compute output layer from input layers thru hidden layers is called forward propagation

→ other network architectures.



\* Non-linear classification example : XOR / XNOR

$x_1$	$x_2$	$y$	$x_1$	$x_2$	$y$
0	0	0	0	0	1
0	1	1	0	1	0
1	0	1	1	0	0
1	1	0	1	1	1

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

$$\begin{aligned} y &= x_1 \text{ XNOR } x_2 \\ y &= \text{NOT}(x_1 \text{ XOR } x_2) \end{aligned}$$

# Start with simple example : AND

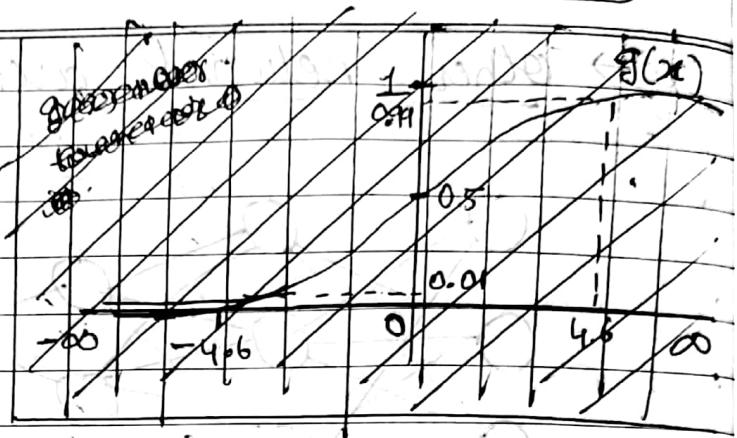
$x_1$	$x_2$	$y$	$x_1, x_2 \in \{0, 1\}$
0	0	0	$y = x_1 \text{ AND } x_2$
0	1	0	
1	0	0	
1	1	1	

$(+1)$       -30      20      20

$h_{\theta}(x) \rightarrow h_{\theta}(x)$

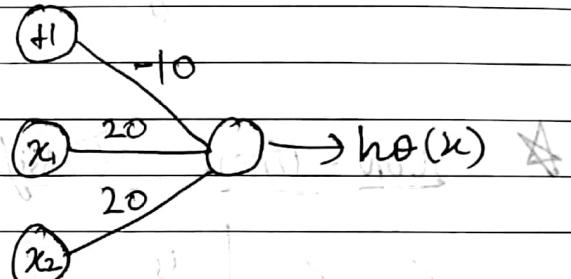
$$h_{\theta}(x) = g(-30 + 20x_1 + 20x_2)$$

$x_1$	$x_2$	$h_\theta(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$



## # Example: DR function

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	1

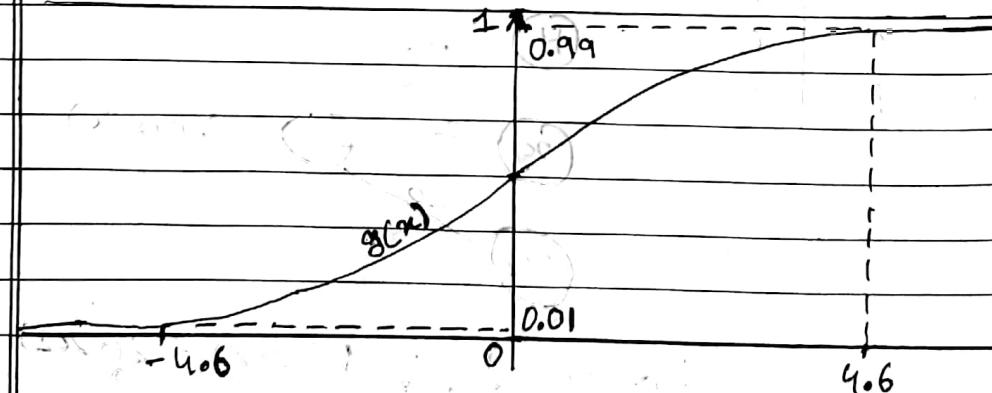


$$h_\theta(x) = g(-10 + 20x_1 + 20x_2)$$

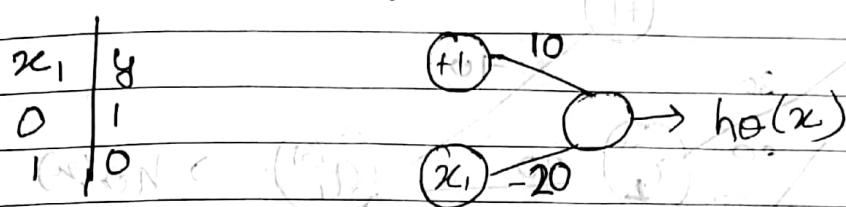
$\uparrow \quad \uparrow \quad \uparrow$   
 $\theta_0 \quad \theta_1 \quad \theta_2$

$x_1$	$x_2$	$h_\theta(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$g(10) \approx 1$
1	1	$g(30) \approx 1$

## # Sigmoid function



## # Example : Not function.



$$h_0(x) = g(10 - 20x_1)$$

$$\begin{array}{c|cc} & 10 & -20 \\ \hline x_1 & 0 & 1 \\ \hline & 10 & -20 \\ & 0 & 1 \\ \hline & 0 & 0 \end{array}$$

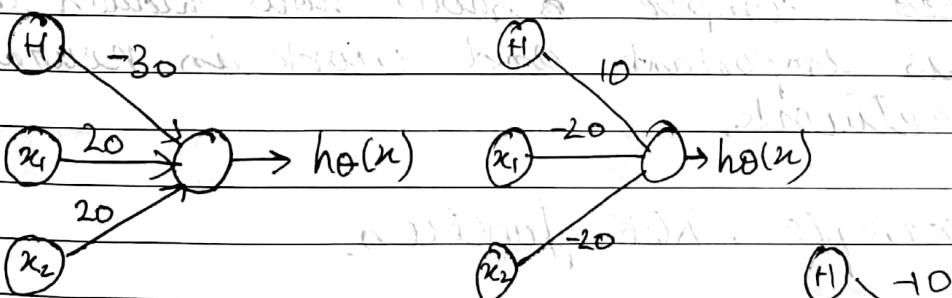
$x_1$	$h_0(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

## # Example : XNOR function.

~~→ putting it together~~

$x_1$	$x_2$	$y_2$
0	0	1 (01)A
0	1	0 (10)B
1	0	0 (01)C
1	1	1 (11)D

putting it together.



$\rightarrow x_1 \text{ AND } x_2$

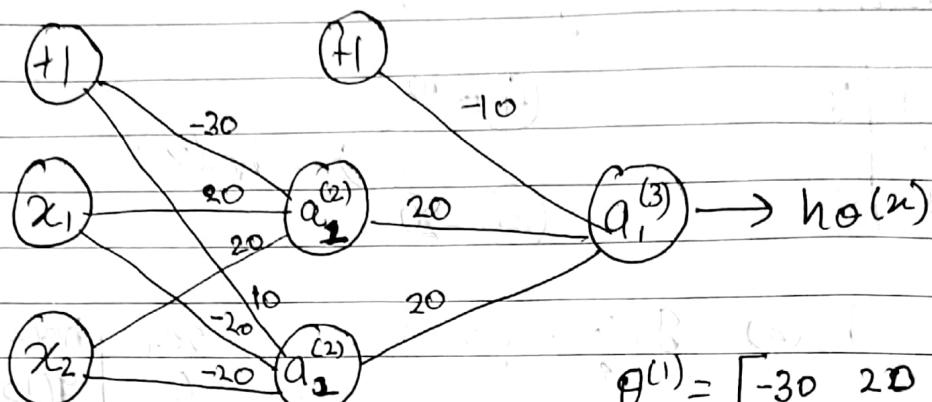
$\rightarrow (\text{Not } x_1) \text{ AND } (\text{Not } x_2)$

$$\Rightarrow \text{Not } x_1 = \text{Not } x_2 = 1$$

if and only if

$$x_1 = x_2 = 0, \text{ or}$$

vice-versa.



$$\theta^{(1)} = \begin{bmatrix} -30 & 20 & 20 \\ 10 & -20 & -20 \end{bmatrix}$$

$$\theta^{(2)} = \begin{bmatrix} -10 & 20 & 20 \end{bmatrix}$$

As

~~X1 AND X2~~ ~~gives~~

~~(NOT X1 AND NOT X2)~~  $a_1^{(2)} = x_1 \text{ AND } x_2$

$a_2^{(2)} = (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$

$h\theta(x) = a_1^{(2)} \text{ OR } a_2^{(2)}$

$x_1 \ x_2 \ a_1^{(2)} \ a_2^{(2)}$

0 0 0 1

0 1 0 0

1 0 0 0

1 1 1 0

$h\theta(x)$

$g(10) \approx 1$

$g(-10) \approx 0$

$g(-10) \approx 0$

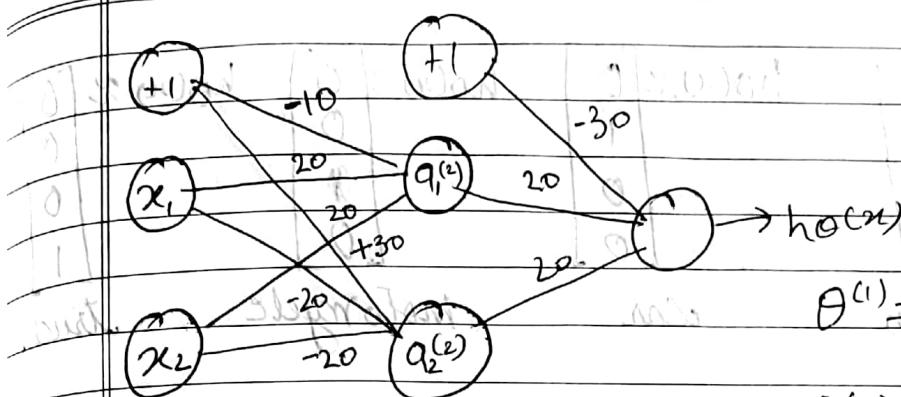
$g(10) \approx 1$

→ this and next

last Example is show how hidden layer is important and work in neural network.

# Example : XOR function.

$x_1$	$x_2$	y
0	0	0
0	1	1
1	0	1
1	1	0



$$\theta^{(1)} = \begin{bmatrix} -10 & 20 & 20 \\ +30 & -20 & -20 \end{bmatrix}$$

$$\theta^{(2)} = \begin{bmatrix} -30 & 20 & 20 \end{bmatrix}$$

$$a_1^{(2)} = x_1 \text{ OR } x_2$$

$$a_2^{(2)} = (\text{NOT } x_1) \text{ OR } (\text{NOT } x_2)$$

$$h(x) = a_1^{(2)} \text{ AND } a_2^{(2)}$$

$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h(x)$
0	0	0	1	$g(-10) \approx 0$
0	1	1	1	$g(10) \approx 1$
1	0	1	0	$g(10) \approx 1$
1	1	1	0	$g(-10) \approx 0$

### \* Multiple Output unit: One vs all

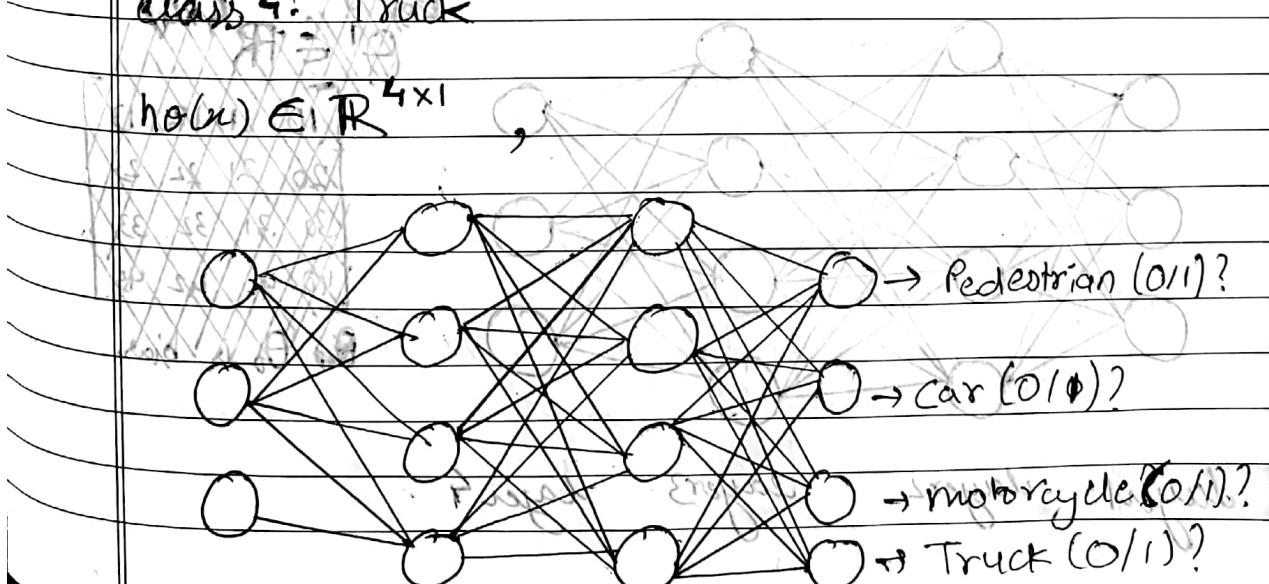
class 1: Pedestrian

class 2: (car, bicycle, truck, motorcycle)

class 3: motorcycle

class 4: Truck

$$h(x) \in \mathbb{R}^{4 \times 1}$$



$\rightarrow$  Pedestrian (0/11)?

$\rightarrow$  car (0/10)?

$\rightarrow$  motorcycle (0/11)?

$\rightarrow$  Truck (0/1)?

$$\text{want } h_{\theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad h_{\theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

when pedestrian car motorcycle truck

Training set :

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}) \dots (x^{(m)}, y^{(m)})$$

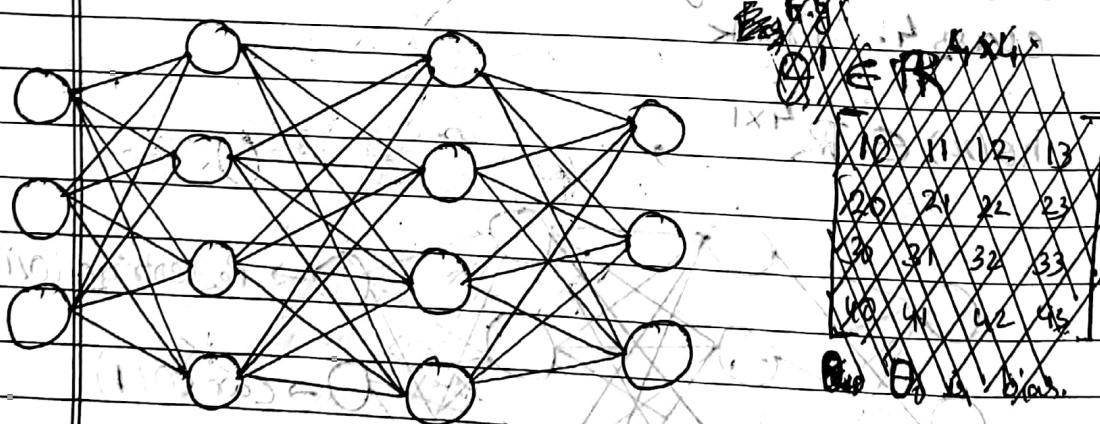
where

$$y^{(i)} \text{ one of } y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

if pedestrian if car if motorcycle if truck

\* cost function and Backpropagation

$\Rightarrow$  Neural Network (Classification)



layer1 layer2 layer3 layer4

## Training Examples:

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

$L =$  total no. of layers in network  
(In this case,  $L=4$ )

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

(i.e.,  $S_1=3$ ,  $S_2=4$ ,  $S_3=4$ ,  $S_4=3$ )

### Binary Classification

$(K=1)$   $k=k$  = number of output units

### Multi-class Classification

$(K \geq 3)$

$$y = 0 \text{ or } 1$$

$$y \in \mathbb{R}^k$$

$$h_\theta(x) \in \mathbb{R}^k$$

$$\text{E.g.: } \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$h_\theta(x) \in \mathbb{R}^k$$

### Cost function

For Logistic regression, we are trying to

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

minimize  $J(\theta)$  with respect to  $\theta$ .

The cost function is convex.

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

## Neural Network:

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

#8 Suppose we want to try to minimize  $J(\theta)$  as a function of  $\theta$ , using one of advanced optimization methods. What do we need to supply code to compute (as a function of  $\theta$ )?

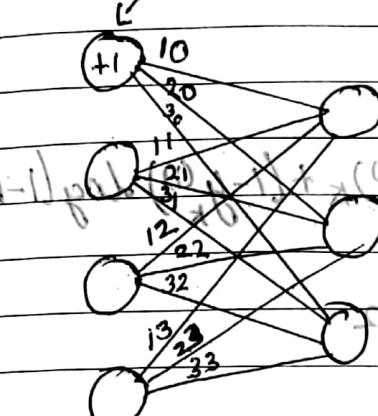
A → We need to supply  $J(\theta)$  and partial derivative term  $\frac{\partial}{\partial \theta_{ij}^{(l)}}$  for every  $i, j, l$ .

→ We have added a few nested summations to account for our multiple output nodes. In the first part of equation, before the square brackets, we have an additional nested summation that loops through the number of output nodes.

In the regularization part, after the square brackets, we must account for multiple theta matrices. The number of columns in our current theta matrix is equal to the number of nodes in our current layer (excluding the bias unit). The number of rows in our current theta matrix is equal to the number of nodes in the next layer (excluding the bias unit). As before with logistic regression, we square every term.

Example \*

bias unit



Theta matrix

$$\theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

$$\begin{matrix} & 10 & 11 & 12 & 13 \\ 20 & & 21 & 22 & 23 \\ 30 & 31 & 32 & 33 \end{matrix}$$

We will not include

$\theta$  values from bias unit in regularization, like we did in regularized logistic regression cost function.

Note 2

$$\theta \in \mathbb{R}^n$$

$$(\theta)^T \leftarrow$$

- the double sum simply adds up the logistic regression costs calculated for each cell in output layer.

The triple sum simply adds up the squares of all the individual  $\theta$ s in the entire network. (if we use  $\theta$  values from bias unit, it would work about the same and doesn't make big difference. But maybe this is convention of not regularizing the bias term is just slightly more common)

The  $i$  in the triple sum does not refer to training example i.

$$((\cdot))B = (\cdot)$$

$$((\cdot))A(b)$$

$\Rightarrow$  Forward propagation algorithm

gradient computation.

$$J(\theta) = -\frac{1}{m} \sum_{k=1}^m \left[ y^{(1)} \log h_\theta(x^{(1)})_k + (1-y^{(1)}) \log (1-h_\theta(x^{(1)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{j=1}^{S_l} \sum_{i=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

shallow tan new job

$\rightarrow$  implement  $J(\theta)$  first without  $\theta$

new  $\theta$ , otherwise you will mi.

designed on basis

existing  $\rightarrow$  need code to compute

$\rightarrow J(\theta)$

$\frac{\partial}{\partial \theta_j^{(1)}} J(\theta)$

$$\theta_{ij}^{(1)} \in R \quad s_{x_i} \times (S_{l+1})$$

$\theta_{ij}^{(1)}$  are those which make old job well  
otherwise else converged visual.  
rough top with me after this

In order to compute  $J(\theta)$  we will use above  
equation and to compute  $\frac{\partial}{\partial \theta_j^{(1)}} J(\theta)$ :

start in  $\theta$  initialization set this to compute  $\frac{\partial}{\partial \theta_j^{(1)}}$

logic given one training example  $(x, y)$ .

and some int words from below to find

# forward propagation:

and with probabilities  $a^{(1)}, a^{(2)}, a^{(3)}, a^{(4)}$   
(normal from input to output)

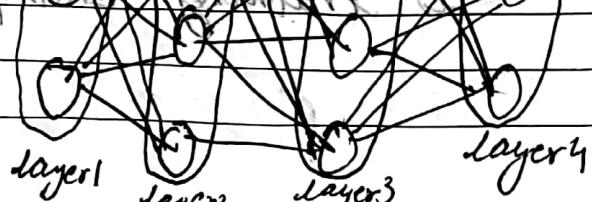
$$c^{(1)} = x$$

now we do this

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

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

(add  $a_0^{(2)}$ )



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

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

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

$$a^{(4)} = h\theta(x) = g(z^{(4)})$$

$\Rightarrow$  Gradient Computation: Backpropagation algorithm

intuition:  $\delta_j^{(l)}$  = "Error" of node  $j$  in layer  $l$ .

for each output unit (layer  $l=4$ )

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

$\Rightarrow$  Vectorized form:

$$\delta^{(4)} = a^{(4)} - y$$

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

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

#  $*$  is element-wise multiplication.

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

for  
brief

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} s_i^{(l+1)} \quad (\text{ignoring } \lambda) \quad \text{Date: } \boxed{\quad / \quad / \quad} \quad * \text{enliv} \quad \text{Page: } \boxed{\quad / \quad}$$

Backpropagation algorithm:

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

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ).

# ( $\Delta$  used to compute  $\delta$ )

$$(\delta) \Delta \theta_{ij}^{(l)} = \delta$$

for  $i=1$  to  $m$

$$\text{set } a^{(1)} = x^{(i)}$$

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

using  $y^{(i)}$ , compute  $s^{(L)} = a^{(L)} - y^{(i)}$

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

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

# vectorized form :  $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

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

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

#  $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = D_{ij}^{(l)}$

$$(\delta) \Delta \theta_{ij}^{(l)} = D_{ij}^{(l)}$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = \frac{\partial J(\theta)}{\partial a_{ij}^{(l)}} \cdot \frac{\partial a_{ij}^{(l)}}{\partial z_{ij}^{(l)}} \cdot \frac{\partial z_{ij}^{(l)}}{\partial \theta_{ij}^{(l)}}$$

← basic idea.

\*enlive  
Date: \_\_\_\_\_  
Page: \_\_\_\_\_

In the capital-letter matrix  $D$  is used as an "accumulator" to add up our values as we went along and eventually compute our derivative. Thus we get  $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}$

## \* Backpropagation Intuition

Recall that cost function for a neural network is:

$$J(\theta) = -\frac{1}{m} \sum_{t=1}^m \sum_{k=1}^K \left[ y^{(t)} \log(h_\theta(x^{(t)}))_k + (1-y^{(t)}) \log(1-h_\theta(x^{(t)}))_k \right] + \frac{\lambda}{2m} \sum_{j=1}^J \sum_{i=1}^I (\theta_{ji}^{(l)})^2$$

(logit intuition) if we consider simple non-multiclass classification ( $K=1$ ) and disregard regularization, the cost is computed with:

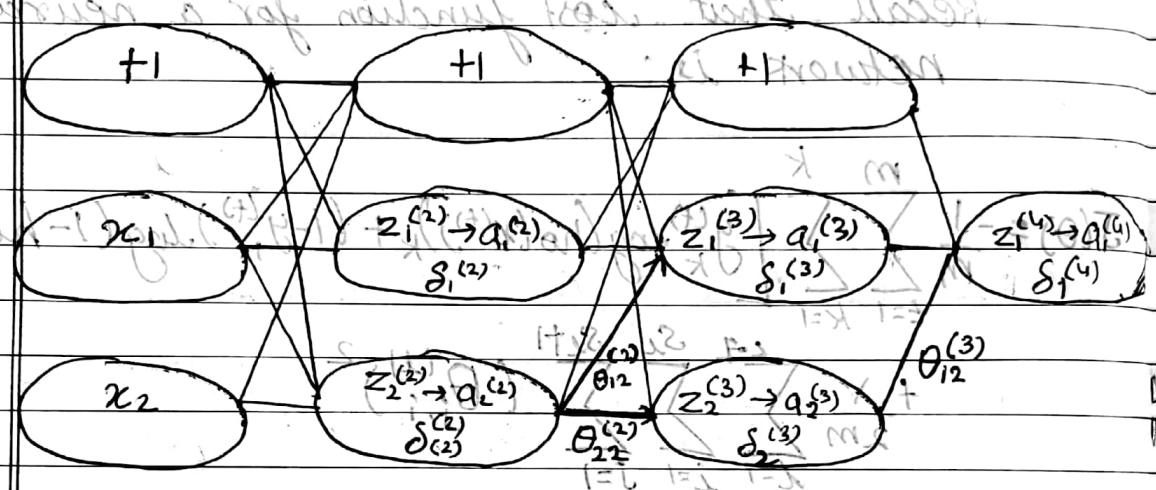
$$\text{Cost}(t) = y^{(t)} \log(h_\theta(x^{(t)})) + (1-y^{(t)}) \log(1-h_\theta(x^{(t)}))$$

Intuitively,  $\delta_j^{(l)}$  is the "error" for  $\theta_{ji}^{(l)}$  (unit  $i$  in layer  $l$ ). More formally, the delta values are actually the derivative of the cost function.

$$\delta_j^{(l)} = \frac{\partial \text{Cost}(t)}{\partial z_j^{(l)}}$$

Recall that our derivative is the slope of a line tangent to the cost function, so the steeper the slope the more incorrect we are. Let us consider the following neural network below and see how we could calculate some  $\delta_j^{(l)}$ .

During a ref adjustment from  $\theta_{12}$  to  $\theta_{12}'$



$\delta_j^{(l)}$  = "error" of cost for  $a_j^{(l)}$  (unit  $j$  in layer  $l$ )

formally,  $\delta_j^{(l)} = \partial_{j,l} \text{cost}(i)$  (for  $j > 0$ ) where

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

in the image, to calculate  $\delta_2^{(2)}$ , we multiply the weights  $\theta_{12}^{(1)}$  and  $\theta_{22}^{(2)}$  by their respective  $\delta$  value found to the right of each edge. So we get  $\delta_2^{(2)} = \theta_{12}^{(1)} * \delta_1^{(1)} + \theta_{22}^{(2)} * \delta_2^{(1)}$ . To calculate every single possible  $\delta_j^{(l)}$ , we could start from the right of our diagram. We

can think of our edges as our  $\theta_{ij}$ . Going from right to left, to calculate the value of  $s_j^{(4)}$ , you can just take the overall sum of each weight times the  $s_i$  it is coming from. Hence, another example could be

~~Want~~

$$s_2^{(3)} = \theta_{12}^{(3)} * s_1^{(4)}$$

Octave/matlab

\* Implementation note: Unrolling parameter.

# Advanced optimization (Theta parameter is vector.)

function [J, Val, gradient] = CostFunction(theta)

% J = cost function (R^n+1) → R

...  
...

: [C] = optimfun([@costfunction], initialTheta, options)

But with neural networks, we working with  
arrays of matrices: need to unroll #

( $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  - matrices) ( $\theta_1, \theta_2, \theta_3$ )  
 (D<sup>(1)</sup>, D<sup>(2)</sup>, D<sup>(3)</sup> - matrices) (D<sub>1</sub>, D<sub>2</sub>, D<sub>3</sub>)  
 (W<sub>1</sub>, W<sub>2</sub>, W<sub>3</sub> - matrices) (W<sub>11</sub>, W<sub>12</sub>, W<sub>13</sub>, W<sub>21</sub>, W<sub>22</sub>, W<sub>23</sub>, W<sub>31</sub>, W<sub>32</sub>, W<sub>33</sub>)  
 "Unroll" into Vectors

: scroll down. o)

(S) (S) (S) A 289 training initini .mat

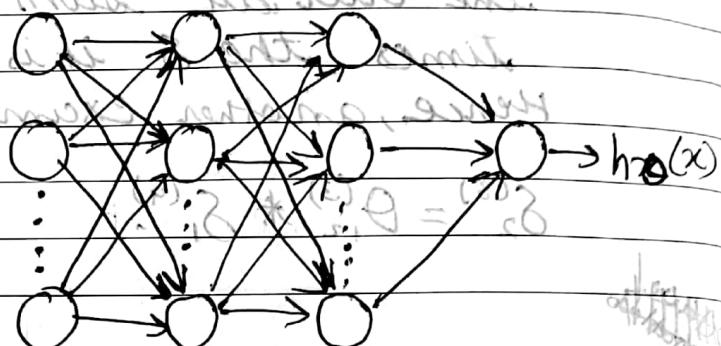
Example: explore uses for third node  
sigmoid at step at about math prior)

$$S_1 = 10, S_2 = 10, S_3 = 10, S_{4,5} \text{ values not}$$

typical uses for math no uses left.  
many priors) and what's next.

get between sigmoids, 2nd, 3rd, 4th

bottom (outputs)



$$\theta^{(1)} \in \mathbb{R}^{10 \times 11}, \theta^{(2)} \in \mathbb{R}^{10 \times 11}, \theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

~~$$D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$$~~

(other) [initialization] = [twisting, doing nothing]

# unroll  $\theta^{(1)}, D^{(2)}$  into vector

# ThetaVector = [Theta1(:, 1); Theta1(:, 2); Theta1(:, 3); ...; Theta1(:, 10); Theta2(:, 1); Theta2(:, 2); Theta2(:, 3); ...; Theta2(:, 10); Theta3(:, 1); Theta3(:, 2); Theta3(:, 3); ...; Theta3(:, 11)];

$$DVec = [D1(:, 1), D2(:, 1), D3(:, 1)];$$

new prior not so strong because this way

# we can get back our original matrices

$$\text{Theta1} = \text{reshape}(\text{ThetaVector}(1:110), 10, 11)$$

$$\text{Theta2} = \text{reshape}(\text{ThetaVector}(111:220), 10, 11)$$

$$\text{Theta3} = \text{reshape}(\text{ThetaVector}(221:231), 1, 11)$$

matrix with "None":

# To Summarize:

→ Have initial parameters  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

→ Unimportant to get "initial Theta's" to pass to

our function. We can use  $\theta = \theta_0$ .

↳ fminunc (@CostFunction, initialTheta, options)

→ function [JVal, gradientVec] = CostFunction(thetaVec)

↳ From thetaVec, get  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ .

↳ Use forward/back propagation to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\theta)$ .

unimportant to unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get "gradientVec".

## \* Gradient Checking (Numerical estimate of gradient)

Gradient checking will assure that our backpropagation works as intended. we can approximate the derivative of our cost function with:

$$\frac{\partial J(\theta)}{\partial \theta_j} \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

way to do this is to take two points on the curve.

one point is at  $\theta_j + \epsilon$  and one is at  $\theta_j - \epsilon$ .

then calculate the slope between these two points.

and divide by  $2\epsilon$  to get the derivative.

With multiple theta matrices, we can approximate the derivative with respect to  $\theta_j$  as follows:

$$\frac{\partial J(\theta)}{\partial \theta_j} \approx \frac{J(\theta_1, \dots, \theta_j + \epsilon, \dots, \theta_n) - J(\theta_1, \dots, \theta_j - \epsilon, \dots, \theta_n)}{2\epsilon}$$

way to do this is to take two points on the curve.

one point is at  $\theta_j + \epsilon$  and one is at  $\theta_j - \epsilon$ .

then calculate the slope between these two points.

and divide by  $2\epsilon$  to get the derivative.

↳ A small value for  $\epsilon$  ( $\epsilon \ll 10^{-4}$ ) such as  $\epsilon = 10^{-4}$ ,

guarantees that math will work out properly.

we are only adding or subtracting epsilon to the  $\theta_j$  matrix. In octave we can (and it's good to do it as you see in the code) summt

$$\text{epsilon} = 1e-4;$$

(3b pts) write for  $i = 1:n$ ,  $\nabla \text{theta}$  for each  $\theta_j$  constant

$$\text{thetaPlus} = \text{theta};$$

$$\text{thetaPlus}(i) += \text{epsilon};$$

$$\text{thetaMinus} = \text{theta};$$

$$\text{thetaMinus}(i) -= \text{epsilon};$$

$$\text{"gradient"} \text{gradApprox}(i) = (\text{J(thetaPlus)} - \text{J(thetaMinus)}) / (2 * \text{epsilon})$$

grad

end;

~~for i = 1:n, print theta(i)~~

So once we compute our "gradApprox" Vector, we can check that  $\text{gradApprox} \approx \text{deltaVector}$

$$\left( \frac{\partial J(\theta)}{\partial \theta_j} \right)$$

$$(\theta - \alpha)T - (\theta + \alpha)T \approx (\theta)T \delta$$

Once you have verified once that your backpropagation algorithm is correct, you don't need to compute gradApprox again.

The code to compute gradApprox can be very slow

~~so it's better to do it once and then use it multiple times~~

~~it's important to trigger this only during init~~

Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of cost function), your code

will be very slow and expensive.

## \* Random Initialization

Initializing all theta weights to zero does NOT work with neural networks. When we backpropagate, all nodes will update to the same value repeatedly. Instead, we can randomly initialize our weights for our  $\Theta$  matrices using the following

syntex of methods:  $\text{rand} = \text{randn}$  (from np.random)

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

### Random initialization: Symmetry breaking

Initialize each  $\Theta_{ij}^{(1)}$  to a random value in  $[-\epsilon, \epsilon]$

(i.e.  $-\epsilon \leq \Theta_{ij}^{(1)} \leq \epsilon$ )

and similarly for  $\Theta_{ij}^{(2)}$ .

dimension of  $\Theta_{ij}^{(1)} \in \mathbb{R}^{10 \times 11}$ ,  $\Theta_{ij}^{(2)} \in \mathbb{R}^{1 \times 11}$

(i.e. Theta1 is 10x11 dimension and Theta2 is 1x11

dimension) if repeat habit 1: thought.

normalized in the next repeat habit 1 now

Example: good note good way to do it.

$$\text{Theta1} = \text{rand}(10, 11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$$

$$\text{Theta2} = \text{rand}(1, 11) * (2 * \text{INIT\_EPSILON}) - \text{INIT\_EPSILON};$$

We initialize each  $\Theta_{ij}^{(1)}$  to a random value

between  $[-\epsilon, \epsilon]$ . Using the above formula guarantees that we get the desired bound. The same procedure applies to all the  $\Theta$ 's.

→  $\text{rand}(x, y)$  is just a function in Octave that will initialize a matrix of random real numbers between 0 & 1.

→ Note: The Epsilon used above is unrelated to the epsilon from gradient checking.

## ★ Putting it Together: Initialization methods ★

→ first, pick a network architecture & choose the layout of your neural network, including how many hidden units in each layer and how many layers in total you want to have.

- Number of input units = dimensions of feature  $x^{(i)}$

• Number of output units = number of classes.

→ Step 1: Number of hidden units per layer = usually more the better (must balance with cost of computation as it increase with more hidden units).

- Default: 1 hidden layer. If you have more than 1 hidden layer, then it is recommended that you have the same number of units in every hidden layer.

## # Training a neural network

1. Randomly initialize the weights

2. Implement forward propagation to get  $h(\theta(x^{(i)})$ )

3. Implement the cost function  $J(\theta, x^{(i)}, y^{(i)})$

4. Implement backpropagation to compute  $\frac{\partial}{\partial \theta_j} J(\theta, x^{(i)}, y^{(i)})$

5. Implement gradient descent algorithm

forward pass). We proceed to perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$

Get activation  $a^{(l)}$  and delta term

$$\delta^{(l)} \text{ for } l=2 \dots L$$

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} \cdot (a^{(l)})^T$$

and;

$$\text{Compute } \frac{\partial J(\theta)}{\partial \theta_{ij}}$$

5. Use gradient checking to compare  $\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}$  computed using backpropagation vs. using numerical estimate of gradient

→ 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  $\theta$

$\rightarrow J(\theta)$  will be non-convex

Initial min. at bottom

and in this situation local. min. can

exist at other minima. Since there

are many local minima, it's hard

to find the global min. without

using some sort of global optimiz.

such as gradient descent, SGD,

# Debugging a learning algorithm.

Suppose you have implemented regularized linear regression to predict housing prices. In a regularization term,

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

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Get more training example.
- Try smaller sets of features
- Try adding polynomial features  
e.g.  $(x_1^2, x_1^3, x_1 x_2, \text{etc.})$
- Try decreasing  $\lambda$
- Try increasing  $\lambda$

# Diagnosis: A test that you can run often to gain insight what is/ isn't working with a learning and gain guidance as to how best to improve its performance

Diagnostic can take time to implement, but doing so can be a very good use of ~~your~~ your time.

## \* Evaluating a hypothesis

we can move on to evaluate our new hypothesis.

A hypothesis may have a low error for the training Examples but still be inaccurate (because of Overfitting). Thus, to evaluate a hypothesis, given a dataset of training Example, we can split up the data into two sets: a training set and a test set.

The new procedure using these two sets is then:

1. Learn  $\theta$  and minimize  $J_{\text{train}}(\theta)$  using the training set.

2. Compute the test set error  $J_{\text{test}}(\theta)$

• The test set error

B) for linear regression:  $J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (\hat{h}_{\theta}(x^{(i)})) - y^{(i)}_i$

2) for classification or Misclassification Error (or misclassification error) (aka 0/1 misclassification error):

$$\text{error}_0(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5 \text{ and } y=0 \text{ or} \\ & \text{if } h_{\theta}(x) < 0.5 \text{ and } y=1 \\ 0 & \text{otherwise} \end{cases}$$

This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is:

$$\text{Test Error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_0(x_i^{(i)}), y_i^{(i)})$$

This gives us the proportion of the test data that was misclassified.

## Model Selection and Train / Validation / Test Sets

just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis. it could overfit and as a result your prediction on the test set would be poor. The error of your hypothesis as measured on the data set which you found the parameters will be lower than the error on any other data set.

Given many models with different polynomial degrees we can use a systematic approach to identify the 'best' function. In order to

choose the model of your hypothesis, to you can test each degree of polynomial and look at the error result.

~~One way to break down our dataset into the three sets this is~~

- Training set: 60%

~~Cross Validation Set : 20%~~

- Test set: 20%

we can now calculate three separate error values for the three different sets using the following method:

1. Optimize the parameters in  $\theta$  using the training set for each polynomial degree.

2. find the polynomial degree  $d$  with the test least error using the cross validation set.

3. Estimate the generalization error using the test set with  $J_{\text{test}}(\theta^{(d)})$  ( $d = \text{the degree of polynomial with lowest error}$ );

This way, the degree of polynomial  $d$  has not been trained using test set.

and now  $(\theta)$  is not trained with test set.

Generalization error depends on  $d$ .

Scanned with CamScanner

# Diagnosing Bias Vs. Variance

In this topic we examine the relationship between the degree of the polynomial  $d$  and the underfitting or overfitting of our hypothesis.

- We need to distinguish whether bias or variance is the problem contributing to bad prediction.

Training High bias is underfitting and high variance is overfitting. Ideally, we need find a golden mean between these two extremes.

→ Training error will tend to decrease as we increase the degree  $d$  of the polynomial.

At the same time, the cross validation error will tend to decrease as we increased  $d$  up to a point and then it will increase, forming a convex curve.

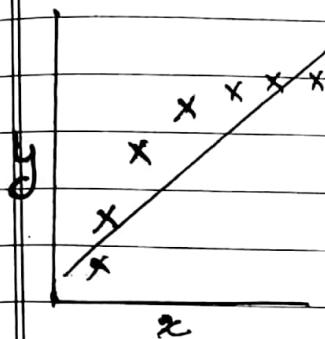
High bias (Underfitting): both  $J_{train}(\theta)$  and  $J_{cv}(\theta)$  will be high.

Also,  $J_{cv}(\theta) \approx J_{train}(\theta)$

High Variance (Overfitting):  $J_{train}(\theta)$  will be low and  $J_{cv}(\theta)$  will be much greater than  $J_{train}(\theta)$ .

Bias/Variance.

(tifrobin) 2018



$$\theta_0 + \theta_1 x$$

(d=1)

High bias  
(Underfitting)

High sd HWS ( $\theta$ ) small  $\leftarrow$

$$(\theta) \text{ small} \leftarrow (\theta) \text{ small}$$

(tifrobin) 2018

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

(d=2)

Just right

$$\theta_0 + \theta_1 x + \theta_2 x^2 +$$

$$\theta_3 x^3 + \theta_4 x^4$$

(d=4)  $\star$

High Variance  
(Overfitting)

$$\theta_0 + \theta_1 x + \theta_2 x^2 +$$

$$\theta_3 x^3 + \theta_4 x^4$$

# Training Error :  $J_{\text{train}} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Gross Validation Error :  $J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$

Underfitting  
(high bias)

$J_{\text{cv}}(\theta)$

Overfitting  
(high Variance)

$J_{\text{train}}(\theta)$

optimal value  
for d

d (Polynome degree)

## Bias (Underfit)

- $J_{\text{train}}(\theta)$  will be high
- $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$

## Variance (Overfit)

- $J_{\text{train}}(\theta)$  will be low
- $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$

$$+ \delta x_i \theta + \dots + \delta x_n \theta + \theta_0$$

$$+ \delta x_{n+1} \theta + \dots + \delta x_{n+i} \theta + \theta_0 \quad (i=1)$$

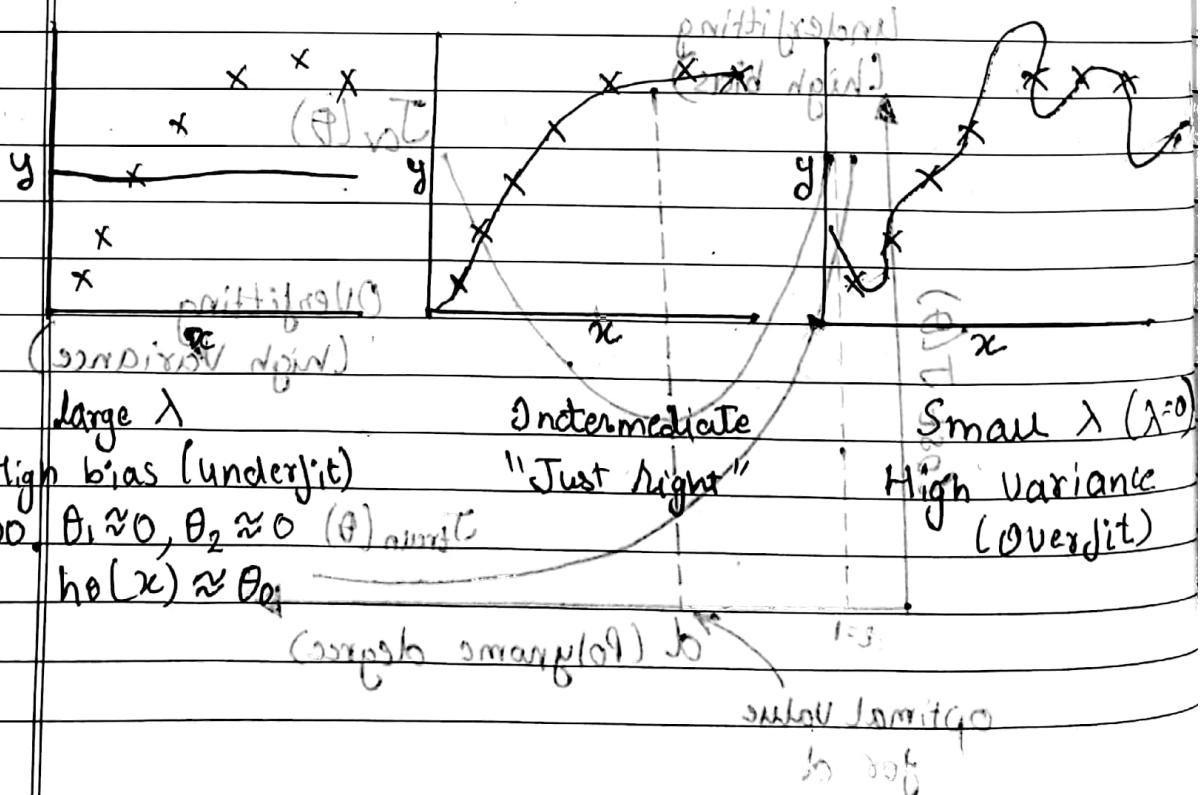
$$\dots + \delta x_{n+1} \theta + \dots + \delta x_{n+i} \theta + \theta_0 \quad (i=n)$$

## Regularization and Bias / Variance

linear regression with regularization.

$$\text{model: } h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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



We see that as  $\lambda$  increases, our fit becomes more rigid. On the other hand, as  $\lambda$  approaches 0, we tend to overfit the data. So how do we choose our parameter  $\lambda$  to get its "just right"? In order to choose the model and the regularization terms, we need to:

1. Create a list of lambdas (i.e.  $\lambda \in [0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24]$ );  $\theta \leftarrow (\theta)^T \text{ min } \leftarrow 10.0 = \theta^*$
  2. Create a set of models with different degrees or any other variants.
  3. Iterate through the  $\lambda$ s and for each  $\lambda$  go through all the models to learn some  $\theta$ .
  4. Compute the cross validation error using the learned  $\theta$  (computed with  $\lambda$ ) on the  $J_{cv}(\theta)$  without regularization or  $\lambda = 0$ .
  5. Select the best combo that produces the lowest error on the cross validation set.  $\theta^* = (\theta)_{\text{min}}$  no fit w/  $\lambda = 0$ .
- Using the best combo  $\theta^*$  and  $\lambda$ , apply it on  $J_{test}(\theta)$  to see if it has a good generalization of the problem.

E.g.

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

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

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Consider if we have  $m_{cv}$  normal fns in total for  $m_{cv}$

$$J_{cv}(\theta) = \frac{1}{m_{cv}} \sum_{i=1}^{m_{cv}} \left( h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)} \right)^2$$

with  $2m_{cv}$  observations &  $\theta$  is fixed

so we can do this with different  $\theta$ . Therefore

"Fwd. step" this at  $m_{test}$  at. & returns  $J_{cv}(\theta)$

$$J_{test}(\theta) = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \left( h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)} \right)^2$$

at each  $m_{test}$ , we have different  $\theta$ .

i.) smallest in tail of  $J_{cv}(\theta)$

1. try  $\lambda = 0.8 \rightarrow \min J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$

2. try  $\lambda = 0.01 \rightarrow \min J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$

3. try  $\lambda = 0.02 \rightarrow \min J(\theta) \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$

4. try  $\lambda = 0.04 \rightarrow \min J(\theta) \rightarrow \theta^{(4)} \rightarrow J_{cv}(\theta^{(4)})$

5. try  $\lambda = 0.08 \rightarrow \min J(\theta) \rightarrow \theta^{(5)} \rightarrow J_{cv}(\theta^{(5)})$

much it seems like its diverging up

12. try  $\lambda = 10 \rightarrow \min J(\theta) \rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$

and now (A new minimum)  $\theta$  having. But

$\theta = \lambda$  Now pick  $\theta$  with lowest  $J_{cv}(\theta)$ . (But

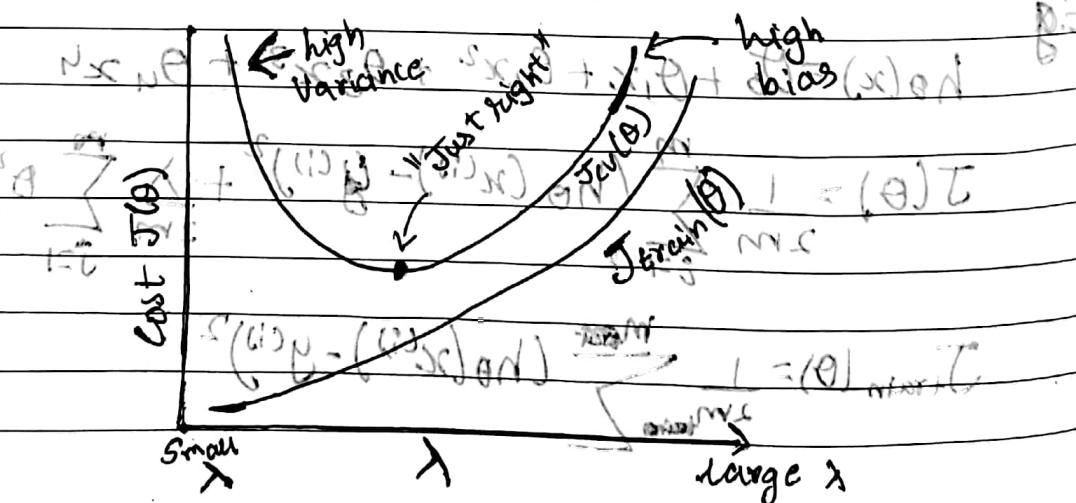
Say,  $\theta^{(5)}$  didn't reach tail yet. So

now it's time to move toward tail.

NOW, first apply it on  $J_{test}(\theta^{(5)})$  to see

Test Error (to see if loss good generalization  
including bias of the problem)

Bias/Variance as a function of the  $\lambda$



## Learning Curves

Training an algorithm on a very few number of data points (such as 1, 2, or 3) will easily have 0 errors because we can always find a quadratic curve that touches exactly those number of points. Hence:

- As the training set gets larger, the error for a quadratic function increases.
- The error value will plateau out after a certain  $m$ , or training set size.

### Experiencing high bias:

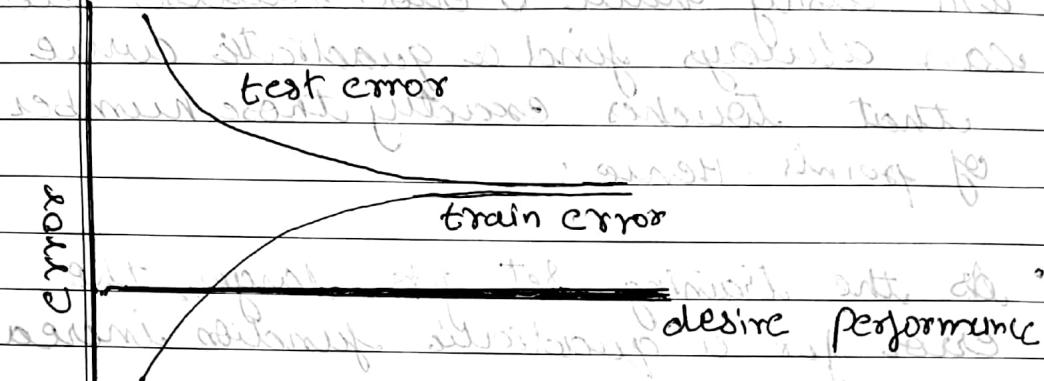
Low Training Set Size: Causes  $J_{train}(\theta)$  to be low and  $J_{cv}(\theta)$  will be high.

Large Training Set Size: Causes  $J_{train}(\theta)$  ~~to be~~ and  $J_{cv}(\theta)$  to be high with  $J_{train}(\theta) \approx J_{cv}(\theta)$ .

If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much.

## More on Bias Vs. Variance.

Typical learning curve for high bias  
(at fixed model complexity)



~~size of training set size~~

### Experiencing high Variance

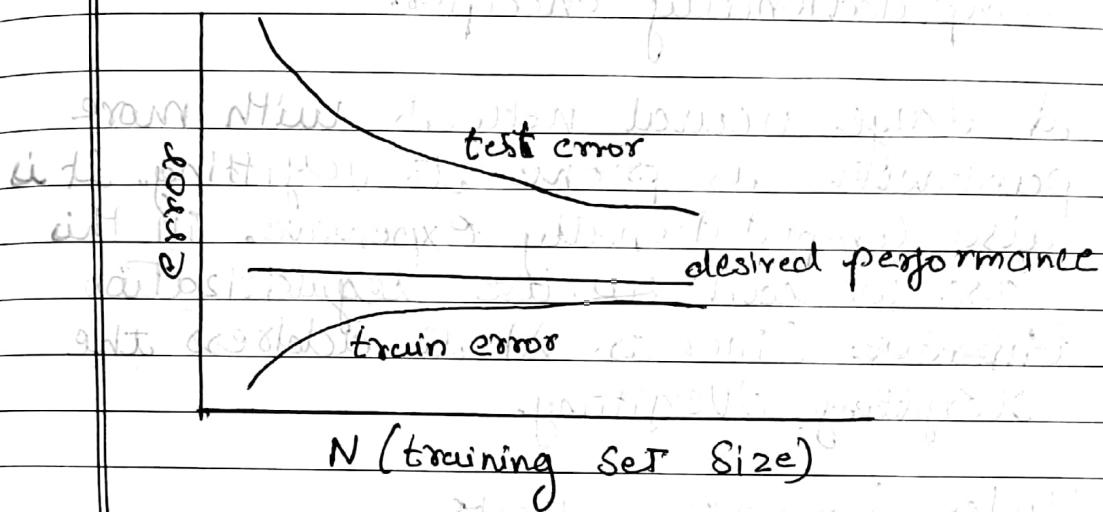
at low training set size:  $J_{\text{train}}(\theta)$  will be low and  $J_{\text{cv}}(\theta)$  will be high.

Large training set size:  $J_{\text{train}}(\theta)$  increases with ~~the~~ training set size and  $J_{\text{cv}}(\theta)$  continues to decrease without leveling off. Also,  $J_{\text{train}}(\theta) < J_{\text{cv}}(\theta)$  but the difference between them becomes significant.

If a ~~bad~~ learning algorithm is suffering from high variance, getting more

training data is likely to help.

Typical Learning curve for high Variance  
(at fixed model complexity)



\* Our decision process can be broken down  
well as follows:

- Getting more training examples : fixes high Variance
- Trying Smaller sets of ~~test~~ features : fixes high Variance
- Adding new features : fixes high bias
- Adding polynomial features : fixes high bias
- Decreasing  $\lambda$  : fixes high bias
- Increasing  $\lambda$  : fixes high Variance

# Diagnosing Neural Networks.

- A neural network with fewer parameters is prone to underfitting. It is also computationally cheaper.
- A large neural network with more parameters is prone to overfitting. It is also computationally expensive. In this case we can use regularization (increase  $\lambda$ ) to address the overfitting.

## Model Complexity Effects:

- Lower-order polynomials (low model complexity) have high bias and low variance. In this case, the model fits poorly (consistently).
- Higher-order polynomials (high model complexity) fit the training data extremely well and the test data extremely poorly. These have low bias on the training data, but very high variance.
- In reality, we would want to choose a model somewhere in between, that can generalize well but also fits the data reasonably well.

# Error analysis

## Recommended approach

(x) Start with simple algorithms that you can implement quickly. Implement it and test it on your cross-validation data.

(b) Plot learning curves to decide if more data, more features, etc. are likely to help.

(c) Error analysis: Manually examine the examples (in cross-validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

for example: Errors

Mcv = 500 examples in cross-validation set  
algorithm misclassifies 100 emails  
manually examine 100 errors, and categorize them based on:

- (i) What type of email is it?
- (ii) What cues/features you think would have helped the algorithm classify them correctly.

# Error metrics for skewed classes

## Cancer Classification Example

- Train a logistic regression model on  $\{x\}$ .  
 → Two labels: 1 if cancer was present. (y=1) if Cancer, y=0 otherwise
- find that you got 1% error on test-set.  
 → (99% correct diagnoses)
- only 0.50% of patients have cancer  
 → Skewed classes (imbalanced data)
- It means it's 99.5% of patients do not have cancer.

## Precision / Recall

$y=1$  in presence of rare class that we want to detect

Predicted class

	True	False	
Positive	True Positive	False Positive	Bad result
0	False	True	Good result
	negative	negative	

→ Precision

(of all patients where we predicted  $y=1$ , what fraction has cancer?)

$$\frac{\text{True positive}}{\text{True positive} + \text{False positive}} = \frac{\text{True Positive}}{\text{Predicted Positive}}$$

→ Recall

(of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positive}}{\text{True positive} + \text{False negative}} = \frac{\text{True positives}}{\text{actual positives.}}$$

→ Suppose we want to predict  $y=1$  (cancer) only if very confident. : Higher precision, lower recall.

→ Suppose we want to avoid missing too many cases of cancer (avoid false negatives):  
Higher recall, lower precision.

\* F-Score (F score)

How to compare precision/recall numbers?

$P = \text{Precision}$

$$\text{f1 Score} = 2 * \frac{PR}{P+R} \quad (\because R = \text{Recall})$$

$$P=0 \quad \text{or} \quad R=0 \Rightarrow \text{f1 Score} = 0$$

$$P=1 \quad \text{or} \quad R=1 \Rightarrow \text{f1 Score} = 1$$

Model

Model selection & Model choosing (in ML)  
 Model selection without tuning  
 Random Forest (with out tuning)

Naive Bayes = Interpreting Bayes

Naive Bayes (without tuning)

Naive Bayes also known as Naive

Naive Bayes is based on Bayes formula

# Anomaly detection

→ problem motivation:

Anomaly detection example.

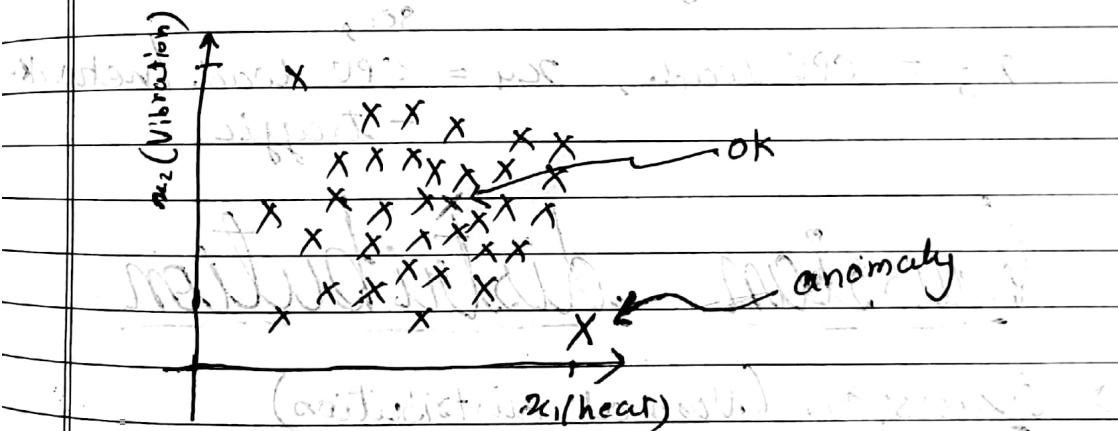
$x_1$  = heat generated

$x_2$  = vibration intensity.

⋮  
⋮  
⋮

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

New data:  $x_{\text{test}}$

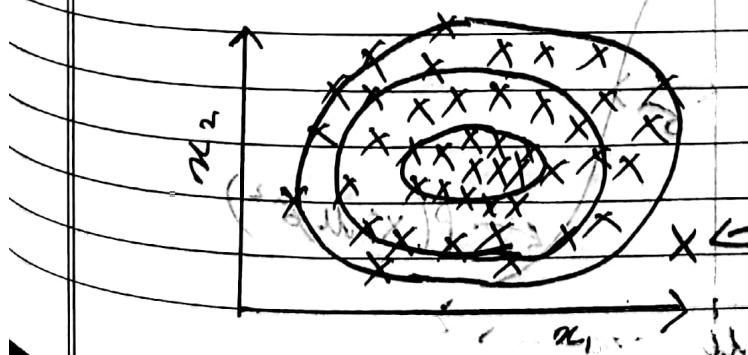


→ Density estimation:

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

→ Is  $x_{\text{test}}$  anomalous?

model:  $P(x)$



$P(x_{\text{test}}) < \epsilon \rightarrow \text{flag anomaly}$

$P(x_{\text{test}}) \geq \epsilon \rightarrow \text{OK}$

## # Anomaly detection example

→ fraud detection:

→  $x^{(i)}$  = feature of user i's activities.

→ model  $P(x)$  from data.

→ Identify unusual users by checking which have  $P(x) < \epsilon$

→ Manufacturing

→ Monitoring computers in a data center.

$x^{(i)}$  = feature of machine i

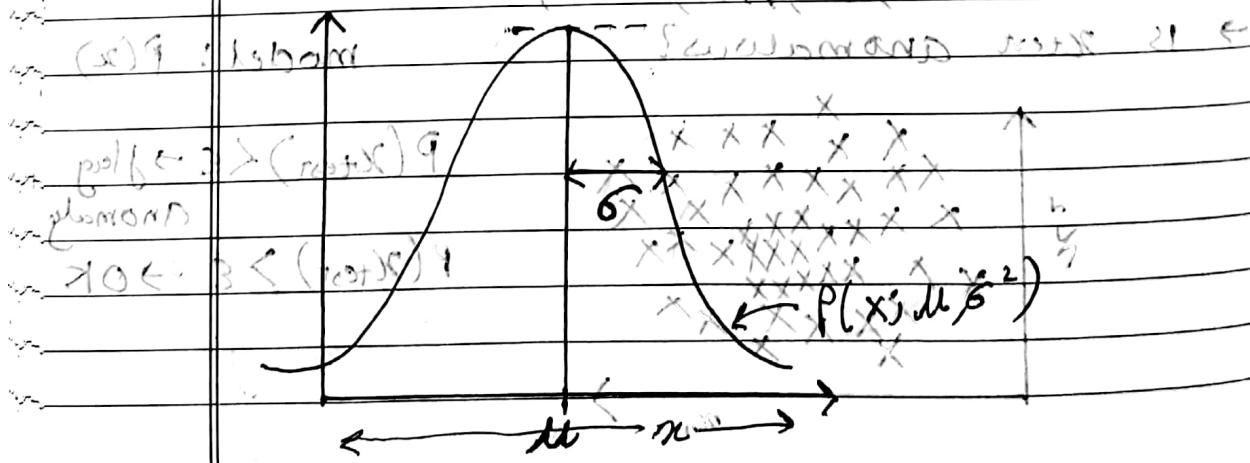
$x_1$  = memory use,  $x_2$  = number of disk accesses/sec,

$x_3$  = CPU load,  $x_4$  = CPU load/network-traffic.

## → Gaussian distribution

→ Gaussian (Normal distribution)

Say  $x \in \mathbb{R}$ . If  $x$  is distributed Gaussian with mean  $\mu$ , variance  $\sigma^2$ .



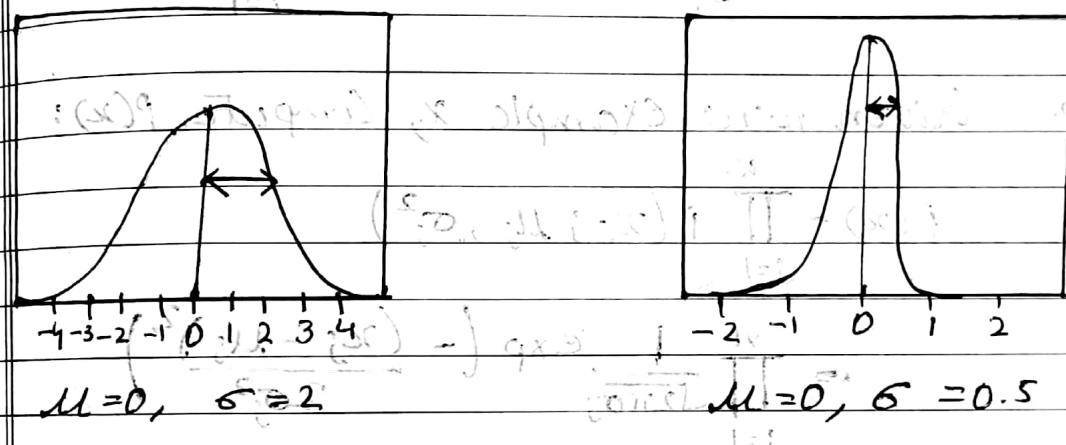
$\sigma$ : standard deviation;  $\sigma^2$ : Variance,  $\mu$ : mean

$$x \sim N(\mu, \sigma^2)$$

all  $x$  according to formula

$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

→ Gaussian distribution example:



→ parameter estimation:

Dataset:  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ ,  $x^{(i)} \in \mathbb{R}$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Algorithm

training set:  $\{x^{(1)}, \dots, x^{(m)}\}$ , Each example is  $x \in \mathbb{R}^n$

$$\begin{aligned} P(x) &= P(x_1; \mu_1, \sigma_1^2) \cdot P(x_2; \mu_2, \sigma_2^2) \cdot P(x_3; \mu_3, \sigma_3^2) \cdots P(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

→ Anomaly detection algorithm:

1. choose features  $x_i$  that you think  $x$  might be indicative of anomalous examples.
2. Fit parameters  $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3. Given new example  $x$ , Compute  $P(x)$ :

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$

$$= \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if  $P(x) < \epsilon$

→ The important of real-number evaluations.

When developing a learning algorithm (choosing features, etc.) making decisions is much easier if we have a way of evaluating our learning algorithm.

Assume we have some labeled data of anomalous and non-anomalous examples ( $y=0$  if normal,  $y=1$  if anomalous).

Training set:  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$  (assume normal examples / not anomalous)

- Cross Validation Set:  $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
- Test set:  $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

→ Aircraft engines motivating example.

- 10000 good (normal) engines
- 20 flawed engines (anomalous)

→ Training set: 6000 good engines

→ CV: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )

→ Test: 2000 good engines ( $y=0$ ), 10 anomalous ( $y=1$ )

→ If we want to find a threshold for a test set, it's sufficient to mark in the CV set.

→ If we want to find a threshold for a test set, it's sufficient to mark in the CV set.

• fit model  $P(x)$  on training set  $\{x^{(1)}, \dots, x^{(m)}\}$

• On a cross-validation / test example  $x_c$ , predict

$$y = \begin{cases} 1 & \text{if } P(x_c) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } P(x_c) \geq \varepsilon \text{ (normal)} \end{cases}$$

possible evaluation metrics:

- True positive, false positive, false negative, True negative

- precision / Recall

-  $F_1 = \frac{2 \cdot \text{precision} \cdot \text{Recall}}{\text{precision} + \text{Recall}}$

Can also use cross-validation set to choose parameter  $\varepsilon$

if necessary, iterate

## Anomaly detection vs Supervised learning

- |   |   |
|---|---|
| <ul style="list-style-type: none"> <li>→ Very Small number of positive examples (<math>y=1</math>). (0-20 is common)</li> <li>→ Large number of negative (<math>y=0</math>) examples.</li> <li>→ many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like. future positive example likely to be similar to ones in training set</li> <li>→ future anomalies may look nothing like any of the anomalous examples we've seen so far.</li> </ul> | <ul style="list-style-type: none"> <li>→ large number of positive and negative examples.</li> <li>→ enough Positive example for algorithm to get a sense of what positive example are like, future positive example likely to be similar to ones in training set</li> <li>→ Example:</li> </ul> |
|---|---|

- |   |  |
|---|--|
| <ul style="list-style-type: none"> <li>• Manufacturing (e.g., aircraft engines)</li> <li>• monitoring machine in a data center</li> </ul> | <ul style="list-style-type: none"> <li>• Email: Spam classification</li> <li>• Weather prediction</li> <li>• Cancer classification.</li> </ul> |
|---|--|