

## Neural Network

cost function

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

(L) = total no. of layers in network

(S<sub>l</sub>) = no. of units (not counting bias unit)  
in layer l  
→ l 層上的神經元有幾個。

J( $\theta$ )

把輸出層加總在一起

$$= -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^{S_l} y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)}))_k \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_j^{(l)})^2$$

regularization 時不要加  $\theta_{0i}$ , 其為 bias  
的 weight

(cost function for classification)

# Backpropagation Algorithm

$J(\theta)$

把輸出層加總在一起

$$= -\frac{1}{m} \left[ \sum_{i=1}^m \left[ \sum_{k=1}^K y_k^{(i)} \log(h_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-h_\theta(x^{(i)}))_k \right] \right]$$

$$+ \frac{\lambda}{2m} \sum_{i=1}^{L-1} \sum_{j=1}^{S_i} \sum_{j'=1}^{S_{i+1}} (\theta_{j,j'}^{(i)})^2$$

Need code to compute:

$$\min_{\theta} J(\theta) \longrightarrow \frac{\partial}{\partial \theta_j^{(i)}} J(\theta)$$

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

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

$$\delta_j^{(4)} = \alpha_j^{(4)} - y_i$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} * g'(z^{(3)}) = \alpha^{(3)} * (1 - \alpha^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} * g'(z^{(2)}) = \alpha^{(2)} * (1 - \alpha^{(2)})$$

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

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

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

for  $i = 1 : m$

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

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

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

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

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

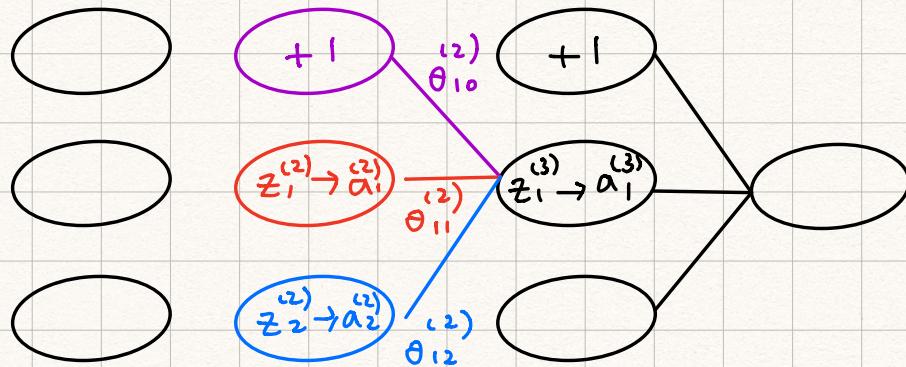
$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}$  if  $j \neq 0$

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

$\nabla_{\theta_{ij}}^{(2)} = \frac{\partial}{\partial \theta_{ij}} J(\theta) *$

## Backpropagation intuition

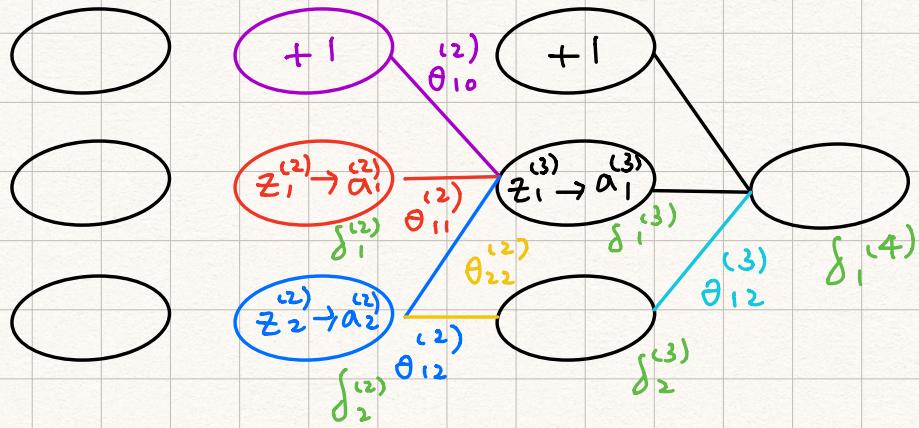
### Forward Propagation



$$z_1^{(3)} = \theta_{10}^{(2)} \times 1 + \theta_{11}^{(2)} \alpha_1^{(2)} + \theta_{12}^{(2)} \cdot \alpha_2^{(2)}$$

$z \xrightarrow{\text{sigmoid}} \alpha$

### Backpropagation Propagation



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

$$\delta_2^{(3)} = \theta_{12}^{(3)} \delta_1^{(4)}$$

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

$$= \frac{\partial}{\partial z_j^{(l)}} \text{cost}(\cdot) \text{, for } j > 0$$

## Back propagation in Practice

### Advanced optimization

function [jVal, gradient]

皆需向量型式

代价函数 偏导数值

= costFunction(theta)

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

matrix size :

$$\theta^{(1)} = [ ]_{10 \times 11} \quad \theta^{(2)} = [ ]_{10 \times 11}$$

$$\theta^{(3)} = [ ]_{1 \times 11}$$

- matrix to vectors ?

thetaVec

$$= [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)];$$

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

- vectors to matrix ?

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

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

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

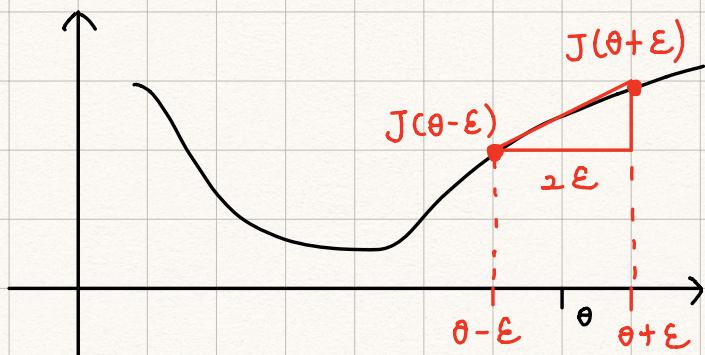
## Summary

$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  合成 - 向量 initialTheta 後，  
丟入 fminunc

從 costFunction(thetaVec) 的 thetaVec 取得  
 $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$  後，用  $f_p, b_p$  計算  
 $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\theta)$

最後，將  $D^{(1)}, D^{(2)}, D^{(3)}$  合成 - 向量 gradientVec

## Gradient Checking 梯度檢查



$$\frac{d}{d\theta} J(\theta) = \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon} \quad \epsilon = 10^{-4}$$

Parameter vector  $\theta$

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \dots, \theta_n)}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \dots, \theta_n)}{2\epsilon}$$

⋮  
⋮

use for loop get grad Approx

for i = 1 : n

thetaPlus = theta

thetaPlus(:, i) = thetaPlus(:, i) + epsilon

thetaMinus = .....

thetaMinus(:, i) = ..... - epsilon

grad Approx(:, i) = (J(thetaPlus) - J(thetaMinus))

end

check that grad Approx  $\approx$  DVec

↑  
From b p

## Summary

1. 以 bp 計算 DVec
2. 應用梯度檢驗法獲得 grad Approx
3. 確保 grad Approx ≈ DVec
4. 關掉梯度檢驗法，以 bp 訓練資料

## Random initialization

### Initial Theta setting

initialTheta = zeros(n, 1) ?

or initialTheta 裡的值都相等

在 logistic regression 可行，但是

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

$$\alpha_1^{(2)} = \alpha_2^{(2)}$$

不論如何計算，所有  $\theta$  得到的更新都一樣

## Symmetry breaking

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

Theta |

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

rand(10, 11) 產生  $[ ]_{10 \times 11}$ , one

Put it together

## Picking a network architecture

- Input: Dimension of feature
- output: Number of classes
- Reasonable default: 1 hidden layer, or  
 $f > 1$  hidden layer, have same no. of  
hidden units in every layer  
隱藏層的神經元越多越好

## Training a neural network

1. 訓練機初始化權重
2. 以  $f_p$  得到  $h_{\theta}(x^{(i)})$  for any  $x^{(i)}$
3. 計算 cost function  $J(\theta)$
4. 以 bp 計算  $\frac{\partial}{\partial \theta_j} J(\theta)$
5. 梯度檢查以確認梯度確實有下降，  
程式無 bug 後，  
關掉梯度檢查
6. 以梯度下降法 or 其他算法來最小化  
cost function

試圖找出某個  $\theta$  使得  $J(\theta)$  為局部  
最小 or 全域最小

for  $i = 1:m$

Perform  $f_p$  and  $b_p$  using example  
 $(x^{(i)}, y^{(i)})$

(Get activation  $a(l)$  and delta terms  
 $d(l)$  for  $l = 2, \dots, L$ )