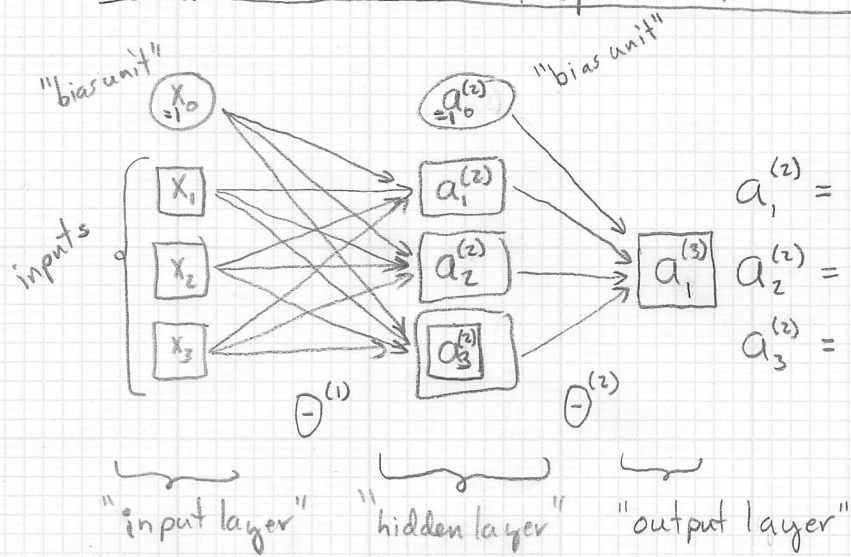


# Neural Networks: Representation

20120214 ABR

①



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

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

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

$$a_1^{(3)} = g(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \dots)$$

$$a^{(j+1)} = g(\theta^T a^{(j)})$$

$$= g(z^{(j+1)})$$

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

$$z^{(j+1)} = \theta^T a^{(j)}$$

Note that input data usually resides in rows of an input matrix:

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

If  $\theta^{(j)}$  is stored in-order, then it must be transposed to conform with  $X$ , and the resulting output vector  $a^{(j+1)}$  is also a row vector:

$$X \cdot \theta^{(j)T} = [x_0 \ x_1 \ x_2 \dots] \begin{bmatrix} \theta_{10}^{(j)} & \theta_{20}^{(j)} & \dots \\ \theta_{11}^{(j)} & \theta_{21}^{(j)} & \dots \\ \theta_{12}^{(j)} & & \ddots \\ \vdots & & \end{bmatrix}$$

$$= [z_1^{(j)} \ z_2^{(j)} \ z_3^{(j)} \dots] = z^{(j)}$$

Then  $a^{(j+1)} = g(z^{(j)})$  is in order for multiplication with  $\theta^{(j+1)T}$  and so on down the chain.

For classification purposes, we select the output layer element with the highest activation level. Octave provides built-in functions for this step.

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

Sum over output units

Where:  $h_{\theta}(x) \in \mathbb{R}^K$ ,  $(h_{\theta}(x))_i = i^{\text{th}}$  output

$K$  = number of classes, number of output units ( $K = S_L$ )

$S_l$  = number of units in layer  $l$  (not counting bias unit)

$L$  = total number of layers in network

$y \in \{0, 1\}^K$

## Back-Propagation: Gradient and Algorithm

need:  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}$  (ignoring  $\lambda$  for now)

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

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

$$\delta_j^{(L-1)} = \Theta^{(L-1)T} \delta_j^{(L)} * g'(z^{(L-1)}) \Rightarrow g'(z^{(L-1)}) = a^{(L-1)} * (1 - a^{(L-1)})$$

$$\delta_j^{(L-2)} = \Theta^{(L-2)T} \delta_j^{(L-1)} * g'(z^{(L-2)})$$

$\vdots$   
No  $\delta_j^{(1)}$

### Algorithm:

-- set  $\Delta_{ij}^{(L)} = 0$

-- for  $i = 1$  to  $m$ :

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

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

-- using  $y^{(i)}$ , compute  $\delta_j^{(L)} = a_j^{(L)} - y_j^{(i)}$

-- for  $l = L-1$  to  $2$ :

$$\delta_j^{(l)} = \Theta^{(l+1)T} \delta_j^{(l+1)} * g'(z^{(l)})$$

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

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



regularization terms are added as follows:

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

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

## Intuition:

consider the case of one training example  $(x, y)$

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

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

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

(don't forget to add  $a_0^{(2)}$ )

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

$$a^{(3)} = g(z^{(3)})$$

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

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

$$h_{\Theta}(x) = \text{or } g(\Theta^{(3)} g(\Theta^{(2)} g(\Theta^{(1)} \cdot x)))$$

## Back-Propagation: Concise Summary

- set  $\Delta_{ij}^{(L)} = 0$

- For  $i = 1$  to  $m$ :

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

- compute all layer activations  $a^{(l)}$

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

- for  $l = L-1$  to  $2$ :

-  $\delta^{(l)} = \Theta^{(l+1)T} \delta^{(l+1)} * (a^{(l)} * (1 - a^{(l)}))$

-  $\Delta_{ij}^{(l)} += a_j^{(l)} \delta_i^{(l+1)} \quad \text{or } \Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)} \cdot a^{(l)T}$

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

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

← These are the components of the gradient

$$\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \left[ y^{(i)} \log(h_{\Theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\Theta}(x^{(i)})) \right]$$

$$\text{Use } \left[ \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) \approx \frac{J(\Theta + \epsilon) - J(\Theta - \epsilon)}{2\epsilon} \right] \text{ to test gradient function from back propagation.}$$

test gradient function from back propagation.

\*\* constant-initialization ("problem of symmetric weights") does not work ( $\frac{\partial}{\partial \Theta} J = 0$ )

## "Problem of Symmetric Weights"

Constant initialization is a "flat" spot on the cost function. Initialize with random deviates to break symmetry.

## Network Architecture Choices

- # Input units } defined by problem
- # Output units }
- # hidden layers  $\rightarrow$  1 is a reasonable default
- # units/hidden layer  $\rightarrow$  generally more is better

## Training

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 back-propagation to compute partial derivatives  $\frac{\partial}{\partial \theta_{jk}^{(l)}} J(\theta)$

for  $i = 1:m$  (this can be vectorized)

Perform forward propagation & backward propagation using example  $(x^{(i)}, y^{(i)})$ , saving activations and deltas  $(a^{(l)}, \delta^{(l)})$  for  $l = 2, \dots, L$

Compute  $\Delta$

Compute  $D(\theta, \lambda)$

5. Use gradient checking to test backpropagation code
6. Use gradient descent or optimized solvers to minimize  $J(\theta)$