

Cost Function:

$L$ : total number of buyers in the network.

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

$k$ : number of output units/classes.

$m$ : number of training samples.

$h(x)_k$ : The hypothesis that results in the  $k^{\text{th}}$  output.

二分类问题: Binary classification.

$$h(x) \in \mathbb{R}$$

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

$S_L$ : 输出层的 node 数量.

Multi-class classification ( $k$  classes).

$$h(x) \in \mathbb{R}^k$$
  
 $S_L = k \quad (k \geq 3)$   
Output:  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$

We use the same cost function as logistic regression:

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

For Neural Networks:

$$J(\theta) = -\frac{1}{m} \sum_i \left[ \sum_k y_k^{(i)} \log(h(\theta)x^{(i)})_k + (1-y_k^{(i)}) \log(1-(h(\theta)x^{(i)})_k) \right]$$
  
 $\Theta^{(l)} \in \mathbb{R}^{(L-1) \times S_l}$   
 $\Theta^{(l)} + \frac{\lambda}{2m} \sum_{j=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{j,i}^{(l)})^2$   
→ 排除第  $0$  行,  $j$  循环所有的行,  $i$  循环所有的列, 由  $S_L$  层的激活单元数量所决定.

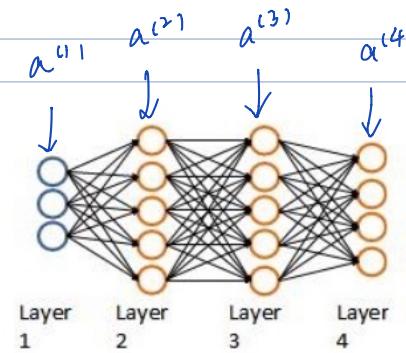
Notice: ① Theta matrix 的列数 = 当前网络层节点数 (包含 bias 单元).  
② Theta matrix 的行数 = 下一层网络层节点数 (不包含 bias 单元)

Backpropagation Algorithms:

$$\min_{\theta} J(\theta)$$

Need code to compute:

$$\frac{\partial}{\partial \theta_{ij}} J(\theta) \leftarrow \text{gradients.}$$



Given one training example  $(x, y)$ :

Forward propagation:

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

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

$$a^{(2)} = g(z^{(2)}) \text{ (add } a_0^{(2)})$$

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

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

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

$$a^{(4)} = g(z^{(4)}) \text{ (add } a_0^{(4)})$$

先对激活函数,

再添加偏置.

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

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

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

$$\underline{\delta^{(3)}} = (\theta^{(3)})^T \underline{\delta^{(4)}} * g'(z^{(3)}) \text{ element-wise}$$

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

$$g(z^{(2)}) * (1 - g(z^{(2)}))$$

BP 时, 依一层的误差对前一层  
的影响.

因为第一层是输入层，不存在误差。当我们有了所有的误差表达式后，便可以计算  $J(\theta)$  的偏导数了：

$$\frac{\partial}{\partial \theta_{ij}} J(\theta) = a_j^{(l)} s_i^{(l+1)} \quad (\text{ignore } x_j \text{ if } i=0)$$

Note:  $l$ : 表示目前所计算的是第几层；

$j$ : 因前计算层中激活单元的下标，也将是下一层的第  $j$  个输入变量的下标；

$i$ : 下一层中误差单元的下标，是受到权重矩阵  $(\theta)$  中第  $i$  行影响的下一层中的误差单元的下标。

Training set  $\{(x^{(1)}, y^{(1)})^T, \dots, (x^{(m)}, y^{(m)})^T\} \leftarrow m \text{ 个训练样本.}$

Set  $b_{ij}^{(l)} = 0$  (for all  $l, i, j$ ) as an accumulator to compute  $\frac{\partial}{\partial \theta_{ij}} J(\theta)$ .  
 $\Delta_{ij}$  表示误差矩阵，表示第  $l$  层的第  $i$  个 unit 受到第  $j$  个参数而导致的误差。

For  $i=1 \text{ to } m$ :

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

↑ 第  $l$  层的误差矩阵。

Forward propagation:  $a^{(l)} = a^{(l-1)}$ ,  $l=2, 3, \dots, L$ .

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

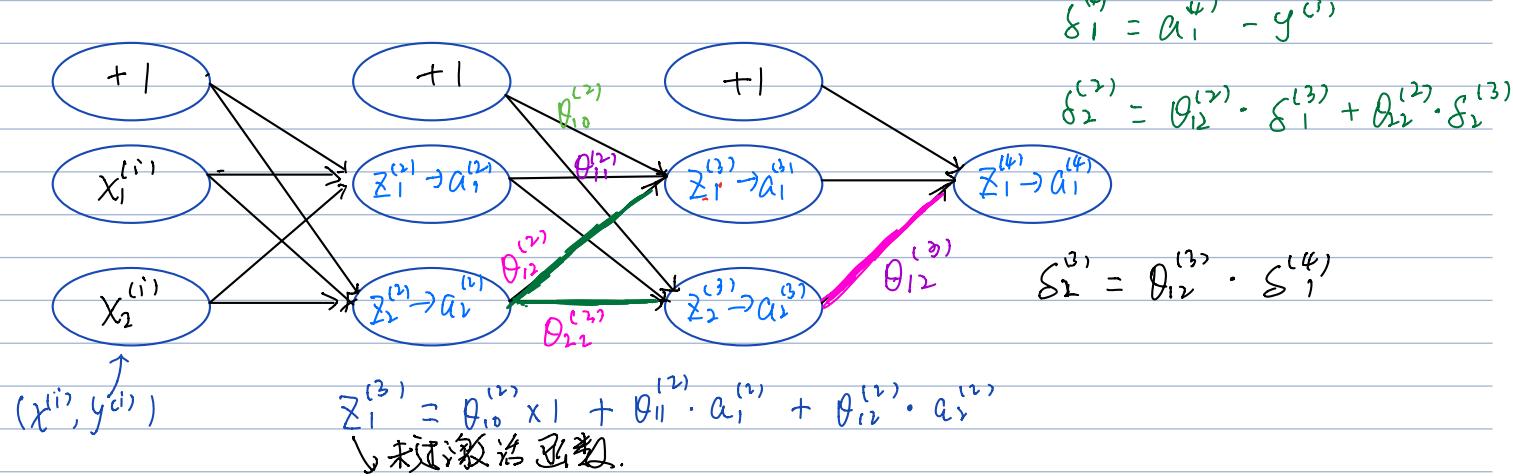
compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$\Delta_{ij} := \Delta_{ij}^{(L-1)} + a_j^{(L-1)} s_i^{(L-1)}$  ( $\Delta^{(L)} := \Delta^{(L-1)} + s^{(L-1)} (a^{(L)})^T$ )

$D_{ij}^{(L)} := \frac{1}{m} b_{ij}^{(L)} + \lambda \theta_{ij}^{(L)}$  if  $j \neq 0 \leftarrow \text{multi-class}$

$D_{ij}^{(L)} := \frac{1}{m} \Delta_{ij}^{(L)}$  if  $j = 0 \leftarrow \text{binary.}$

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



$$J(\theta) = -\frac{1}{m} \left[ \sum y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{s_l+1} (\theta_{ij}^{(l)})^2$$

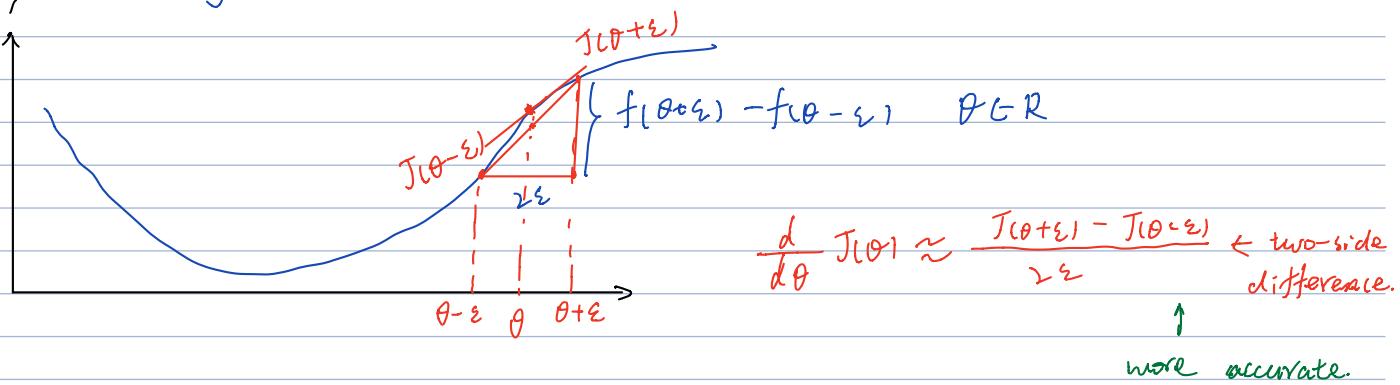
We focus on single example  $(x^{(i)}, y^{(i)})$ ,  $\lambda = 0$ , 1 output unit

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

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

$\delta_j^{(l)}$  = "error" of cost for  $a_j^{(l)}$ , Formally,  $\delta_j^{(l)} = \frac{\partial}{\partial \theta_{ij}} \text{Cost}(i)$  (for  $j > 0$ )

## Gradient checking:



$\theta \in \mathbb{R}^n$  (e.g.  $\theta$  is "unrolled" version of  $\theta^{(1)}$ ,  $\theta^{(2)}$ ,  $\theta^{(3)}$ ).

$$\theta = [\theta_1, \theta_2, \dots, \theta_n]$$

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

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

⋮

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2\epsilon}$$

numerically estimate gradient  
of  $J(\theta)$  w.r.t to  $\theta_j$ .

In practical = gradApprox  $\approx$  prec  
 ↓  
 estimation From Back Prop --

Implementation Note:

- Implement backprop to compute prec
- Implement numerical gradient check to compute gradApprox.
- Make sure they give similar values.
- Turn off gradient checking before using backprop code for training.

Random Initialization:

如果令所有的参数都为0, hidden layer所有nodes都会得到一样的值.

Random Initialization: Symmetry breaking

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

$$\text{eg. } \theta_{ij} = \text{rand}(10 * 1) * (2\epsilon) - \epsilon.$$

## Training a Neural Network:

1. Randomly initialize the weights;
2. Implement forward propagation to get  $h(\theta)(x^{(i)})$  for any  $x^{(i)}$ .
3. Implement the cost function  $J(\theta)$
4. Implement Backprop to compute partial derivatives
5. Use gradient checking to confirm that your Backprop works, then disable it
6. Use gradient descent or other built-in optimization function to minimize the cost function with the weights in theta.

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