

# Knowledges:

BP

反向传播法  $\Rightarrow$  传统神经网 络训练方法

: 随机设置初值，计算当前网络 output，与标签残差去改变前面各层参数，直至收敛。

本质即为梯度下降法求得  $J(\theta)$  的偏导数，确定下降的方向。

梯度扩散：在深层网络结构中，残差传递至最前面层已变得太小，即梯度扩散了。

# Neural Network:

## I. Cost function:

来源  $\Rightarrow$  Logistic Regression:

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

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

Neural Network:

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

## II. Forward Propagation : (FP)

$$a^{(1)} = x \quad (\text{No need to add bias unit})$$

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

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

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

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

⋮

$$z^{(k+1)} = \theta^{(k)} \cdot a^{(k)}$$

$$a^{(k+1)} = g(z^{(k+1)}) \quad (a^{(k+1)} \text{ add } a_0^{(k+1)})$$

⋮

How to get the expression of  $g^{(L)} \Rightarrow$

$$\textcircled{1} \quad \delta^l = (\theta^l)^T \cdot \frac{\partial E^{(l+1)}}{\partial x^{(l+1)}}$$

$$E = \frac{1}{2} \sum_{i=1}^m (y_{\text{output}} - y_{\text{real}})^2$$

## 反向传播

现在我们获得了  $\delta^l$ , 接下来便可以根据链式法则得出

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} = y_{\text{output}} - y_{\text{real}}$$

$$\frac{\partial E}{\partial x} = \frac{dy}{dx} \frac{\partial E}{\partial y} = \frac{d}{dx} f(x) \frac{\partial E}{\partial y}$$

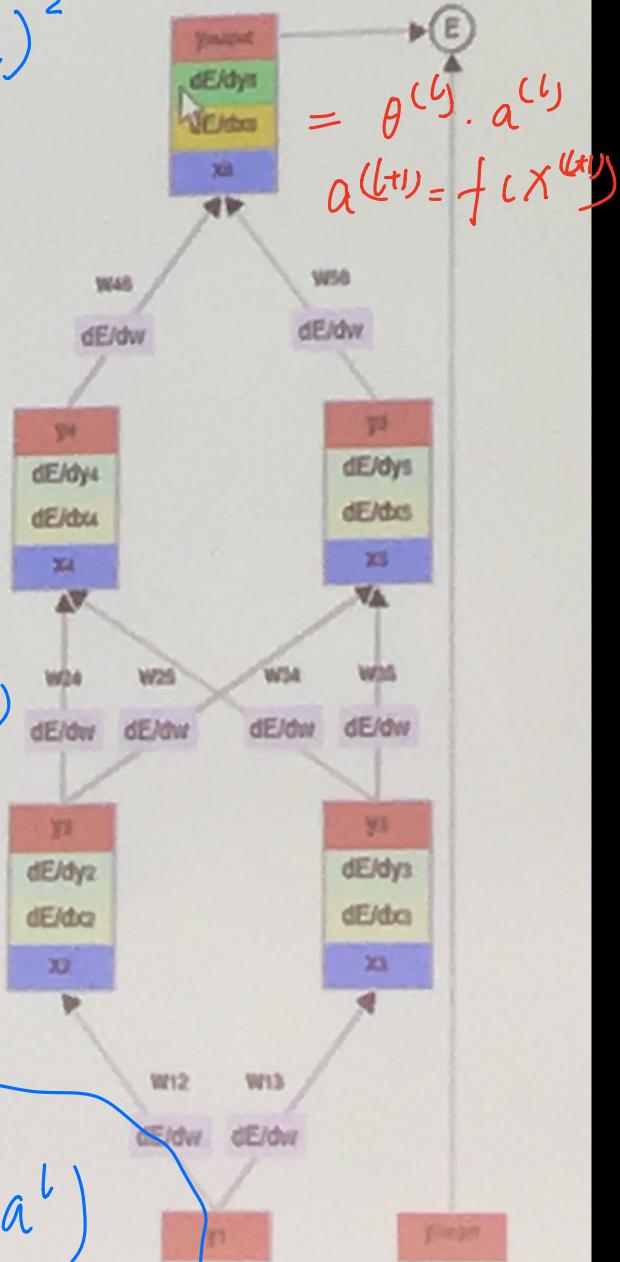
$$\text{其中, 当 } f(x) \text{ 是 S 型激活函数} \Rightarrow f(x) = \frac{1}{1 + e^{-x}}$$

$$\text{时, } \frac{d}{dx} f(x) = f(x)(1 - f(x))$$

$$\begin{aligned} f'(x) &= f(x)(1 - f(x)) \\ &= a^l \cdot (1 - a^l) \end{aligned}$$

$$\delta^l = (\theta^l)^T \cdot \frac{\partial E^{(l+1)}}{\partial x^{(l+1)}}$$

$$\begin{aligned} \frac{\partial E}{\partial y^{(l+1)}} \cdot a^l \cdot (1 - a^l) \\ \delta^{(l+1)} \end{aligned}$$



# III - Back Propagation (BP):

I. Calculate

$$\delta_j^{(L)}$$

for each  $j \in \{1, \dots, S_L\}$

$$l \in \{2, \dots, L\}$$

$\Rightarrow$  error for node  $j$  in layer  $l$ .

the first layer is input layer, there is no error for every node  $j$ .

How to calculate  $\Rightarrow$  right to the left (or back)

$$\text{error of layer } L \leftarrow \delta^{(L)} = a^{(L)} - y$$

$$\frac{\partial E}{\partial a}$$

$$\frac{\partial a}{\partial z} = \frac{\partial g(z)}{\partial z}$$

$$\text{error of layer } (L-1) \leftarrow \delta^{(L-1)} = (\theta^{(L)})^T \cdot \delta^{(L)} \cdot g'(z^{(L-1)})$$

( $L-1$ ), calculate using  $\delta^{(L)}$

$$= (\theta^{(L-1)})^T \cdot \delta_j^{(L)} \cdot a^{(L-1)} \cdot * (1 - a^{(L-1)})$$

II. Calculate  $\Delta_{ij}^{(l)}$

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

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

$$= g(z) \cdot (1 - g(z))$$

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

$$\emptyset \quad \Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \cdot \delta_i^{(l+1)}$$

$$\Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + \underbrace{\delta_i^{(l+1)} \cdot \alpha_j^{(l)}}_{\sim}$$

IN calculate  $\Delta_{ij}^{(l)}$   $\Leftrightarrow \frac{\partial J(\theta)}{\partial \theta_{ij}}$  . BP cost function  
对第  $l$  层的  $i$  行的参数

$$\frac{\partial J(\theta)}{\partial \theta_{ij}} = D_{ij}^{(l)} = \begin{cases} \frac{1}{m} \cdot \Delta_{ij}^{(l)} + \lambda \cdot \theta_{ij}^{(l)}, & j \neq 0 \Rightarrow \text{other unit} \\ \frac{1}{m} \cdot \Delta_{ij}^{(l)} - \text{if } j = 0 \Rightarrow \text{bias unit} \end{cases}$$

IV. Gradient checking  $\Rightarrow$  numerical approximation

$$\text{of } \frac{\partial J(\theta)}{\partial \theta}$$

$$\frac{\partial J(\theta)}{\partial \theta_i} \approx \text{gradApprox} =$$

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

$\triangle \epsilon$  is usually  $10^{-3}, 10^{-4}$  or even smaller.

$\triangle$  Note that the most important thing about

gradient checking is that to calculate each

$$\frac{\partial J(\theta)}{\partial \theta_i}, \text{ we remain other } \theta_j (j \neq i) \text{ unchanged}$$

△ After checking - turn off the gradient checking

and just use backpropagation code for learning -

because gradient checking is very slow !

## V. Random initialization

Reason : We can't use the all zero initialization to init  $\overrightarrow{\theta}$  as we've done in Logistic Regression - because this will cause symmetric problems.

# Application:

OCR: Optical Character Recognition.

Pipeline process:

