

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 δ^l , 接下来便可以根据链式法则得出

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

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

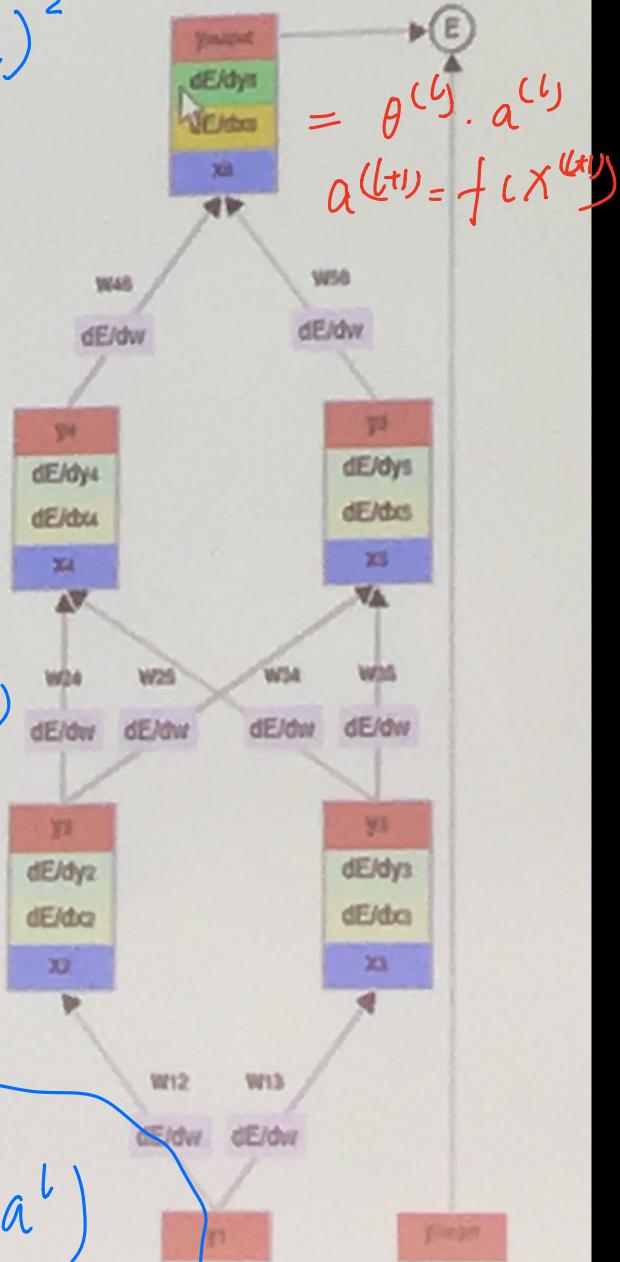
$$\text{其中, 当 } f(x) \text{ 是 S 型激活函数} \Rightarrow f(x) = \frac{1}{1 + e^{-x}}$$

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

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

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

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



III - Back Propagation (BP):

I. Calculate

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

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

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

error for node j in layer l .

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

$(l-1)$ th layer (hidden layer)

input:	$a^{(l-1)}$	$w^{(l-1)}/\theta^{(l-1)}, b^{(l-1)}$
		$z^{(l-1)} = w^{(l-1)} \cdot a^{(l-1)} + b^{(l-1)}$
		$\delta^{(l-1)} = \frac{\partial J}{\partial z^{(l-1)}}, \Delta^{(l-1)} = \frac{\partial J}{\partial w^{(l-1)}}$
		$= \frac{\partial J}{\partial z^{(l)}} \cdot \frac{\partial z^{(l)}}{\partial w^{(l-1)}}$

output layer

$$z^{(l)} = a^{(l)} \cdot w^{(l)} + b^{(l)}$$

$$a^{(l)} = f^{(l)}(z^{(l)})$$

$$y$$

$$\hat{y} = z^{(l)}$$

$$J^{(l)} = (y - \hat{y})$$

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

error of layer L $\leftarrow \delta^{(L)} = z^{(L-1)} - y$

$$\delta^{(L-1)} = \frac{\partial J}{\partial z^{(L-1)}} = \underbrace{\frac{\partial J}{\partial z^{(L)}}}_{\text{是 } l} \cdot \frac{\partial z^{(L)}}{\partial z^{(L-1)}}$$

$$g(z^{(L-1)}). w^L$$

$$\frac{\partial z^{(L)}}{\partial z^{(L-1)}} = \delta^{(L)} \cdot g'(z^{(L-1)}) \cdot w^L$$

error of layer $(L-1)$ \leftarrow

$$= \delta^{(L)} \cdot$$

$$z^{(L-1)} \cdot (1 - z^{(L-1)}) (w^L)^T.$$

第 L 层为输出层，
认为没有参数。

II. Calculate $\Delta_{ij}^{(l-1)}$

$$\begin{aligned} g(z) &= \frac{1}{1+e^{-z}} \\ g'(z) &= \left(\frac{1}{1+e^{-z}}\right) \cdot (1 - 1 + e^{-z}) \\ &= g(z) \cdot (1 - g(z)) \end{aligned}$$

How \Rightarrow ① Set $W_{ij}^{(l-1)} = 0$ for all l, i, j

$$a^{(l+1)} \rightarrow 0 \quad \text{②} \quad W_{ij}^{(l-1)} = W_{ij}^{(l-1)} + \underbrace{\alpha_j^{(l-1)} \cdot \delta_i^{(l-1)}}_{\left(\frac{\partial J}{\partial w_{ij}^{(l-1)}} \right) \left(\frac{\partial z^{(l+1)}}{\partial w_{ij}^{(l-1)}} \right) \rightarrow \alpha^{(l-1)}}$$

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

$$\left(\Rightarrow \Delta^{(l-1)} = \delta^{(l-1)} \cdot a^{(l-1)} \right)$$

III calculate $D_{ij}^{(l)}$ $\Leftrightarrow \underbrace{\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}}}_{\text{BP cost function}} \cdot \text{对第 } l \text{ 层的 } j \text{ 行的参数的偏导数}$

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

IV. Gradient checking \Rightarrow numerical approximation
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 ϵ is usually $10^{-3}, 10^{-4}$ or even smaller.

\triangle Note that the most important thing about gradient checking is that to calculate each

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

\triangle 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 $\vec{\theta}$ as we've done in Logistic Regression - because this will cause symmetric problems.

Application:

OCR: Optical Character Recognition.

Pipeline process:

