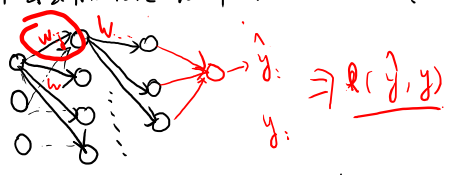


反向传播之即算法的推导：

1. 为什么会有反向传播算法产生。大家看图：



现：想求得参数权重matrix的更新。  
就应搞清楚，参数关于输出的梯度。

但是：请大家观察一下，参数关于output的gradient可以直接一箭到地的计算出来吗？  
答案显然是不可以的！

所以传播之前定义几个符号以便理解：

1° 选第l层的参数为  $W^{(l)}$   
 $W_{ij}^{(l)}$  表示：第l层的第i个神经元与第(l+1)层的第j个神经元的连接权重！

2°  $l(y, \hat{y})$  表示输出层的损失函数 loss.

3° 且：
$$\begin{cases} z^{(l)} = W^{(l)} a^{(l-1)} \\ z^{(l+1)} = W^{(l+1)} a^{(l)} \end{cases}$$
  

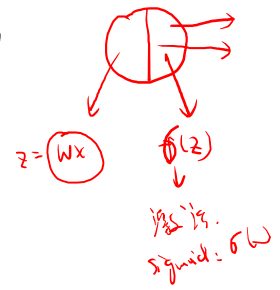
$$\begin{cases} a^{(l)} = \sigma(z^{(l)}) \\ a^{(l+1)} = \sigma(z^{(l+1)}) \end{cases}$$

其中： $\sigma(\cdot)$  为 sigmoid 函数。

$$\frac{\partial l(y, \hat{y})}{\partial W_{ij}^{(l)}} = \frac{\partial l(y, \hat{y})}{\partial z^{(l)}} \cdot \frac{\partial z^{(l)}}{\partial W_{ij}^{(l)}}$$

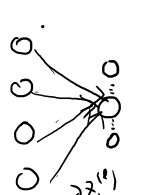
↓  
先不管这该怎么得的！

前项： $\frac{\partial z^{(l)}}{\partial W_{ij}^{(l)}}$  怎么求



$$\frac{\partial z^{(l)}}{\partial w_{ij}^{(l)}}$$

大家想清楚。  
z 是啥？



$$z = [z_1, z_2, \dots, z_n]$$

故:  $\frac{\partial z^{(l)}}{\partial w_{ij}^{(l)}} = \left[ \frac{\partial z_1^{(l)}}{\partial w_{ij}^{(l)}}, \frac{\partial z_2^{(l)}}{\partial w_{ij}^{(l)}}, \dots, \frac{\partial z_n^{(l)}}{\partial w_{ij}^{(l)}} \right]$

$$= [0, 0, \dots, \frac{\partial (w_{ij}^{(l)} \cdot a^{(l-1)})}{\partial w_{ij}^{(l)}}, 0, \dots]$$

$$= [0, 0, \dots, a^{(l-1)}, \dots, 0, 0]$$

$$\text{即: } \nabla a_j^{(l-1)}$$

$$z^{(l)} = w^{(l)} \cdot \underbrace{a^{(l-1)}}$$

再推导:  $\frac{\partial \ell(y, y)}{\partial z^{(l)}} = \text{即: } \frac{\partial \ell(y, y)}{\partial z^{(l)}} \stackrel{?}{=} \delta^{(l)}$

$$\frac{\partial \ell(y, y)}{\partial z^{(l)}} = \frac{\partial \ell(y, y)}{\partial z^{(l+1)}} \cdot \frac{\partial z^{(l+1)}}{\partial z^{(l)}} \stackrel{?}{=} \frac{\partial \ell(y, y)}{\partial z^{(l+1)}} \cdot \frac{\partial a^{(l)}}{\partial z^{(l)}}$$

$\delta^{(l)}$

$$\delta^{(l+1)} \quad w^{(l+1)} \quad \sigma'(\cdot)$$

$$a^{(l)} = \sigma(z^{(l)})$$

$$z^{(l+1)} = w^{(l+1)} \cdot a^{(l)}$$

$$z^{(l)} = [z_1^{(l)}, \dots, z_n^{(l)}]$$

故:  $\frac{\partial \ell(y, y)}{\partial w_{ij}^{(l)}} = \nabla a_j^{(l-1)} \cdot \delta^{(l+1)} \cdot w^{(l+1)} \cdot \sigma'(\cdot)$

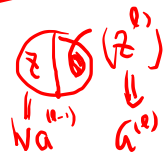
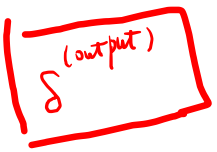
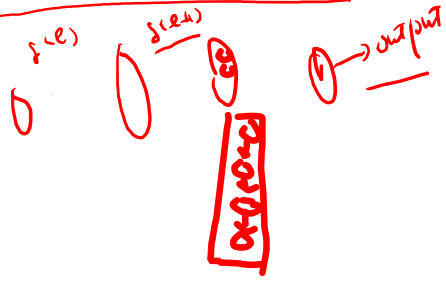
$$= \nabla a_j^{(l-1)} \cdot \delta^{(l)}$$

$$\sigma(1-\sigma)$$

$$z_n^{(l)}$$

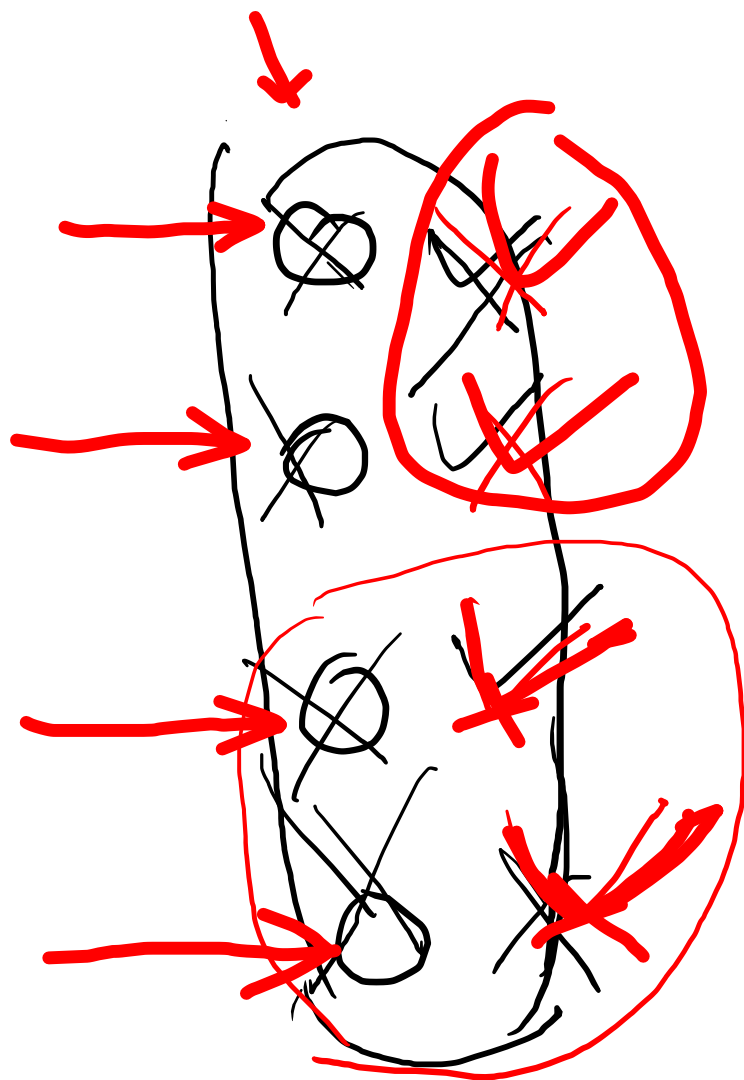
故:  $\delta^{(l)} = \delta^{(l+1)} \cdot w^{(l+1)} \cdot \sigma'(\cdot)$

$$\delta^{(l)} = \delta^{(l+1)} \cdot w^{(l+1)} \cdot \sigma'(\cdot)$$



1. ~~the~~ dropout

2. BN : Batch - Normalization



0

0

0

0

0.5  
—

x 0.5

↓

归一化 × 手动

每个信息的归一化

$y = y_1 + y_2 \Rightarrow \dot{y} = \dot{y}_1 + \dot{y}_2$   
 $y' = y_1' + y_2'$

$$\frac{y_1 = Wx}{y_2 = W^2 x}$$

1.0

$$\{x_1, x_2, \dots, x_j, \dots, x_{n_{size}}\}$$

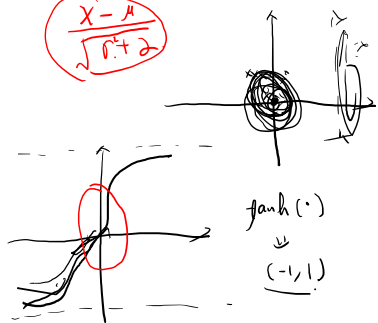

---

$\mu, \sigma^2$

$$\mu = \frac{1}{N} \sum_{j=1}^N x_j$$

$$\sigma^2 = \frac{1}{N} \sum_{j=1}^N (x_j - \mu)^2$$

$$\frac{x - \mu}{\sqrt{\sigma^2 + 2}}$$



1. 梯度 (缓斜)
2. 流速 速度 快
3. 初始方式  $\uparrow$
4. peening - vete 可猜不

$y = y_{in}$   
 $\text{NAN} \rightarrow \text{er} \downarrow \checkmark$

$$\frac{x-\mu}{\sqrt{\sigma^2 + 2}}$$

$$y_i = a \cdot \frac{x_i - \mu}{\sqrt{\sigma^2 + 2}} + b$$

Diagram illustrating the components of the equation:
 

- $a$  is labeled  $l_3$ .
- $x_i - \mu$  is labeled  $l_2$ .
- $\sqrt{\sigma^2 + 2}$  is labeled  $l_1$ .
- $b$  is labeled  $l_4$ .
- The entire expression is labeled  $meq_1$ .
- The expression is also labeled triviale.

RUN, L5 tri X

80% 反效果

## Layer - Norm

Sequence  $\mu, \sigma^2$

















