

一. 对 softmax 求导

设  $z_j$  是输出

$$y_j = \text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$

$$\text{Loss}(\hat{y}, y) = -\sum y_j \log(\hat{y}_j)$$

设正确类别为  $t$ ,  $y_t = 1$ ,  $y_{j \neq t} = 0$   $\text{Loss}(\hat{y}, y) = -\log(\hat{y}_t)$   $\frac{\partial L}{\partial \hat{y}_t} = -\frac{1}{\hat{y}_t}$   $j=t$  时

$$\frac{\partial \hat{y}_t}{\partial z_t} = \frac{\partial \left( \frac{e^{z_t}}{\sum_k e^{z_k}} \right)}{\partial z_t} = \frac{e^{z_t} \cdot \sum_k e^{z_k} - e^{z_t} \cdot e^{z_t}}{\sum_k e^{z_k}^2} = \frac{e^{z_t}}{\sum_k e^{z_k}} - \left( \frac{e^{z_t}}{\sum_k e^{z_k}} \right)^2 = \hat{y}_t (1 - \hat{y}_t)$$

$j \neq t$  时

$$\frac{\partial \hat{y}_t}{\partial z_t} = -\hat{y}_t \hat{y}_j$$

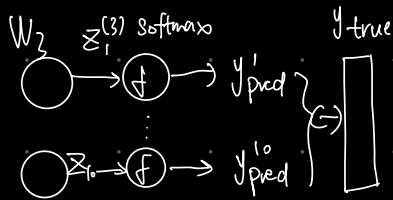
则气链式:  $\frac{\partial L}{\partial z_j} = \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_t} = -\frac{1}{\hat{y}_t} \cdot \hat{y}_t (1 - \hat{y}_t) = \hat{y}_t - 1 = \hat{y}_t - y_t$

$j=t$

$j \neq t$   $\frac{\partial L}{\partial z_j} = \frac{\partial L}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial z_j} = -\frac{1}{\hat{y}_t} (-\hat{y}_t \hat{y}_j) = \hat{y}_j = \hat{y}_j - y_j = \hat{y}_j - 0$

$$\frac{\partial L}{\partial z_j} = \hat{y}_j - y_j$$

$\text{grad-} y_{\text{pred}} = y_{\text{pred}} - y_{\text{true}}$   
 $\in \mathbb{R}^{m \times 10}$



$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial w_3^{(1)}}$$

$z^{(3)} = w_3 a^{(2)} + b_3$   
 $64 \times 10 \quad m \times 64$

$\frac{\partial z^{(1)}}{\partial w_3} = a^{(2)} \in \mathbb{R}^{m \times 64}$

$\frac{\partial L}{\partial w_3} \in (64 \times m) \cdot (m \times 10)$   
 $= \mathbb{R}^{64 \times 10}$

$\therefore \frac{\partial L}{\partial w_3} = \frac{1}{m} \cdot \text{np.dot}(a_2^T, \text{grad-} y_{\text{true}})$

$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial b_3} = \frac{1}{m} \cdot \text{np.sum}(\text{grad-} y_{\text{pred}})$   
 $\mathbb{R}^{m \times 10}, \text{axis}=0 \Rightarrow \mathbb{R}^{10}$

$\mathbb{R}^{m \times 64} \Leftarrow \frac{\partial L}{\partial a^{(2)}} = \frac{\partial L}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial a^{(2)}} = \frac{1}{m} \cdot \text{np.dot}(\text{grad-} y_{\text{true}}, w_3^T) \in \mathbb{R}^{m \times 64}$   
 $\text{grad-} y_{\text{true}} \in \mathbb{R}^{m \times 10}$   
 $w_3 \in \mathbb{R}^{64 \times 10}$   
 矩阵乘法其实就是 sum 的过程

$$\frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial a^{(1)}} \cdot \underbrace{\frac{\partial a^{(2)}}{\partial z^{(2)}}}_{\text{grad}}$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial W_2} \rightarrow a^{(1)} \in \mathbb{R}^{m \times 128}, W_2 \in \mathbb{R}^{128 \times 64}$$

$$\text{grad-}z_2 \in \mathbb{R}^{m \times 64} \quad \text{---} \quad \frac{1}{m} \text{np.dot}(a_1^T, \text{grad-}z_2)$$

$$\text{总归: } \frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W_3} = \text{grad-}y_{\text{pred}} \cdot a^{(2)} = \frac{1}{m} \text{np.dot}(a_2^T, \text{grad-}y_{\text{pred}})$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial W_2} = \text{grad-}z_2 \cdot a^{(1)} = \frac{1}{m} (a_1^T, \text{grad-}z_2)$$

$$\text{第} l \text{层: } \frac{\partial L}{\partial W_l} = \frac{\partial L}{\partial z^{(l)}} \cdot \frac{\partial z^{(l)}}{\partial W_l} = \text{grad-}z_l \cdot a^{(l-1)}$$

$$\text{第} l \text{层: } \text{grad-}z_l \text{ 用 } \frac{\partial L}{\partial a^{(l)}}, \text{ 即 } \frac{\partial L}{\partial a^{(l)}} = \frac{\partial L}{\partial z^{(l+1)}} \cdot \frac{\partial z^{(l+1)}}{\partial a^{(l)}}$$

$$= \text{grad-}z_{l+1} \cdot W_{l+1}^{(l+1)} = \text{np.dot}(\text{grad-}z_{l+1}, W_{l+1}^T)$$

$\downarrow$   $\mathbb{R}^{m, d_{l+1}}$        $\downarrow$   $\mathbb{R}^{d_l, d_{l+1}}$

二. 为什么矩阵的  $\frac{\partial L}{\partial W_{ij}}$  要乘  $1/m$ ?

在 SGD 中, 一个 batch 有  $m$  个样本, 乘  $1/m$  是这一批样本的梯度均值。

以  $W_3 \in \mathbb{R}^{64 \times 10}$  为例。

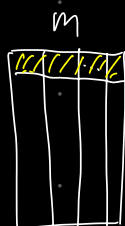
$$\frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W_3} = \text{grad-}y_{\text{pred}} \cdot a^{(2)}$$

$\downarrow$   $\mathbb{R}^{m \times 10}$        $\downarrow$   $\mathbb{R}^{m \times 64}$

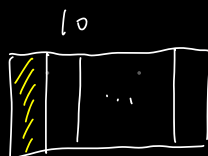
$$\mathbb{R}^{64 \times 10}$$

$$a_2^T$$

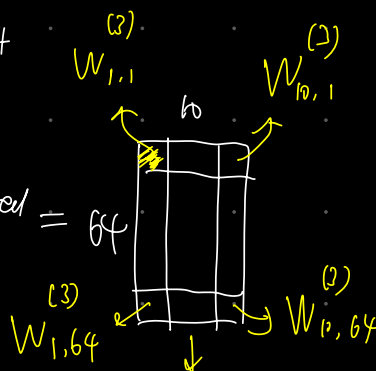
$$64$$



$$x m$$



$$\text{grad-}y_{\text{pred}} = 64$$



为什么  $\frac{\partial L}{\partial a^{(i)}}$  不用  $\times \frac{1}{m}$ ?

这个矩阵中每个元素是  $m$   
样本梯度的累加,  $\therefore \times \frac{1}{m}$

$$\frac{\partial L}{\partial a^{(i)}} = \frac{\partial L}{\partial z^{(i+1)}} \cdot \frac{\partial z^{(i+1)}}{\partial a^{(i)}} = \text{grad-}z^{(i+1)} \cdot W^{(i+1)}$$

$\downarrow$                        $\downarrow$   
 $z^{(i)} \in \mathbb{R}^{m \times 10}$      $W^{(i)} \in \mathbb{R}^{64 \times 10}$

以  $a^{(i)}$  为例,  $t \in \mathbb{R}$

$$z^{(i)} \cdot W^{(i)T} \Rightarrow \text{mp.dot}(z^{(i)}, W^{(i)T})$$

$$m \times 10, 10 \times 64 \Rightarrow m \times 64$$

$m$  的维度本来就要  
留下来, 不用  $\times \frac{1}{m}$