

Network structure :

$$a = \text{input} \cdot W_1 + b_1$$

$$b = \text{ReLU}(a)$$

$$c = b \cdot W_2 + b_2$$

$$d = \text{softmax}(c)$$

$$\text{Loss} = -\frac{1}{N} \sum \log(d_k) + \frac{\alpha}{2} (\|W_1\|^2 + \|W_2\|^2)$$

shape :

$$\text{input.shape} : (N_{\text{batch}}, N_{\text{in}})$$

$$W_1.\text{shape} : (N_{\text{in}}, N_{\text{hid}})$$

$$a.\text{shape} : (N_{\text{batch}}, N_{\text{hid}})$$

$$W_2.\text{shape} : (N_{\text{hid}}, N_{\text{out}})$$

$$c.\text{shape} : (N_{\text{batch}}, N_{\text{out}})$$

$$d.\text{shape} : (N_{\text{batch}}, N_{\text{out}})$$

$$s_k.\text{shape} : (N_{\text{batch}}, N_{\text{out}}) \rightarrow \text{每一行仅第 } k \text{ 个} = 1, \text{ ground truth}$$

对应  $d$  中的  $d_k$  ( $d_k$  为一个数值)

Hint:

$$Y = XW, \quad \frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \frac{\partial Y}{\partial X} = \frac{\partial L}{\partial Y} W^T$$

$$Y = XW, \quad \frac{\partial L}{\partial W} = X^T \frac{\partial L}{\partial Y}$$

Deduction:

$$\frac{\partial \text{Loss}}{\partial w_1} = \frac{1}{N} (\text{input}^T * \frac{\partial \text{Loss}}{\partial a}) + \alpha \|w_1\|$$

$$= \frac{1}{N} \left[ \text{input}^T * \left( \frac{\partial \text{Loss}}{\partial c} * \frac{\partial c}{\partial b} \odot \frac{\partial b}{\partial a} \right) \right] + \alpha \|w_1\|$$

$$= \frac{1}{N} \left[ \text{input}^T * \left( (d - \delta_k) * W_2^T \odot I(a > 0) \right) \right] + \alpha \|w_1\|$$

已经应用了 softmax 梯度, 注意到第一层对  $w_1$  求导, ReLU 的导数  
是大于 0 而不是大于等于 0

$$\frac{\partial \text{Loss}}{\partial b_1} = \frac{\partial \text{Loss}}{\partial a} = \frac{1}{N} \left( (d - \delta_k) * W_2^T \odot I(a > 0) \right)$$

$$\frac{\partial \text{Loss}}{\partial w_2} = \frac{1}{N} (b^T * (d - \delta_k)) + \alpha \|w_2\|$$

$$\frac{\partial \text{Loss}}{\partial b_2} = \frac{\partial \text{Loss}}{\partial c} = \frac{1}{N} (d - \delta_k)$$