

## 1 Linear Module

### 1.a

$$\left[ \frac{\partial L}{\partial \mathbf{W}} \right]_{ij} = \frac{\partial L}{\partial W_{ij}} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial W_{ij}} \quad (1)$$

$$\frac{\partial Y_{sn}}{\partial W_{ij}} = \frac{\partial}{\partial W_{ij}} \left( \sum_m [\mathbf{X}]_{sm} [\mathbf{W}^\top]_{mn} + [\mathbf{B}]_{sn} \right) = \sum_m X_{sm} \frac{\partial W_{nm}}{\partial W_{ij}} + \frac{\partial B_{sn}}{\partial W_{ij}} \quad (2)$$

$$= \sum_m X_{sm} \delta_{ni} \delta_{mj} + 0 = X_{sj} \delta_{ni} \quad (3)$$

$$\sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial W_{ij}} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} X_{sj} \delta_{ni} = \sum_s \frac{\partial L}{\partial Y_{si}} X_{sj} \quad (4)$$

$$\therefore \frac{\partial L}{\partial \mathbf{W}} = \left( \frac{\partial L}{\partial \mathbf{Y}} \right)^\top \mathbf{X} \in \mathbb{R}^{N \times M} \quad (5)$$

### 1.b

$$\left[ \frac{\partial L}{\partial \mathbf{b}} \right]_j = \frac{\partial L}{\partial b_j} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial b_j} \quad (6)$$

$$\frac{\partial Y_{sn}}{\partial b_j} = \frac{\partial}{\partial b_j} \left( \sum_m [\mathbf{X}]_{sm} [\mathbf{W}^\top]_{mn} + [\mathbf{B}]_{sn} \right) = \sum_m \frac{\partial X_{sm} W_{nm}}{\partial b_j} + \frac{\partial B_{sn}}{\partial b_j} \quad (7)$$

$$= 0 + \delta_{nj} = \delta_{nj} \quad (8)$$

$$\sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial b_j} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \delta_{nj} = \sum_s \frac{\partial L}{\partial Y_{sj}} \quad (9)$$

$$\therefore \frac{\partial L}{\partial \mathbf{b}} = \sum_s \frac{\partial L}{\partial \mathbf{Y}_s} \in \mathbb{R}^{1 \times N} \quad (10)$$

For clarification, in equation (10) we sum over the rows  $s \in S$  of  $\mathbf{Y}$ , so the  $j$ -th element of  $\mathbf{b}$  is the sum of all elements in position  $j$  of the rows of  $\mathbf{Y}$ .

### 1.c

$$\left[ \frac{\partial L}{\partial \mathbf{X}} \right]_{ij} = \frac{\partial L}{\partial X_{ij}} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial X_{ij}} \quad (11)$$

$$\frac{\partial Y_{sn}}{\partial X_{ij}} = \frac{\partial}{\partial X_{ij}} \left( \sum_m [\mathbf{X}]_{sm} [\mathbf{W}^\top]_{mn} + [\mathbf{B}]_{sn} \right) = \sum_m \frac{\partial X_{sm} W_{nm}}{\partial X_{ij}} + \frac{\partial B_{sn}}{\partial X_{ij}} \quad (12)$$

$$= \sum_m \delta_{si} \delta_{mj} W_{nm} + 0 = \delta_{si} W_{nj} \quad (13)$$

$$\sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \frac{\partial Y_{sn}}{\partial X_{ij}} = \sum_{s,n} \frac{\partial L}{\partial Y_{sn}} \delta_{si} W_{nj} = \sum_n \frac{\partial L}{\partial Y_{in}} W_{nj} \quad (14)$$

$$\therefore \frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \mathbf{Y}} \mathbf{W} \in \mathbb{R}^{S \times M} \quad (15)$$

## 1.d

$$\left[ \frac{\partial L}{\partial \mathbf{X}} \right]_{ij} = \frac{\partial L}{\partial X_{ij}} = \sum_{s,m} \frac{\partial L}{\partial Y_{sm}} \frac{\partial Y_{sm}}{\partial X_{ij}} \quad (16)$$

$$= \sum_{s,m} \frac{\partial L}{\partial Y_{sm}} \frac{\partial h(X_{sm})}{\partial X_{ij}} \quad (17)$$

$$= \sum_{s,m} \frac{\partial L}{\partial Y_{sm}} \delta_{si} \delta_{mj} = \frac{\partial L}{\partial Y_{ij}} \quad (18)$$

$$\therefore \frac{\partial L}{\partial \mathbf{X}} = \frac{\partial L}{\partial \mathbf{Y}} \circ h'(\mathbf{X}) \in \mathbb{R}^{S \times M} \quad (19)$$

From equation (17) to (18) we get two Kronecker deltas because the derivative of the element-wise activation function is zero for all elements except the one we are taking the derivative of. This makes intuitive sense, given we are differentiating an element-wise function. The final derivative of the loss w.r.t. the input  $\mathbf{X}$  is the Hadamard product of the derivative of the loss w.r.t. the output  $\mathbf{Y}$  and the element-wise derivative of the activation function  $h$  w.r.t. the input  $\mathbf{X}$ . We can assume the shapes of  $\mathbf{X}$  and  $\mathbf{Y}$  are compatible, since  $\mathbf{Y} = h(\mathbf{X})$