

1. Derivative of Sigmoid

$$S(x) = \frac{1}{1 + e^x}$$

$$\frac{d}{dx} S(x) = \frac{1}{(1 + e^x)^2} (e^{-x}) \quad (\text{chain rule})$$

$$= \frac{1 + e^{-x} - 1}{(1 + e^x)^2} = \frac{1 + e^{-x}}{(1 + e^x)^2} - \frac{1}{(1 + e^x)^2}$$

$$= \frac{1}{1 + e^x} - \left(\frac{1}{1 + e^x} \right)^2$$

$$= S(x) - S(x)^2$$

$$= \underline{\underline{S(x) (1 - S(x))}}$$

→ The properties of the sigmoid function's derivative is convenient for neural networks, because the gradients for the layer can be calculated using simple multiplication and subtraction rather than performing the exponentials of sigmoid function.

2.1 derivative of Softmax.

$$y_c = \frac{e^{z_c}}{\sum_{d=1}^C e^{z_d}}$$

(Quotient rule)

$$\text{if } i = j : \frac{\partial y_i}{\partial z_i} = \frac{\partial e^{z_i} / z_c}{\partial z_i} = \frac{e^{z_i} z_c - e^{z_i} e^{z_i}}{z_c^2}$$

$$= \frac{e^{z_i}}{z_c} \cdot \frac{z_c - e^{z_i}}{z_c} = y_i (1 - y_i)$$

$$\text{if } \boxed{i \neq j} \quad \frac{\partial y_i}{\partial z_j} = \frac{\partial e^{z_i} / \zeta_c}{\partial z_j} = \frac{0 - e^{z_i} e^{z_j}}{\zeta_c^2}$$

$$= -\frac{e^{z_i}}{\zeta_c} \cdot \frac{e^{z_j}}{\zeta_c} = \underline{\underline{-y_i y_j}}$$

2.2 Cross Entropy loss derivative

$$L = - \sum_i t_i \log(y_i)$$

$$\frac{\partial L}{\partial z_i} = - \sum_{j=1}^C \frac{\partial t_j \log(y_j)}{\partial z_i} = - \sum_{j=1}^C t_j \frac{\partial}{\partial z_i} \log(y_j)$$

$$= - \sum_{j=1}^C \frac{t_j}{y_j} \frac{\partial y_j}{\partial z_i} \quad (\text{chain rule})$$

$$= -\frac{t_i}{y_i} \frac{\partial y_i}{\partial z_i} - \sum_{j \neq i}^C \frac{t_j}{y_j} \frac{\partial y_j}{\partial z_i}$$

$$= -\frac{t_i}{y_i} \log(y_i (1 - y_i)) - \sum_{j \neq i}^C \frac{t_j}{y_j} (-y_j y_i)$$

$$= -t_i + t_i y_i + \sum_{j \neq i}^C t_j y_i$$

$$= -t_i + \sum_{j=1}^C t_j y_i$$

$$= -t_i + y_i \sum_{j=1}^C t_j$$

$$= \underline{\underline{y_i - t_i}}$$