

$$\hat{y}_i = \text{softmax}(z) = \frac{e^{z_i}}{\sum e^{z_i}} \quad [\text{logits of } c \text{ classes}]$$

$$z = W_L^T h + b_L$$

$$h = g(W_1^T x + b_1)$$

Loss function

$$L(y, \hat{y}) = - \sum_{i=1}^c y_i \cdot \log(\hat{y}_i)$$

$y \rightarrow$ one hot encoding

$y_i \rightarrow$ predicted prob

compute derivative

$$\frac{dL}{d\hat{y}_i} = - \frac{y_i}{\hat{y}_i}$$

$$\frac{dL}{dz_i} = \sum_{j=1}^c \frac{\partial L}{\partial \hat{y}_j} \cdot \frac{\partial \hat{y}_j}{\partial z_i}$$

we will derivate $\hat{y}_j = \frac{e^{z_j}}{\sum_k e^{z_k}}$

$$\frac{\partial \hat{y}_j}{\partial z_i} = \begin{cases} \hat{y}_j(1 - \hat{y}_j) & i=j \\ -\hat{y}_i \hat{y}_j & i \neq j \end{cases}$$

combine $\frac{dL}{d\hat{y}_j}$, $\frac{\partial \hat{y}_j}{\partial z_i}$ $\left[\frac{dL}{d\hat{y}_j} = - \frac{y_j}{\hat{y}_j} \right]$

$$i=j \rightarrow \frac{dL}{dz_i} = \hat{y}_i(1 - \hat{y}_i) \cdot \left(- \frac{y_i}{\hat{y}_i} \right) = \hat{y}_i - y_i$$

$$i \neq j \rightarrow \frac{dL}{dz_i} = -\hat{y}_i y_j \cdot \left(- \frac{y_j}{\hat{y}_j} \right) = 0$$