

Categorical cross-entropy loss

$$\mathcal{L}_{\text{CCE}} = - \sum_{i=1}^k y_i \log(\hat{y}_i)$$

where \hat{y}_i is obtained using
SOFTMAX function

$$\hat{y}_i = \frac{e^{x_i}}{\sum_{j=1}^k e^{x_j}}, \quad \hat{y}_1 + \hat{y}_2 + \dots + \hat{y}_k = 1$$

Compute the derivative of the
loss function.

$$x_j = [x_1, x_2, x_3, \dots, x_k]$$

$$\hat{y}_i = [y_1, y_2, y_3, \dots, y_k]$$

$$\hat{y}_1 = \frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}}$$

$$\hat{y}_2 = \frac{e^{x_2}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}}$$

⋮

$$\hat{y}_k = \frac{e^{x_k}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}}$$

$$\Rightarrow \frac{\partial \hat{y}_i}{\partial x_j} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial x_1} & \frac{\partial \hat{y}_1}{\partial x_2} & \dots & \frac{\partial \hat{y}_1}{\partial x_k} \\ \frac{\partial \hat{y}_2}{\partial x_1} & \frac{\partial \hat{y}_2}{\partial x_2} & \dots & \frac{\partial \hat{y}_2}{\partial x_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_k}{\partial x_1} & \frac{\partial \hat{y}_k}{\partial x_2} & \dots & \frac{\partial \hat{y}_k}{\partial x_k} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial \hat{y}_1}{\partial x_1} & \frac{\partial \hat{y}_1}{\partial x_2} & \dots & \frac{\partial \hat{y}_1}{\partial x_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_k}{\partial x_1} & \frac{\partial \hat{y}_k}{\partial x_2} & \dots & \frac{\partial \hat{y}_k}{\partial x_k} \end{bmatrix}$$

$$\textcircled{1} \quad h(x) = \frac{f(x)}{g(x)}$$

$$\frac{\partial h(x)}{\partial x} = \frac{f'(x)g(x) - g'(x)f(x)}{g(x)^2}$$

$$\textcircled{2} \quad \frac{\partial (f(x) + h(x))}{\partial x} = \frac{\partial f(x)}{\partial x} + \frac{\partial h(x)}{\partial x}$$

$$\begin{aligned} \frac{\partial \hat{y}_1}{\partial x_1} &= \frac{\partial \left(\frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \right)}{\partial x_1} \\ &= \frac{e^{x_1} (e^{x_1} + e^{x_2} + \dots + e^{x_k}) - e^{x_1} \times e^{x_1}}{(e^{x_1} + e^{x_2} + \dots + e^{x_k})^2} \\ &= \frac{e^{x_1} \times (\Sigma - e^{x_1})}{\Sigma \times \Sigma} \end{aligned}$$

$$\boxed{\frac{\partial \hat{y}_1}{\partial x_1} = \hat{y}_1 (1 - \hat{y}_1)}$$

$$\boxed{\frac{\partial y_1}{\partial u_1} = y_1 (1 - y_1)}$$

$$\|y\| \quad \frac{\partial \hat{y}_2}{\partial u_2} = \hat{y}_2 (1 - \hat{y}_2)$$

$$\frac{\partial \hat{y}_k}{\partial u_k} = \hat{y}_k (1 - \hat{y}_k)$$

$$\frac{\partial \hat{y}_2}{\partial u_1} = \frac{\partial \left(\frac{e^{u_2}}{e^{u_1} + e^{u_2} + \dots + e^{u_k}} \right)}{\partial u_1}$$

$$= \frac{0 - e^{u_1} e^{u_2}}{\sum x \sum}$$

$$\boxed{\frac{\partial \hat{y}_2}{\partial u_1} = -\hat{y}_1 \times \hat{y}_2}$$

$$\boxed{\frac{\partial \hat{y}_k}{\partial u_1} = -\hat{y}_1 \times \hat{y}_k}$$

$$\frac{\partial \hat{y}_i}{\partial x_j} = \begin{bmatrix} \hat{y}_1(1-\hat{y}_1) & -\hat{y}_1 \cdot \hat{y}_2 & \dots & -\hat{y}_1 \cdot \hat{y}_k \\ -\hat{y}_2 \cdot \hat{y}_1 & \hat{y}_2(1-\hat{y}_2) & \dots & -\hat{y}_2 \cdot \hat{y}_k \\ \vdots & \vdots & \ddots & \vdots \\ -\hat{y}_k \cdot \hat{y}_1 & -\hat{y}_k \cdot \hat{y}_2 & \dots & \hat{y}_k(1-\hat{y}_k) \end{bmatrix}$$

⇒ In general,

$$\frac{\partial \hat{y}_i}{\partial x_j} = \begin{cases} y_i(1-y_i) & \forall i=j \\ -y_i \times y_j & \forall i \neq j \end{cases}$$

$$\begin{aligned} \frac{\partial \mathcal{L}_{\text{CCE}}}{\partial x_j} &= \frac{\partial}{\partial x_j} \left(-\sum_{i=1}^k y_i \log(\hat{y}_i) \right) \\ &= -\sum_{i=1}^k y_i \left(\frac{\partial (\log \hat{y}_i)}{\partial x_j} \right) \\ &= -\sum_{i=1}^k y_i \times \frac{1}{\hat{y}_i} \left(\frac{\partial \hat{y}_i}{\partial x_j} \right) \\ &= -\sum_{i=1}^k \frac{y_i}{\hat{y}_i} \times \left(\hat{y}_i (1 \{i=j\} - \hat{y}_j) \right) \\ &\dots \rightarrow \sum_{i=1}^k u_i (1 \{i=j\} - \hat{y}_i) \end{aligned}$$

$$= - \sum_{i=1}^k y_i (\mathbb{1}_{\{i=j\}} - \hat{y}_j)$$

$$= \sum_{i=1}^k y_i \cdot \hat{y}_j - \sum_{i=1}^k y_i \cdot \mathbb{1}_{\{i=j\}}$$

$$= \sum_{i=1}^k y_i \cdot \hat{y}_j - y_j$$

$$= \hat{y}_j \times \sum_{i=1}^k y_i - y_j$$

$$\boxed{\frac{\partial L_{\text{CE}}}{\partial x_j} = \hat{y}_j - y_j}$$