

$$\text{Softmax } y_k = P(l=k, \vec{x})$$

$$= \frac{e^{z_k}}{\sum_j e^{z_j}}$$

Note $y_k = P(l=k, z_1, \dots, z_K)$

Thus, y_k depends on z_1, \dots, z_K .

$$\frac{\partial y_k}{\partial z_j} = \frac{\partial}{\partial z_j} \left(\frac{e^{z_k}}{e^{z_1} + \dots + e^{z_K}} \right)$$

If $\boxed{j=k}$...

Two indices, so two cases ✓

$$\frac{\partial y_k}{\partial z_k} = \frac{e^{z_k} z - (e^{z_k}) e^{z_k}}{(z)^2}$$

$$= \frac{e^{z_k}}{z} - \left(\frac{e^{z_k}}{z} \right)^2 = y_k (1 - y_k)$$

$\leftarrow z = \sum_j e^{z_j}$

If $\boxed{j \neq k}$

$$\frac{\partial y_k}{\partial z_j} = \frac{0 - (e^{z_k})(e^{z_j})}{(z)^2} = -y_k y_j$$

Putting those together

$$\boxed{\frac{\partial y_k}{\partial z_j} = y_k (\delta_{jk} - y_j)}$$

✓
 δ_{jk} is the Kronecker delta
 $\delta_{jk} = \begin{cases} 1 & \text{if } j=k \\ 0 & \text{if } j \neq k \end{cases}$

For cross-entropy, $E(\vec{y}, \vec{t}) = - \sum_{k=1}^K t_k \ln y_k$

$$\frac{\partial E}{\partial z_j} = \sum_{k=1}^K \frac{\partial E}{\partial y_k} \frac{\partial y_k}{\partial z_j}$$

✓ $\frac{\partial E}{\partial y_k} = \frac{-t_k}{y_k}$

$$\therefore \frac{\partial E}{\partial z_j} = \sum_{k=1}^K \frac{-t_k}{y_k} y_k (\delta_{jk} - y_j)$$

$$= - \sum_{k=1}^K t_k (\delta_{jk} - y_j)$$

Simplify ✓

$$= - (t_j - y_j)$$

$$= y_j - t_j$$

since $t_k = 1$ only for the correct class, which we call j here.

Or, stated as a vector

$$\frac{\partial E}{\partial \vec{z}} = \vec{y} - \vec{t}$$

Alternative Solution

Softmax Gradient

We substitute the softmax activation function for the output node

$$y_k = \frac{e^{z_k}}{\sum_{i=1}^K e^{z_i}}$$

inside the loss function

$$E(\vec{y}, \vec{t}) = - \sum_k^K t_k \ln(y_k)$$

to get

$$E(\vec{y}, \vec{t}) = - \sum_k^K t_k \ln \left(\frac{e^{z_k}}{\sum_{i=1}^K e^{z_i}} \right)$$

and evaluate the gradient with respect to the input current to the output layer

$$\begin{aligned} \frac{\partial E}{\partial z_j} &= \frac{\partial}{\partial z_j} \left[- \sum_k^K t_k \ln \left(\frac{e^{z_k}}{\sum_{i=1}^K e^{z_i}} \right) \right] \\ &= \frac{\partial}{\partial z_j} \left[- \sum_k^K t_k \left(z_k - \ln \left(\sum_{i=1}^K e^{z_i} \right) \right) \right] \\ &= \frac{\partial}{\partial z_j} \left[\sum_k^K t_k \left(\ln \left(\sum_{i=1}^K e^{z_i} \right) - z_k \right) \right] \\ &= \frac{\partial}{\partial z_j} \left[\sum_k^K t_k \ln \left(\sum_{i=1}^K e^{z_i} \right) - t_k z_k \right] \\ &= \frac{\partial}{\partial z_j} \left[\sum_k^K t_k \ln \left(\sum_{i=1}^K e^{z_i} \right) \right] - t_j \\ &= \sum_k^K \left(\frac{t_k e^{z_j}}{\sum_{i=1}^K e^{z_i}} \right) - t_j \end{aligned}$$

$$= \sum_k t_k y_j - t_j$$

$$= y_j \left(\sum_k t_k \right) - t_j$$

$$= y_j - t_j$$

$$= \begin{cases} y_j & \text{if } j \neq \gamma \\ y_j - 1 & \text{if } j = \gamma \end{cases}$$

As a vector: $\frac{\partial E}{\partial \vec{z}} = \vec{y} - \vec{t}$.