Derive Gradient ~~Descent~~ For | → softmax → cross entropy →|

1. $Z →$ | softmax function | $→$ softmax_out

$$\begin{bmatrix} z_0 \\ z_1 \\ z_2 \\ \vdots \\ z_k \end{bmatrix} \rightarrow \text{softmax\_out}_i = \frac{e^{z_i}}{\sum\limits_{j=0}^{k} e^{z_j}} \rightarrow \begin{bmatrix} \text{softmax\_out}_0 \\ \text{softmax\_out}_1 \\ \text{softmax\_out}_2 \\ \vdots \\ \text{softmax\_out}_k \end{bmatrix}$$
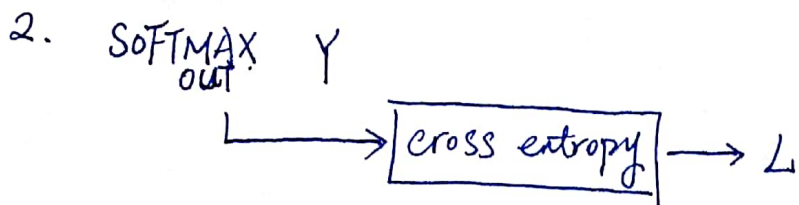
$$\frac{\partial \text{softmax\_out}_n}{\partial z_m} = \frac{\partial}{\partial z_m} e^{z_n} \left( e^{z_0} + e^{z_1} + \cdots + e^{z_k} \right)^{-1}$$

Let
$e^{z_0} + e^{z_1} + \cdots + e^{z_k} = S$

$$= \begin{cases} e^{z_m}(S)^{-1} + e^{z_m} \cdot (-1)(S)^{-2} e^{z_m} & , \text{ when } m=n \\ (-1) e^{z_n}(S)^{-2} e^{z_m} & , \text{ when } m \neq n \end{cases}$$

$$= \begin{cases} e^{z_m}(S)^{-1} \left( 1 - e^{z_m}(S)^{-1} \right) & , \text{ when } m=n \\ -\left( e^{z_n} S^{-1} \right)\left( e^{z_m} S^{-1} \right) & \text{ when } m \neq n \end{cases}$$

$$= \begin{cases} \text{softmax\_out}_m \left( 1 - \text{softmax\_out}_m \right) & \text{ when } m=n \\ -(\text{softmax\_out}_n)(\text{softmax\_out}_m) & \text{ when } m \neq n \end{cases}$$

2. SOFTMAX Y
   out

$$\longrightarrow \boxed{\text{cross entropy}} \longrightarrow L$$

$$\begin{bmatrix} \text{softmax\_out }_0 \\ \text{softmax\_out}_1 \\ \vdots \\ \text{softmax\_out}_k \end{bmatrix} \begin{bmatrix} y_0 \\ y_1 \\ \vdots \\ y_k \end{bmatrix}$$

$$\longrightarrow \boxed{- \sum_{j=0}^{k} y_i \log(\text{softmax\_out}_i)} \longrightarrow L$$

$$\frac{\partial L}{\partial \, \text{softmax\_out}_m} = -(y_m)\left(\frac{1}{\text{softmax\_out}_m}\right)$$

3. Combine 1. and 2.

$$\frac{\partial L}{\partial z_m} = \frac{\partial L}{\partial \, \text{SOFTMAX\_OUT}} * \frac{\partial \, \text{SOFTMAX\_OUT}}{\partial z_m}$$

$$= \sum_{n=0}^{k} \frac{\partial L}{\partial \, \text{softmax\_out}_n} * \frac{\partial \, \text{softmax\_out}_n}{\partial z_m}$$

assume $y_j = 1$ for $j = \ell$ else $y_j = 0$ recall that Y is a one hot vector

$$= - y_\ell \left(\frac{1}{\text{softmax\_out}_\ell}\right) \times \frac{\partial \, \text{softmax\_out}_\ell}{\partial z_m}$$

$$= \begin{cases} - y_m \left(\frac{1}{\text{softmax\_out}_m}\right)(\text{softmax\_out}_m)(1 - \text{softmax\_out}_m) & \text{when } \ell = m \\ + y_\ell \left(\frac{1}{\text{softmax\_out}_\ell}\right)(\text{softmax\_out}_\ell)(\text{softmax\_out}_m) & \text{when } \ell \neq m \end{cases}$$

$$= \text{softmax out}_m - y_m$$