



$$\frac{\partial l}{\partial w} = ? \quad \rightarrow \quad \frac{\partial l}{\partial w} = \frac{\partial l}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial x} \cdot \frac{\partial x}{\partial w}$$

* $\frac{\partial l}{\partial \sigma} \rightarrow$ Derivative of l w.r.t softmax output

$$\frac{\partial l}{\partial b} = ? \quad \rightarrow \quad \frac{\partial l}{\partial b} = \frac{\partial l}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial x} \cdot \frac{\partial x}{\partial b}$$

* $\frac{\partial \sigma}{\partial x} \rightarrow$ Derivative of softmax w.r.t its inputs (x_i)

* Let's start with computing $\frac{\partial l}{\partial \sigma}$

$$\frac{\partial l}{\partial \sigma} = \left[\frac{\partial l}{\partial \sigma_1}, \dots, \frac{\partial l}{\partial \sigma_N} \right]$$

$$= \left[\frac{-t_1}{\sigma_1}, \dots, \frac{-t_N}{\sigma_N} \right]$$

* Then compute $\frac{\partial \sigma}{\partial x}$

$$\frac{\partial \sigma}{\partial x} = \begin{bmatrix} \sigma_1(1-\sigma_1(x)) & \dots & -\sigma_N(x)\sigma_1(x) \\ \vdots & \ddots & \vdots \\ -\sigma_1(x)\sigma_N(x) & \dots & -\sigma_N(x)(1-\sigma_N(x)) \end{bmatrix}$$

$$* \frac{\partial \ell}{\partial w} = \frac{\partial \ell}{\partial \sigma} \frac{\partial \sigma}{\partial x} \left(\frac{\partial x}{\partial w} \right) \rightarrow \text{Input}$$

$$* \frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial \sigma} \frac{\partial \sigma}{\partial x} \left(\frac{\partial x}{\partial b} \right) \rightarrow 1$$

$$\hookrightarrow \frac{\partial \ell}{\partial x} ??$$

$$* \frac{\partial \ell}{\partial x} = \frac{\partial \ell}{\partial \sigma} \cdot \frac{\partial \sigma}{\partial x} = \begin{bmatrix} \frac{-t_1}{\sigma_1} & \dots & \frac{-t_N}{\sigma_N} \end{bmatrix} \begin{bmatrix} \sigma_1 (1 - \sigma_1(x)) & \dots & -\sigma_N(x) \sigma_1(x) \\ \vdots & \ddots & \vdots \\ -\sigma_1(x) \sigma_N(x) & \dots & -\sigma_N(x) (1 - \sigma_N(x)) \end{bmatrix}$$

$1 \times N \quad N \times N$

$$= \sigma_i(x) (u_i - \sigma(x)^T)$$

* Let's assume that we have two linear layers and (w_1, b_1) is parameters of Linear 1 and (w_2, b_2) is parameters of Linear 2. (So Input \rightarrow Linear 1 \rightarrow Linear 2 \rightarrow Softmax \rightarrow Loss)

$$\boxed{\frac{\partial \ell}{\partial w_2} = \frac{\partial \ell}{\partial \sigma} \frac{\partial \sigma}{\partial x} \frac{\partial x}{\partial w_2}} \quad | \quad \boxed{\frac{\partial \ell}{\partial w_1} = \frac{\partial \ell}{\partial \sigma} \frac{\partial \sigma}{\partial x} \frac{\partial x}{\partial \text{Input2}} \frac{\partial \text{Input2}}{\partial w_1}}$$