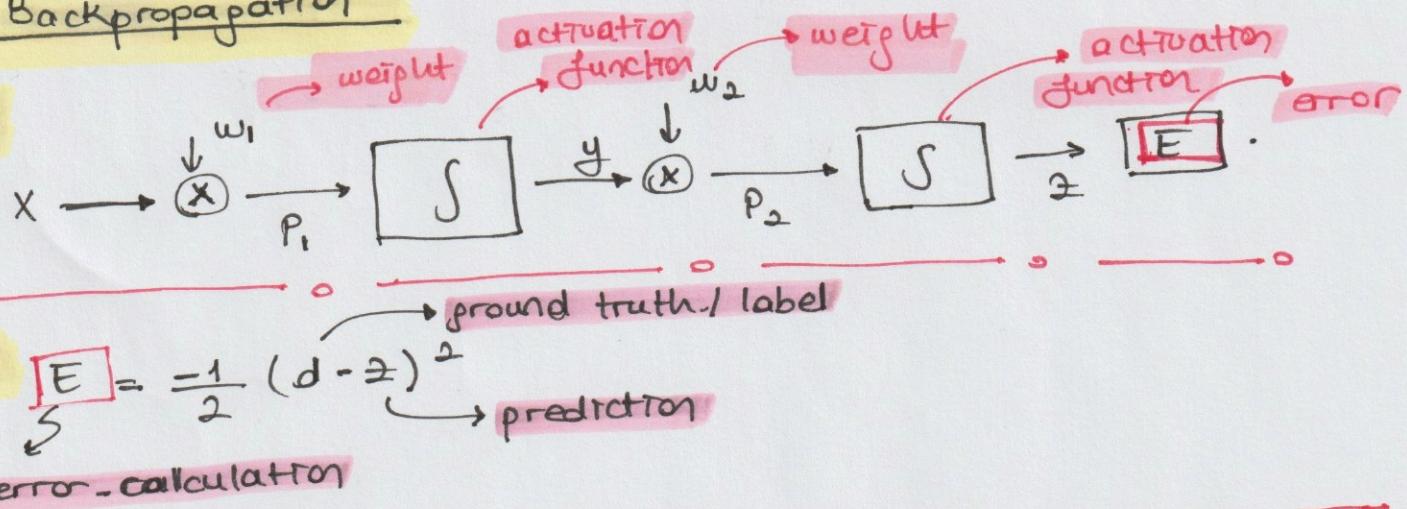


Backpropagation

①



②

$$E = \frac{1}{2} (d - p_2)^2$$

ground truth / label

prediction

③

$$\frac{\partial E}{\partial w_1} = ?$$

$$\frac{\partial E}{\partial w_2} = ?$$

④

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial z} \cdot \frac{\partial z}{\partial p_2} \cdot \frac{\partial p_2}{\partial w_2}$$

$$\frac{\partial E}{\partial z} = \frac{-1}{2} \cdot 1 \cdot (d - z)$$

$$= (d - z)$$

$$\frac{\partial p_2}{\partial w_2} = y$$

calculated from label and output of the current layer.

calculated from output of the current layer.

input of the current layer.

* let's say we use sigmoid activation function, so $z = \text{Act}(p_2)$

Sigmoid

$$\frac{\partial z}{\partial p_2} = z(1-z)$$

⑤

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial z} \cdot \frac{\partial z}{\partial p_2} \cdot \frac{\partial p_2}{\partial y} \cdot \frac{\partial y}{\partial p_1} \cdot \frac{\partial p_1}{\partial w_1}$$

$$\frac{\partial E}{\partial z} = (d - z)$$

$$\frac{\partial z}{\partial p_2} = z(1-z)$$

$$\frac{\partial p_2}{\partial y} = w_2$$

$$\frac{\partial y}{\partial p_1} = y(1-y)$$

$$\frac{\partial p_1}{\partial w_1} = x$$

①

Known (calculated) weight of the next layer (already calculated in fprop) of the current layer

Question - 1

⇒ The last layer of the network:

$$\begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \rightarrow \begin{bmatrix} \sigma_1(x) \\ \vdots \\ \sigma_N(x) \end{bmatrix}$$

Softmax function takes a N dimensional vector, and produces a N dimensional vector $\rightarrow R^N \rightarrow R^N$

$$x = [x_1, \dots, x_N] \in R^{1 \times N}$$

$$\sigma(x) = \text{softmax}(x) = [\sigma_1(x), \dots, \sigma_N(x)]$$

$$\sigma_i(x) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

(a)

$$D_x \sigma(x) \in R^{N \times N} = ?$$

Reminder : $J = \left[\frac{\partial f}{\partial x_1} \cdots \frac{\partial f}{\partial x_N} \right]$
 Jacobian Matrix

$$\Rightarrow \begin{bmatrix} \frac{\partial \sigma_1(x)}{\partial x_1} & \cdots & \frac{\partial \sigma_N(x)}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial \sigma_1(x)}{\partial x_N} & \cdots & \frac{\partial \sigma_N(x)}{\partial x_N} \end{bmatrix}$$

⇒ How we're gonna compute $\frac{\partial \sigma_i(x)}{\partial x_i}$ or $\frac{\partial \sigma_i(x)}{\partial x_j}$ values??

⇒ We're gonna use two tricks:

① calculate the diagonal entries

② calculate off-diagonal entries.

① Calculate diagonal entries

$$\frac{\partial \sigma_i(x)}{\partial x_i} = \frac{1}{\sum_{j=1}^N e^{x_j}} \left[\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \right] \Rightarrow \text{TRICK-1}$$

$$\left(\frac{1}{\sum_{j=1}^N e^{x_j}} \right)' = \frac{f' f - f f'}{f^2}$$

②

$$f \rightarrow e^{xi}$$

$$g \rightarrow \sum_{j=1}^N e^{x_j}$$

$$f' \rightarrow (e^{xi})' \rightarrow e^{xi}$$

$$g' \rightarrow \left(\sum_{j=1}^N e^{x_j} \right)' \rightarrow e^{xi}$$

or

$$g' \rightarrow \frac{\partial (e^{x_1} + e^{x_2} + \dots + e^{x_i} + \dots + e^{x_N})}{\partial x_i}$$

$$\rightarrow e^{xi}$$

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

$$= \frac{e^{xi} \left(\sum_{j=1}^N e^{x_j} \right) - e^{xi} \cdot e^{xi}}{\left(\sum_{j=1}^N e^{x_j} \right)^2}$$

② Calculate off-diagonal entries:

$$\frac{\partial \sigma_i(x)}{\partial x_j} = \frac{\partial}{\partial x_j} \left[\frac{e^{xi}}{\sum_{j=1}^N e^{x_j}} \right] \Rightarrow \text{Trick - 2}$$

$$(c \cdot f^{-1})' = -c \cdot f^{-2} \cdot f'$$

where $c = \text{constant}$

$$= \frac{\partial}{\partial x_j} \left[e^{xi} \cdot \frac{1}{\sum_{j=1}^N e^{x_j}} \right]$$

c

$$\rightarrow f^{-1} \text{ where } f = \sum_{j=1}^N e^{x_j}$$

$$f \rightarrow \sum_{j=1}^N e^{x_j}$$

$$c \rightarrow e^{xi}$$

$$f' \rightarrow \left(\sum_{j=1}^N e^{x_j} \right)' \rightarrow e^{x_j}$$

$$f' \rightarrow \frac{\partial (e^{x_1} + e^{x_2} + \dots + e^{x_j} + \dots + e^{x_N})}{\partial x_j}$$

$$f' \rightarrow e^{x_j}$$

$$\frac{\partial \sigma_i(x)}{\partial x_j} = \frac{-e^{x_i} \cdot e^{x_j}}{\left(\sum_{j=1}^N e^{x_j}\right)^2}$$

①

$$\frac{\partial \sigma_i(x)}{\partial x_i} = \frac{e^{x_i} \left(\sum_{j=1}^N e^{x_j} \right) - (e^{x_i})^2}{\left(\sum_{j=1}^N e^{x_j}\right)^2}$$

$$\begin{aligned} &= \frac{e^{x_i} \left(\sum_{j=1}^N e^{x_j} \right)}{\left(\sum_{j=1}^N e^{x_j}\right)^2} - \left(\frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}} \right)^2 \\ &\quad \sum_{j=1}^N e^{x_j} \\ &= \sigma_i(x) - \sigma_i^2(x) \end{aligned}$$

$$= \sigma_i(1 - \sigma_i(x)) \quad \text{diagonal}$$

②

$$\frac{\partial \sigma_i(x)}{\partial x_j} = \underbrace{\frac{e^{x_j}}{\sum_{j=1}^N e^{x_j}}} \cdot \underbrace{\frac{e^{x_j}}{\sum_{j=1}^N e^{x_j}}} \quad \sigma_i(x) \quad \sigma_j(x)$$

$$= -\sigma_i(x) \sigma_j(x) \quad \text{off-diagonal}$$

finally, $D_x \sigma(x) =$

$\sigma_1(x)(1 - \sigma_1(x))$	\dots	$\sigma_N(x)\sigma_1(x)$
$-\sigma_1(x)\sigma_2(x)$	\dots	$-\sigma_N(x)\sigma_2(x)$
\vdots	\ddots	\vdots
$-\sigma_1(x)\sigma_N(x)$	\dots	$\sigma_N(x)(1 - \sigma_N(x))$

diagonal

$$(b) \quad z = [v_1, \dots, v_N]_{1 \times N} \left[\begin{array}{c} D_x \sigma(x) \\ \vdots \\ D_x \sigma(x) \end{array} \right]_{N \times N}$$

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

$$= [z_1, \dots, z_N]_{1 \times N}$$

\Rightarrow Let's calculate z_i , specifically $i=1$. For z_i :

$$z_1 = v_1 \cdot \text{col}_1(D_x \sigma(x))$$

$$= (v_1, \sigma_1(x)) (1 - \sigma_1(x)) + v_2 (-\sigma_1(x) \sigma_2(x)) + \dots + v_N \frac{(-\sigma_1(x))}{\sigma_N(x)}$$

$$= v_1 \sigma_1(x) - v_1 \sigma_1^2(x) + v_2 (-\sigma_1(x) \sigma_2(x)) + \dots + v_N (-\sigma_1(x) \sigma_N(x))$$

$$= \sigma_1(x) (v_1 - v_1 \sigma_1(x) - v_2 \sigma_2(x) - \dots - v_N \sigma_N(x)) \quad \text{Trick}$$

$$= \sigma_1(x) (v_1 - v_1 \sigma_1(x) - v_2 \sigma_2(x) - \dots - v_N \sigma_N(x) + v_1 \sigma_1(x) - v_1 \sigma_1(x))$$

$$= - \sum_{j=1}^N v_j \sigma_j(x)$$

$$= \sigma_1(x) (v_1 - v_1 \sigma_1(x) - \sum_{j=1}^N v_j \sigma_j(x) + v_1 \sigma_1(x))$$

$$= \sigma_1(x) (v_1 - v \cdot \sigma(x)^T)$$

$$= \sigma_1(x) (v_1 - v \cdot \sigma(x)^T)$$

Generally,
$$z_i = \sigma_i(x) (v_i - v \cdot \sigma(x)^T)$$

$$(c) \ell(\mathbf{z}, \mathbf{t}) = -\sum_{i=1}^N t_i \ln(z_i) \text{ with } \mathbf{t} \in [0, 1]^{1 \times N}$$

for example $\mathbf{t} = [0 \ 0 \ 1 \ 0]$

$$D_{\mathbf{z}} \ell(\mathbf{z}, \mathbf{t}) \in \mathbb{R}^{1 \times N} = ?$$

$$D_{\mathbf{z}} \ell(\mathbf{z}, \mathbf{t}) = \left[\frac{\partial \ell(\mathbf{z}, \mathbf{t})}{\partial z_1}, \dots, \frac{\partial \ell(\mathbf{z}, \mathbf{t})}{\partial z_N} \right]$$

$$\frac{\partial \ell(\mathbf{z}, \mathbf{t})}{\partial z_i} = ?$$

$$\frac{\partial \ell(\mathbf{z}, \mathbf{t})}{\partial z_i} = \frac{\partial}{\partial z_i} (-t_i \ln(z_i))$$

$$\text{or } \frac{\partial}{\partial z_i} (-t_1 \ln(z_1) - t_2 \ln(z_2) - \dots - t_i \ln(z_i) - t_N \ln(z_N))$$

$$= \frac{\partial}{\partial z_i} (-t_i \ln(z_i)) \quad \text{Trick } \ln(x)' = \frac{1}{x}$$

$$x = z_i$$

$$= \frac{\partial}{\partial z_i} (-t_i \ln(z_i)) = -t_i \frac{1}{z_i}$$

$$D_{\mathbf{z}} \ell(\mathbf{z}, \mathbf{t}) = \left[\frac{-t_1}{z_1}, \dots, \frac{-t_N}{z_N} \right]$$

- (d) Division by zero when $z_i = 0$ or any class gets 0 probability.