

Softmax

$$X \in \mathbb{R}^{k \times n}$$

k = batch size

n = input dim

$$Y = \text{Softmax}(X)$$

$$Y \in \mathbb{R}^{k \times n}$$

Forward Propagation

Suppose $X \in \mathbb{R}^{1 \times n}$, $Y \in \mathbb{R}^{1 \times n}$ then

$$Y = \text{Softmax}(X)$$

$$Y_i = \frac{e^{x_i}}{\sum_j e^{x_j}}, \quad \text{here } 1 \leq i \leq n, 1 \leq j \leq n$$

①

Now Suppose $X \in \mathbb{R}^{k \times n}$, $Y \in \mathbb{R}^{k \times n}$

here $k = \text{batch size}$
 $n = \text{input dim}$

then by equation ①

$$Y_{ij} = \frac{e^{x_{ij}}}{\sum_d e^{x_{id}}}, \quad \begin{array}{l} \text{where} \\ 1 \leq i \leq k \\ 1 \leq j \leq n \\ 1 \leq d \leq n \end{array}$$

————— ②

Derivative

Suppose $X \in \mathbb{R}^{1 \times n}$, $Y \in \mathbb{R}^{1 \times n}$

$$\text{then } \frac{dL}{dx} = \frac{dL}{dY} \frac{dY}{dx}$$

$$\text{where } \frac{dL}{dY} \in \mathbb{R}^{1 \times n}$$

Now by equation (1)

$$Y_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\left(\frac{\partial Y_i}{\partial x} \right)_j = \left\{ \right.$$

To solve above suppose $n=3$
and

$$X = [x_1, x_2, x_3]$$

$$Y = [Y_1, Y_2, Y_3]$$

then,

$$Y_i = \frac{e^{x_i}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

Now there are two cases when

$$\frac{\partial Y_i}{\partial x_i} \text{ and } \frac{\partial Y_i}{\partial x_j}$$

So for first

$$\frac{dy_1}{dx_1} = \frac{d}{dx_1} \left(\frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \right)$$

$$\frac{dy_1}{dx_1} = e^{x_1} \frac{d}{dx_1} \left(\frac{1}{e^{x_1} + e^{x_2} + e^{x_3}} \right) +$$

$$\frac{d e^{x_1}}{dx_1} \left(\frac{1}{e^{x_1} + e^{x_2} + e^{x_3}} \right)$$

$$\frac{dy_1}{dx_1} = \frac{-e^{x_1} e^{x_1}}{(e^{x_1} + e^{x_2} + e^{x_3})^{-2}} + \frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}}$$

Now Suppose

$$p_1 = \frac{e^{x_1}}{(e^{x_1} + e^{x_2} + e^{x_3})^{-2}}$$

then

$$\frac{\partial Y_1}{\partial x_1} = -P_1^2 + P_1$$

$$\frac{\partial Y_1}{\partial x_1} = P_1(1 - P_1)$$

3

now for $\frac{\partial Y_i}{\partial x_j}$

$$\frac{\partial Y_1}{\partial x_2} = \frac{d}{dx_2} \left(\frac{e^{x_1}}{e^{x_1} + e^{x_2} + e^{x_3}} \right)$$

$$\frac{\partial Y_1}{\partial x_2} = \frac{-e^{x_1}e^{x_2}}{(e^{x_1} + e^{x_2} + e^{x_3})^2}$$

Now Suppose

$$P_1 = \frac{e^{x_1}}{(e^{x_1} + e^{x_2} + e^{x_3})}$$

$$P_2 = \left(\frac{e^{x_2}}{e^{x_1} + e^{x_2} + e^{x_3}} \right)$$

then

$$\frac{\partial y_1}{\partial x_2} = -p_1 p_2$$

(4)

Now we can generalize (3) and (4)

$$\left(\frac{\partial y}{\partial x}\right)_{i,j} = \left(\frac{\partial y_i}{\partial x}\right)_j = \begin{cases} p_i(1-p_i) & i=j \\ -p_i p_j & i \neq j \end{cases}$$

where, $\frac{\partial y}{\partial x} \in \mathbb{R}^{n \times n}$

(5)

Now

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

hence $\frac{\partial L}{\partial y} \in \mathbb{R}^{1 \times n}$, $\frac{\partial y}{\partial x} \in \mathbb{R}^{n \times n}$

then by equation (5)

$$\left(\frac{\partial L}{\partial x}\right)_i = \left(\frac{\partial L}{\partial y}\right)_1 \left(\frac{\partial y}{\partial x}\right)_{i1} + \dots + \left(\frac{\partial L}{\partial y}\right)_n \left(\frac{\partial y}{\partial x}\right)_{in}$$
$$\frac{\partial L}{\partial x} \in \mathbb{R}^{1 \times n}$$

————— (6)

Now Suppose $x \in \mathbb{R}^{k \times n}$, $y \in \mathbb{R}^{k \times n}$

=> Hence take each input independent
and once find put it in matrix
So the output from (6) will be

$$\frac{\partial L}{\partial y} \in \mathbb{R}^{k \times n}, \quad \frac{\partial y}{\partial x} \in \mathbb{R}^{k \times n \times n}, \quad \frac{\partial L}{\partial x} \in \mathbb{R}^{k \times n}$$

=> Means we have to apply (6) k
times.

