

# 1 Derivative of Binary Cross-Entropy cost function with respect to a weight w

The Binary Cross-Entropy cost function is defined as:

$$C = \frac{1}{n} \sum_{i=1}^n h_i^r \log(h_i^p) + (1 - h_i^r) \log(1 - h_i^p) \quad (1)$$

Taking the derivative of  $C$  with respect to a weight  $w$ :

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n h_i^r \frac{1}{h_i^p} \frac{\partial h_i^p}{\partial w} - (1 - h_i^r) \frac{1}{1 - h_i^p} \frac{\partial h_i^p}{\partial w} \quad (2)$$

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n h_i^r \frac{1}{h_i^p} \frac{\partial h_i^p}{\partial w} - (1 - h_i^r) \frac{1}{1 - h_i^p} \frac{\partial h_i^p}{\partial w} \quad (3)$$

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n \frac{\partial h_i^p}{\partial w} \frac{h_i^r(1 - h_i^p) - (1 - h_i^r)h_i^p}{h_i^p(1 - h_i^p)} \quad (4)$$

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n \frac{\partial h_i^p}{\partial w} \frac{h_i^r - y_i^p}{h_i^p(1 - h_i^p)} \quad (5)$$

## 2 Binary Cross-Entropy cost function with sigmoid activation

$$h_i^p = \sigma(z_i) \quad (6)$$

$$\frac{\partial h_i^p}{\partial w} = \frac{\partial z_i}{\partial w} \sigma'(z_i) \quad (7)$$

$$= \frac{\partial z_i}{\partial w} \sigma(z_i)(1 - \sigma(z_i)) \quad (8)$$

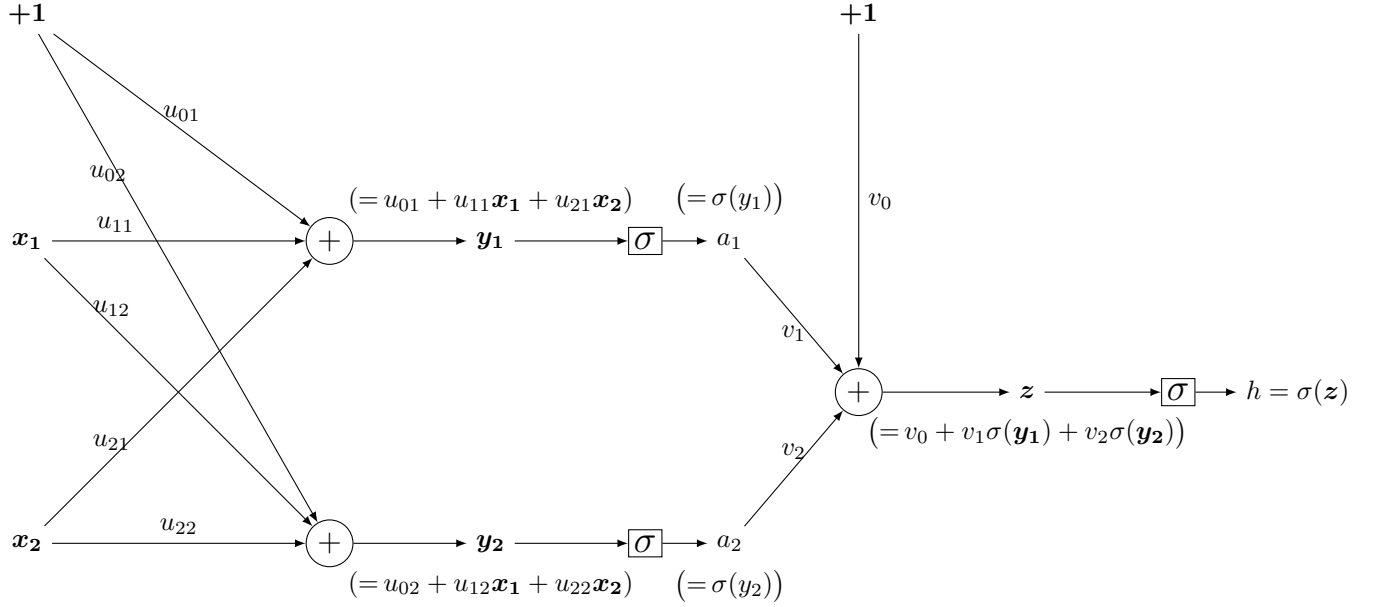
$$= \frac{\partial z_i}{\partial w} h_i^p(1 - h_i^p) \quad (9)$$

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n \frac{\partial z_i}{\partial w} h_i^p(1 - h_i^p) \frac{h_i^r - h_i^p}{h_i^p(1 - h_i^p)} \quad (10)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial z_i}{\partial w} (h_i^p - h_i^p) \quad (11)$$

## 3 2-layer architecture with sigmoid activation

Let us suppose we want to use the following linear neural network architecture, with a sigmoid activation function  $\sigma$ :



Let us define the following matrices and operations which help model the above network:

$$\begin{aligned}
 \mathbf{X}_0 &= \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \\ x_{41} & x_{42} \end{bmatrix} & \mathbf{X} = \begin{bmatrix} 1 & \mathbf{X}_0 \end{bmatrix} &= \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{32} \\ 1 & x_{41} & x_{42} \end{bmatrix} & \mathbf{u} = \begin{bmatrix} u_{01} & u_{02} \\ u_{11} & u_{12} \\ u_{21} & u_{22} \end{bmatrix} \\
 \mathbf{Y}_0 = \mathbf{X}\mathbf{u} &= \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ y_{31} & y_{32} \\ y_{41} & y_{42} \end{bmatrix} & \mathbf{A}_0 = \sigma(\mathbf{Y}_0) &= \begin{bmatrix} \sigma(y_{11}) & \sigma(y_{12}) \\ \sigma(y_{21}) & \sigma(y_{22}) \\ \sigma(y_{31}) & \sigma(y_{32}) \\ \sigma(y_{41}) & \sigma(y_{42}) \end{bmatrix} & \begin{bmatrix} 1 & \mathbf{A}_0 \end{bmatrix} = \begin{bmatrix} 1 & a_{11} & a_{12} \\ 1 & a_{21} & a_{22} \\ 1 & a_{31} & a_{32} \\ 1 & a_{41} & a_{42} \end{bmatrix} & \mathbf{v} = \begin{bmatrix} v_0 \\ v_1 \\ v_2 \end{bmatrix} \\
 \mathbf{z} = \mathbf{A}\mathbf{v} &= \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} & \mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \end{bmatrix} &= \sigma(\mathbf{z}) = \begin{bmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \sigma(z_3) \\ \sigma(z_4) \end{bmatrix}
 \end{aligned}$$

The final output of the network can also be written as:

$$\begin{aligned}
 h_i^{pred} &= \sigma(v_0 + v_1a_{i1} + v_2a_{i2}) \\
 &= \sigma\left(v_0 + v_1\sigma(u_{01} + u_{11}x_{i1} + u_{21}x_{i2}) + v_2\sigma(u_{02} + u_{12}x_{i1} + u_{22}x_{i2})\right)
 \end{aligned}$$

Let us use the binary cross-entropy cost function as defined in section 1:

$$C = \frac{1}{n} \sum_{i=1}^n h_i^r \log(h_i^p) + (1 - h_i^r) \log(1 - h_i^p) \quad (12)$$

Remembering that  $h_i^p = \sigma(z_i)$ , the derivative with respect to weight  $w$  is:

$$\frac{\partial C}{\partial w} = \frac{1}{n} \sum_{i=1}^n \frac{\partial z_i}{\partial w} (h_i^p - h_i^r) \quad (13)$$

$z_i$  is given by:

$$z_i = v_0 + v_1a_{i1} + v_2a_{i2} \quad (14)$$

$$= v_0 + v_1\sigma(u_{01} + u_{11}x_{i1} + u_{21}x_{i2}) + v_2\sigma(u_{02} + u_{12}x_{i1} + u_{22}x_{i2}) \quad (15)$$

Using (??) we can calculate the derivative of  $C$  with respect to each weight  $u_{ij}$  and  $v_i$  of the network and write them in vector form:

$$\frac{\partial z_i^{pred}}{\partial \mathbf{u}} = \begin{bmatrix} \frac{\partial z_i^{pred}}{\partial u_{01}} & \frac{\partial z_i^{pred}}{\partial u_{02}} \\ \frac{\partial z_i^{pred}}{\partial u_{11}} & \frac{\partial z_i^{pred}}{\partial u_{12}} \\ \frac{\partial z_i^{pred}}{\partial u_{21}} & \frac{\partial z_i^{pred}}{\partial u_{22}} \end{bmatrix} = \begin{bmatrix} v_1 a_{i1}(1 - a_{i1}) & v_2 a_{i2}(1 - a_{i2}) \\ v_1 x_{i1} a_{i1}(1 - a_{i1}) & v_2 x_{i1} a_{i2}(1 - a_{i2}) \\ v_1 x_{i2} a_{i1}(1 - a_{i1}) & v_2 x_{i2} a_{i2}(1 - a_{i2}) \end{bmatrix}$$

$$\frac{\partial z_i^{pred}}{\partial \mathbf{v}} = \begin{bmatrix} \frac{\partial z_i^{pred}}{\partial v_0} \\ \frac{\partial z_i^{pred}}{\partial v_1} \\ \frac{\partial z_i^{pred}}{\partial v_2} \end{bmatrix} = \begin{bmatrix} 1 \\ a_{i1} \\ a_{i2} \end{bmatrix}$$

Let us define  $d_i$  such that

$$d_i = z_i^{reality} - z_i^{pred} \quad \mathbf{d} = \begin{bmatrix} d_1 \\ d_2 \\ d_3 \\ d_4 \end{bmatrix}$$

By plugging in the corresponding derivatives found above in to (dC), we obtain:

$$\begin{aligned} \frac{\partial C}{\partial u_{01}} &= v_1 \sum_{i=1}^n a_{i1}(1 - a_{i1})d_i & \frac{\partial C}{\partial u_{02}} &= v_2 \sum_{i=1}^n a_{i2}(1 - a_{i2})d_i \\ \frac{\partial C}{\partial u_{11}} &= v_1 \sum_{i=1}^n x_{i1} a_{i1}(1 - a_{i1})d_i & \frac{\partial C}{\partial u_{12}} &= v_2 \sum_{i=1}^n x_{i1} a_{i2}(1 - a_{i2})d_i \\ \frac{\partial C}{\partial u_{21}} &= v_1 \sum_{i=1}^n x_{i2} a_{i1}(1 - a_{i1})d_i & \frac{\partial C}{\partial u_{22}} &= v_2 \sum_{i=1}^n x_{i2} a_{i2}(1 - a_{i2})d_i \end{aligned}$$

$$\begin{aligned} \frac{\partial C}{\partial v_0} &= \sum_{i=1}^n d_i \\ \frac{\partial C}{\partial v_1} &= \sum_{i=1}^n a_{i1} d_i \\ \frac{\partial C}{\partial v_2} &= \sum_{i=1}^n a_{i2} d_i \end{aligned}$$

Which can be cast in matrix form as:

$$\begin{aligned} \frac{\partial C}{\partial \mathbf{u}} &= [\text{diag}(\mathbf{d}) \times \mathbf{X}]^T [(1 - \mathbf{A}) \odot \mathbf{A}] \times \text{diag} \left( \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \right) \\ \frac{\partial C}{\partial \mathbf{v}} &= \mathbf{A}^T \mathbf{d} \end{aligned}$$

Where  $\text{diag}(\mathbf{u})$  denotes the diagonal matrix whose diagonal entries are the entries of vector  $u$ .