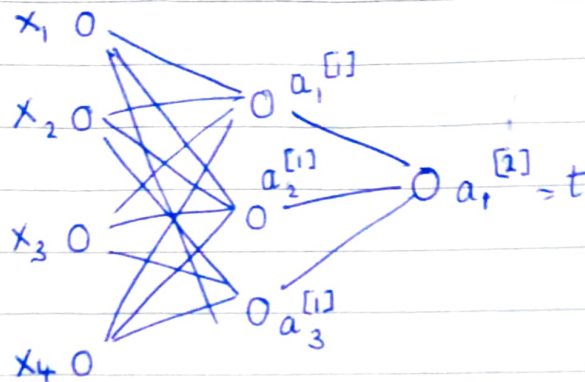


Given network :



Loss function : $\mathcal{L}^{(i)}(t^{(i)}, y^{(i)}) = -y^{(i)} \log t^{(i)} - (1-y^{(i)}) \log (1-t^{(i)})$
 Cost : $J(W, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}^{(i)}(t^{(i)}, y^{(i)})$

Now, to generate predictions, apply forward propagation :

For instance i , we have :

$$z^{[1](i)} = W^{[1]} x^{(i)} + b^{[1]} \quad - (1)$$

$$a^{[1](i)} = g^{[1]}(z^{[1](i)}), \text{ where } g^{[1]} = \text{ReLU} \quad - (2)$$

$$= \begin{cases} z_j^{[1](i)}, & z_j^{[1]} > 0 \\ 0, & \text{otherwise} \end{cases} \quad \forall j$$

$$z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]} \quad - (3)$$

$$a^{[2](i)} = t^{(i)} = g^{[2]}(z^{[2](i)}), \text{ where } g^{[2]} = \sigma \quad - (4)$$

$$= \frac{1}{1 + e^{-z^{[2](i)}}}$$

$$\mathcal{L}^{(i)}(t^{(i)}, y^{(i)}) = -y^{(i)} \log t^{(i)} - (1-y^{(i)}) \log (1-t^{(i)})$$

Vectorizing across all instances, we have :

$$Z^{[1]} = W^{[1]} X + b^{[1]} \quad - (5)$$

$$A^{[1]} = \text{ReLU}(Z^{[1]}) \quad - (6)$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]} \quad - (7)$$

$$A^{[2]} = \sigma(Z^{[2]}) = T \quad - (8)$$

$$J(W, b) = (-Y \log^T(T) - (1-Y) \log^T(1-T)) / m.$$



Now, parameters are updated via gradient descent:

$$W^{[2]} = W^{[2]} - \frac{\partial J}{\partial W^{[2]}} \cdot \alpha$$

$$b^{[2]} = b^{[2]} - \frac{\partial J}{\partial b^{[2]}} \cdot \alpha$$

To compute derivatives, use backpropagation:

For layer 2, we have j for each instance i :

$$W_{jk}^{[2]} = W_{jk}^{[2]} - \frac{\partial \mathcal{L}^{(i)}}{\partial W_{jk}^{[2]}} \cdot \alpha \Rightarrow W_k^{[2]} = W_k^{[2]} - \frac{\partial \mathcal{L}^{(i)}}{\partial W_k^{[2]}} \alpha$$

$$b_j^{[2]} = b_j^{[2]} - \frac{\partial \mathcal{L}^{(i)}}{\partial b_j^{[2]}} \cdot \alpha \Rightarrow b^{[2]} = b^{[2]} - \frac{\partial \mathcal{L}^{(i)}}{\partial b_j^{[2]}} \alpha$$

$$\begin{aligned} \text{Now, } \frac{\partial \mathcal{L}^{(i)}}{\partial W_k^{[2]}} &= \frac{\partial \mathcal{L}^{(i)}}{\partial t^{(i)}} \frac{\partial t^{(i)}}{\partial z^{[2](i)}} \frac{\partial z^{[2](i)}}{\partial W_k^{[2]}} \\ &= \frac{t^{(i)} - y^{(i)}}{t^{(i)}(1-t^{(i)})} t^{(i)}(1-t^{(i)}) a_k^{[1](i)} \\ &= (t^{(i)} - y^{(i)}) a_k^{[1](i)} \\ \frac{\partial \mathcal{L}^{(i)}}{\partial b^{[2]}} &= \frac{\partial \mathcal{L}^{(i)}}{\partial t^{(i)}} \frac{\partial t^{(i)}}{\partial z^{[2](i)}} \frac{\partial z^{[2](i)}}{\partial b^{[2]}} \\ &= (t^{(i)} - y^{(i)}) \cdot 1 \end{aligned}$$

Vectorizing across instances & k :

$$\begin{aligned} \frac{\partial J}{\partial W_k^{[2]}} &= \frac{1}{m} \left[\frac{\partial \mathcal{L}^{(1)}}{\partial W_k^{[2]}} \quad \frac{\partial \mathcal{L}^{(2)}}{\partial W_k^{[2]}} \quad \frac{\partial \mathcal{L}^{(3)}}{\partial W_k^{[2]}} \quad \dots \right] \\ &= \frac{1}{m} \begin{bmatrix} t^{(1)} - y^{(1)} & t^{(2)} - y^{(2)} & \dots & t^{(m)} - y^{(m)} \end{bmatrix} \begin{bmatrix} a_k^{1} \\ a_k^{[1](2)} \\ a_k^{[1](3)} \\ \vdots \\ a_k^{[1](m)} \end{bmatrix} \\ &= \frac{1}{m} (T - Y) A_k^{[1]T} \end{aligned}$$

$$\frac{\partial J}{\partial W^{[2]}} = \frac{1}{m} (T - Y) A^{[1]T}$$

$$\frac{\partial J}{\partial b^{[2]}} = \frac{1}{m} \sum_{i=1}^m t^{(i)} - y^{(i)} = \text{np.sum}(T - Y, \text{axis}=1, \text{keepdims}=\text{True})$$

Now, for layer 1, for each instance i we have:

$$\begin{aligned} \frac{\partial \mathcal{L}^{(i)}}{\partial W_{jk}^{[1]}} &= \frac{\partial \mathcal{L}^{(i)}}{\partial t^{(i)}} \frac{\partial t^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial a_j^{[1](i)}} \frac{\partial a_j^{[1](i)}}{\partial z_j^{[1](i)}} \frac{\partial z_j^{[1](i)}}{\partial W_{jk}^{[1]}} \\ &= \frac{t^{(i)} - y^{(i)}}{t^{(i)} - t^{(i)}} \cdot t^{(i)} (1 - t^{(i)}) W_j^{[2]} \text{ReLU}'(z_j^{[1](i)}) x_k^{(i)} \\ &= (t^{(i)} - y^{(i)}) W_j^{[2]} \text{ReLU}'(z_j^{[1](i)}) x_k^{(i)} \end{aligned}$$

$$\text{where } \text{ReLU}'(x) = \begin{cases} 1 & x > 0 \\ 0 & x = 0 \\ 0 & x < 0 \end{cases}$$

$$\begin{aligned} \frac{\partial \mathcal{L}^{(i)}}{\partial b_j^{[1]}} &= \frac{\partial \mathcal{L}^{(i)}}{\partial t^{(i)}} \frac{\partial t^{(i)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial a_j^{[1](i)}} \frac{\partial a_j^{[1](i)}}{\partial z_j^{[1](i)}} \frac{\partial z_j^{[1](i)}}{\partial b_j^{[1]}} \\ &= (t^{(i)} - y^{(i)}) W_j^{[2]} \text{ReLU}'(z_j^{[1](i)}) \cdot 1 \end{aligned}$$

Vectorizing across instances:

$$\begin{aligned} \frac{\partial J}{\partial W_{jk}^{[1]}} &= \frac{1}{m} \left[\frac{\partial \mathcal{L}^{(1)}}{\partial W_{jk}^{[1]}} \quad \frac{\partial \mathcal{L}^{(2)}}{\partial W_{jk}^{[1]}} \quad \dots \quad \frac{\partial \mathcal{L}^{(m)}}{\partial W_{jk}^{[1]}} \right] \\ &= \frac{1}{m} \begin{bmatrix} t^{(1)} - y^{(1)} & W_j^{[2]} & \text{ReLU}'(z_j^{1}) \\ t^{(2)} - y^{(2)} & W_j^{[2]} & \text{ReLU}'(z_j^{[1](2)}) \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} x_k^{(1)} \\ x_k^{(2)} \\ \vdots \\ x_k^{(m)} \end{bmatrix} \end{aligned}$$

$$= \frac{1}{m} (T - Y) (W_j^{[2]} \odot \text{ReLU}'(z_j^{[1]})) x_k^T$$

$$\begin{aligned} \frac{\partial J}{\partial W^{[1]}} &= \frac{1}{m} (T - Y) (W^{[2]} \odot \text{ReLU}'(z^{[1]})) x^T \\ &= \frac{1}{m} (T - Y) \odot \text{ReLU}'(z^{[1]}) x^T \end{aligned}$$

$$\frac{\partial J}{\partial W^{[1]}} = W^{[2]T} (T - Y) \odot \text{ReLU}'(z^{[1]}) x^T$$

$$\begin{aligned}\frac{\partial J}{\partial b_j^{[1]}} &= \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}^{(i)}}{\partial b_j^{(i)}} \\ &= \text{np.sum}(W_j^{[2]} (T-Y) \odot \text{ReLU}'(z_j^{[1]}), \text{axis}=1, \text{keepdims}=T) \\ \frac{\partial J}{\partial b^{[1]}} &= \frac{1}{m} \text{np.sum}(W^{[2]T} (T-Y) \odot \text{ReLU}'(Z^{[1]}), \text{axis}=1, \text{keepdims}=T)\end{aligned}$$

Final update equations for gradient descent:

$$W^{[2]} = W^{[2]} - \frac{\alpha}{m} (T-Y) A^{[1]T}$$

$$b^{[2]} = b^{[2]} - \frac{\alpha}{m} \sum_{\text{axis}=1} (T-Y)$$

$$W^{[1]} = W^{[1]} - \frac{\alpha}{m} W^{[2]T} (T-Y) \odot \text{ReLU}'(Z^{[1]}) X^T$$

$$b^{[1]} = b^{[1]} - \frac{\alpha}{m} \sum_{\text{axis}=1} (W^{[2]T} (T-Y) \odot \text{ReLU}'(Z^{[1]}))$$