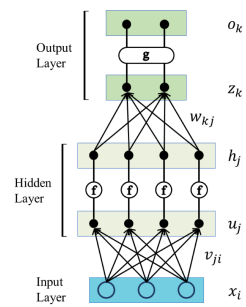


Problem 1. (60points)

Consider a neural network which receives an input $x = [x_1 \ x_2 \ x_3]^T$ and generates an output $o = [o_1 \ o_2]^T$. It consists of one input layer, one hidden layer, and one output layer. The hidden layer contains 4 neurons and the output layer has 2 neurons. The input x is transformed by a set of weights $\{v_{ji} \mid j = 1, 2, 3, 4, i = 1, 2, 3\}$ and the output of the hidden layer is transformed by weights $\{w_{kj} \mid k = 1, 2, j = 1, 2, 3, 4\}$. f is a *Softmax* activation function:

$$\text{ELU}(u) = \begin{cases} u & \geq 0 \\ e^u - 1 & < 0 \end{cases}$$

g is a *Softmax* activation function.



1) **(10points)** Derive the expressions for u_j , h_j , z_k , and o_k using the input, activation function, and weights.

$$u_j = \sum_i v_{ji} x_i = \sum_{i=1}^3 v_{ji} x_i = v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3$$

$$h_j = f(u_j) = \text{ELU}(u_j) = \text{ELU}(v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3)$$

$$z_k = \sum_j w_{kj} h_j = \sum_j w_{kj} \text{ELU}(u_j) = \sum_{j=1}^4 w_{kj} \text{ELU}(v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3)$$

$$o_k = g(z_k) = \frac{e^{z_k}}{e^{z_1} + e^{z_2}} = \frac{e^{\sum_{j=1}^4 w_{kj} \text{ELU}(v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3)}}{e^{\sum_{j=1}^4 w_{1j} \text{ELU}(v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3)} + e^{\sum_{j=1}^4 w_{2j} \text{ELU}(v_{j1} x_1 + v_{j2} x_2 + v_{j3} x_3)}}$$

2) **(10points)** Suppose that $x = [1.0 \ 1.5 \ -0.5]^T$ and the ground truth $o^* = [1 \ 0]^T$. And, the loss function \mathcal{L} and the initial weights are given as follows:

$$\mathcal{L} = - \sum_{k=1}^2 o_k^* \log o_k$$

$$v = \begin{bmatrix} -1.0 & 2.0 & 0.0 \\ 0.5 & -1.0 & 1.0 \\ 1.5 & 2.0 & -0.4 \\ 0.2 & 0.1 & 0.3 \end{bmatrix}$$

$$w = \begin{bmatrix} 1.0 & 1.0 & 0.5 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.1 \end{bmatrix}$$

Compute u_j , h_j , z_k , and \mathcal{L} ($j = 1, 2, 3, 4$ and $k = 1, 2$).

$$u = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{bmatrix} = \vec{v} \cdot \vec{x} = \begin{bmatrix} -1.0 & 2.0 & 0.0 \\ 0.5 & -1.0 & 1.0 \\ 1.5 & 2.0 & -0.4 \\ 0.2 & 0.1 & 0.3 \end{bmatrix} \begin{bmatrix} 1.0 \\ 1.5 \\ -0.5 \end{bmatrix} = \begin{bmatrix} 2 \\ -0.75 \\ 1.75 \\ 0.2 \end{bmatrix}$$

$$h = \begin{bmatrix} \text{ELU}(u_1) \\ \text{ELU}(u_2) \\ \text{ELU}(u_3) \\ \text{ELU}(u_4) \end{bmatrix} = \begin{bmatrix} 2 \\ -0.75 \\ 1.75 \\ 0.2 \end{bmatrix}$$

$$z = \vec{w} \cdot h = \begin{bmatrix} 1.0 & 1.0 & 0.5 & 0.5 \\ 0.3 & 0.4 & 0.5 & 0.1 \end{bmatrix} \begin{bmatrix} 2 \\ -0.75 \\ 1.75 \\ 0.2 \end{bmatrix} = \begin{bmatrix} 3.075 \\ 2.659 \end{bmatrix}$$

$$o = g(z) = \begin{bmatrix} g(z_1) \\ g(z_2) \end{bmatrix} = \begin{bmatrix} 0.9118 \\ 0.2662 \end{bmatrix}$$

$$\mathcal{L} = -(o_1^* \log o_1 + o_2^* \log o_2) = -(1 \times \log(0.9118)) = 0.1344$$

3) **(10points)** Derive the expressions for $\frac{\partial \mathcal{L}}{\partial w_{kj}}$ and $\frac{\partial \mathcal{L}}{\partial v_{ji}}$. (Hint: $\frac{\partial \mathcal{L}}{\partial z_k} = o_k - o_k^*$)

$$\frac{\partial \mathcal{L}}{\partial w_{kj}} = \underbrace{\frac{\partial \mathcal{L}}{\partial o_k}}_{o_k - o_k^*} \underbrace{\frac{\partial o_k}{\partial z_k}}_{h_j} \frac{\partial z_k}{\partial w_{kj}} = (o_k - o_k^*) h_j$$

$$\frac{\partial \mathcal{L}}{\partial v_{ji}} = \frac{\partial \mathcal{L}}{\partial h_j} \cdot \frac{\partial h_j}{\partial v_{ji}} = \frac{\partial \mathcal{L}}{\partial h_j} \cdot \frac{\partial h_j}{\partial u_j} \cdot \frac{\partial u_j}{\partial v_{ji}} = f'(u_j) x_i \sum_k (o_k - o_k^*) w_{kj}$$

$$\frac{\partial \mathcal{L}}{\partial h_j} = \sum_k \frac{\partial \mathcal{L}}{\partial o_k} \cdot \frac{\partial o_k}{\partial z_k} \cdot \frac{\partial z_k}{\partial h_j} = \sum_k (o_k - o_k^*) w_{kj}$$

$$f(u) = \begin{cases} 1 & (u \geq 0) \\ e^u & (u < 0) \end{cases}$$

4) **(15points)** Given $x = [1.0 \ 1.5 \ -0.5]^T$, compute the values of $\frac{\partial \mathcal{L}}{\partial w}$ and $\frac{\partial \mathcal{L}}{\partial v}$.

$$\frac{\partial \mathcal{L}}{\partial w} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w_{11}} & \frac{\partial \mathcal{L}}{\partial w_{12}} & \frac{\partial \mathcal{L}}{\partial w_{13}} & \frac{\partial \mathcal{L}}{\partial w_{14}} \\ \frac{\partial \mathcal{L}}{\partial w_{21}} & \frac{\partial \mathcal{L}}{\partial w_{22}} & \frac{\partial \mathcal{L}}{\partial w_{23}} & \frac{\partial \mathcal{L}}{\partial w_{24}} \end{bmatrix}, \quad \frac{\partial \mathcal{L}}{\partial v} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial v_{11}} \cdot \frac{\partial \mathcal{L}}{\partial v_{12}} \cdot \frac{\partial \mathcal{L}}{\partial v_{13}} \\ \frac{\partial \mathcal{L}}{\partial v_{21}} \cdot \frac{\partial \mathcal{L}}{\partial v_{22}} \cdot \frac{\partial \mathcal{L}}{\partial v_{23}} \\ \frac{\partial \mathcal{L}}{\partial v_{31}} \cdot \frac{\partial \mathcal{L}}{\partial v_{32}} \cdot \frac{\partial \mathcal{L}}{\partial v_{33}} \\ \frac{\partial \mathcal{L}}{\partial v_{41}} \cdot \frac{\partial \mathcal{L}}{\partial v_{42}} \cdot \frac{\partial \mathcal{L}}{\partial v_{43}} \end{bmatrix}$$

given \vec{x} and \vec{o}^* is

$$\vec{o} = [0.9118 \ 0.2662]^T$$

$$h = [2, -0.75, 1.75, 0.2]^T$$

$$\therefore \vec{o} - \vec{o}^* = [-0.2662 \ 0.2662]^T$$

$$\therefore \frac{\partial \mathcal{L}}{\partial \vec{w}} = [-0.9118, 0.0495] \cdot [2, -0.75, 1.75, 0.2]^T$$

$$= \begin{bmatrix} -0.5374 & 0.2068 & -1.251 & -0.0532 \\ 0.5374 & -0.2068 & 1.251 & 0.0532 \end{bmatrix}$$

$$\therefore \frac{\partial \mathcal{L}}{\partial \vec{v}} = \begin{bmatrix} -0.1863 & -0.2795 & 0.0932 \\ -0.0354 & -0.0535 & 0.0198 \\ 0 & 0 & 0 \\ -0.1065 & -0.1547 & 0.0532 \end{bmatrix}$$

5) **(15points)** Using the values of $\frac{\partial \mathcal{L}}{\partial w}$ and $\frac{\partial \mathcal{L}}{\partial v}$ from 4) and a learning rate $\eta = 0.1$, update the weights \vec{v} and \vec{w} via a gradient descent update rule.

$$\vec{w}_{\text{new}} = \vec{w}_{\text{old}} - 0.1 \frac{\partial \mathcal{L}}{\partial \vec{w}}$$

$$w_i^{\text{new}} = w_i^{\text{old}} - 0.1 \frac{\partial \mathcal{L}}{\partial w_i}$$

$$= \begin{bmatrix} 1.0532 & 0.9793 & 0.681 & 0.5053 \\ 0.7468 & 0.4209 & 0.3747 & 0.0467 \end{bmatrix}$$