

## Problem 1

Consider the following multi-layer neural network (Figure 1) which includes an input layer, one hidden layer, and an output layer, containing  $d$ ,  $n$ ,  $q$  neurons respectively. The parameters between the input layer and the hidden layer are  $\mathbf{W}_1 \in \mathbb{R}^{d \times n}$ ,  $\mathbf{b}_1 \in \mathbb{R}^n$ , and the parameters between the hidden layer and the output layer are  $\mathbf{W}_2 \in \mathbb{R}^{n \times q}$ ,  $\mathbf{b}_2 \in \mathbb{R}^q$ . Where  $\mathbf{W}_1$ ,  $\mathbf{W}_2$  are the weight matrices and  $\mathbf{b}_1$ ,  $\mathbf{b}_2$  are the bias vectors.

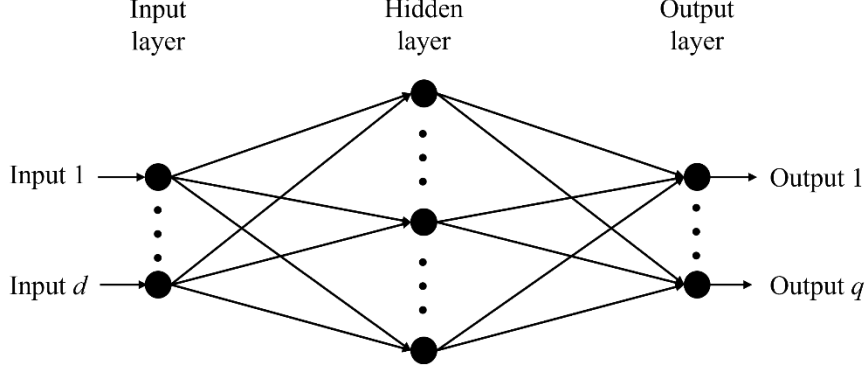


Figure 1: Neural Network

Let us first compute the forward propagation. Let  $\mathbf{x} \in \mathbb{R}^d$  be the input. The hidden layer is computed as follows:

$$\mathbf{h} = \mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1 \in \mathbb{R}^n \quad (1-1)$$

Then the ReLU activation function is applied to Eq.(1-1):

$$\mathbf{a} = \text{ReLU}(\mathbf{h}) \in \mathbb{R}^n \quad (1-2)$$

The output layers' activation is obtained using the following transformation

$$\mathbf{z} = \mathbf{W}_2^\top \mathbf{a} + \mathbf{b}_2 \in \mathbb{R}^q \quad (1-3)$$

Finally, the soft-max function is applied to Eq.(1-3) to obtain the probability for each category.

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_q \end{bmatrix} = \text{Softmax}(\mathbf{z}) = \begin{bmatrix} \text{Softmax}(z_1) \\ \text{Softmax}(z_2) \\ \vdots \\ \text{Softmax}(z_q) \end{bmatrix} = \begin{bmatrix} \frac{\exp(z_1)}{\sum_{i=1}^q \exp(z_i)} \\ \frac{\exp(z_2)}{\sum_{i=1}^q \exp(z_i)} \\ \vdots \\ \frac{\exp(z_q)}{\sum_{i=1}^q \exp(z_i)} \end{bmatrix} \in \mathbb{R}^q \quad (1-4)$$

Here,  $\hat{y}_i$ ,  $z_i$  is the  $i$ -th element of  $\hat{\mathbf{y}}$ ,  $\mathbf{z}$ ,  $\hat{\mathbf{y}}$  is the predicted output by the feed-forward

neural network.

Your task is to compute the derivatives of the Cross-Entropy loss function given in Eq.(1-5) with respect to  $\mathbf{W}_1$ ,  $\mathbf{W}_2$ ,  $\mathbf{b}_1$ ,  $\mathbf{b}_2$  by hand, i.e.,

$$\frac{\partial Loss}{\partial \mathbf{W}_1}, \frac{\partial Loss}{\partial \mathbf{W}_2}, \frac{\partial Loss}{\partial \mathbf{b}_1}, \frac{\partial Loss}{\partial \mathbf{b}_2}.$$
$$Loss = - \sum_{i=1}^q y_i \ln(\hat{y}_i) \quad (1-5)$$

Here,  $y_i \in \{0, 1\}$  is the ground-truth indicator,  $\hat{y}_i$  is the  $i$ -th element of  $\hat{\mathbf{y}}$ .

Show all the intermediate derivative computation steps. You might benefit from making a rough schematic of the back-propagation process.