

# Machine Learning (DS4023) Assignment 4

Deadline: Nov. 25, 2024.

## Problem 1: Neural Networks (20 pts)

Consider a 3-layer fully connected neural network with the following architecture:

- **Input layer:**  $n = 4$  neurons.
- **Hidden layer:**  $m = 3$  neurons using a custom activation function  $f(x) = \text{ReLU}(x) + \sin(x)$ .
- **Output layer:**  $k = 2$  neurons using the softmax activation function  $\sigma(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}}$ .

The network parameters (**weights and biases**) are given as:

- $\mathbf{W}_1 \in \mathbb{R}^{3 \times 4}$  and  $\mathbf{b}_1 \in \mathbb{R}^3$  for the hidden layer.
- $\mathbf{W}_2 \in \mathbb{R}^{2 \times 3}$  and  $\mathbf{b}_2 \in \mathbb{R}^2$  for the output layer.

Given the input vector  $\mathbf{x} \in \mathbb{R}^4$  and target output  $\mathbf{y} \in \mathbb{R}^2$ . Define the loss function as cross-entropy loss:

$$\text{Loss} = - \sum_{i=1}^k y_i \log(\hat{y}_i),$$

where  $\hat{y}$  is the output after the softmax activation.

### Your tasks:

- 1) Derive the equations for the forward pass through the network, including both the hidden and output layers. (3 pts)
- 2) Calculate the outputs  $\mathbf{Z}_1$ ,  $\mathbf{H}$ ,  $\mathbf{Z}_2$ , and  $\hat{\mathbf{y}}$  explicitly for a given input  $\mathbf{x} = [1, -1, 0.5, 2]^T$  and the following initial weights and biases:

$$\mathbf{W}_1 = \begin{pmatrix} 0.1 & -0.2 & 0.3 & 0.4 \\ 0.5 & -0.3 & 0.1 & -0.2 \\ 0.4 & 0.2 & -0.5 & 0.3 \end{pmatrix}, \quad \mathbf{b}_1 = \begin{pmatrix} 0.1 \\ -0.1 \\ 0.05 \end{pmatrix}$$

$$\mathbf{W}_2 = \begin{pmatrix} -0.3 & 0.2 & 0.1 \\ 0.4 & -0.5 & 0.3 \end{pmatrix}, \quad \mathbf{b}_2 = \begin{pmatrix} 0.05 \\ -0.05 \end{pmatrix}.$$

Note that  $\mathbf{Z}_1$  is the net input to the hidden layer,  $\mathbf{H}$  is the activation output of the hidden layer, and  $\mathbf{Z}_2$  is the net input to the output layer. (3 pts)

2. Derive the gradient of the loss with respect to each parameter ( $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$ ) in the network and obtain the gradient values using results from the first question. Use matrix calculus to express the gradients. Hint: You can first calculate the error terms  $\delta_2$  and  $\delta_1$  for each layer and use them to express the gradients. (10 pts)
3. Suppose the learning rate  $\alpha = 0.001$ . Please calculate the updated parameter values after one back propagation process. (4 pts)

## Problem 2: Programming (80 pts)

Complete the jupyter notebook attached on programming for CNN and RNN. Submit the completed file.

### To submit:

1. A file containing the written answer to the Problems 1.
2. The Jupyter notebook with solutions to Problem 2.

## Solution (Full marks can be given if it is correct for the first three decimal points.)

1. The net input to the hidden layer,  $\mathbf{Z}_1$ , is calculated as:

$$\mathbf{Z}_1 = \mathbf{W}_1 \mathbf{x} + \mathbf{b}_1 \quad (1 \text{ pt})$$

Expanding the matrix-vector multiplication and addition:

$$\mathbf{Z}_1 = \begin{pmatrix} 0.1 & -0.2 & 0.3 & 0.4 \\ 0.5 & -0.3 & 0.1 & -0.2 \\ 0.4 & 0.2 & -0.5 & 0.3 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \\ 0.5 \\ 2 \end{pmatrix} + \begin{pmatrix} 0.1 \\ -0.1 \\ 0.05 \end{pmatrix},$$

which is

$$\mathbf{Z}_1 = \begin{pmatrix} 1.25 \\ 0.45 \\ 0.55 \end{pmatrix} + \begin{pmatrix} 0.1 \\ -0.1 \\ 0.05 \end{pmatrix} = \begin{pmatrix} 1.35 \\ 0.35 \\ 0.6 \end{pmatrix}. \quad (0.5 \text{ pt})$$

The activation of the hidden layer  $\mathbf{H}$  is computed by applying the custom activation function  $f(x) = \text{ReLU}(x) + \sin(x)$  to each element of  $\mathbf{Z}_1$ .

$$\mathbf{H} = f(\mathbf{Z}_1) = \begin{pmatrix} \text{ReLU}(1.35) + \sin(1.35) \\ \text{ReLU}(0.35) + \sin(0.35) \\ \text{ReLU}(0.6) + \sin(0.6) \end{pmatrix} = \begin{pmatrix} 2.3257 \\ 0.6929 \\ 1.1646 \end{pmatrix}. \quad (1 \text{ pt})$$

The net input to the output layer  $\mathbf{Z}_2$  is calculated as:

$$\mathbf{Z}_2 = \mathbf{W}_2 \mathbf{H} + \mathbf{b}_2 \quad (1 \text{ pt})$$

Expanding the matrix-vector multiplication and addition:

$$\mathbf{Z}_2 = \begin{pmatrix} -0.3 & 0.2 & 0.1 \\ 0.4 & -0.5 & 0.3 \end{pmatrix} \begin{pmatrix} 2.3257 \\ 0.6929 \\ 1.1646 \end{pmatrix} + \begin{pmatrix} 0.05 \\ -0.05 \end{pmatrix} = \begin{pmatrix} -0.3927 \\ 0.8832 \end{pmatrix} \quad (0.5 \text{ pt})$$

Applying the softmax function to  $\mathbf{Z}_2$ :

$$\hat{y}_i = \frac{e^{\mathbf{Z}_{2,i}}}{\sum_{j=1}^2 e^{\mathbf{Z}_{2,j}}}, \quad (1 \text{ pt})$$

We have

$$\begin{aligned} \hat{y}_1 &= \frac{e^{-0.3927}}{e^{-0.3927} + e^{0.8832}} = \frac{0.6752}{0.6752 + 2.4182} = \frac{0.6752}{3.0934} \approx 0.2182, \\ \hat{y}_2 &= \frac{e^{0.8832}}{e^{-0.3927} + e^{0.8832}} = \frac{2.4182}{3.0934} \approx 0.7818. \end{aligned}$$

Thus:

$$\hat{\mathbf{y}} = \begin{pmatrix} 0.2182 \\ 0.7818 \end{pmatrix} \quad (1 \text{ pt})$$

## 2. Gradient calculation

We will compute gradients with respect to  $\mathbf{W}_2$ ,  $\mathbf{b}_2$ ,  $\mathbf{W}_1$ , and  $\mathbf{b}_1$  by applying the chain rule in backpropagation.

### Gradient of Loss with Respect to $\mathbf{W}_2$ and $\mathbf{b}_2$

Since we are using the softmax activation with cross-entropy loss, we have:

$$\boldsymbol{\delta}_2 = \hat{\mathbf{y}} - \mathbf{y}, \quad (1 \text{ pt})$$

where  $\boldsymbol{\delta}_2$  represents the error at the output layer.

The gradient of the loss with respect to  $\mathbf{W}_2$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_2} = \boldsymbol{\delta}_2 \mathbf{H}^T \quad (1 \text{ pt})$$

where  $\mathbf{H}^T$  is the transpose of the hidden layer output.

The gradient of the loss with respect to  $\mathbf{b}_2$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{b}_2} = \boldsymbol{\delta}_2 \quad (1 \text{ pt})$$

These expressions are derived from the chain rule applied to the softmax and cross-entropy, making the derivations for  $\mathbf{W}_2$  and  $\mathbf{b}_2$  correct.

### Gradient of Loss with Respect to $\mathbf{W}_1$ and $\mathbf{b}_1$

To find the gradients with respect to  $\mathbf{W}_1$  and  $\mathbf{b}_1$ , we need the hidden layer error term  $\delta_1$ , which backpropagates from the output layer through the hidden layer.

#### Calculation of $\delta_1$

The error term  $\delta_1$  for the hidden layer is influenced by the output layer error term  $\delta_2$  backpropagated through  $\mathbf{W}_2$ :

$$\delta_1 = (\mathbf{W}_2^T \delta_2) \odot f'(\mathbf{Z}_1) \quad (1 \text{ pt})$$

where  $f'(\mathbf{Z}_1)$  represents the element-wise derivative of the activation function applied to  $\mathbf{Z}_1$ .

For the activation function  $f(x) = \text{ReLU}(x) + \sin(x)$ , the derivative  $f'(x)$  is:

$$f'(x) = \begin{cases} 1 + \cos(x) & \text{if } x > 0 \\ \cos(x) & \text{if } x \leq 0 \end{cases} \quad (1 \text{ pt})$$

Thus,  $f'(\mathbf{Z}_1)$  is applied element-wise to each component of  $\mathbf{Z}_1$  to calculate  $\delta_1$ .

### Gradients with Respect to $\mathbf{W}_1$ and $\mathbf{b}_1$

Using  $\delta_1$ , the gradients for  $\mathbf{W}_1$  and  $\mathbf{b}_1$  are as follows:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_1} = \delta_1 \mathbf{x}^T \quad (1 \text{ pt})$$

where  $\mathbf{x}^T$  is the transpose of the input vector  $\mathbf{x}$ .

The gradient of the loss with respect to  $\mathbf{b}_1$  is simply:

$$\frac{\partial \text{Loss}}{\partial \mathbf{b}_1} = \delta_1 \quad (1 \text{ pt})$$

These gradients are derived by applying the chain rule from the hidden layer output  $\mathbf{H}$  back to  $\mathbf{W}_1$  and  $\mathbf{b}_1$  through  $\mathbf{Z}_1$ , taking into account the influence of the custom activation function's derivative.

### Summary of Gradient Expressions

Gradients for Output Layer:

$$\begin{aligned} \frac{\partial \text{Loss}}{\partial \mathbf{W}_2} &= \delta_2 \mathbf{H}^T \\ \frac{\partial \text{Loss}}{\partial \mathbf{b}_2} &= \delta_2 \end{aligned}$$

Gradients for Hidden Layer:

$$\begin{aligned} \frac{\partial \text{Loss}}{\partial \mathbf{W}_1} &= \delta_1 \mathbf{x}^T \\ \frac{\partial \text{Loss}}{\partial \mathbf{b}_1} &= \delta_1 \end{aligned}$$

This completes the derivation process for the gradients of each parameter in the network, considering the custom activation function  $f(x) = \text{ReLU}(x) + \sin(x)$  in the hidden layer.

**Substitution:**

Output layer error term  $\delta_2$ :

$$\delta_2 = \hat{\mathbf{y}} - \mathbf{y} = \begin{pmatrix} 0.2182 \\ 0.7818 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -0.7818 \\ 0.7818 \end{pmatrix}.$$

To compute  $\delta_1$ , we need  $\mathbf{W}_2^T \delta_2$  and the derivative of the activation function,  $F'(\mathbf{Z}_1)$ .

$$\mathbf{W}_2^T = \begin{pmatrix} -0.3 & 0.4 \\ 0.2 & -0.5 \\ 0.1 & 0.3 \end{pmatrix}$$

$$\mathbf{W}_2^T \delta_2 = \begin{pmatrix} -0.3 & 0.4 \\ 0.2 & -0.5 \\ 0.1 & 0.3 \end{pmatrix} \begin{pmatrix} -0.7818 \\ 0.7818 \end{pmatrix} = \begin{pmatrix} 0.5473 \\ -0.5473 \\ 0.1716 \end{pmatrix}. \quad (\mathbf{0.5 \ pt})$$

For  $F'(\mathbf{Z}_1)$ :

$$F'(\mathbf{Z}_1) = \begin{pmatrix} 1 + \cos(1.35) \\ 1 + \cos(0.35) \\ 1 + \cos(0.6) \end{pmatrix} \approx \begin{pmatrix} 1.219 \\ 1.9394 \\ 1.8253 \end{pmatrix}.$$

Thus:

$$\delta_1 = (\mathbf{W}_2^T \delta_2) \odot F'(\mathbf{Z}_1) = \begin{pmatrix} 0.5473 \\ -0.5473 \\ 0.1546 \end{pmatrix} \odot \begin{pmatrix} 1.219 \\ 1.9394 \\ 1.8253 \end{pmatrix} = \begin{pmatrix} 0.6672 \\ -1.0614 \\ 0.2855 \end{pmatrix}. \quad (\mathbf{0.5 \ pt})$$

The gradient of the loss with respect to  $\mathbf{W}_2$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_2} = \delta_2 H^T = \begin{pmatrix} -0.7818 \\ 0.7818 \end{pmatrix} \begin{pmatrix} 2.3257 & 0.6929 & 1.1646 \end{pmatrix}.$$

Calculating each term:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_2} = \begin{pmatrix} -1.8182 & -0.5417 & -0.9105 \\ 1.8182 & 0.5417 & 0.9105 \end{pmatrix}. \quad (\mathbf{0.5 \ pt})$$

The gradient of the loss with respect to  $\mathbf{b}_2$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{b}_2} = \delta_2 = \begin{pmatrix} -0.7818 \\ 0.7818 \end{pmatrix}. \quad (\mathbf{0.5 \ pt})$$

The gradient of the loss with respect to  $\mathbf{W}_1$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_1} = \delta_1 \mathbf{x}^T = \begin{pmatrix} 0.6672 \\ -1.0614 \\ 0.2855 \end{pmatrix} \begin{pmatrix} 1 & -1 & 0.5 & 2 \end{pmatrix}.$$

Calculating each term:

$$\frac{\partial \text{Loss}}{\partial \mathbf{W}_1} = \begin{pmatrix} 0.6672 & -0.6672 & 0.3336 & 1.3344 \\ -1.0614 & 1.0614 & -0.5307 & -2.1228 \\ 0.2855 & -0.2855 & 0.1428 & 0.5710 \end{pmatrix}. \quad (\mathbf{0.5 \ pt})$$

The gradient of the loss with respect to  $\mathbf{b}_1$  is:

$$\frac{\partial \text{Loss}}{\partial \mathbf{b}_1} = \boldsymbol{\delta}_1 = \begin{pmatrix} 0.6672 \\ -1.0614 \\ 0.2855 \end{pmatrix}. \quad (0.5 \text{ pt})$$

**3. Using the update rule:**

New parameter = Old parameter  $- \alpha \times$  Gradient.

Updated  $\mathbf{W}_2$ :

$$\begin{aligned} \mathbf{W}_2' &= \mathbf{W}_2 - 0.001 \times \frac{\partial \text{Loss}}{\partial \mathbf{W}_2} \\ \mathbf{W}_2' &\approx \begin{pmatrix} -0.2982 & 0.2005 & 0.1009 \\ 0.3982 & -0.5005 & 0.2991 \end{pmatrix}. \quad (1 \text{ pt}) \end{aligned}$$

Updated  $\mathbf{b}_2$ :

$$\mathbf{b}_2' \approx \begin{pmatrix} 0.0508 \\ -0.0508 \end{pmatrix}. \quad (1 \text{ pt})$$

Updated  $\mathbf{W}_1$ :

$$\mathbf{W}_1' \approx \begin{pmatrix} 0.0993 & -0.1993 & 0.2997 & 0.3987 \\ 0.5011 & -0.3011 & 0.1005 & -0.1979 \\ 0.3997 & 0.2003 & -0.5001 & 0.2994 \end{pmatrix}. \quad (1 \text{ pt})$$

Updated  $\mathbf{b}_1$ :

$$\mathbf{b}_1' \approx \begin{pmatrix} 0.0993 \\ -0.0989 \\ 0.0497 \end{pmatrix}. \quad (1 \text{ pt})$$