

HW4

15 Points

1 Instructions

You may choose whatever software (overleaf, Latex, Word, etc.) to write your answers. But make sure to save your file as PDF (the only acceptable format) and upload to Canvas by the deadline.

This homework is based on Neural Network lectures and primarily focuses on math and logic questions with a small programming component for the last question. To better familiarize yourself with the questions and the lecture materials, first go through the description section. The five HW problems are posted afterwards.

2 Description of Neural Networks

- Neural networks refer to a general class of models for making predictions on data. The key feature of neural networks is that the input data is processed in multiple stages, called *layers*. The parameters in each layer are *trained* or *learned* from examples. This multi-layer processing with trainable parameters in each layer is loosely inspired by biological neural systems.
- To describe the neural network, let \mathbf{x} denote the input vector and let y denote the target variable that we wish to predict from \mathbf{x} .
- We consider networks for both **regression** or **classification** problems:
 - For regression problems, y is a scalar or vector and is typically continuous-valued. In this case, the neural network produces a prediction \hat{y} of y of the same dimension as y .
 - For classification problems, the target variable $y \in \{1, \dots, K\}$ and the neural network typically provides a soft prediction of the class label. Specifically, the networks makes a prediction of the probability $P(y = k|\mathbf{x})$ for each class label k given the input \mathbf{x} .
- For both regression and classification problems, we assume the input \mathbf{x} is a vector of dimension N_i so that

$$\mathbf{x} = (x_1, \dots, x_{N_i}).$$

- In this note, we look at neural networks with one hidden layer. In such a network, the neural network mapping is performed in two steps – each step is called a *layer*:
 - *Hidden layer* produces outputs \mathbf{z}^h and \mathbf{u}^h of dimension N_h .

- *Output layer* produces outputs \mathbf{z}^o and \mathbf{u}^o of dimension N_o .

The dimensions N_h and N_o will be discussed below.

- The equations for the two layers with a single input \mathbf{x} are:

$$\text{Hidden layer: } z_j^H = \sum_{k=1}^{N_i} W_{jk}^H x_k + b_j^H, \quad u_j^H = g_{\text{act}}(z_j^H), \quad j = 1, \dots, N_h \quad (1a)$$

$$\text{Output layer: } z_j^O = \sum_{k=1}^{N_h} W_{jk}^O u_k^H + b_j^O, \quad u^O = g_{\text{out}}(\mathbf{z}^O). \quad j = 1, \dots, N_o. \quad (1b)$$

- In the hidden layer, the function $g_{\text{act}}(z)$ is called the *activation function*. There are three common choices:

- Hard threshold:

$$g_{\text{act}}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{if } z < 0. \end{cases} \quad (2)$$

- Sigmoid: $g_{\text{act}}(z) = 1/(1 + e^{-z})$.

- Rectified linear unit (ReLU): $g(z) = \max\{0, z\}$.

- The output map $g_{\text{out}}(z)$ and dimension N_o depends on the type of prediction problem we are using the neural network for:

- Binary classification: In this case, the response is $y = \{0, 1\}$. To predict the response, we typically take $N_o = 1$ and the output u^o is the class probability:

$$P(y = 1|\mathbf{x}) = u^O = g_{\text{out}}(z^O) = \frac{1}{1 + e^{-z^O}}. \quad (3)$$

The variable z^O is called the *logit* and the output mapping (3) is a sigmoid. We can also use z^O to make a hard decision by selecting the most likely class:

$$\hat{y} = \begin{cases} 1 & z^O \geq 0 \\ 0 & z^O < 0. \end{cases} \quad (4)$$

- K -class classification: In this case, the target is a class label $y = 1, \dots, K$. We typically take $N_o = K$ and use the soft-max function for the class probability:

$$P(y = k|\mathbf{x}) = u_k^O = g_{\text{out},k}(z^O) = \frac{e^{z_k^O}}{\sum_{\ell=1}^K e^{z_\ell^O}}. \quad (5)$$

Again, we can make a hard decision on the class label by selecting the highest logit score:

$$\hat{y} = \arg \max_{k=1, \dots, K} z_k^O. \quad (6)$$

- Regression: In this case, one is trying to predict $\mathbf{y} \in \mathbb{R}^d$, where d is the number of variables to predict and each component y_j is typically continuous-valued. For this problem, we take $N_o = d$ and \mathbf{u}^o is the prediction of \mathbf{y} ,

$$\hat{\mathbf{y}} = \mathbf{u}^o = g_{\text{out}}(\mathbf{z}^o) = \mathbf{z}^o. \quad (7)$$

In (7), we will either say there is no activation or an *identity* activation.

- The number, N_h of hidden variables (also called the number of *hidden units*) represents the model complexity. Higher values of N_h can express more complex mappings, but will need more data to train.
- When $N_o = 1$ (single output unit), we drop the subscript j in (1b) and write

$$\text{Output layer: } z^o = \sum_{k=1}^{N_h} W_k^H u_k^H + b^o, \quad \hat{y} = g_{\text{out}}(z^o). \quad (8)$$

- Matrix form of the neural network with single sample input \mathbf{x} :

$$\begin{aligned} \mathbf{z}^H &= \mathbf{W}^H \mathbf{x} + \mathbf{b}^H, & \mathbf{u}^H &= g_{\text{act}}(\mathbf{z}^H), \\ \mathbf{z}^O &= \mathbf{W}^O \mathbf{u}^H + \mathbf{b}^O, & \mathbf{u}^O &= g_{\text{out}}(\mathbf{z}^O). \end{aligned}$$

- Batch form: Often we have a batch of input-output samples, (\mathbf{x}_i, y_i) , $i = 1, \dots, N$ (e.g. a mini-batch in training or test). In this case, we need to add an index i to each sample and rewrite (1) as:

$$\text{Hidden layer: } z_{ij}^H = \sum_{k=1}^{N_h} W_{jk}^H x_{ik} + b_j^H, \quad u_{ij}^H = g_{\text{act}}(z_{ij}^H), \quad j = 1, \dots, N_h \quad (9a)$$

$$\text{Output layer: } z_{ij}^O = \sum_{k=1}^{N_h} W_{jk}^O u_{ik}^H + b_j^O, \quad u_j^O = g_{\text{out}}(z_i^O), \quad j = 1, \dots, N_o. \quad (9b)$$

- Matrix form of batch processing in neural networks: Define the matrices,

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1,N_h} \\ \vdots & \vdots & \vdots \\ x_{N1} & \cdots & x_{N,N_h} \end{bmatrix}, \quad \mathbf{Z}^H = \begin{bmatrix} (\mathbf{z}_1^H)^\top \\ \vdots \\ (\mathbf{z}_N^H)^\top \end{bmatrix} = \begin{bmatrix} z_{11}^H & \cdots & z_{1,N_h}^H \\ \vdots & \vdots & \vdots \\ z_{N1}^H & \cdots & z_{N,N_h}^H \end{bmatrix},$$

Similarly define

$$\mathbf{U}^H = \begin{bmatrix} (\mathbf{u}_1^H)^\top \\ \vdots \\ (\mathbf{u}_N^H)^\top \end{bmatrix} = \begin{bmatrix} u_{11}^H & \cdots & u_{1,N_h}^H \\ \vdots & \vdots & \vdots \\ u_{N1}^H & \cdots & u_{N,N_h}^H \end{bmatrix}, \quad \mathbf{U}^O = \begin{bmatrix} (\mathbf{u}_1^O)^\top \\ \vdots \\ (\mathbf{u}_N^O)^\top \end{bmatrix} = \begin{bmatrix} u_{11}^O & \cdots & u_{1,N_o}^O \\ \vdots & \vdots & \vdots \\ u_{N1}^O & \cdots & u_{N,N_o}^O \end{bmatrix}.$$

Hence each row contains all the components for the i -th sample. Then, we can write

$$\mathbf{Z}^H = \mathbf{X} [\mathbf{W}^H]^\top + \mathbf{B}^H, \quad \mathbf{U}^H = g_{\text{act}}(\mathbf{Z}^H) \quad (10a)$$

$$\mathbf{Z}^O = \mathbf{U}^H [\mathbf{W}^O]^\top + \mathbf{B}^O, \quad \mathbf{U}^O = g_{\text{out}}(\mathbf{Z}^O), \quad (10b)$$

where \mathbf{B}^H and \mathbf{B}^O repeat the bias vectors on the N rows

$$\mathbf{B}^H = \begin{bmatrix} (\mathbf{b}^H)^\top \\ \vdots \\ (\mathbf{b}^H)^\top \end{bmatrix} \quad \mathbf{B}^O = \begin{bmatrix} (\mathbf{b}^O)^\top \\ \vdots \\ (\mathbf{b}^O)^\top \end{bmatrix} \quad (11)$$

Problems

1. Consider the neural network (1) with a scalar input x and parameters

$$W^H = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad b^H = \begin{bmatrix} -1 \\ -3 \end{bmatrix} \quad W^O = [-1, 2], \quad b^O = 0.5,$$

using a hard threshold activation function (2) and threshold output function (4).

- (a) What is N_h , the number of hidden units? What is N_o , the number of output units?
 - (b) Write \mathbf{z}^H in terms of x . Show the functions for each component z_j^H .
 - (c) Write \mathbf{u}^H in terms of x . Show the functions for each component u_j^H .
 - (d) Write z^O in terms of x .
 - (e) What is \hat{y} in terms of x ?
2. Consider the data set for four points with scalar features x_i and binary class labels $y_i = 0, 1$.

x_i	0	1	3	5
y_i	0	0	1	0

- (a) Find a neural network with $N_h = 2$ units, $N_o = 1$ output units such that $\hat{y}_i = y_i$ on all four data points. Use a network similar in structure to the previous problem. Also, you want to find features that can distinguish between $x = 3$ and $x = \{0, 1, 5\}$. Since there are many features, use two features: whether $x \geq 2$ and $x \geq 4$. Use a hard threshold activation function (2) and threshold output function (4). State the weights and biases used in each layer.
 - (b) Compute the values of \hat{y}_i and all the intermediate variables \mathbf{z}_i^H , \mathbf{u}_i^H and z_i^O for each sample $x = x_i$.
 - (c) Now suppose we are given a new sample, $x = 3.5$. What does the network predict as \hat{y} ?
3. *Two-dimensional example:* Consider a neural network on a 2-dimensional input $\mathbf{x} = (x_1, x_2)$ with weights and biases:

$$W^H = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}, \quad b^H = \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix} \quad W^O = [1, 1, -1], \quad b^O = -1.5.$$

Assume the network uses hard a threshold activation function (2) and threshold output function (4).

- (a) Write the components of \mathbf{z}^H and \mathbf{u}^H as a function of (x_1, x_2) . For each component j , indicate where in the (x_1, x_2) plane $u_j^H = 1$.
- (b) Write z^O as a function of (x_1, x_2) . In what region is $\hat{y} = 1$?
4. *Architecture choices for different problems:* For each problem, state possible selections for the dimensions N_i , N_h , N_o and the functions $g_{\text{act}}(\cdot)$ and $g_{\text{out}}(\cdot)$. Indicate which parameters are free to choose.
- (a) One wants a neural network to take as an input a 20×20 gray scale image and determine which letter ('a' to 'z') the image is of.
- (b) One extracts 120 features of a sample of a speech recording (like the MFCCs). Based on the audio samples, the network is to determine if the speech is male or female.
- (c) One wants a neural network to predict the stock price based on the average stock price of the last five days.
5. *Implementation in python:* Write python code for implementing the following steps for a batch of samples:
- (a) The hidden layer step (10a).
- (b) The output layer step (10b) for binary classification with a sigmoid output (3).
- (c) The output layer step (10b) for K -class classification with a softmax output (5).

For all examples, avoid for-loops and instead use Python broadcasting.

Solutions

1. (a) Since \mathbf{W}^H has 2 outputs, $N_h = 2$. Since W^O has 1 output, $N_o = 1$.
- (b) The components of the linear variables \mathbf{z}^H in the hidden layer are:

$$\mathbf{z}^H = \mathbf{W}^H \mathbf{x} + \mathbf{b}^H = \begin{bmatrix} 1 \\ 1 \end{bmatrix} x + \begin{bmatrix} -1 \\ -3 \end{bmatrix} = \begin{bmatrix} x - 1 \\ x - 3 \end{bmatrix}.$$

- (c) The activation outputs in the hidden layer are

$$\mathbf{u}^H = g_{\text{act}}(\mathbf{z}^H) = \begin{bmatrix} g_{\text{act}}(x - 1) \\ g_{\text{act}}(x - 3) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\{x \geq 1\}} \\ \mathbb{1}_{\{x \geq 3\}} \end{bmatrix}.$$

- (d) The linear variable in the output layer is

$$\begin{aligned} z^O &= W^O \mathbf{u}^H + b^O = [-1, 2] \begin{bmatrix} \mathbb{1}_{\{x \geq 1\}} \\ \mathbb{1}_{\{x \geq 3\}} \end{bmatrix} + 0.5 \\ &= -\mathbb{1}_{\{x \geq 1\}} + 2\mathbb{1}_{\{x \geq 3\}} + 0.5 = \begin{cases} 0.5 & \text{when } x < 1 \\ -0.5 & \text{when } x \in [1, 3) \\ 1.5 & \text{when } x \geq 3 \end{cases} \end{aligned}$$

(e) The predicted output \hat{y} is

$$\hat{y} = \mathbb{1}_{\{z^O \geq 0\}} = \begin{cases} 0 & \text{if } x \in [1, 3) \\ 1 & \text{else} \end{cases}$$

2. (a) We can use a network similar in structure to the previous problem. In the hidden layer, we want to find features that can distinguish between $x = 3$ which belongs to class $y = 1$, and $x = \{0, 1, 5\}$ which belong to class $y = 0$. There are lots of choices, but we will extract two features: whether $x \geq 2$ and $x \geq 4$. We can do this with the linear transforms:

$$\mathbf{W}^H = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad b^H = \begin{bmatrix} -2 \\ -4 \end{bmatrix},$$

so that

$$\mathbf{z}^H = \mathbf{W}^H \mathbf{x} + \mathbf{b}^H = \begin{bmatrix} 1 \\ 1 \end{bmatrix} x + \begin{bmatrix} -2 \\ -4 \end{bmatrix} = \begin{bmatrix} x - 2 \\ x - 4 \end{bmatrix}.$$

Then, the activation outputs in the hidden layer are

$$\mathbf{u}^H = g_{\text{act}}(\mathbf{z}^H) = \begin{bmatrix} g_{\text{act}}(x - 2) \\ g_{\text{act}}(x - 4) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\{x \geq 2\}} \\ \mathbb{1}_{\{x \geq 4\}} \end{bmatrix}.$$

Now we need to combine these functions in a way that they are negative for samples x_i when $y_i = 0$ and positive on sample x_i when $y_i = 1$. Similar to the previous problem, take

$$W^O = [1, -2], \quad b^O = -0.5.$$

Then,

$$\begin{aligned} z^O &= W^O \mathbf{u}^H + b^O = [1, -2] \begin{bmatrix} \mathbb{1}_{\{x \geq 2\}} \\ \mathbb{1}_{\{x \geq 4\}} \end{bmatrix} - 0.5 \\ &= \mathbb{1}_{\{x \geq 2\}} - 2\mathbb{1}_{\{x \geq 4\}} - 0.5 = \begin{cases} -0.5 & \text{when } x < 2 \\ 0.5 & \text{when } x \in [2, 4) \\ -1.5 & \text{when } x \geq 4. \end{cases} \end{aligned}$$

The predicted output \hat{y} is

$$\hat{y} = \mathbb{1}_{\{z^O \geq 0\}} = \begin{cases} 1 & \text{if } x \in [2, 4) \\ 0 & \text{else} \end{cases}$$

This matches the training data.

- (b) Table 1 computes the values for all intermediate variables for the four training data points. We can see that $\hat{y}_i = y_i$ for all four points.
- (c) The value for $x = 3.5$ is shown in the final column of Table 1. For this value of $x = 3.5$, we get $\hat{y} = 1$.

	Training				Test
x_i	0	1	3	5	3.5
$z_{i1}^H = x - 2$	-2	-1	1	3	1.5
$z_{i2}^H = x - 2$	-4	-3	-1	1	0.5
$u_{i1}^H = g_{\text{act}}(z_{i1}^H)$	0	0	1	1	1
$u_{i1}^H = g_{\text{act}}(z_{i2}^H)$	0	0	0	1	0
$z_i^O = u_{i1}^H - 2u_{i1}^H - 0.5$	-0.5	-0.5	0.5	-1.5	0.5
$\hat{y}_i = g_{\text{out}}(z_i^O)$	0	0	1	0	1
y_i	0	0	1	0	

Table 1: Computations for the neural network in Problem 2

3. (a) The linear functions in the hidden layer are:

$$\mathbf{z}^H = \mathbf{W}^H \mathbf{x} + \mathbf{b}^H = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ -1 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_1 + x_2 - 1 \end{bmatrix}$$

Hence, the activation functions are

$$\mathbf{u}^H = g_{\text{act}}(\mathbf{z}^H) = \begin{bmatrix} g_{\text{act}}(x_1) \\ g_{\text{act}}(x_2) \\ g_{\text{act}}(x_1 + x_2 - 1) \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{\{x_1 \geq 0\}} \\ \mathbb{1}_{\{x_2 \geq 0\}} \\ \mathbb{1}_{\{x_1 + x_2 \geq 1\}} \end{bmatrix}.$$

- (b) The output z^O is

$$\begin{aligned} z^O &= W^O \mathbf{u}^H + b^O = [1, 1, -1] \begin{bmatrix} \mathbb{1}_{\{x_1 \geq 0\}} \\ \mathbb{1}_{\{x_2 \geq 0\}} \\ \mathbb{1}_{\{x_1 + x_2 \geq 1\}} \end{bmatrix} - 1.5 \\ &= \mathbb{1}_{\{x_1 \geq 0\}} + \mathbb{1}_{\{x_2 \geq 0\}} - \mathbb{1}_{\{x_1 + x_2 \geq 1\}} - 1.5. \end{aligned}$$

It is best to draw this. If you do that (I will add the figures later), you will see that in the triangle defined by the lines

$$x_1 \geq 0, \quad x_2 \geq 0, \quad x_1 + x_2 < 1, \tag{12}$$

we have

$$z^O = (1) + (1) - (0) - 1.5 = 0.5.$$

Outside this region, $z^O < 0$. Hence, $z^O \geq 0$ only in the region (12). Therefore,

$$\hat{y} = \mathbb{1}_{\{z^O \geq 0\}} = \begin{cases} 1 & \text{if } x_1 \geq 0, x_2 \geq 0, x_1 + x_2 < 1 \\ 0 & \text{else} \end{cases}$$

4. (a) Assuming we represent the image pixels as a vector, the input size is $N_i = 20 \times 20 = 400$. Since this is a multi-class classification problem with 26 classes, one would generally take an output size of $N_o = 26$ and use a soft-max for the output map $g_{\text{out}}(\cdot)$. The number of hidden units N_h and the activation $g_{\text{act}}(\cdot)$ are free to choose

- (b) The input size is $N_i = 120$. Since this is a binary classification problem, we would generally take $N_o = 1$ and use a sigmoid output for $g_{\text{out}}(\cdot)$. The number of hidden units N_h and the activation $g_{\text{act}}(\cdot)$ are free to choose.
 - (c) We would take $N_i = 5$ for the stock prices for the last five days. Since this is a regression problem predicting a scalar value, we would take $N_o = 1$ and use an identity activation, $g_{\text{out}}(z) = z$ (i.e. there is no activation). Again, the number of hidden units N_h and the activation $g_{\text{act}}(\cdot)$ are free to choose.
5. (a) The operations in the hidden layer in (10) can be implemented as:

```
Zh = X.dot(Wh.T) + bh[None,:]
Uh = 1/(1+np.exp(-Zh))
```

where **Wh** and **bh** are the weight matrix and bias vector. Note that we have used python broadcasting on the bias vector **bh**

- (b) For a binary classification problem, $N_o = 1$ and we can represent \mathbf{Z}^o as an N -dimensional vector, **zo**. Also, we can represent the weight \mathbf{W}^o as an N_h -dimensional vector **wo** and the bias b^o as a scalar **bo**. In this case,

```
Zo = Uh.dot(Wo.T) + bo
Uo = 1/(1+np.exp(-Zo))
```

Note that when you add a scalar to a vector, Python automatically broadcasts the scalar to match the vector size.

- (c) For K -class classification, we perform

```
Zo = Uh.dot(Wo.T) + bo[None,:]
expZo = np.exp(Zo)
Uo = expZo / np.sum(expZo,axis=1)[: ,None]
```

Note the method for normalizing the outputs in a softmax classifier.