

# Homework 1: Backpropagation

CSCI-GA 2572 Deep Learning

Spring 2022

The goal of homework 1 is to help you understand the common techniques used in Deep Learning and how to update network parameters by the using backpropagation algorithm.

Part 1 has two sub-parts, 1.1, 1.2, 1.3 majorly deal with the theory of backpropagation algorithm whereas 1.4 is to test conceptual knowledge on deep learning. For part 1.2 and 1.3, you need to answer the questions with mathematical equations. You should put all your answers in a PDF file and we will not accept any scanned hand-written answers. It is recommended to use  $\LaTeX$ .

For part 2, you need to program in Python. It requires you to implement your own forward and backward pass without using autograd. You need to submit your `mlp.py` file for this part.

The due date of homework 1 is 23:55 EST of 02/17. Submit the following files in a zip file `your_net_id.zip` through NYU Brightspace:

- `theory.pdf`
- `mlp.py`

The following behaviors will result in penalty of your final score:

1. 5% penalty for submitting your files without using the correct format. (including naming the zip file, PDF file or python file wrong, or adding extra files in the zip folder, like the testing scripts from part 2).
2. 20% penalty for late submission within the first 24 hours. We will not accept any late submission after the first 24 hours.
3. 20% penalty for code submission that cannot be executed using the steps we mentioned in part 2. So please test your code before submit it.

# 1 Theory (50pt)

To answer questions in this part, you need some basic knowledge of linear algebra and matrix calculus. Also, you need to follow the instructions:

1. Every vector is treated as column vector.
2. You need to use the numerator-layout notation for matrix calculus. Please refer to [Wikipedia](#) about the notation.
3. You are only allowed to use vector and matrix. You cannot use tensor in any of your answer.
4. Missing transpose are considered as wrong answer.

## 1.1 Two-Layer Neural Nets

You are given the following neural net architecture:

$$\text{Linear}_1 \rightarrow f \rightarrow \text{Linear}_2 \rightarrow g$$

where  $\text{Linear}_i(\mathbf{x}) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$  is the  $i$ -th affine transformation, and  $f, g$  are element-wise nonlinear activation functions. When an input  $\mathbf{x} \in \mathbb{R}^n$  is fed to the network,  $\hat{\mathbf{y}} \in \mathbb{R}^K$  is obtained as the output.

## 1.2 Regression Task

We would like to perform regression task. We choose  $f(\cdot) = (\cdot)^+ = \text{ReLU}(\cdot)$  and  $g$  to be the identity function. To train this network, we choose MSE loss function  $\ell_{\text{MSE}}(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$ , where  $\mathbf{y}$  is the target output.

- (1pt) Name and mathematically describe the 5 programming steps you would take to train this model with PyTorch using SGD on a single batch of data.
- (5pt) For a single data point  $(x, y)$ , write down all inputs and outputs for forward pass of each layer. You can only use variable  $\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}$  in your answer. (note that  $\text{Linear}_i(\mathbf{x}) = \mathbf{W}^{(i)}\mathbf{x} + \mathbf{b}^{(i)}$ ).

Layer	Input	Output
Linear <sub>1</sub>		
$f$		
Linear <sub>2</sub>		
$g$		
Loss		

- (c) (8pt) Write down the gradient calculated from the backward pass. You can only use the following variables:  $\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \frac{\partial \ell}{\partial \hat{\mathbf{y}}}, \frac{\partial z_2}{\partial \hat{\mathbf{y}}}, \frac{\partial \hat{\mathbf{y}}}{\partial z_3}$  in your answer, where  $z_1, z_2, z_3, \hat{\mathbf{y}}$  are the outputs of  $\text{Linear}_1, f, \text{Linear}_2, g$ .

Parameter	Gradient
$\mathbf{W}^{(1)}$	
$\mathbf{b}^{(1)}$	
$\mathbf{W}^{(2)}$	
$\mathbf{b}^{(2)}$	

- (d) (3pt) Show us the elements of  $\frac{\partial z_2}{\partial z_1}, \frac{\partial \hat{\mathbf{y}}}{\partial z_3}$  and  $\frac{\partial \ell}{\partial \hat{\mathbf{y}}}$  (be careful about the dimensionality)?

**Solution:**

(a)

1. Generate a prediction ( $\hat{\mathbf{y}} = \text{model}(\mathbf{x})$ )
2. Compute the loss function (loss = criterion( $\hat{\mathbf{y}}, \mathbf{y}$ ))
3. zero  $\nabla \text{params}$ , clear the gradients of all parameters. (optimizer.zero\_grad())
4. Compute and Accumulate  $\nabla \text{params}$ : Backpropagation to compute the gradients of the loss function with respect to the parameters. (loss.backward())
5. Step in towards  $-\nabla \text{params}$ : Update the parameters and update the optimizer status. (optimizer.step())

(b)

$\mathbf{x}, \mathbf{y}, \mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}$

For a single data point  $(\mathbf{x}, \mathbf{y})$ , we have:

Layer	Input	Output
$\text{Linear}_1$	$\mathbf{x}$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
$f$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$	$f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
$\text{Linear}_2$	$f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$	$\mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$
$g$	$\mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$	$\mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$
Loss	$\mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$	$\ \mathbf{W}^{(2)}f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)} - \mathbf{y}\ ^2$

The input dim is n, the output dim is k, we define the output dim of middle layer as L1, L2=K. We have B = batch-size = 1 here, so we have:

$\mathbf{x}.\text{shape} = (\text{feature-in}, \text{batch-size}) = (n, B)$

$\mathbf{W1}.\text{shape} = (L1, n)$

$\mathbf{b1}.\text{shape} = (L1)$

$\mathbf{z1}.\text{shape} = (\mathbf{W1} \mathbf{x} + \mathbf{b1}).\text{shape} = (L1, B)$

$\mathbf{z2}.\text{shape} = f(\mathbf{z1}).\text{shape} = (L1, B)$

$\mathbf{W2}.\text{shape} = (K, L1)$

$\mathbf{b2}.\text{shape} = (K)$

$\mathbf{z3}.\text{shape} = (\mathbf{W2} \mathbf{x} + \mathbf{b2}).\text{shape} = (K, B)$

$\mathbf{y\_hat}.\text{shape} = g(\mathbf{z3}).\text{shape} = (K, B)$

(c)

Parameter	Gradient
$\mathbf{W}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \mathbf{x}^T$
$\mathbf{b}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$
$\mathbf{W}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} (f(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}))^T$
$\mathbf{b}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}}$

**We use matrix multiplication in the activation gradient backpropagation, the gradient of activation function is a diagonal matrix.**

For  $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}}$ , which is equivalent to  $\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \cdot f'(z_3)$ , which is

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \cdot g'(z_3) = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \cdot \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad (1)$$

$$\equiv \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} = I_K \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \quad (2)$$

$$\equiv \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \begin{bmatrix} 1 & 0 & \vdots & 0 \\ 0 & 1 & \vdots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \vdots & 1 \end{bmatrix} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \quad (3)$$

$$\mathbf{y} = [y_1, y_2, \dots, y_n]^T \quad (4)$$

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = 2(\hat{\mathbf{y}} - \mathbf{y}) \quad (5)$$

$$\mathbf{y} \in \mathbb{R}^{K \times 1}$$

$$\frac{\partial \ell}{\partial \mathbf{z}_3} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} = 2(\hat{\mathbf{y}} - \mathbf{y}) \quad (6)$$

$$\frac{\partial \ell}{\partial \mathbf{W}^{(2)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{W}^{(2)}} \quad (7)$$

$$= \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \mathbf{z}_2^T \quad (8)$$

$$= \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} (f(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}))^T \quad (9)$$

$$(10)$$

$$\frac{\partial \ell}{\partial \mathbf{b}^{(2)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{b}^{(2)}} \quad (11)$$

$$= \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \quad (12)$$

$$\frac{\partial \ell}{\partial \mathbf{W}^{(1)}} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}^{(1)}} \quad (13)$$

$$= \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \mathbf{x}^T \quad (14)$$

$$(15)$$

$$\hat{\mathbf{y}} \in \mathbb{R}^{K \times 1}, \mathbf{z}_3 \in \mathbb{R}^{K \times 1}, \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \in \mathbb{R}^{K \times K}.$$

$$\frac{\partial \ell}{\partial \mathbf{b}^{(1)}} = (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}}) \frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_2} \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \frac{\partial \mathbf{z}_1}{\partial \mathbf{b}^{(1)}} = \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \quad (16)$$

(d)

If we only consider one data point,

$$\mathbf{z}_1 = \mathbf{w}_1 * \mathbf{x} + \mathbf{b}_1$$

$$\mathbf{z}_2 = f(\mathbf{z}_1)$$

$$\mathbf{z}_3 = \mathbf{W}_2 * \mathbf{z}_2 + \mathbf{b}_2$$

$$\mathbf{y}_{\text{hat}} = g(\mathbf{z}_3)$$

$\mathbf{z}_2 \in \mathbb{R}^{L_1}$ ,  $\mathbf{z}_1 \in \mathbb{R}^{L_1}$ ,  $\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} \in \mathbb{R}^{L_1 \times L_1}$ , vector by vector derivative.

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} \frac{\partial z_2[1]}{\partial z_1[1]} & \frac{\partial z_2[2]}{\partial z_1[1]} & \dots & \frac{\partial z_2[L_1]}{\partial z_1[1]} \\ \frac{\partial z_2[1]}{\partial z_1[2]} & \frac{\partial z_2[2]}{\partial z_1[2]} & \dots & \frac{\partial z_2[L_1]}{\partial z_1[2]} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_2[1]}{\partial z_1[L_1]} & \frac{\partial z_2[2]}{\partial z_1[L_1]} & \dots & \frac{\partial z_2[L_1]}{\partial z_1[L_1]} \end{bmatrix} = \begin{bmatrix} \mathbb{1}_{z_{21}>0} & 0 & \dots & 0 \\ 0 & \mathbb{1}_{z_{22}>0} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \mathbb{1}_{z_{2L_1}>0} \end{bmatrix} \quad (17)$$

For the element in the matrix, if the corresponding element in the  $\mathbf{z}_2$  is greater than 0, then the element in the matrix is 1, otherwise it is 0.

$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$ , vector-by-vector derivative

$$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} = I_K \quad (18)$$

$\mathbf{z}_3 \in \mathbb{R}^K$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^K$ ,  $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \in \mathbb{R}^{K \times K}$ . Only the diagonal elements are one.

$\frac{\partial \ell}{\partial \hat{\mathbf{y}}}$ , scalar-by-vector derivative.

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \begin{bmatrix} \frac{\partial \ell}{\partial \hat{y}_1} & \frac{\partial \ell}{\partial \hat{y}_2} & \dots & \frac{\partial \ell}{\partial \hat{y}_K} \end{bmatrix}^T = 2(\hat{\mathbf{y}} - \mathbf{y}) = 2[\hat{y}_0 - y_0, \hat{y}_1 - y_1, \dots]^T \quad (19)$$

where  $\mathbf{y} = [y_0, y_1, \dots, y_K]^T$ .

$\ell \in \mathbb{R}$ ,  $\hat{\mathbf{y}} \in \mathbb{R}^K$ ,  $\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \in \mathbb{R}^K$ .

### 1.3 Classification Task

We would like to perform multi-class classification task, so we set both  $f, g = \sigma$ , the logistic sigmoid function  $\sigma(z) \doteq (1 + \exp(-z))^{-1}$ .

- (2pt + 6pt + 2pt) If you want to train this network, what do you need to change in the equations of (b), (c) and (d), assuming we are using the same MSE loss function.
- (2pt + 6pt + 2pt) Now you think you can do a better job by using a *Binary Cross Entropy* (BCE) loss function  $\ell_{\text{BCE}}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{K} \sum_{i=1}^K -[\hat{y}_i \log(\hat{y}_i) + (1 - \hat{y}_i) \log(1 - \hat{y}_i)]$ . What do you need to change in the equations of (b), (c) and (d)?
- (1pt) Things are getting better. You realize that not all intermediate hidden activations need to be binary (or soft version of binary). You decide to use  $f(\cdot) = (\cdot)^+$  but keep  $g$  as  $\sigma$ . Explain why this choice of  $f$  can be beneficial for training a (deeper) network.

**Solution:**

### 1.3.1 (a)

Layer	Input	Output
Linear <sub>1</sub>	$\mathbf{x}$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
$f$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
Linear <sub>2</sub>	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$
$g$	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
Loss	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$	$\ \sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}) - \mathbf{y}\ ^2$

Here we do not expand the  $\sigma$  to prevent the fomula from becoming too long. Forward table do not change much, only change the f and g to  $\sigma$

MSE Loss and f, g =  $\sigma$ . the gradient table do not changed, but the term  $\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$ ,  $\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$  are changed.

Parameter	Gradient
$\mathbf{W}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \mathbf{x}^T$
$\mathbf{b}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$
$\mathbf{W}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} (f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}))^T$
$\mathbf{b}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}}$

For the (d) part,

$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$  becomes element-wise  $\sigma(z_1)(1 - \sigma(z_1)) = z_2(1 - z_2)$

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} \frac{\partial z_2[1]}{\partial z_1[1]} & \frac{\partial z_2[2]}{\partial z_1[1]} & \dots & \frac{\partial z_2[L1]}{\partial z_1[1]} \\ \frac{\partial z_2[1]}{\partial z_1[2]} & \frac{\partial z_2[2]}{\partial z_1[2]} & \dots & \frac{\partial z_2[L1]}{\partial z_1[2]} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_2[1]}{\partial z_1[L1]} & \frac{\partial z_2[2]}{\partial z_1[L1]} & \dots & \frac{\partial z_2[L1]}{\partial z_1[L1]} \end{bmatrix} = \begin{bmatrix} (\mathbf{z}_{2[1]})(1 - (\mathbf{z}_{2[1]})) & 0 & \dots & 0 \\ 0 & (\mathbf{z}_{2[2]})(1 - (\mathbf{z}_{2[2]})) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & (\mathbf{z}_{2[L1]})(1 - (\mathbf{z}_{2[L1]})) \end{bmatrix} \quad (20)$$

$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$  is changed to element-wise  $\sigma(z_3)(1 - \sigma(z_3)) = \hat{y}(1 - \hat{y})$ .

$$\frac{\partial \hat{\mathbf{y}}}{\partial \hat{\mathbf{z}}_3} = \begin{bmatrix} \hat{y}_1(1 - \hat{y}_1) & 0 & \dots & 0 \\ 0 & \hat{y}_2(1 - \hat{y}_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \hat{y}_K(1 - \hat{y}_K) \end{bmatrix} \quad (21)$$

$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = 2(\hat{\mathbf{y}} - \mathbf{y})$  do not change.

### 1.3.2 (b)

BCE Loss and  $f, g = \sigma$ . the gradient table representations do not changed, but the terms  $\frac{\partial \ell}{\partial \hat{\mathbf{y}}}, \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}, \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$  are changed.

Layer	Input	Output
Linear <sub>1</sub>	$\mathbf{x}$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$
$f$	$\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$
Linear <sub>2</sub>	$\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$
$g$	$\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}$	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
Loss	$\sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$	$\frac{1}{K} \sum_{i=1}^K -[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$

Parameter	Gradient
$\mathbf{W}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \mathbf{x}^T$
$\mathbf{b}^{(1)}$	$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} ((\mathbf{W}^{(2)})^T (\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}})) \frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$
$\mathbf{W}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}} (f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}))^T$
$\mathbf{b}^{(2)}$	$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \frac{\partial \ell}{\partial \hat{\mathbf{y}}}$

where

$$\hat{\mathbf{y}} = \sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)}), \quad \mathbf{y} \in \mathbb{R}^K$$

$\frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3}$  is changed to the diagonal  $\sigma(z_3)(1 - \sigma(z_3))$ .

$$\frac{\partial \hat{\mathbf{y}}}{\partial \hat{\mathbf{z}}_3} = \begin{bmatrix} \hat{y}_1(1 - \hat{y}_1) & 0 & \cdots & 0 \\ 0 & \hat{y}_2(1 - \hat{y}_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \hat{y}_K(1 - \hat{y}_K) \end{bmatrix} \quad (22)$$

$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1}$  is changed to the diagonal  $\sigma(z_1)(1 - \sigma(z_1))$ .

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} (\mathbf{z}_{2[1]})(1 - (\mathbf{z}_{2[1]})) & 0 & \cdots & 0 \\ 0 & (\mathbf{z}_{2[2]})(1 - (\mathbf{z}_{2[2]})) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & (\mathbf{z}_{2[L1]})(1 - (\mathbf{z}_{2[L1]})) \end{bmatrix} \quad (23)$$

$$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} = \frac{1}{K} \left[ \frac{\hat{y}_1 - y_1}{\hat{y}_1 * (1 - \hat{y}_1)}, \frac{\hat{y}_2 - y_2}{\hat{y}_2 * (1 - \hat{y}_2)}, \dots, \frac{\hat{y}_K - y_K}{\hat{y}_K * (1 - \hat{y}_K)} \right]^T \quad (24)$$

$\frac{\partial \ell}{\partial \hat{\mathbf{y}}} \in \mathbb{R}^K$ .



$$\frac{\partial \ell}{\partial \mathbf{z}_3} = \frac{\partial \ell}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_3} \quad (25)$$

$$= \frac{1}{K} \begin{bmatrix} \hat{y}_1 - y_1 \\ \hat{y}_2 - y_2 \\ \vdots \\ \hat{y}_K - y_K \end{bmatrix} \quad (26)$$

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1},$$

$$\frac{\partial \mathbf{z}_2}{\partial \mathbf{z}_1} = \begin{bmatrix} z_{2_1}(1-z_{2_1}) & 0 & \cdots & 0 \\ 0 & z_{2_2}(1-z_{2_2}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & z_{2_K}(1-z_{2_K}) \end{bmatrix} \quad (27)$$

### 1.3.3 (c)

The advantage of using the ReLU over Sigmoid in the middle layers:

- Preventing Small Derivatives: When the value is too large or too small, the derivatives of sigmoid is close to 0, but relu is a non-saturated activation function. This phenomenon does not exist.
- Increasing the sparsity of the network: when the value is less than 0, the derivative of relu is 0, and the derivative of sigmoid is close to 0.

## 1.4 Conceptual Questions

- (1pt) Why is softmax actually  $\text{soft}(\arg)\max$ ?
- (1pt) In what situations,  $\text{soft}(\arg)\max$  can become unstable?
- (1pt) Should we have two consecutive linear layers in a neural network? Why or why not?
- (4pt) Can you draw the graph of the derivative for the following functions?
  - `ReLU()`
  - `LeakyReLU(negative_slope=0.01)`
  - `Softplus(beta=1)`
  - `GELU()`

- (e) (4pt) What are 4 different types of linear transformations? What is the role of linear transformation and non linear transformation in a neural network?
- (f) (1pt) How should we adjust the learning rate as we increase or decrease the batch size?

**Solution:**

(a)

Soft-{Function-Name} is the name of the softened version of the function.

The argmax function takes a vector as input and returns a one-hot vector of the maximum value.

The softmax (softargmax) function is a softened version of this function, which return a distribution but not the value of the maximum value.

(b)

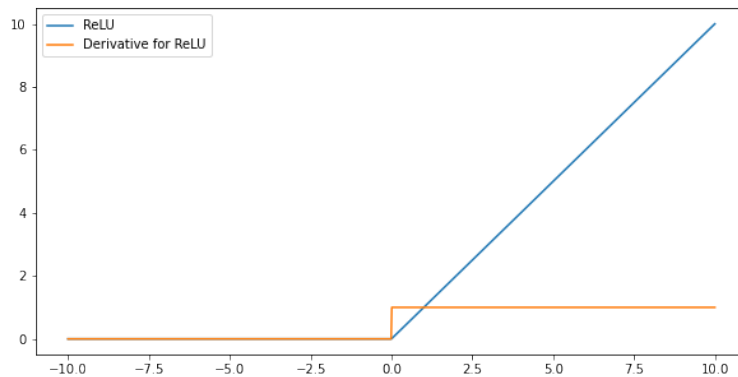
If the maximum value and the minimum value are in quite different scales, the softmax function will become unstable (numerically unstable for values with a very large range).

(c)

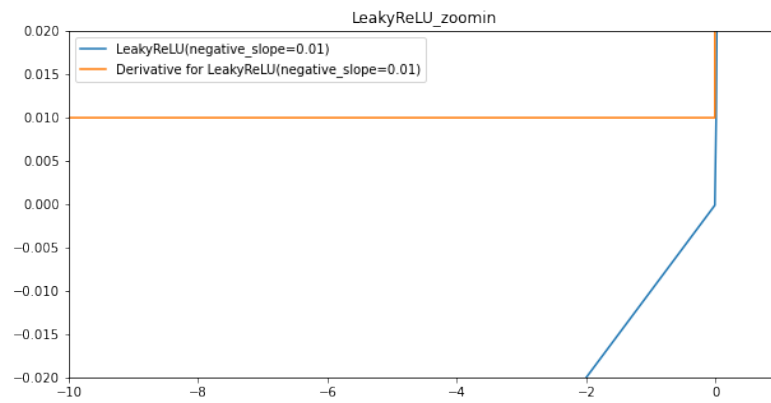
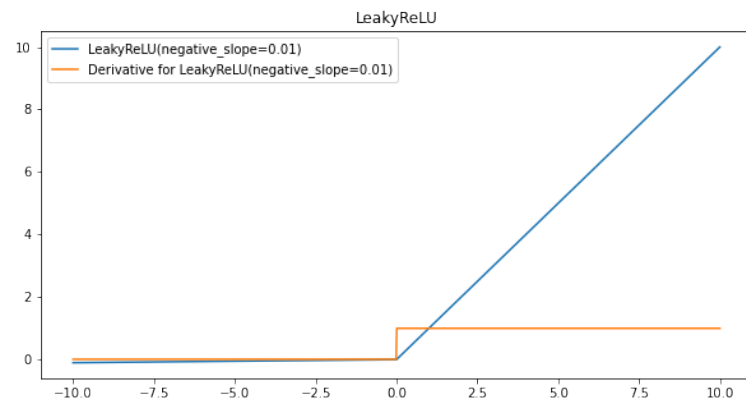
We should not have two consecutive linear layers in a neural network. Two consecutive linear layers can be transfered/merged to a single linear layer.

(d)

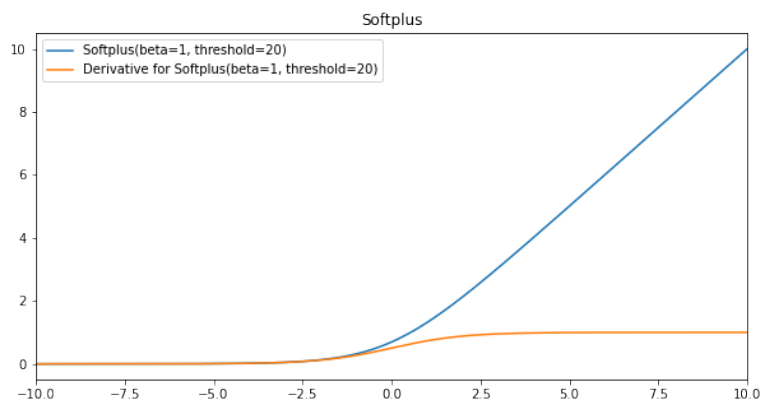
ReLU



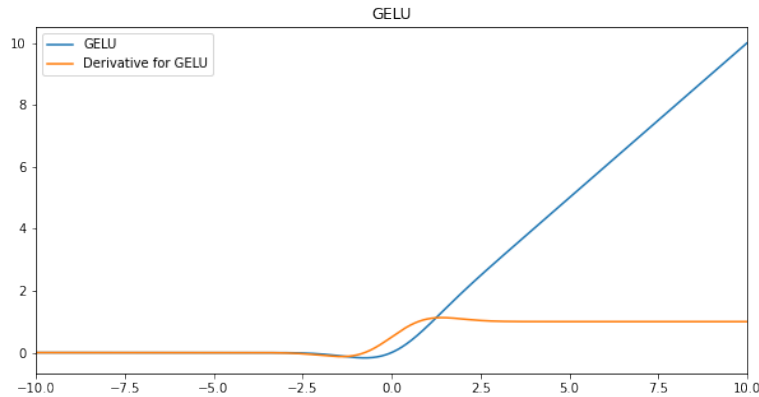
LeakyReLU



Softplus(beta=1)



## GELU



(d) The four types of linear transformations are: rotations, zooming, flipping, and shearing.

The linear transformation can remap or change the input space or representation space, which help us get better representation and feature mapping.

For non linear transformation, we do not want a combined linear layer if all the layers are linear. For increasing the non-linearity and solve the non-linear tasks, we should use a non-linear layer. From one perspective, nonlinear transformation can map the task to a linear seperable task, just like we usually take linear layer as last layer. (before logit, for classification, and we do not consider the softmax as a layer)

(e)

When we increasing the batch size, the learning rate should be decreased. When we decrease the batch size, the learning rate should be increased.

## 2 Implementation (50pt)

You need to implement the forward pass and backward pass for Linear, ReLU, Sigmoid, MSE loss, and BCE loss in the attached `mlp.py` file. We provide three example test cases `test1.py`, `test2.py`, `test3.py`. We will test your implementation with other hidden test cases, so please create your own test cases to make sure your implementation is correct.

**Recommendation:** Go through this [Pytorch tutorial](#) to have a thorough understanding of Tensors.

Extra instructions:

1. Please use Python version  $\geq 3.7$  and PyTorch version 1.7.1. We recommend

you to use Miniconda to manage your virtual environment.

2. We will put your `mlp.py` file under the same directory of the hidden test scripts and use the command `python hiddenTestScriptName.py` to check your implementation. So please make sure the file name is `mlp.py` and it can be executed with the example test scripts we provided.
3. You are not allowed to use PyTorch autograd functionality in your implementation.
4. Be careful about the dimensionality of the vector and matrix in PyTorch. It is not necessarily follow the the Math you got from part 1.