

# 705.641.81: Natural Language Processing

## Self-Supervised Models

### Homework 3: Building Your Neural Network!

For homework deadline and collaboration policy, please see our Canvas page.

Name: \_\_\_\_Breanna Burd\_\_\_\_

Collaborators, if any: \_\_\_\_\_N/A\_\_\_\_\_

Sources used for your homework, if any: \_\_\_\_\_

<https://akshay-a.medium.com/activation-functions-with-derivative-and-python-code-sigmoid-vs-tanh-vs-relu-1e9a5b5e3134>

<https://mmuratarat.github.io/2019-01-27/derivation-of-softmax-function>

This assignment is focusing on understanding the fundamental properties of neural networks and their training.

**Homework goals:** After completing this homework, you should be comfortable with:

- thinking about neural networks
- key implementation details of NNs, particularly in PyTorch,
- explaining and deriving Backpropagation,
- debugging your neural network in case it faces any failures.

## Concepts, intuitions and big picture

1. Suppose you have built a neural network. You decide to initialize the weights and biases to be zero. Which of the following statements are True? (Check all that apply)
  - Each neuron in the first hidden layer will perform the same computation. So even after multiple iterations of gradient descent each neuron will be computing the same thing as other neurons in the same layer.
    - Each neuron in the first hidden layer will perform the same computation in the first iteration. But after one iteration of gradient descent they will learn to compute different things because we have “broken symmetry”.
    - Each neuron in the first hidden layer will compute the same thing, but neurons in different layers will compute different things, thus we have accomplished “symmetry breaking” as described in lecture.
    - The first hidden layer’s neurons will perform different computations from each other even in the first iteration; their parameters will thus keep evolving in their own way.
2. Vectorization allows you to compute forward propagation in an  $L$ -layer neural network without an explicit for-loop (or any other explicit iterative loop) over the layers  $l = 1, 2, \dots, L$ . True/False?
  - True
  - False
3. The tanh activation usually works better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data better for the next layer. True/False?
  - True
  - False
4. Which of the following techniques does NOT prevent a model from overfitting?
  - Data augmentation ▫ Dropout ▫ Early stopping ■ None of the above
5. Why should dropout be applied during training? Why should dropout NOT be applied during evaluation?

During training, dropout should be applied to apply regularization by randomly “dropping” neurons. Dropout should not be applied during evaluation because the model should already be regularized and each neuron has its own weights that should be used in the computations for the prediction.

6. Explain why initializing the parameters of a neural net with a constant is a bad idea.  
Initializing the parameters with a constant causes a similar problem as #1 because each neuron in each layer would be performing the same calculations. This means the model would be performing many of the same calculations but no additional relationships would be learned, and the model would perform like there was only one neuron in each layer.
7. You design a fully connected neural network architecture where all activations are sigmoids. You initialize the weights with large positive numbers. Is this a good idea? Explain your answer.  
This is not a good idea because the sigmoid function tends to cause exploding/vanishing gradients. This case would most likely cause vanishing gradients since the sigmoid function converges at 0 and 1.
8. Explain what is the importance of “residual connections”.  
The importance of residual connections support the flow of gradients. Additionally, residual connections help prevent the vanishing gradients problem, allowing models to become deeper.
9. What is cached (“memoized”) in the implementation of forward propagation and backward propagation?
  - Variables computed during forward propagation are cached and passed on to the corresponding backward propagation step to compute derivatives.
    - Caching is used to keep track of the hyperparameters that we are searching over, to speed up computation.
    - Caching is used to pass variables computed during backward propagation to the corresponding forward propagation step. It contains useful values for forward propagation to compute activations.
10. Which of the following statements is true?
  - The deeper layers of a neural network are typically computing more complex features of the input than the earlier layers.
    - The earlier layers of a neural network are typically computing more complex features of the input than the deeper layers.

## Revisiting Jacobians

Recall that Jacobians are generalizations of multi-variate derivatives and are extremely useful in denoting the gradient computations in computation graph and Backpropagation. A potentially confusing aspect of using Jacobians is their dimensions and so, here we're going to focus on understanding Jacobian dimensions.

### Recap:

Let's first recap the formal definition of Jacobian. Suppose  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a function that takes a point  $\mathbf{x} \in \mathbb{R}^n$  as input and produces the vector  $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$  as output. Then the Jacobian matrix of  $\mathbf{f}$  is defined to be an  $m \times n$  matrix, denoted by  $\mathbf{J}_{\mathbf{f}}(\mathbf{x})$ , whose  $(i, j)$ th entry is  $\mathbf{J}_{ij} = \frac{\partial f_i}{\partial x_j}$ , or:

$$\mathbf{J} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial x_1} & \cdots & \frac{\partial \mathbf{f}}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla^\top f_1 \\ \vdots \\ \nabla^\top f_m \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

### Examples:

The shape of a Jacobian is an important notion to note. A Jacobian can be a vector, a matrix, or a tensor of arbitrary ranks. Consider the following special cases:

- If  $f$  is a scalar and  $\mathbf{w}$  is a  $d \times 1$  column vector, the Jacobian of  $f$  with respect to  $\mathbf{w}$  is a row vector with  $1 \times d$  dimensions.
- If  $\mathbf{y}$  is a  $n \times 1$  column vector and  $\mathbf{z}$  is a  $d \times 1$  column vector, the Jacobian of  $\mathbf{z}$  with respect to  $\mathbf{y}$ , or  $\mathbf{J}_{\mathbf{z}}(\mathbf{y})$  is a  $d \times n$  matrix.
- Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  and  $\mathbf{B} \in \mathbb{R}^{l \times p \times q}$ . Then the Jacobian  $\mathbf{J}_{\mathbf{A}}(\mathbf{B})$  is a tensor of shape  $(m \times n) \times (l \times p \times q)$ . More broadly, the shape of the Jacobian is determined as (shape of the output)  $\times$  (shape of the input).

### Problem setup:

Suppose we have:

- $\mathbf{X}$ , an  $n \times d$  matrix,  $x_i \in \mathbb{R}^{d \times 1}$  correspond to the rows of  $\mathbf{X} = [x_1, \dots, x_n]^\top$
- $\mathbf{Y}$ , a  $n \times k$  matrix
- $\mathbf{W}$ , a  $k \times d$  matrix and  $\mathbf{w}$ , a  $d \times 1$  vector

For the following items, compute (1) the shape of each Jacobian, and (2) an expression for each Jacobian:

1.  $f(\mathbf{w}) = c$  (constant)

1)  $\mathbf{J}_{\mathbf{f}}(\mathbf{w})$  shape is a  $1 \times d$  vector.

2)  $\mathbf{J} = \mathbf{0}_{1 \times d}$  (This is a  $1 \times d$  vector with all zeros because the function doesn't have any variables)

2.  $f(w) = \|w\|^2$  (squared L2-norm)  
 $f(w) = w^T w \rightarrow$  returns a scalar value  $\rightarrow (1 \times d) \times (d \times 1) = (1 \times 1)$  output  
  - 1)  $J_f(w)$  shape is a  $1 \times d$  vector.
  - 2)  $J = 2w^T$  (The derivative for each element is  $2w_j$ )
3.  $f(w) = w^T x_i$  (vector dot product)  
 $(1 \times d) \times (d \times 1) = (1 \times 1)$  output  
  - 1)  $J_f(w)$  shape is a  $1 \times d$  vector.
  - 2)  $J = x_i^T$  (The derivative for each element in the  $w$  vector is  $x_i$ )
4.  $f(w) = Xw$  (matrix-vector product)  
 $(n \times d) \times (d \times 1) = (n \times 1)$  output  
  - 1)  $J_f(w)$  shape is a  $n \times d$  matrix with respect to  $w$
  - 2)  $J = X$
5.  $f(w) = w$  (vector identity function)  
  - 1)  $J_f(w)$  shape is a  $d \times d$  vector.
  - 2)  $J = I_d$  (This is the identity matrix of diagonal 1s within a  $d \times d$  matrix)
6.  $f(w) = w^2$  (element-wise power)  
 $(d \times 1)$  output  
  - 1)  $J_f(w)$  shape is a  $d \times d$  matrix
  - 2)  $J = 2w$ , excluding the diagonal within the matrix where  $i=j$  (the diagonal values are 1)
7. **Extra Credit:**  $f(W) = XW^T$  (matrix multiplication)  
 $(n \times d) \times (d \times k) = (n \times k)$  output  
  - 1)  $J_f(W)$  shape is a  $(n \times k) \times (k \times d)$  matrix with respect to  $W$
  - 2) I'm not exactly sure what the equation would be for this Jacobian...

## Activations Per Layer, Keeps Linearity Away!

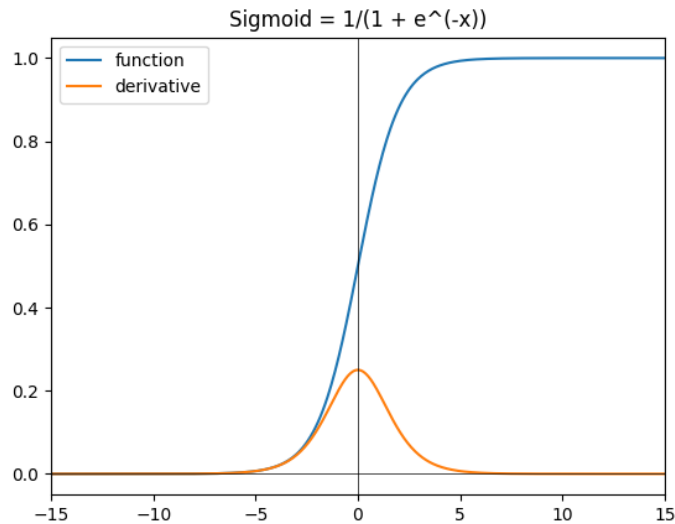
Based on the content we saw at the class lectures, answer the following:

1. Why are activation functions used in neural networks?

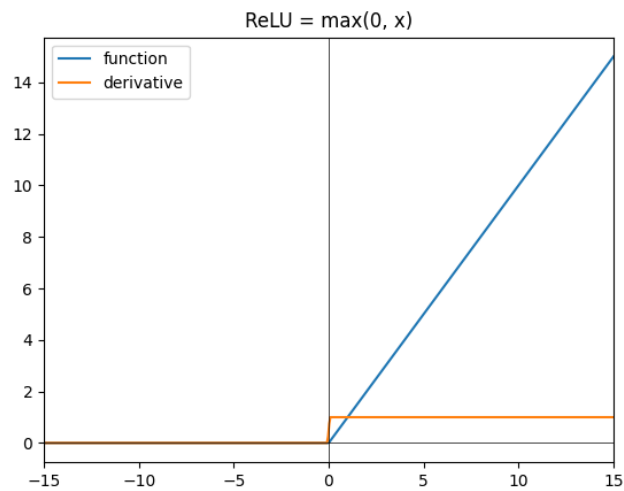
Activation functions are used to integrate non-linearity in the neural network to allow the model to learn complex relationships.

2. Write down the formula for three common action functions (sigmoid, ReLU, Tanh) and their derivatives (assume scalar input/output). Plot these activation functions and their derivatives on  $(-\infty, +\infty)$ .

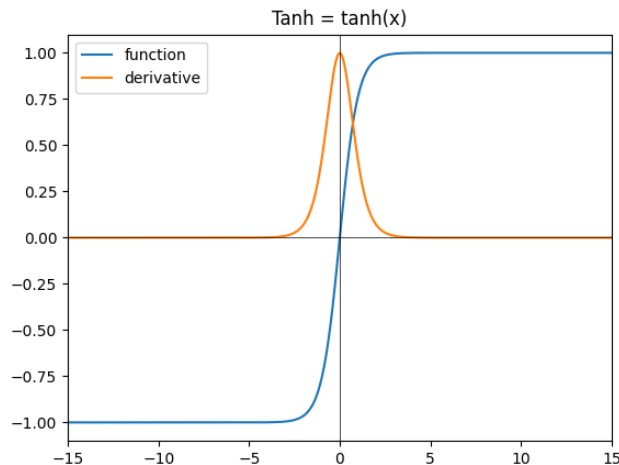
Sigmoid:  $f(x) = \frac{1}{1+e^{-x}}$   $f'(x) = f(x) * (1 - f(x))$



ReLU:  $f(x) = \max(0, x)$   $f'(x) = \begin{cases} 0, & x < 0 \\ x, & x > 0 \end{cases}$  (undefined at  $x=0$  but is often treated as 0 within the derivative in practice)



Tanh:  $f(x) = \tanh(x)$   $f'(x) = 1 - \tanh^2(x)$



**\*\*Note that I attempted to plot from  $[-\infty, \infty]$  but ran into errors and even from  $[-100, 100]$  it was difficult to see the derivative. I settled on  $[-15, 15]$  but we can see each of these functions and derivatives converge. The function plots match the ones provided in the Module 2 notes slide 19. \*\***

3. What is the “vanishing gradient” problem? (respond in no more than 3 sentences)  
Which activation functions are subject to this issue and why? (respond in no more than 3 sentences).

The vanishing gradient problem is when the derivatives during backpropagation fall to 0 for most/all of the weights, causing very little to be learned and poor model accuracy. The tanh and sigmoid activation functions are subject to this issue because the majority of the derivatives calculated in backpropagation fall to 0 (as shown in the graphs above, the derivatives are only non-zero between about  $[-5, 5]$ ).

4. Why zero-centered activation functions impact the results of Backprop?  
They impact the results of backprop because the calculations are able to capture more complex relationships by allowing both positive and negative gradients, which can also lead to faster convergence.
5. Remember the Softmax function  $\sigma(\mathbf{z})$  and how it extends sigmoid to multiple dimensions? Let’s compute the derivative of Softmax for each dimension. Prove that:

$$\frac{d\sigma(\mathbf{z})_i}{dz_j} = \sigma(\mathbf{z})_i (\delta_{ij} - \sigma(\mathbf{z})_j)$$

where  $\delta_{ij}$  is the Kronecker delta function.<sup>1</sup>

---

<sup>1</sup> [https://en.wikipedia.org/wiki/Kronecker\\_delta](https://en.wikipedia.org/wiki/Kronecker_delta)

Softmax:  $\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$  for  $i=1, \dots, K$

$\frac{\partial \sigma_i}{\partial z_j} = \frac{\partial}{\partial z_j} \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$  then use the quotient rule for  $f(x) = g(x)/h(x)$  where  $g$  is the numerator ( $g_i = e^{z_i}$ ) and  $h$  is the denominator ( $h_i = \sum_{j=1}^K e^{z_j}$ ).

$\frac{\partial h_i}{\partial z_j} = e^{z_j}$  for  $k=j$  and 0 otherwise. Thus, applying the quotient rule gives us:

$$\frac{e^{z_i} * \sum_{j=1}^K e^{z_k} - e^{z_i} * e^{z_j}}{[\sum_{j=1}^K e^{z_k}]^2} = \frac{e^{z_i} (\sum_{j=1}^K e^{z_k} - e^{z_j})}{[\sum_{j=1}^K e^{z_k}]^2} = \sigma(\mathbf{z})_i (1 - \sigma(\mathbf{z})_j) \text{ for } i=j \text{ and}$$

$\frac{0 * \sum_{j=1}^K e^{z_k} - e^{z_i} * e^{z_j}}{[\sum_{j=1}^K e^{z_k}]^2} = -\frac{e^{z_j}}{\sum_{j=1}^K e^{z_k}} * \frac{e^{z_i}}{\sum_{j=1}^K e^{z_k}} = -\sigma(\mathbf{z})_j \sigma(\mathbf{z})_i$ . The piecewise function for the derivative is:

$\begin{cases} \sigma(\mathbf{z})_i (1 - \sigma(\mathbf{z})_j), & \text{if } i = j \\ -\sigma(\mathbf{z})_j \sigma(\mathbf{z})_i, & \text{if } i \neq j \end{cases}$ . Now apply the Kronecker delta where  $\delta_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$ , and we get  $\frac{\partial \sigma(\mathbf{z})_i}{\partial z_j} = \sigma(\mathbf{z})_i (\delta_{ij} - \sigma(\mathbf{z})_j)$ .

6. Use the above point to prove that the Jacobian of the Softmax function is the following:

$$\mathbf{J}_\sigma(\mathbf{z}) = \text{diag}(\sigma(\mathbf{z})) - \sigma(\mathbf{z})\sigma(\mathbf{z})^\top$$

where  $\text{diag}(\cdot)$  turns a vector into a diagonal matrix. Also, note that  $\mathbf{J}_\sigma(\mathbf{z}) \in \mathbb{R}^{K \times K}$ .

$$\frac{\partial \sigma_i}{\partial z_j} = \sigma_i (\delta_{ij} - \sigma_j) = \sigma_i \delta_{ij} - \sigma_i \sigma_j = (\text{diag}(\sigma))_{ij} - (\sigma \sigma^\top)_{ij} = \text{diag}(\sigma(\mathbf{z})) - \sigma(\mathbf{z})\sigma(\mathbf{z})^\top$$

## Simulating XOR

1. Can a single-layer network simulate (represent) an XOR function on  $\mathbf{x} = [x_1, x_2]$ ?

$$y = \text{XOR}(\mathbf{x}) = \begin{cases} 1, & \text{if } \mathbf{x} = (0,1) \text{ or } \mathbf{x} = (1,0) \\ 0, & \text{if } \mathbf{x} = (1,1) \text{ or } \mathbf{x} = (0,0). \end{cases}$$

Explain your reasoning using the following single-layer network definition:  $\hat{y} = \text{ReLU}(W \cdot \mathbf{x} + b)$

As we discovered in HW2, ReLU within a single layer is not able to learn complex relationships because it is a linear activation function, but XOR is not linearly separable.

2. Repeat (1) with a two-layer network:

$$\begin{aligned} \mathbf{h} &= \text{ReLU}(W_1 \cdot \mathbf{x} + \mathbf{b}_1) \\ \hat{y} &= W_2 \cdot \mathbf{h} + b_2 \end{aligned}$$



Note that this model has an additional layer compared to the earlier question: an input layer  $\mathbf{x} \in \mathbb{R}^2$ , a hidden layer  $\mathbf{h}$  with ReLU activation functions that are applied component-wise, and a linear output layer, resulting in scalar prediction  $\hat{y}$ . Provide a set of weights  $W_1$  and  $W_2$  and biases  $b_1$  and  $b_2$  such that this model can accurately model the XOR problem.

$\mathbf{x} = [x_1, x_2]$ . We want to output 1 if  $x_1$  does not equal  $x_2$  and 0 if they are equal for the XOR function. Based on this, we can apply  $W_1 = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$  and  $b_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ . Thus,  $h_1 = \text{ReLU}(x_1 - x_2)$  and  $h_2 = \text{ReLU}(x_2 - x_1)$ . Now we choose  $W_2 = \begin{bmatrix} 1 & 1 \end{bmatrix}$  and  $b_2 = 0$ . To give us the expected results for all four input combinations.

3. Consider the same network as above (with ReLU activations for the hidden layer), with an arbitrary differentiable loss function  $\ell: \{0,1\} \times \{0,1\} \rightarrow \mathbb{R}$  which takes as input  $\hat{y}$  and  $y$ , our prediction and ground truth labels, respectively. Suppose all weights and biases are initialized to zero. Show that a model trained using standard gradient descent will not learn the XOR function given this initialization.

ReLU(0) is always 0 because  $\text{ReLU} = \max(0, 0)$ . So forward pass is 0. When we perform gradients, the cached hidden layer output is still 0 so the derivative is 0, thus the weights are not updated. The  $b_2$  bias is the only element able to be trained with the weights and biases initialized to zero. Thus the model will not be able to learn the non-linear XOR function.

4. **Extra Credit:** Now let's consider a more general case than the previous question: we have the same network with an arbitrary hidden layer activation function:

$$\mathbf{h} = f(W_1 \cdot \mathbf{x} + \mathbf{b}_1)$$

$$\hat{y} = W_2 \cdot \mathbf{h} + b_2$$

Show that if the initial weights are any uniform constant, then gradient descent will not learn the XOR function from this initialization.

*A computation graph, so elegantly designed  
With nodes and edges, so easily combined  
It starts with inputs, a simple array  
And ends with outputs, in a computationally fair way*

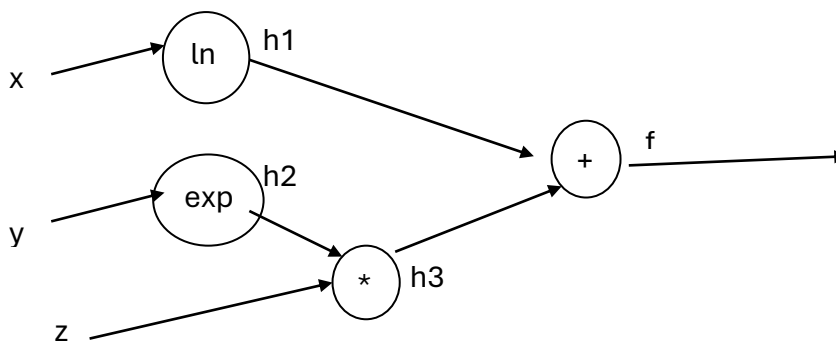
*Each node performs, an operation with care  
And passes its results, to those waiting to share  
The edges connect, each node with its peers  
And flow of information, they smoothly steer*

*It's used to calculate, complex models so grand  
And trains neural networks, with ease at hand  
Backpropagation, it enables with grace  
Making deep learning, a beautiful race*

## Neural Nets and Backpropagation

Draw the computation graph for  $f(x, y, z) = \ln x + \exp(y) \cdot z$ . Each node in the graph should correspond to only one simple operation (addition, multiplication, exponentiation, etc.). Then we will follow the forward and backward propagation described in class to estimate the value of  $f$  and partial derivatives  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$  at  $[x, y, z] = [1, 3, 2]$ . For each step, show your work.

1. Draw the computation graph for  $f(x, y, z) = \ln x + \exp(y) \cdot z$ . The graph should have three input nodes for  $x, y, z$  and one output node  $f$ . Label each intermediate node  $h_i$ .



2. Run the forward propagation and evaluate  $f$  and  $h_i$  ( $i = 1, 2, \dots$ ) at  $[x, y, z] = [1, 3, 2]$ .  
 $f(1, 3, 2) = \ln(1) + 2e^3 = 0 + 2e^3 = 2e^3$

3. Run the backward propagation and give partial derivatives for each intermediate operation, i.e.,  $\frac{\partial h_i}{\partial x}$ ,  $\frac{\partial h_j}{\partial h_i}$ , and  $\frac{\partial f}{\partial h_i}$ . Evaluate the partial derivatives at  $[x, y, z] = [1, 3, 2]$ .

$$\frac{\partial h_1}{\partial x} = \frac{1}{x} @ x=1 \rightarrow 1/1 = 1$$

$$\frac{\partial h_2}{\partial x} = 0, \frac{\partial h_3}{\partial x} = 0$$

$$\frac{\partial h_2}{\partial y} = e^y @ y=3 \rightarrow e^3$$

$$\frac{\partial h_3}{\partial h_2} = z @ z=2 \rightarrow 2, \frac{\partial h_3}{\partial z} = h_2 = e^3$$

$$f = h_1 + h_3, \frac{\partial f}{\partial h_1} = 1, \frac{\partial f}{\partial h_2} = 0, \frac{\partial f}{\partial h_3} = 1$$

4. Aggregate the results in (c) and evaluate the partial derivatives  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$  with chain rule. Show your work.

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial h_1} \frac{\partial h_1}{\partial x} = 1 * 1 = 1$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial h_3} \frac{\partial h_3}{\partial h_2} \frac{\partial h_2}{\partial y} = 1 * 2 * e^3 = 2e^3$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial h_3} \frac{\partial h_3}{\partial z} = 1 * e^3 = e^3$$

## Programming

In this programming homework, we will

- implement MLP-based classifiers for the sentiment classification task of homework 1.

### *Skeleton Code and Structure:*

The code base for this homework can be found at [this GitHub repo](#) under the hw3 directory. Your task is to fill in the missing parts in the skeleton code, following the requirements, guidance, and tips provided in this pdf and the comments in the corresponding .py files. The code base has the following structure:

- `mlp.py` reuse the sentiment classifier on movie reviews you implemented in homework 1, with additional requirements to implement MLP-based classifier architectures and forward pass .
- `main.py` provides the entry point to run your implementations `mlp.py`
- `hw3.md` provides instructions on how to setup the environment and run each part of the homework in `main.py`

**TODOs** — Your tasks include 1) generate plots and/or write short answers based on the results of running the code; 2) fill in the blanks in the skeleton to complete the code. We will explicitly mark these plotting, written answer, and filling-in-the-blank tasks as **TODOs** in the following descriptions, as well as a `# TODO` at the corresponding blank in the code. **TODOs** (Copy from your HW1). We are reusing most of the `model.py` from homework 1 as the starting point for the `mlp.py` - you will see in the skeleton that they look very similar. Moreover, in order to make the skeleton complete, for all the `# TODO` (Copy from your

**HW1**), please fill in the blank below them by copying and pasting the corresponding implementations you wrote for homework 1 (i.e. the corresponding # **TODO** in homework 1.)

### *Submission:*

Your submission should contain two parts: 1) plots and short answers under the corresponding questions below; and 2) your completion of the skeleton code base, in a .zip file

## MLP-based Sentiment Classifier

In both homework 1 & 2, our implementation of the `SentimentClassifier` is essentially a single-layer feedforward neural network that maps input features directly to 2-dimensional output logits. In this part of the programming homework, we will expand the architecture of our classifier to multi-layer perceptron (MLP).

### Reuse Your HW1 Implementation

**TODOs** (Copy from your HW1): for all the # **TODO** (**Copy from your HW1**) in `mlp.py`, please fill in the blank below them by copying and pasting the corresponding implementations you wrote for homework 1 (i.e. the corresponding # **TODO** in the `model.py` in homework 1).

### Build MLPs

Remember from the lecture that MLP is a multi-layer feedforward network with perceptrons as its nodes. A perceptron consists of non-linear activation of the affine (linear) transformation of inputs.

**TODOs:** Complete the `__init__` and forward function of the `SentimentClassifier` class in `mlp.py` to build MLP classifiers that supports custom specification of architecture (i.e. number and dimension of hidden layers)

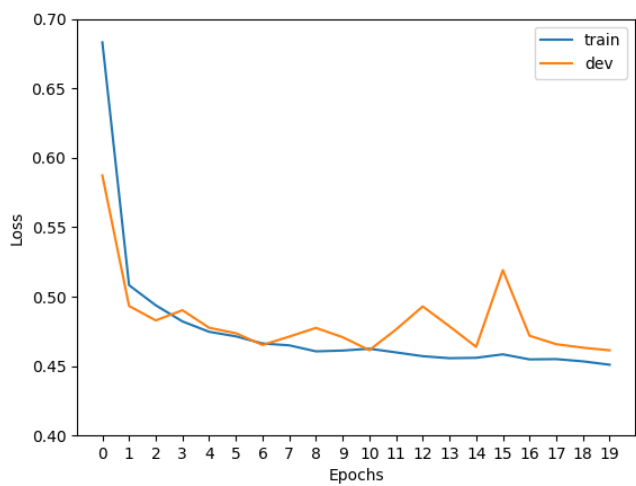
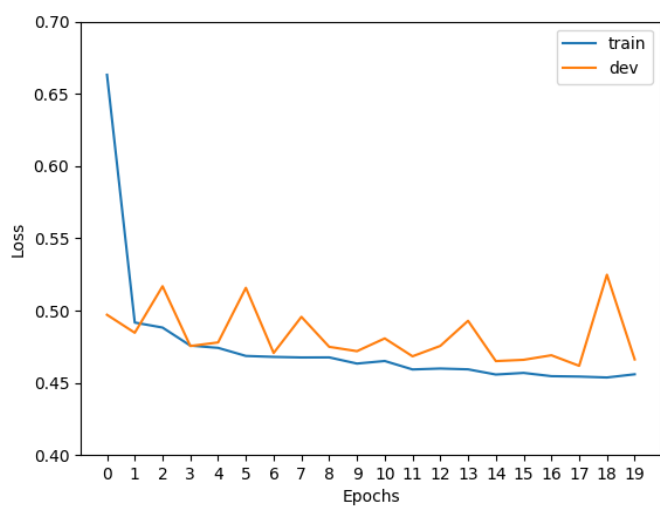
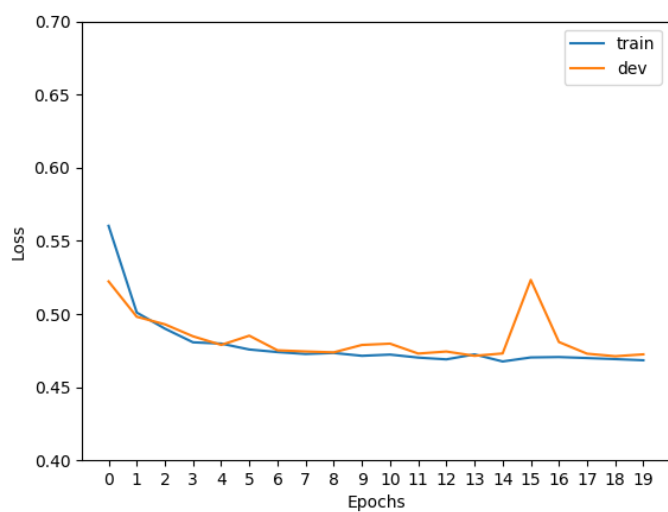
**Hint:** check the comments in the code for specific requirements about input, output, and implementation. Also, check out the document of [nn.ModuleList](#) about how to define and implement forward pass of MLPs as a stack of layers.

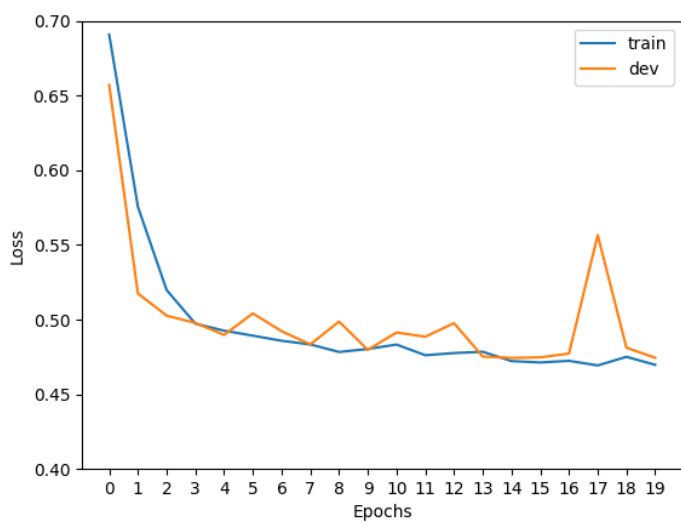
### Train and Evaluate MLPs

We provide in `main.py` several MLP configurations and corresponding recipes for training them.

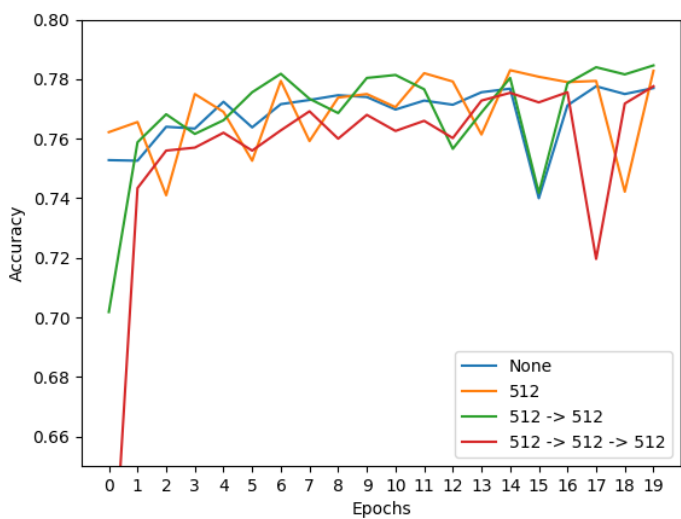
**TODOs** Once you finished 6.1.2, you can run `load_data_mlp` and `explore_mlp_structures` to train and evaluate these MLPs and paste two sets of plots here:

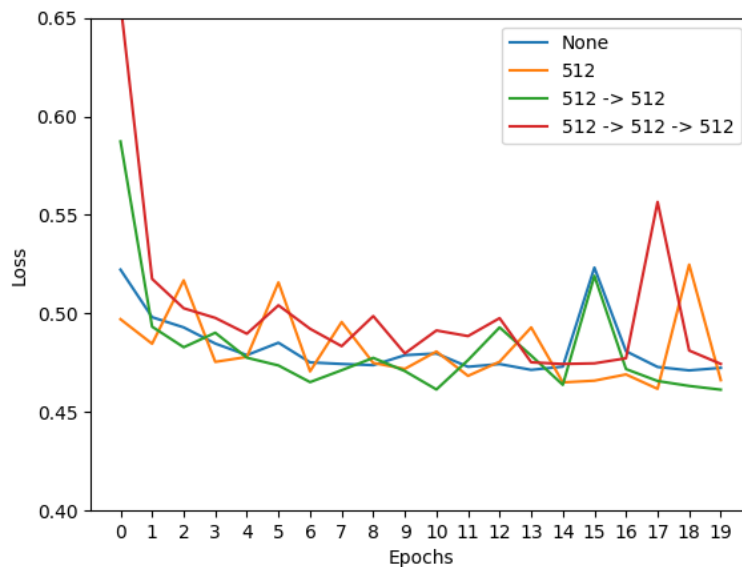
- 4 plots of train & dev loss for each MLP configuration





- 2 plots of dev losses and accuracies across MLP configurations





and describe in 2-3 sentences your findings.

**Hint:** what are the trends of train & dev loss and are they consistent across different configurations? Are deeper models always better? Why?

The 512->512 example appears to perform the best overall. The 512->512->512 actually performed the worst overall, which shows that deeper models are not always better. This could be caused by vanishing/exploding gradients and overfitting.

## Embrace Non-linearity: The Activation Functions

Remember we have learned why adding non-linearity is useful in neural nets and gotten familiar with several non-linear activation functions both in the class and 3. Now it is time to try them out in our MLPs!

**Note:** for the following TODO and the TODO in 6.1.5, we fix the MLP structure to be with a single 512-dimension hidden layer, as specified in the code. You only need to run experiments on this architecture.

**TODOs:** Read and complete the missing lines of the two following functions:

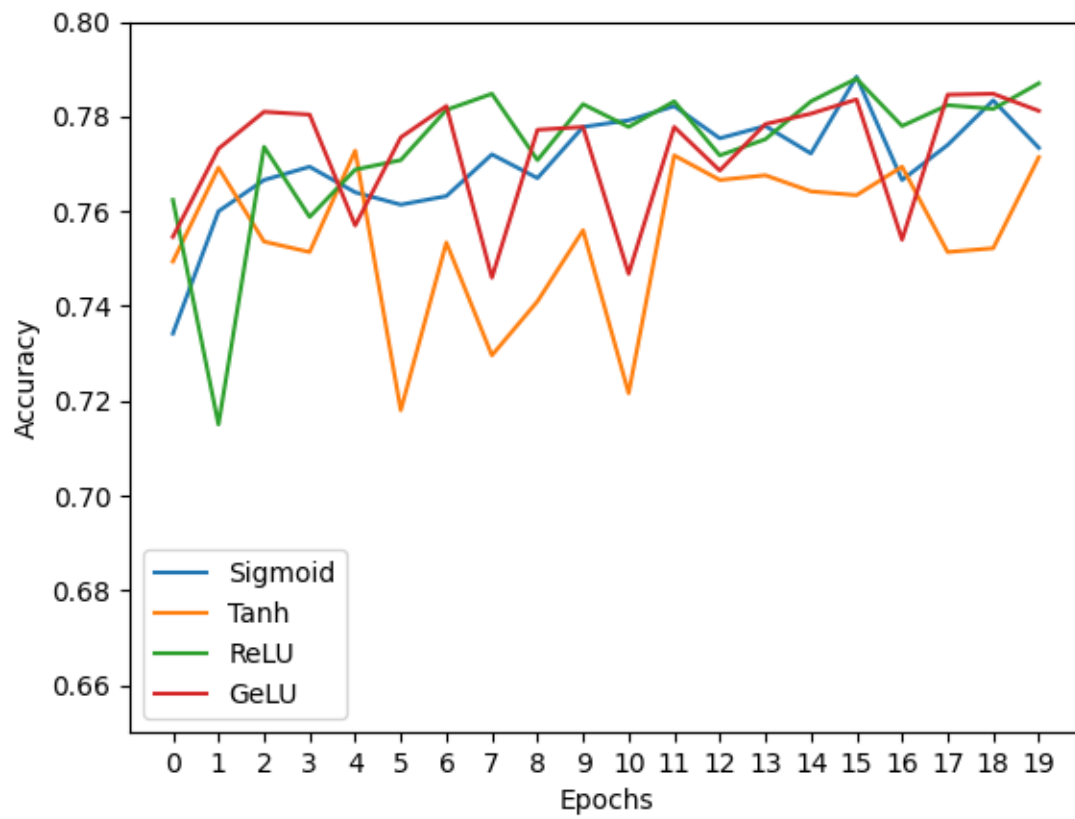
- `__init__` function of the `SentimentClassifier` class: define different activation functions given the input activation type.  
**Hint:** we have provided you with a demonstration of defining the Sigmoid activation, you can search for the other `nn.<activation>` in PyTorch documentation.
- `explore_mlp_activations` in `main.py`: iterate over the activation options, define the corresponding training configurations, train and evaluate the model, and

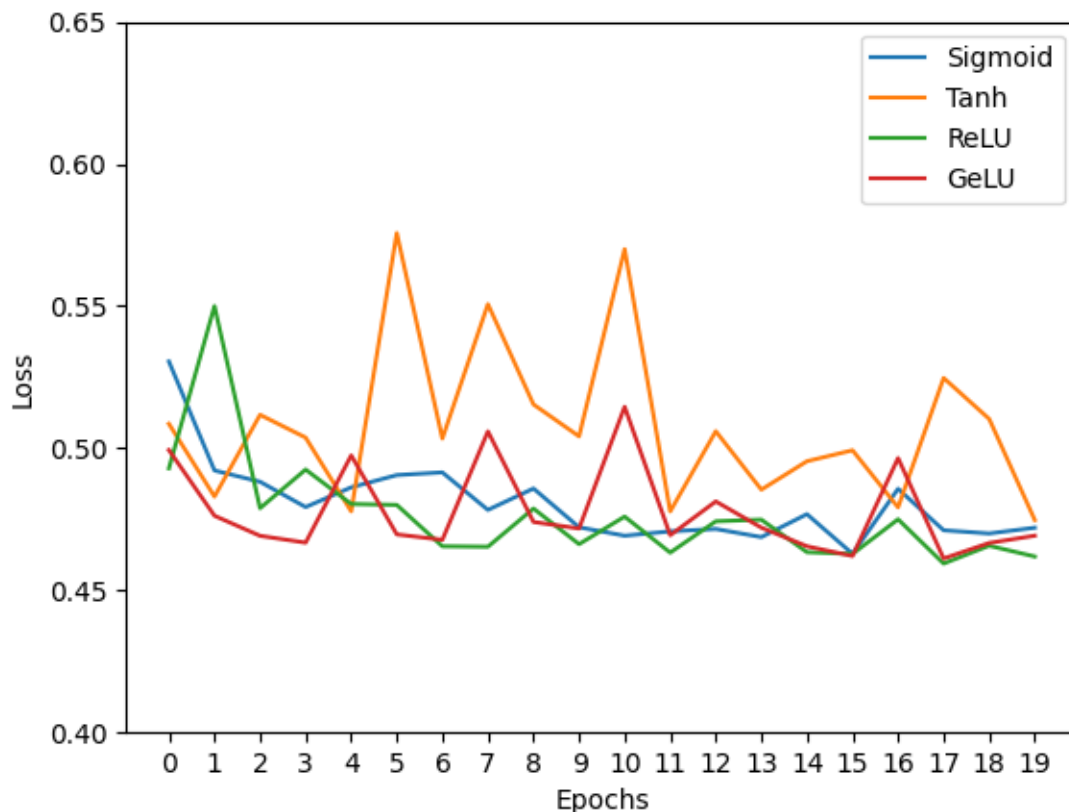
visualize the results. Note: you only need to generate the plots of dev loss and dev acc across different configurations, by calling `visualize_configs`, you **do not** need to plot the train-dev loss curves for each configuration (i.e. no need to call `visualize_epochs`). We provide you with a few choices of common activation functions, but feel free to try out the others.

**Hint:** You can refer to `explore_mlp_structure` as a demonstration of how to define training configurations with fixed hyper-parameters & iterate over hyper-parameters/design choices of interests (e.g. hidden dimensions, choice of activation), and plot the evaluation results across configurations.

Once you complete the above functions, run `explore_mlp_activations` and paste the two generated plots here. Describe in 2-3 sentences your findings.







The Tanh activation function performed the worst overall. The ReLU appears to perform very consistently meaning that it handle the noise better than some of the other activation functions, which showcases its ability to perform faster convergence.

### Hyper-parameter Tuning: Learning Rate

The training process mostly involves learning model parameters, which are automatically performed by gradient-based methods. However, certain parameters are “unlearnable” through gradient optimization while playing a crucial role in affecting model performance, for example, learning rate and batch size. We typically refer to these parameters as *Hyper-parameters*.

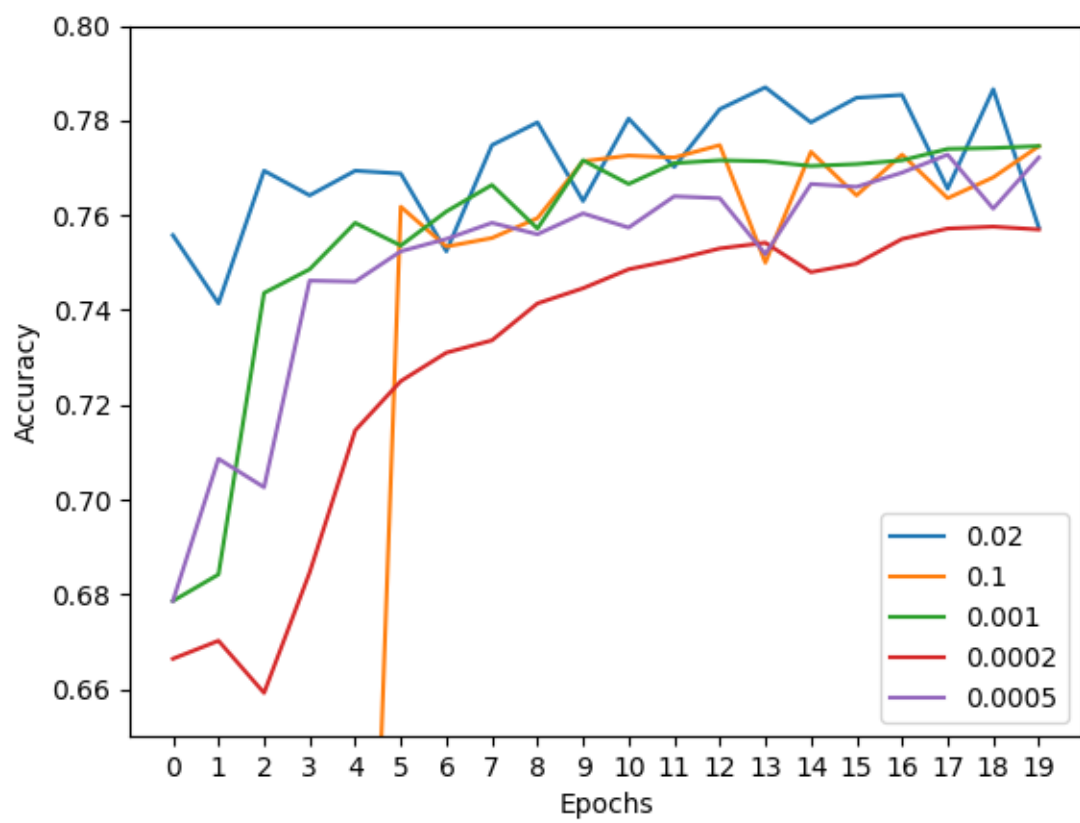
We will now take the first step to tune these hyper-parameters by exploring the choices of one of the most important one - learning rate, on our MLP. (There are lots of tutorials on how to tune the learning rate manually or automatically in practice, for example [this note](#) can serve as a starting point.)

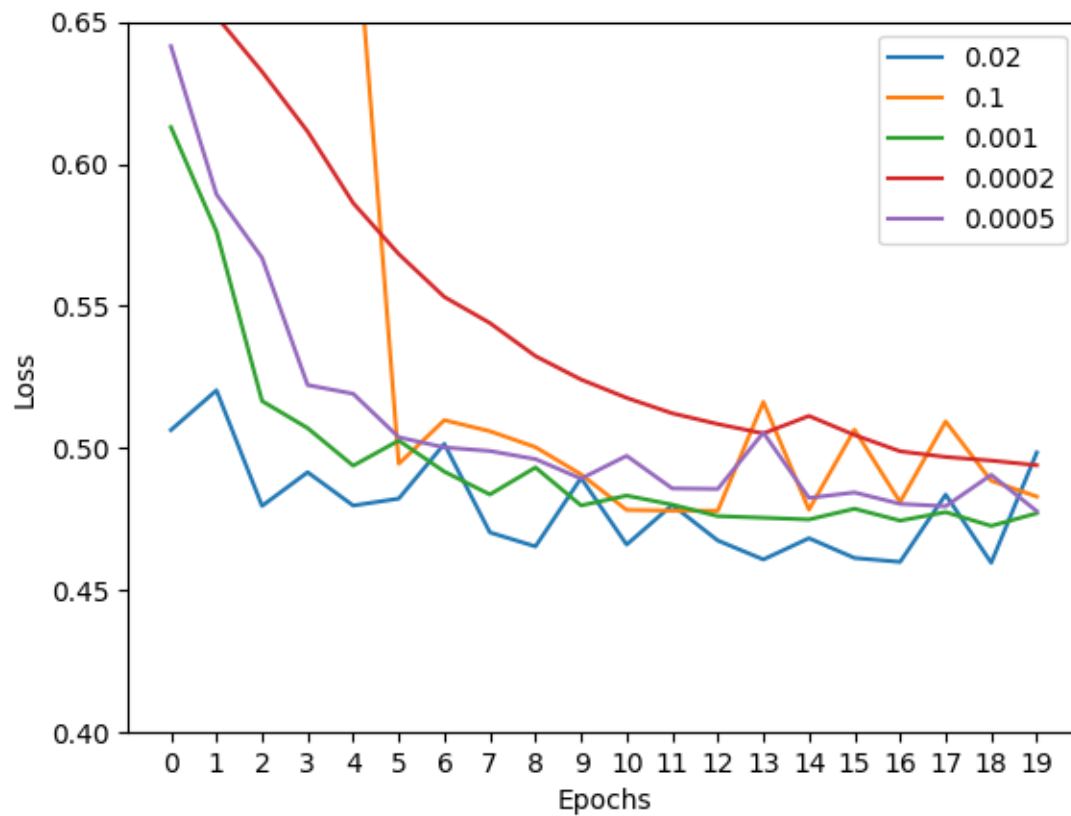
**TODOs:** Read and complete the missing lines in `explore_mlp_learning_rates` in `main.py` to iterate over different learning rate values, define the training configurations, train and evaluate the model, and visualize the results. Note: same as above, you only need to generate the plots of dev loss and dev acc across different configurations, by calling

visualize\_configs, you **do not** need to plot the train-dev loss curves for each configuration (i.e. no need to call visualize\_epochs). We provide you with the default learning rate we set to start with, and we encourage you to add more learning rate values to explore and include in your final plots curves of **at least 4 different representative learning rates**.

**Hint:** again, you can checkout explore\_mlp\_structure as a demonstration for how to perform hyper-parameter search.

Once you complete the above functions, run explore\_mlp\_learning\_rates and paste the two generated plots here. Describe in 2-3 sentences your findings.





As shown in these plots, the learning rate is crucial for model training. Choosing a learning rate that is too high can overshoot minima and cause sporadic results (shown by 0.1). When the learning rate is too small, the model can take a very long time to find the minima and converge.