

# Naïve Bayes and Neural Networks

## Solutions - DO NOT DISTRIBUTE

**Overview and Objectives.** In this homework, we are going to do some exercises to understand Naïve Bayes a bit better and then get some hand-on experience with simple Neural Networks in a multi-class classification setting.

### How to Do This Assignment.

- Each question that you need to respond to is in a blue "Task Box" with its corresponding point-value listed.
- We prefer typeset solutions (L<sup>A</sup>T<sub>E</sub>X / Word) but will accept scanned written work if it is legible. If a TA can't read your work, they can't give you credit.
- Programming should be done in Python and numpy. If you don't have Python installed, you can install it from [here](#). This is also the link showing [how to install numpy](#). You can also search through the internet for numpy tutorials if you haven't used it before. Google and APIs are your friends!

### You are **NOT** allowed to...

- Use machine learning package such as `sklearn`.
- Use data analysis package such as `panda` or `seaborn`.
- Discuss low-level details or share code / solutions with other students.

**Advice.** Start early. There are two sections to this assignment – one involving working with math (20% of grade) and another focused more on programming (80% of the grade). Read the whole document before deciding where to start.

**How to submit.** Submit a zip file to Canvas. Inside, you will need to have all your working code and `hw3-report.pdf`. You will also submit test set predictions to a [class-wide Kaggle competition](#).

## 1 Written Exercises: Analyzing Naïve Bayes [5pts]

### 1.1 Bernoulli Naïve Bayes As A Linear Classifier

Consider a Naïve Bayes model for binary classification where the features  $X_1, \dots, X_d$  are also binary variables. As part of training this Naïve Bayes model, we would estimate the conditional distributions  $P(X_i|y=c)$  for  $c = \{0, 1\}$  and  $i = 1, \dots, d$  as well as the prior distribution  $P(y=c)$  for  $c = \{0, 1\}$ . For notational simplicity, let's denote  $P(X_i=1|y=c)$  as  $\theta_{ic}$  and  $P(X_i=0|y=c)$  as  $1 - \theta_{ic}$ . Likewise, let  $\theta_1$  be  $P(y=1)$  and  $\theta_0 = 1 - \theta_1$  be  $P(y=0)$ .

If we write out the posterior of this classifier (leaving out the normalizing constant), we would have:

$$P(y = 1|x_1, \dots, x_d) = \frac{P(y = 1) \prod_{i=1}^d P(x_i|y = 1)}{P(x_1, \dots, x_d)} \propto \theta_1 \prod_{i=1}^d \theta_{i1}^{x_i} (1 - \theta_{i1})^{1-x_i} \quad (1)$$

$$P(y = 0|x_1, \dots, x_d) = \frac{P(y = 0) \prod_{i=1}^d P(x_i|y = 0)}{P(x_1, \dots, x_d)} \propto \theta_0 \prod_{i=1}^d \theta_{i0}^{x_i} (1 - \theta_{i0})^{1-x_i} \quad (2)$$

The classification rule in Naïve Bayes is to choose the class with the highest posterior probability, that is to say by predicting  $y = \operatorname{argmax}_c P(y = c|x_1, \dots, x_d)$  for a new example  $\mathbf{x} = [x_1, \dots, x_d]$ . In order for the prediction to be class 1,  $P(y = 1|x_1, \dots, x_d)$  must be greater than  $P(y = 0|x_1, \dots, x_d)$ , or equivalently we could check if:

$$\frac{P(y = 1|x_1, \dots, x_d)}{P(y = 0|x_1, \dots, x_d)} > 1 \quad (3)$$

In this setting with binary inputs, we will show that this is a linear decision boundary. This also true for continuous inputs if they are modelled with a member of the exponential family (including Gaussian) – see [proof here in §3.1](#).

► **Q1 Prove Bernoulli Naïve Bayes has a linear decision boundary [4pts]**. Prove the Bernoulli Naïve Bayes model described above has a linear decision boundary. Specifically, you'll need to show that class 1 will be predicted only if  $b + \mathbf{w}^T \mathbf{x} > 0$  for some parameters  $b$  and  $\mathbf{w}$ . To do so, show that:

$$\frac{P(y=1|x_1, \dots, x_d)}{P(y=0|x_1, \dots, x_d)} > 1 \implies b + \sum_{i=1}^d w_i x_i > 0 \quad (4)$$

As part of your report, explicitly write expressions for the bias  $b$  and each weight  $w_i$ .

**Hints:** Expand Eq. (3) by substituting the posterior expressions from Eq. (1) & (2) into Eq. (3). Take the log of that expression and combine like-terms.

► **Q1 Solution.** It can be useful to rewrite Eq. 3 a bit to make life easier for ourselves:

$$\frac{P(y=1|x_1, \dots, x_d)}{P(y=0|x_1, \dots, x_d)} = \frac{\theta_1 \prod_{i=1}^d \theta_{i1}^{x_i} (1 - \theta_{i1})^{1-x_i}}{\theta_0 \prod_{i=1}^d \theta_{i0}^{x_i} (1 - \theta_{i0})^{1-x_i}} \quad (5)$$

$$= \frac{\theta_1}{\theta_0} \prod_{i=1}^d \left( \frac{\theta_{i1}}{\theta_{i0}} \right)^{x_i} \left( \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right)^{1-x_i} \quad (6)$$

Substitute this expression into the inequality in Eq. 4 and take log of both sides:

$$\frac{\theta_1}{\theta_0} \prod_{i=1}^d \left( \frac{\theta_{i1}}{\theta_{i0}} \right)^{x_i} \left( \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right)^{1-x_i} > 1 \quad (7)$$

$$\log \frac{\theta_1}{\theta_0} + \sum_{i=1}^d \left( x_i \log \frac{\theta_{i1}}{\theta_{i0}} + (1 - x_i) \log \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right) > 0 \quad (8)$$

$$\log \frac{\theta_1}{\theta_0} + \sum_{i=1}^d \left( x_i \log \frac{\theta_{i1}}{\theta_{i0}} - x_i \log \frac{1 - \theta_{i1}}{1 - \theta_{i0}} + \log \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right) > 0 \quad (9)$$

$$\log \frac{\theta_1}{\theta_0} + \sum_{i=1}^d \log \frac{1 - \theta_{i1}}{1 - \theta_{i0}} + \sum_{i=1}^d x_i \left( \log \frac{\theta_{i1}}{\theta_{i0}} - \log \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right) > 0 \quad (10)$$

$$\log \left( \frac{\theta_1}{\theta_0} \prod_{i=1}^d \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right) + \sum_{i=1}^d x_i \log \left( \frac{\theta_{i1}(1 - \theta_{i0})}{\theta_{i0}(1 - \theta_{i1})} \right) > 0 \quad (11)$$

$$\implies b + \sum_{i=1}^d w_i x_i > 0 \quad (12)$$

with the parameter  $b$  and  $w_i$  as

$$b = \log \left( \frac{\theta_1}{\theta_0} \prod_{i=1}^d \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \right) \quad w_i = \log \left( \frac{\theta_{i1}(1 - \theta_{i0})}{\theta_{i0}(1 - \theta_{i1})} \right) \quad (13)$$

## 1.2 Duplicated Features in Naïve Bayes

Naïve Bayes classifiers assume features are conditionally independent given the class label. When this assumption is violated, Naïve Bayes classifiers may be over-confident or under-confident in the outputs. To examine this more closely, we'll consider the situation where an input feature is duplicated. These duplicated features are maximally dependent – making this a strong violation of conditional independence. Here, we'll examine a case where the confidence increases when the duplicated features are included *despite no new information actually being added as input*.

► **Q2 Duplicate Features in Naïve Bayes [1pts]**. Consider a Naïve Bayes model with a single feature  $X_1$  for a binary classification problem. Assume a uniform prior such that  $P(y = 0) = P(y = 1)$ . Suppose the model predicts class 1 for an example  $\mathbf{x}$  then we know:

$$P(y = 1|X_1 = x_1) > P(y = 0|X_1 = x_1) \quad (14)$$

Now suppose we make a mistake and duplicate the  $X_1$  feature in our data – now we have two identical inputs  $X_1$  and  $X_2$ . Show that the predicted probability for class 1 is higher than it was before; that is, prove:

$$P(y = 1|X_1 = x_1, X_2 = x_2) > P(y = 1|X_1 = x_1) \quad (15)$$

**Hints:** Use the assumption that  $P(y = 0) = P(y = 1)$  to simplify the posterior expressions in Eq. (15) and Eq. (14). As  $X_2$  is an exact copy of  $X_1$ ,  $P(X_2 = x_2|y)$  is the same as  $P(X_1 = x_1|y)$  for any example.

► **Q2 Solution.** Let's start by seeing what we can learn from Eq. 14.

$$P(y = 1|X_1 = x_1) > P(y = 0|X_1 = x_1) \quad (16)$$

$$\frac{P(X_1 = x_1|y = 1)P(y = 1)}{P(X_1 = x_1)} > \frac{P(X_1 = x_1|y = 0)P(y = 0)}{P(X_1 = x_1)} \quad (17)$$

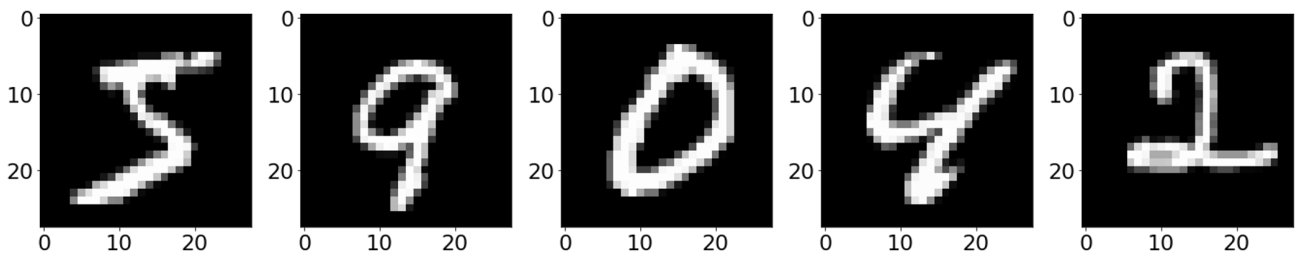
$$P(y = 0) = P(y = 1) \implies P(X_1 = x_1|y = 1) > P(X_1 = x_1|y = 0) \quad (18)$$

Now let's start manipulating  $P(y = 1|X_1 = x_1, X_2 = x_2)$  to see if we can derive Eq. 15.

$$\begin{aligned} P(y = 1|x_1, x_2) &= \frac{P(x_1, x_2|y = 1)P(y = 1)}{P(x_1, x_2|y = 1)P(y = 1) + P(x_1, x_2|y = 0)P(y = 0)} \\ &= \frac{P(x_1|y = 1)P(x_2|y = 1)P(y = 1)}{P(x_1|y = 1)P(x_2|y = 1)P(y = 1) + P(x_1|y = 0)P(x_2|y = 0)P(y = 0)} \\ &> \frac{P(x_1|y = 1)P(x_2|y = 1)P(y = 1)}{P(x_1|y = 1)P(x_2|y = 1)P(y = 1) + P(x_1|y = 1)P(x_2|y = 0)P(y = 0)} \\ &= \frac{P(x_2|y = 1)P(y = 1)}{P(x_2|y = 1)P(y = 1) + P(x_2|y = 0)P(y = 0)} \\ &= \frac{P(x_1|y = 1)P(y = 1)}{P(x_1|y = 1)P(y = 1) + P(x_1|y = 0)P(y = 0)} \\ &= P(y = 1|x_1) \\ \implies P(y = 1|x_1, x_2) &> P(y = 1|x_1) \end{aligned}$$

## 2 Implementing a Neural Network For Digit Identification [20pts]

**Small MNIST.** In this section, we will implement a feed-forward neural network model for predicting the value of a drawn digit. We are using a subset of the **MNIST** dataset commonly used in machine learning research papers. A few example of these handwritten-then-digitized digits from the dataset are shown below:



Each digit is a  $28 \times 28$  greyscale image with values ranging from 0 to 256. We represent an image as a row vector  $x \in \mathbb{R}^{1 \times 784}$  where the image has been serialized into one long vector. Each digit has an associated class label from  $0, 1, 2, \dots, 9$  corresponding to its value. We provide three dataset splits for this homework – a training set containing 5000 examples, a validation set containing 1000, and our test set containing 4220 (no labels). These datasets can be downloaded from the [class Kaggle competition for this homework](#).

## 2.1 Cross-Entropy Loss for Multiclass Classification

Unlike the previous classification tasks we've examined, we have 10 different possible class labels here. How do we measure error of our model? Let's formalize this a little and say we have a dataset  $D = \{\mathbf{x}_i, y_i\}_{i=1}^N$  with  $y_i \in \{0, 1, 2, \dots, 9\}$ . Assume we have a model  $f(x; \theta)$  parameterized by a set of parameters  $\theta$  that predicts  $P(Y|X = x)$  (a distribution over our labels given an input). Let's refer to  $P(Y = c|X = x)$  predicted from this model as  $p_{c|x}$  for compactness. We can write this output as a categorical distribution:

$$P(Y = y|X = x) = \begin{cases} p_{0|x} & \text{if } y = 0, \\ p_{1|x} & \text{if } y = 1 \\ \vdots & \\ p_{9|x} & \text{if } y = 9 \end{cases} = \prod_{c=0}^9 p_{c|x}^{\mathbb{I}[y==c]} \quad (19)$$

where  $\mathbb{I}[\text{condition}]$  is the indicator function that is 1 if the condition is true and 0 otherwise. Using this, we can write our negative log-likelihood of a single example as as:

$$-\log P(D|\theta) = -\sum_{c=0}^9 \mathbb{I}[y_i == c] \log p_{c|x_i} = -\log p_{y_i|x_i} \quad (20)$$

This loss function is also often referred to as a Cross-Entropy loss. In this homework, we will minimize this negative log-likelihood by stochastic gradient descent. In the following, we will refer to this negative log-likelihood as

$$\ell_i(\theta) = -\log p_{y_i|x_i} \quad (21)$$

Note that we write  $\ell_i$  as a function of  $\theta$  because each  $p_{y_i|x_i}$  is produced by our model  $f(x_i; \theta)$ . Summing this over an entire dataset (or a batch of examples as we will do later) of size  $N$  would yield the overall cross-entropy loss:

$$\mathcal{L}(\theta) = -\sum_{i=1}^N \log p_{y_i|x_i} \quad (22)$$

## 2.2 Implementing Backpropagation for Feed-forward Neural Network

In this homework, we'll consider feed-forward neural networks composed of a sequence of linear layers  $\mathbf{x}W_1 + \mathbf{b}_1$  and non-linear activation functions  $g_1(\cdot)$ . As such, a network with 3 of these layers stacked together can be written

$$\mathbf{b}_3 + g_2(\mathbf{b}_2 + g_1(\mathbf{b}_1 + \mathbf{x} * W_1) * W_2) * W_3 \quad (23)$$

Note how that this is a series of nested functions, reflecting the sequential feed-forward nature of the computation. To make our notation easier in the future, I want to give a name to the intermediate outputs at each stage so will expand this to write:

$$\mathbf{z}_1 = \mathbf{x} * W_1 + \mathbf{b}_1 \quad (24)$$

$$\mathbf{a}_1 = g_1(\mathbf{z}_1) \quad (25)$$

$$\mathbf{z}_2 = \mathbf{a}_1 * W_2 + \mathbf{b}_2 \quad (26)$$

$$\mathbf{a}_2 = g_2(\mathbf{z}_2) \quad (27)$$

$$\mathbf{z}_3 = \mathbf{a}_2 * W_3 + \mathbf{b}_3 \quad (28)$$

where  $\mathbf{z}$ 's are intermediate outputs from the linear layers and  $\mathbf{a}$ 's are post-activation function outputs. In the case of our MNIST experiments,  $\mathbf{z}_3$  will have 10 dimensions – a score for each of the possible labels. Finally, the output vector  $\mathbf{z}_3$  is not yet a probability distribution so we apply the softmax function:

$$p_{j|x} = \frac{e^{\mathbf{z}_{3j}}}{\sum_{c=0}^9 e^{\mathbf{z}_{3c}}} \quad (29)$$

and let  $\mathbf{p}_{\cdot|x}$  be the vector of these predicted probability values.

**Gradient Descent for Neural Networks.** Considering this simple 3-layer neural network, there are quite a few parameters spread out through the function – weight matrices  $W_3, W_2, W_1$  and biases vectors  $\mathbf{b}_3, \mathbf{b}_2, \mathbf{b}_1$ . Suppose we would like to find parameters that minimize our loss  $\mathcal{L}$  that measures our error in the network's prediction.

How can we update the weights to reduce this error? Let's use gradient descent and start by writing out the chain rule for the gradient of each of these. I'll work backwards from  $W_3$  to  $W_1$  to expose some structure here.

$$\frac{\delta \mathcal{L}}{\delta W_3} = \frac{\delta \mathcal{L}}{\delta \mathbf{p}_{\cdot|x}} \frac{\delta \mathbf{p}_{\cdot|x}}{\delta \mathbf{z}_3} \frac{\delta \mathbf{z}_3}{\delta W_3} \quad (30)$$

$$\frac{\delta \mathcal{L}}{\delta W_2} = \frac{\delta \mathcal{L}}{\delta \mathbf{p}_{\cdot|x}} \frac{\delta \mathbf{p}_{\cdot|x}}{\delta \mathbf{z}_3} \frac{\delta \mathbf{z}_3}{\delta \mathbf{a}_2} \frac{\delta \mathbf{a}_2}{\delta \mathbf{z}_2} \frac{\delta \mathbf{z}_2}{\delta W_2} \quad (31)$$

$$\frac{\delta \mathcal{L}}{\delta W_1} = \frac{\delta \mathcal{L}}{\delta \mathbf{p}_{\cdot|x}} \frac{\delta \mathbf{p}_{\cdot|x}}{\delta \mathbf{z}_3} \frac{\delta \mathbf{z}_3}{\delta \mathbf{a}_2} \frac{\delta \mathbf{a}_2}{\delta \mathbf{z}_2} \frac{\delta \mathbf{z}_2}{\delta \mathbf{a}_1} \frac{\delta \mathbf{a}_1}{\delta \mathbf{z}_1} \frac{\delta \mathbf{z}_1}{\delta W_1} \quad (32)$$

As I've highlighted in color above, we end up reusing the same intermediate terms over and over as we compute derivatives for weights further and further from the output in our network.<sup>1</sup> As discussed in class, this suggests the straight-forward *backpropagation* algorithm for computing these efficiently. Specifically, we will compute these intermediate colored terms starting from the output and working backwards.

**Forward-Backward Pass in Backpropagation.** One convenient way to implement backpropagation is to consider each layer (or operation)  $f$  as having a forward pass that computes the function output normally as

$$output = f_{forward}(input) \quad (33)$$

and a backward pass that takes in the gradient up to this point in our backward pass and then outputs the gradient of the loss with respect to its input:

$$\frac{\delta \mathcal{L}}{\delta input} = f_{backward} \left( \frac{\delta \mathcal{L}}{\delta output} \right) = \frac{\delta \mathcal{L}}{\delta output} \frac{\delta output}{\delta input} \quad (34)$$

The backward operator will also compute any gradients with respect to parameters of  $f$  and store them to be used in a gradient descent update step after the backwards pass. The starter code implements this sort of framework.

See the snippet on the following page that defines a neural network like the one we've described here, except it allows for a configurable number of linear layers. Please read the comments and code below before continuing reading this document. To give concrete examples of the forward-backward steps for an operator, consider the Sigmoid (aka the logistic) activation function below:

$$Sigmoid(x) = \frac{1}{1 + e^{-x}} \quad (35)$$

The implementation for forward and backward for the Sigmoid is below – in forward it computes Eq.(35), in backward it computes and returns

$$\frac{\delta \mathcal{L}}{\delta input} = \frac{\delta \mathcal{L}}{\delta output} \frac{\delta output}{\delta input} = \frac{\delta \mathcal{L}}{\delta output} Sigmoid(input)(1 - Sigmoid(input)) \quad (36)$$

It has no parameters so does nothing during the "step" function.

```
1 class Sigmoid:
2
3     # Given the input, apply the sigmoid function
4     # store the output value for use in the backwards pass
5     def forward(self, input):
6         self.act = 1/(1+np.exp(-input))
7         return self.act
8
9     # Compute the gradient of the output with respect to the input
10    # self.act*(1-self.act) and then multiply by the loss gradient with
11    # respect to the output to produce the loss gradient with respect to the input
12    def backward(self, grad):
13        return grad * self.act * (1-self.act)
14
15    # The Sigmoid has no parameters so nothing to do during a gradient descent step
16    def step(self, step_size):
17        return
```

<sup>1</sup>I don't repeat this for the bias vectors, as it would simply involve changing the final terms from  $\frac{\delta \mathbf{z}_i}{\delta W_i}$  to  $\frac{\delta \mathbf{z}_i}{\delta \mathbf{b}_i}$ .

```

1 class FeedForwardNeuralNetwork:
2
3     # Builds a network of linear layers separated by non-linear activations
4     # Either ReLU or Sigmoid. Each internal layer has hidden_dim dimensions.
5     def __init__(self, input_dim, output_dim, hidden_dim, num_layers, activation="ReLU"):
6
7         if num_layers == 1: #Just a linear mapping from input to output
8
9             self.layers = [LinearLayer(input_dim, output_dim)]
10
11         else: # At least two layers
12
13             #Layer to map input to hidden dimension size
14             self.layers = [LinearLayer(input_dim, hidden_dim)]
15             self.layers.append(Sigmoid() if activation=="Sigmoid" else ReLU())
16
17             # Hidden layers
18             for i in range(num_layers-2):
19                 self.layers.append(LinearLayer(hidden_dim, hidden_dim))
20                 self.layers.append(Sigmoid() if activation=="Sigmoid" else ReLU())
21
22             # Layer to map hidden dimension to output size
23             self.layers.append(LinearLayer(hidden_dim, output_dim))
24
25     # Given an input, call the forward function of each of our layers
26     # Pass the output of each layer to the next one
27     def forward(self, X):
28         for layer in self.layers:
29             X = layer.forward(X)
30         return X
31
32     # Given an gradient with respect to the network output, call
33     #the backward function of each of our layers. Pass the output of each layer to the
34     #one before it
35     def backward(self, grad):
36         for layer in reversed(self.layers):
37             grad = layer.backward(grad)
38
39     # Tell each layer to update its weights based on the gradient computed in the
40     #backward pass
41     def step(self, step_size=0.001):
42         for layer in self.layers:
43             layer.step(step_size)

```

**Operating on Batches.** The network described in the equations earlier in this section is operating on a single input at a time. In practice, we will want to operate on sets of  $n$  examples at once such that the layer actually computes  $Z = XW + \mathbf{b}$  for  $X \in \mathbb{R}^{n \times \text{input\_dim}}$  and  $Z \in \mathbb{R}^{n \times \text{output\_dim}}$  – call this a batched operation. It is straightforward to change the forward pass to operate on these all at once. For example, a linear layer can be rewritten as  $Z = XW + b$  where the  $+b$  is a broadcasted addition – this is already done in the code above. On the backward pass, we simply need to aggregate the gradient of the loss of each data point with respect to our parameters. For example,

$$\frac{\delta L}{\delta W_1} = \sum_{i=1}^n \frac{\delta L_i}{\delta W_1} \quad (37)$$

where  $L_i$  is the loss of the  $i$ 'th datapoint and  $L$  is the overall loss.

**Deriving the Backward Pass for a Linear Layer.** In this homework, we'll implement the backward pass of a linear layer. To do so, we'll need to be able to compute  $dZ/db$ ,  $dZ/dW$ , and  $dZ/dX$ . For each, we'll start by considering the problem for a single training example  $x$  (i.e. a single row of  $X$ ) and then generalize to the batch setting. In this single-example setting,  $z = xW + b$  such that  $z, b \in \mathbb{R}^{1 \times c}$ ,  $x \in \mathbb{R}^{1 \times d}$ , and  $W \in \mathbb{R}^{d \times c}$ . Once we solve this case, extending to the batch setting just requires summing over the gradient terms for each example.

**$dZ/db$ .** Considering just the  $i$ 'th element of  $z$ , we can write  $z_i = xw_{\cdot,i} + b_i$  where  $w_{\cdot,i}$  is the  $i$ 'th column of  $W$ . From this equation, it is straightforward to observe that element  $b_i$  only effects the corresponding output  $z_i$  such that

$$\frac{dz_i}{db_j} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (38)$$

This suggests that the Jacobian  $dZ/db$  is an identity matrix  $I$  of dimension  $c \times c$ . Applying chain rule and summing

over all our datapoints, we see  $dL/db$  can be computed as a sum of the rows of  $dL/dZ$ :

$$\frac{dL}{db} = \sum_{k=1}^n \frac{dL_k}{dZ_k} \frac{dZ_k}{db} = \sum_{k=1}^n \frac{dL_k}{dZ_k} I = \sum_{k=1}^n \frac{dL_k}{dZ_k} \quad (39)$$

**dZ/dW.** Following the same process of reasoning from the single-example case, we can again write the  $i$ 'th element of  $z$  as  $z_i = xw_{.,i} + b_i$  where  $w_{.,i}$  is the  $i$ 'th column of  $W$ . When considering the derivative of  $z_i$  with respect to the columns of  $W$ , we see that it is just  $x$  for  $w_{.,i}$  and 0 for other columns as they don't contribute to  $z_i$  – that is to say:

$$\frac{\delta z_i}{\delta w_{.,j}} = \begin{cases} x & \text{if } i = j \\ 0 & \text{otherwise} \end{cases} \quad (40)$$

Considering the loss gradient  $\delta L/\delta w_{.,i}$  for a single example, we can write:

$$\frac{\delta L}{\delta w_{.,i}} = \frac{\delta L}{\delta z_i} \frac{\delta z_i}{\delta w_{.,i}} = \frac{\delta L}{\delta z_i} x \quad (41)$$

That is to say, each column  $i$  of  $\frac{\delta L}{\delta W}$  is the input  $x$  scaled by the loss gradient of  $z_i$ . As such, we can compute the gradient for the entire  $W$  as the product:

$$\frac{\delta L}{\delta W} = \mathbf{x}^T \frac{\delta L}{\delta z} \quad (42)$$

Notice that  $\mathbf{x}^T$  is  $d \times 1$  and  $\frac{\delta L}{\delta z}$  is  $1 \times c$  – resulting in a  $d \times c$  gradient that matches the dimension of  $W$ .

Now let's consider if we have multiple datapoints  $x_1, \dots, x_n$  as the matrix  $X$  and likewise multiple activation vectors  $z_1, \dots, z_n$  as the matrix  $Z$ . As our loss simply sums each datapoint's loss, the gradient also decomposes into a sum of  $\frac{\delta L}{\delta z_i}$  terms.

$$\frac{\delta L}{\delta W} = \sum_{k=1}^n \frac{\delta L}{\delta Z_k} \frac{\delta Z_k}{\delta W} = \sum_{k=1}^n X_k^T \frac{\delta L_k}{\delta Z_k} \quad (43)$$

We can write this even more compactly as:

$$\frac{\delta L}{\delta W} = X^T \frac{\delta L}{\delta Z} \quad (44)$$

**dZ/dX.** This follows a very similar path as dZ/dW. We again consider the  $i$ 'th element of  $z$  as  $z_i = xw_{.,i} + b_i$  where  $w_{.,i}$  is the  $i$ 'th column of  $W$ . Taking the derivative with respect to  $x$  it is clear that for  $z_i$  the result will be  $w_{.,i}$ .

$$\frac{\delta z_i}{\delta x} = w_{.,i} \quad (45)$$

This suggests that the rows of  $dZ/dx$  are simply the columns of  $W$  such that  $dZ/dx = W^T$  and we can write

$$\frac{\delta L}{\delta \mathbf{x}} = \frac{\delta L}{\delta z} \frac{\delta z}{\delta \mathbf{x}} = \frac{\delta L}{\delta z} W^T \quad (46)$$

Moving to the multiple example setting, the above expression gives each row of  $dL/dX$  and the entire matrix can be computed efficiently as

$$\frac{dL}{dX} = \frac{dL}{dZ} W^T \quad (47)$$

► **Q3 Implementing the Backward Pass for a Linear Layer [6pt]**. Implement the backward pass function of the linear layer in the skeleton code. The function takes in the matrix  $dL/dZ$  as the variable `grad` and you must compute  $dL/dW$ ,  $dL/db$ , and  $dL/dX$ . The first two are stored as `self.grad_weights` and `self.grad_bias` and the third is returned. The expressions for these can be found above in Eq.39 ( $dL/db$ ), Eq.44 ( $dL/dW$ ), and Eq.47 ( $dL/dX$ ).

```

1 class LinearLayer:
2
3     # Initialize our layer with (input_dim, output_dim) weight matrix and a (1,
4     # output_dim) bias vector
5     def __init__(self, input_dim, output_dim):
6         self.weights = np.random.randn(input_dim, output_dim).astype(np.float64)* np
7         .sqrt(2. / input_dim)
8         self.bias = np.ones( (1,output_dim) ).astype(np.float64)*0.5
9
10    # During the forward pass, we simply compute Xw+b
11    def forward(self, input):
12        self.input = input
13        return self.input@self.weights + self.bias
14
15    #####
16    # Q3 Implementing Backward Pass for Linear
17    #####
18    # Inputs:
19    #
20    # grad dL/dZ -- For a batch size of n, grad is a (n x output_dim) matrix where
21    #                 the i'th row is the gradient of the loss of example i with respect
22    #                 to z_i (the output of this layer for example i)
23    #
24    # Computes and stores:
25    #
26    # self.grad_weights dL/dW -- A (input_dim x output_dim) matrix storing the
27    #                             gradient of the loss with respect to the weights of
28    #                             this layer.
29    #
30    # self.grad_bias dL/db--      A (1 x output_dim) matrix storing the gradient
31    #                             of the loss with respect to the bias of this layer.
32    #
33    # Return Value:
34    #
35    # grad_input dL/dX -- For a batch size of n, grad_input is a (n x input_dim)
36    #                     matrix where the i'th row is the gradient of the loss of
37    #                     example i with respect to x_i (the input of this
38    #                     layer for example i)
39    #####
40
41    def backward(self, grad): # grad is dL/dZ
42        self.grad_weights = ? # Compute dL/dW as in Eq. 32
43        self.grad_bias = ?   # Compute dL/db as in Eq. 27
44        return ?             # Compute dL/dX as in Eq. 35
45
46    # During the gradient descent step, update the weights and biases based on the
47    # stored gradients from the backward pass
48    def step(self, step_size):
49        self.weights -= step_size*self.grad_weights
50        self.bias -= step_size*self.grad_bias

```

1

► **Q3 Solution**

```

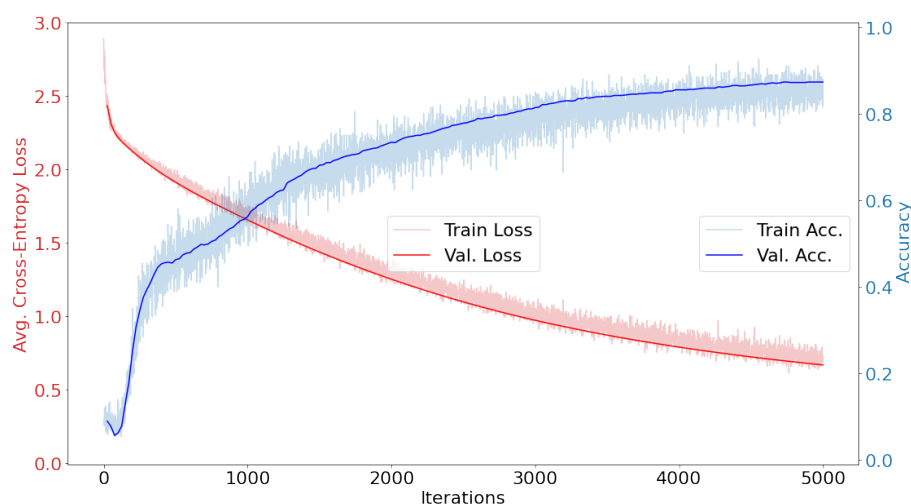
1 def backward(self, grad):
2     self.grad_weights = self.input.T @ grad
3     self.grad_bias = np.sum(grad,axis=0)
4     return grad @ self.weights.T

```

Once you've completed the above task, running the skeleton code should load the digit data and train a 2-layer neural network with hidden dimension of 16 and Sigmoid activations. This model is trained on the training set and evaluated once per epoch on the validation data. After training, it will produce a plot of your results that should look like the one below. This curve plots training and validation loss (cross-entropy in this case) over training iterations (in



red and measured on the left vertical axis). It also plots training and validation accuracy (in blue and measures on the right vertical axis). As you can see, this model achieves between 80% and 90% accuracy on the validation set.



## 2.3 Analyzing Hyperparameter Choices

Neural networks have many hyperparameters. These range from architectural choices (*How many layers? How wide should each layer be? What activation function should be used?*) to optimization parameters (*What batch size for stochastic gradient descent? What step size (aka learning rate)? How many epochs should I train?*). This section has you modify many of these to examine their effect. The default parameters are below for easy reference.

```
1 # GLOBAL PARAMETERS FOR STOCHASTIC GRADIENT DESCENT
2 np.random.seed(102)
3 step_size = 0.01
4 batch_size = 200
5 max_epochs = 200
6
7 # GLOBAL PARAMETERS FOR NETWORK ARCHITECTURE
8 number_of_layers = 2
9 width_of_layers = 16 # only matters if number of layers > 1
10 activation = "ReLU" if False else "Sigmoid"
```

**Optimization Parameters.** Optimization parameters in Stochastic Gradient Descent are very inter-related. Large batch sizes mean less noisy estimates of the gradient, so larger step sizes could be used. But larger batch sizes also mean fewer gradient updates per epoch, so we might need to increase the max epochs. Getting a good set of parameters that work well can be tricky and requires checking the validation set performance. Further, these “good parameters” will vary model-to-model.

► **Q4 Learning Rate [2pts]**. The learning rate (or step size) in stochastic gradient descent controls how large of a step in the direction of the loss gradient we take our parameters at each iteration. The batch size determines how many data points we use to estimate the gradient. Modify the hyperparameters to run the following experiments:

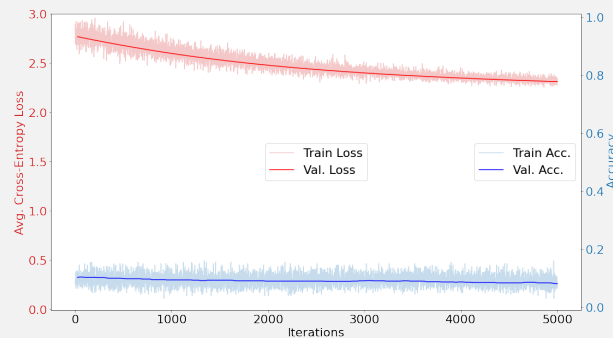
1. Step size of 0.0001 (leave default values for other hyperparameters)
2. Step size of 5 (leave default values for other hyperparameters)
3. Step size of 10 (leave default values for other hyperparameters)

Include these plot in your report and answer the following questions:

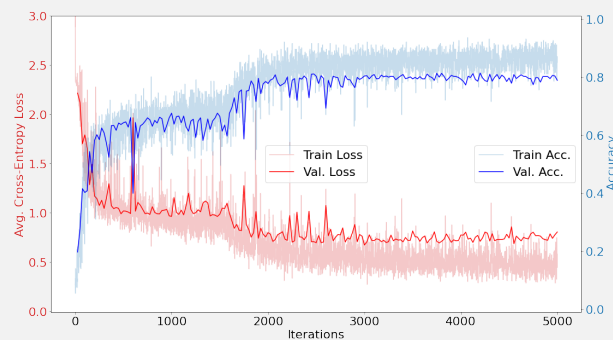
- a) Compare and contrast the learning curves with your curve using the default parameters. What do you observe in terms of smoothness, shape, and what performance they reach?
- b) For (a), what would you expect to happen if the max epochs were increased?

## ► Q4 Solution.

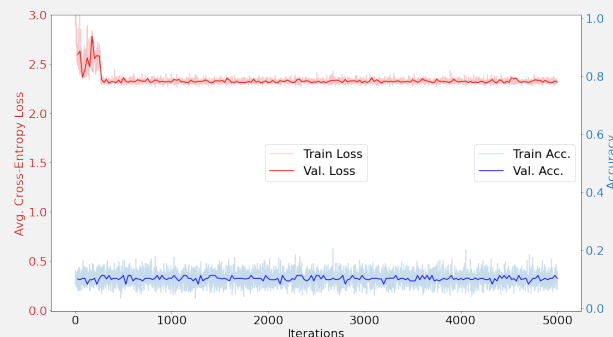
1) Step size of 0.0001 (leave default values for other hyperparameters)



2) Step size of 5 (leave default values for other hyperparameters)



3) Step size of 10 (leave default values for other hyperparameters)



a) Learning curves for (1) are smooth but don't seem to be converging quickly compared to the default parameters. For (2), the model converges quickly but is jumpy and seems to converge to a worse solution for validation accuracy. The divergence between training loss and validation loss suggests some overfitting. For (3), practically no learning occurs at all.

b) If given more training time by extending max\_epochs, the curve for (1) would likely eventually converge whereas neither (2) or (3) would likely benefit.

**Activation Function and Depth.** As networks get deeper (or have more layers) they tend to become able to fit more complex functions (though this also may lead to overfitting). However, this also means the backpropagated gradient has many product terms before reaching lower levels – resulting in the magnitude of the gradients being relatively small. This has the effect of making learning slower. Certain activation functions make this better or worse depending on the shape of their derivative. One popular choice is to use a Rectified Linear Unit or ReLU activation that computes:

$$\text{ReLU}(x) = \max(0, x) \quad (48)$$

This is especially common in very deep networks. In the next question, we'll see why experimentally.

► Q5 ReLU's and Vanishing Gradients [3pts]. Modify the hyperparameters to run the following experiments:

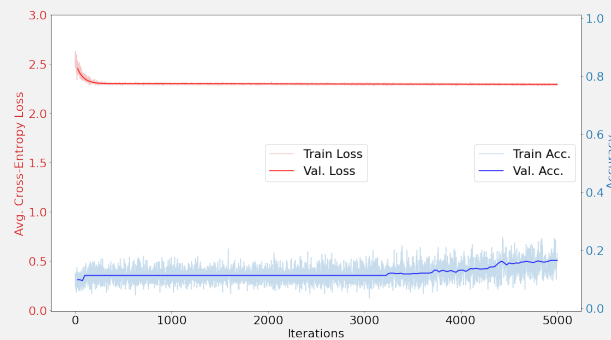
1. 5-layer with Sigmoid Activation (leave default values for other hyperparameters)
2. 5-layer with Sigmoid Activation with 0.1 step size (leave default values for other hyperparameters)
3. 5-layer with ReLU Activation (leave default values for other hyperparameters)

Include these plot in your report and answer the following questions:

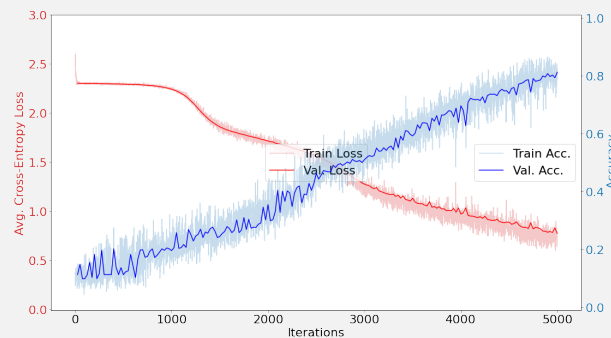
- a) Compare and contrast the learning curves you observe and the curve for the default parameters in terms of smoothness, shape, and what performance they reach. Do you notice any differences in the relationship between the train and validation curves in each plot?
- b) If you observed increasing the learning rate in (2) improves over (1), why might that be?
- c) If (3) outperformed (1), why might that be? Consider the derivative of the sigmoid and ReLU functions.

► Q4 Solution.

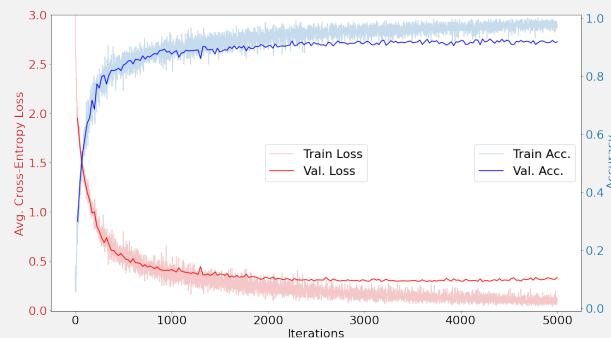
1) 5-layer with Sigmoid Activation (leave default values for other hyperparameters)



2) 5-layer with Sigmoid Activation with 0.1 step size (leave default values for other hyperparameters)



3) 5-layer with ReLU Activation (leave default values for other hyperparameters)



a) Learning curves for (1) fail to learn for thousands of iterations and only start to improve on training error towards the end. For (2), the model learns slowly as well but the large learning rate induces noise. For (3), the ReLU network learns quite quickly and converges to a high accuracy.

b) The first network was converging very slowly by taking too small of steps (low learning rate). Increasing the learning rate led to quicker convergence and better score.

c) The sigmoid activation function's gradient decays to zero outside of a narrow range near zero. The ReLU on the other hand has a constant derivative for any positive input. This enables ReLU networks to train more quickly.

## 2.4 Randomness in Training

There is also a good deal of randomness in training a neural network with stochastic gradient descent – network weights are randomly initialized and the batches are randomly ordered. This can make a non-trivial difference to outcomes.

► **Q6 Measuring Randomness [1pt]**. Using the default hyperparameters, set the random seed to 5 different values and report the validation accuracies you observe after training. What impact does this randomness have on the certainty of your conclusions in the previous questions?

► **Q6 Solution**. Running the default 5 times yields the following validation accuracies: 88.7, 87.2, 87.6, 88.5, 88.3. This variance isn't high enough to change any of the conclusions above, but it is notable that almost a point can vary in either direction though.

## 2.5 Make Your Kaggle Submission

Great work getting here. In this section, you'll submit the predictions of your best model to the **class-wide Kaggle competition**. You are free to make any modification to your neural network or the optimization procedure to improve performance; however, it must remain a feed-forward neural network! For example, you can change any of the optimization hyperparameters or add momentum / weight decay, vary the number of layers or width, add dropout or residual connections, etc.

► **Q7 Kaggle Submission [8pt]**. Submit a set of predictions to Kaggle that outperforms the baseline on the public leaderboard. To make a valid submission, use the train set to train your neural network classifier and then apply it to the test instances in `mnist_small_test.csv` available from Kaggle's Data tab. Format your output as a two-column CSV as below:

```
id,digit
0,3
1,9
2,4
3,1
.
.
```

where the id is just the row index in `mnist_small_test.csv`. You may submit up to 10 times a day. In your report, tell us what modifications you made for your final submission.

**Extra Credit and Bragging Rights [1.25pt Extra Credit]**. The TA has made a submission to the leaderboard. Any submission outperforming the TA on the *private* leaderboard at the end of the homework period will receive 1.25 extra credit points on this assignment. Further, the top 5 ranked submissions will “win HW3” and receive bragging rights.

## 3 Debriefing (required in your report)

1. Approximately how many hours did you spend on this assignment?
2. Would you rate it as easy, moderate, or difficult?
3. Did you work on it mostly alone or did you discuss the problems with others?
4. How deeply do you feel you understand the material it covers (0%–100%)?
5. Any other comments?