



PURDUE UNIVERSITY

OPT4DL

---

## HW 4

---

October 3, 2024

33678229 John Yeary

## Contents

<b>1</b>	<b>Question 1</b>	<b>1</b>
<b>2</b>	<b>Question 2</b>	<b>2</b>
<b>3</b>	<b>Question 3</b>	<b>3</b>
3.1	Part A . . . . .	3
3.1.1	Deriving Hessian . . . . .	3
3.1.2	Proving Convex and L Smooth . . . . .	4
3.2	Part B . . . . .	4
<b>4</b>	<b>Question 4</b>	<b>6</b>
4.1	Explanation . . . . .	6
4.2	Code . . . . .	6

## 1 Question 1

This paper presents a benchmark of various zero order methods for fine tuning a large language model (LLM). This is in stark contrast to first order methods, which requires a backpropagation step. These steps are memory and computationally efficient. As models increase in size, memory requirements scale exponentially. To remedy this, ZO methods offer the opportunity to estimate the gradient without requiring backpropagation. These methods have the potential to reduce memory requirements and computational overhead.

One major finding of the paper was the importance of task alignment during fine tuning. During pre-training, the model is fed massive amounts of data, and training is performed using next token prediction, or similar methods. This informs the model's parameters to be within a 'ballpark' of a final solution. Once the fine tuning step begins, it is important to use a similar format to the pre-training data that was fed to the model. In doing so, the model starts at a reasonable place within its parameter space, so fine tuning is more straightforward.

This starting point is important since zero order methods will be sensitive to starting location. Since these methods are approximating the gradient with finite differences in values, small changes in function value can cause the model's convergence to vary widely. By identifying a better starting point, convergence becomes easier with a ZO model. This may struggle with generalization as the model will become specialized for certain data types.

The paper suggests that there is a tradeoff in memory efficiency and model accuracy. There are a few suggestions from the paper as well. The authors suggest using block ZO, hybrid FO and ZO, and gradient pruning to improve ZO performance.

Block-wise ZO solves the issue of high gradient estimation variance typical of ZO models. Typically the gradient for the entire model is estimated at once, leading to reduced accuracy across the model. By breaking the model into smaller blocks and calculating a gradient approximation for each block, you can limit the variance of each block's variance, effectively improving the training for that block. Since we aren't requiring a backpropagation, each block can be self contained. One drawback to this approach is that it requires multiple forward passes per iteration, since each block will require a forward pass.

Hybrid approaches aim to combine the accuracy of FO methods with the memory and computational efficiency of ZO methods. The idea is to use FO methods for some layers, and ZO for the rest. Typically, FO methods are used for deeper model layers, and critical layers within the model. The remaining layers apply ZO approximations to reduce computational requirements. The difficulty in this approach is identifying which layers to treat as FO and which to treat as ZO. This is likely application specific.

Some ZO methods, such as random gradient estimation, attempt to estimate the gradient by randomly perturbing the model's parameters. This can introduce variance to the gradient estimate, especially with large numbers of parameters. To resolve this issue, gradient pruning selects a subset of parameters to perturb in order to approximate the gradient. The selection criteria for which parameters to include/remove from this perturbation can have a large impact on model accuracy.

There are clear memory advantages of using a zero order approximation, but there are still some questions about performance stability and generalization when applying these methods. In general the performance gap between ZO and FO methods increases with task complexity.

The paper does a good job of providing an overview and benchmark for ZO methods. The authors provide a solid starting point for future investigation into ZO methods.

Future work should investigate the scalability issues with ZO methods. This would likely involve some adaptive methods to reduce variance of the gradient estimate.

More work could also focus on establishing a more theoretical reasoning for why the proposed methods work. This would involve evaluating conditions for when and why to apply each method. This could tie into some of the previous reading, where a decision metric was used to identify when it was fine to use 'easier' samples instead of 'harder' samples for training. Similarly, this would apply to deciding to use FO or ZO methods depending on some external criteria.

## 2 Question 2

A backdoor attack is performed by adding faulty data with a hidden trigger to the training data. This will activate some malicious function within the model. Once trained, the model will behave normally until the trigger appears, causing the model to misbehave and output whatever the malicious functionality is.

These attacks have typically been easy to spot, since any time the trigger appears, the model will behave differently. These were fairly simple to catch since the trigger remains static.

The paper proposes an input aware attack, where multiple triggers are hidden across multiple different training samples. The faulty behavior would then only occur when the unique data and trigger are presented. This would allow the trigger to be unique to each image, and prevent the trigger from occurring on other samples.

The authors proposed a system for generating such an attack: First, they propose creating a trigger generator network. This will encode a unique trigger pattern for each image. They proposed training the model with clean mode, attack mode, and cross trigger mode. Clean mode is simply training, attack mode is training the network to output a predefined 'attack' label when a trigger corresponding to the input is present, and cross trigger mode will train the network to output the correct result even if a trigger for a different image is present.

The authors then tested their attack using various datasets. They tested its effectiveness against multiple forms of defense, and found their attack was able to circumnavigate all of the defenses they tried. They showed that the regions in the input image the network focused on when faced with a trigger was more diverse than in standard backdoor attacks, making it more difficult to identify problematic data.

This paper makes great strides at introducing a key vulnerability to many networks. This exposes vulnerabilities that were previously undocumented or not theorized.

One weakness in the paper is that it focuses solely on image classification tasks. Another shortcoming briefly mentioned by the authors was that the resulting trigger patterns could be easily recognized by a human.

Future work could include creating trigger patterns that appear more realistic, or to expand the domain into other data types. In general the paper itself is a call to researchers to find ways to remove this vulnerability.

### 3 Question 3

#### 3.1 Part A

##### 3.1.1 Deriving Hessian

Start with the loss function:

$$\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N y_i \log(1 + e^{-\langle x_i, w \rangle}) + (1 - y_i) \langle x_i, w \rangle \quad (1)$$

Let:

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & \dots & x_{1d} \\ x_{21} & . & & & x_{2d} \\ \dots & & . & & \\ x_{N1} & & & . & x_{Nd} \end{bmatrix} \quad (2)$$

Where N is the number of training samples, and d is the number of features. Let:

$$W = (w_1, w_2, \dots, w_d) \quad (3)$$

And let  $j \in [1, d]$  Now, the partial derivative of  $\mathcal{L}(w)$  with respect to  $w_j$  is:

$$\begin{aligned} \frac{\partial \mathcal{L}(w)}{\partial w_j} &= \frac{1}{N} \frac{\partial}{\partial w_j} \sum_{i=1}^N y_i \log(1 + e^{-\langle x_i, w \rangle}) + (1 - y_i) \langle x_i, w \rangle \\ &= \frac{1}{N} \sum_{i=1}^N y_i \frac{\partial}{\partial w_j} \log(1 + e^{-\langle x_i, w \rangle}) + \frac{\partial}{\partial w_j} (1 - y_i) \langle x_i, w \rangle \\ &= \frac{1}{N} \sum_{i=1}^N -x_{ij} \frac{1}{1 + e^{W^T x_i}} y_i + (1 - y_i) x_{ij} \end{aligned} \quad (4)$$

Next, take the derivative w.r.t  $w_k$ , where  $k \in [1, d]$

$$\begin{aligned} \frac{\partial^2 \mathcal{L}(w)}{\partial w_j \partial w_k} &= \frac{1}{N} \sum_{i=1}^N -x_{ij} \frac{\partial}{\partial w_k} \left( \frac{1}{1 + e^{W^T x_i}} \right) y_i + 0 \\ &= \frac{1}{N} \sum_{i=1}^N x_{ij} x_{ik} \frac{e^{W^T x_i}}{(1 + e^{W^T x_i})^2} y_i \end{aligned} \quad (5)$$

Let:

$$t_i = \frac{e^{W^T x_i} y_i}{(1 + e^{W^T x_i})^2} \quad (6)$$

Note that since  $y \in [0, 1]$ ,  $t_i \geq 0$  So:

$$T = \begin{bmatrix} t_1 & 0 & 0 & \dots & 0 \\ 0 & t_2 & 0 & \dots & 0 \\ 0 & 0 & t_3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & t_n \end{bmatrix} \quad (7)$$

And:

$$q_i = (x_{1i}, x_{2i} \dots x_{ni})^T \quad (8)$$

For any  $i \in [0, d]$ . These are the columns of the X matrix.

Now:

$$\begin{aligned}\frac{\partial^2 \mathcal{L}(w)}{\partial w_j \partial w_k} &= \frac{1}{N} \sum_{i=1}^N q_j^T T q_k \\ \nabla_w^2 \mathcal{L}(w) &= \frac{1}{N} X^T T X\end{aligned}\tag{9}$$

### 3.1.2 Proving Convex and L Smooth

Since  $t_i \geq 0$  (due to  $y_i = [0, 1]$ ), the eigenvalues for  $T$  must also be  $\geq 0$ . Let  $P \in \mathbf{R}^d$ :

$$\begin{aligned}P^T \nabla \mathcal{L}_w(w)^2 P &= P^T \left( \frac{1}{N} X^T T X \right) P \\ &= \frac{1}{N} P^T X T (P^T X)^T \\ &= \frac{1}{N} \|P^T T X\|^2 \geq 0\end{aligned}\tag{10}$$

For any  $P$  This is positive semi-definite, meaning it is convex, but not strongly.

Next, to prove  $L$  smoothness, show that the maximum and minimum values for  $T$  are bounded:

$$\begin{aligned}\lim_{W^T x_i \rightarrow +\infty} &\approx \frac{e^{W^T x_i}}{e^{2W^T x_i}} \approx 0 \\ \lim_{W^T x_i \rightarrow -\infty} &\approx \frac{e^{W^T x_i}}{1} \approx 0\end{aligned}\tag{11}$$

Therefore the slope must be bounded; i.e.

$$\nabla_w \mathcal{L}(w)^2 \leq L I_d\tag{12}$$

So the loss is  $L$ -Smooth, and convex.

## 3.2 Part B

```

1  # -*- coding: utf-8 -*-
2  """OPT4DL_Hw4_Problem3
3
4  Automatically generated by Colab.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1ITnkKN_lKsvAPFeTOVntvNvVBx_uxN8Q
8  """
9
10 import numpy as np                # advanced math library
11 import matplotlib.pyplot as plt   # MATLAB like plotting routines
12 import random                     # for generating random numbers
13
14 from keras.datasets import mnist  # MNIST dataset is included in Keras
15 from keras.models import Sequential # Model type to be used
16
17 from keras.layers import Dense, Dropout, Activation # Types of layers to be used in our
18     model
19
20 # The MNIST data is split between 60,000 28 x 28 pixel training images and 10,000 28 x 28
21     pixel images
22 (X_train, y_train), (X_test, y_test) = mnist.load_data()
23
24 X_train = X_train.reshape(60000, 784) # reshape 60,000 28 x 28 matrices into 60,000 784-
25     length vectors.

```

```

23 y_train = np.atleast_2d(y_train.T)      # 2d saves headache when multiplying
24 y_train[y_train>=5] = 1                # binary encoding
25 y_train[y_train<5] = 0
26
27 X_test = X_test.reshape(10000, 784)    # reshape 10,000 28 x 28 matrices into 10,000 784-
    length vectors.
28 y_test = np.atleast_2d(y_test.T)      # 2d saves headache when multiplying
29 y_test[y_test>=5] = 1                # binary encoding
30 y_test[y_test<5] = 0
31
32 X_train = X_train.astype('float32')    # change integers to 32-bit floating point numbers
33 X_test = X_test.astype('float32')
34
35 X_train /= 255                        # normalize each value for each pixel for the entire
    vector for each input
36 X_test /= 255
37
38 print("Training matrix shape", X_train.shape)
39 print("Testing matrix shape", X_test.shape)
40 print("Label Matrix Shape:", y_train.shape)
41
42 # Constants and initialization
43 B = 32      # batch size
44 T = 100     # number of epochs
45 eta = 0.0001 # placeholder
46 losses = np.zeros((T,1)) # loss for epoch t
47 w = np.random.normal(0,1.5,size = [784,1])
48 idx = np.arange(60000)
49
50 def sigmoid(x):
51     return 1/(1+np.exp(-x))
52
53 def lossFn(x,w,y):
54     prod = x@w
55     prod = np.clip(prod,0,1)
56     return (y.T@np.log(1+np.exp(-prod)) +(1-y).T@prod)
57
58 for t in range(T):
59     #first, select the batch for this iteration:
60     batchIdx = np.random.choice(X_train.shape[0],size = X_train.shape[0],replace=False)
61     print(f'Epoch ={t}')
62     epochLoss = 0
63     for b in range(batchIdx.shape[0]):
64         thisBatch = batchIdx[b*B:(b+1)*B]
65         X_batch = X_train[thisBatch]
66         y_batch = y_train[thisBatch]
67         # gradient of loss function w.r.t w
68         gradL = X_batch.T@(1+y_batch*sigmoid(X_batch@w)-2*y_batch)
69         w = w-eta*gradL #update rule
70         epochLoss += lossFn(X_batch,w,y_batch)
71     losses[t] = epochLoss/X_train.shape[0] #avg loss/sample
72
73 plt.plot(losses)
74 plt.xlabel('Epoch')
75 plt.ylabel('Loss')
76 plt.title(f'Loss vs Epoch: eta = {eta}')
77 plt.show()
78
79 pred = sigmoid(X_test @ w)
80 pred_labels = (pred >=0.5).astype(int)
81 accuracy = np.mean(pred_labels == y_test)
82 print(f'Accuracy: {accuracy}')

```

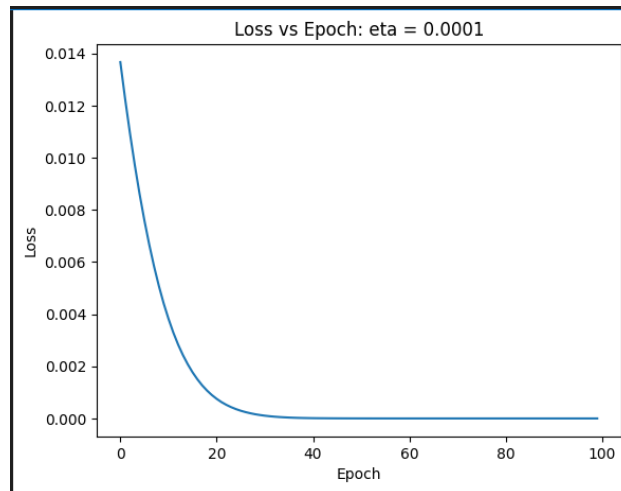


Figure 1: Loss

## 4 Question 4

### 4.1 Explanation

Adding more hidden units to the hidden layer will improve the fitting and allow the network to learn more complex patterns. This comes at the risk of overfitting to the training data and not generalizing well. During training, this would be observed by comparing the training loss/accuracy to the validation loss/accuracy. If the training loss continues to improve while the validation loss stagnates or increases, overfitting has occurred. More neurons may also require more training data/epochs to effectively fit the data.

Decreasing the hidden units to 128 has the opposite problem. The model is more simplistic so struggles to learn complex patterns. This may cause the accuracy to be quite low as the model may never be able to approximate a complex pattern. If the model is underfitting, training and test loss will be high, and accuracy will be low. Few neurons typically require less data/epochs to effectively fit the data, assuming the distribution is simple enough for the reduced set to learn it.

By adding a dropout layer, you effectively combat over fitting by regularizing the network. This prevents a single neuron/ path from dominating the model. This typically results in improved test loss/accuracy since the model is less prone to overfitting. This can reduce training performance, since it will effectively limit the learning capacity of the model. What you gain from this limitation is a more robust model.

### 4.2 Code

My code is really split into two classes: Layer and MLP. The layer class holds members for weights, activation functions, biases, masks (for dropout), and activation input / output for backpropagation. The function 'forward' will handle a single forward pass from this layer, returning the activation output from this layer. The function 'backward' will handle the backpropagation step, which will perform a chain rule and apply the weight and bias update for this layer.

The MLP class holds all of the requested layer objects, and sets up a framework for forward passes and backward passes between layers. This class will also run small batch SGD, calculate cross entropy loss and accuracy for each epoch.

With those classes established, the rest of the code/implementation for each question just relies on instantiating an MLP object and adding layers to it.



```

1  # -*- coding: utf-8 -*-
2  """OPT4DL_HW4_Problem4
3
4  Automatically generated by Colab.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1h9jtgs6dBfbqmKMHQmoa8jXsNORLKP5x
8  """
9
10 import tensorflow as tf
11 import numpy as np
12 import matplotlib.pyplot as plt
13 import pdb
14 np.random.seed(42)
15 #Load fasion MNIST dataset
16 (x_train, y_train), (x_test, y_test) = tf.keras.datasets.fashion_mnist.load_data()
17
18 print(x_train.shape)
19 print(y_train.shape)
20 print(x_test.shape)
21 print(y_test.shape)
22
23 x_train = x_train.reshape(60000,784) # flatten images
24 x_test = x_test.reshape(10000,784)
25 y_train = np.atleast_2d(y_train.T)      # 2d saves headache when multiplying
26 y_test = np.atleast_2d(y_test.T)        # 2d saves headache when multiplying
27
28 x_train = x_train.astype('float32')     # change integers to 32-bit floating point numbers
29 x_test = x_test.astype('float32')
30
31 x_train /= 255                          # normalize each value for each pixel for the entire
    vector for each input
32 x_test /= 255
33
34
35 print(x_train.shape)
36 print(y_train.shape)
37 print(x_test.shape)
38 print(y_test.shape)
39
40 # Parameters
41 B = 256 #batch size
42 num_epochs = 10
43 learning_rate = 0.001
44
45 # Make single layer perceptron
46 class Layer:
47
48     def __init__(self, input_size, output_size, activation = 'relu', dropout_prob=None):
49         #self.weights = np.random.normal(size=(input_size,output_size),loc = 0, scale = 0.01)
50         self.weights = np.random.randn(input_size,output_size)*np.sqrt(2./input_size)
51         self.bias = np.zeros(shape = (1,output_size))
52         self.activation = activation
53         self.activation_output = None
54         self.mask = None
55         self.dropout_prob = dropout_prob
56         self.layer_input = None
57         self.activation_input = None
58
59     def dropout(self,x,p):
60         #hold the mask for backprop
61         self.mask = np.random.binomial(1,1-p,size=x.shape)
62         return x*self.mask / (1-p)
63
64     def relu(self,x):
65         return np.maximum(0,x)
66

```

```

67
68 def forward(self,x,training=True):
69     z = np.dot(x,self.weights) + self.bias
70
71     self.activation_input = z
72     self.layer_input = x
73
74     if self.activation == 'relu':
75         self.activation_output = self.relu(z)
76     elif self.activation == 'dropout' and training:
77         self.activation_output = self.dropout(z,self.dropout_prob)
78     else:
79         self.activation_output = z
80     print(f'Number of dead neurons: {np.sum(self.activation_output == 0)}')
81
82     return self.activation_output
83
84 def backward(self,grad_output,learning_rate,training=True):
85
86     if self.activation == 'relu':
87         dz = grad_output * (self.activation_input > 0).astype(float)
88     elif self.activation == 'dropout' and training:
89         dz = grad_output * self.mask / (1-self.dropout_prob)
90     else:
91         dz = grad_output
92     #compute gradient wrt weights and biases
93     dw = 1/self.layer_input.shape[0]*np.dot(self.layer_input.T,dz)
94     db = np.mean(dz,axis=0,keepdims=True)
95
96     print(f'dw max: {np.max(dw)}')
97     print(f'dw min: {np.min(dw)}')
98     print(f'db max: {np.max(db)}')
99     print(f'db min: {np.min(db)}')
100    #compute gradient wrt input to layer
101    self.weights -= learning_rate * dw
102    self.bias -= learning_rate * db
103    return np.dot(dz,self.weights.T) # Gradient w.r.t. input to the layer
104
105 # Combine single layers into MLP
106
107 class MLP:
108
109     def __init__(self, input_size, output_size, dropout_prob = None):
110         self.layers = []
111         self.input_size = input_size
112         self.output_size = output_size
113         self.dropout_prob = dropout_prob
114         self.epochTrainingLoss = [] #list to store epoch losses for SGD
115         self.epochTrainingAccuracy = [] #list to store epoch accuracies for SGD
116         self.epochValidationLoss = [] # list to store validation losses
117         self.epochValidationAccuracy = [] #list to store validation accuracies for SGD
118
119     def add_layer(self,layer):
120         self.layers.append(layer)
121
122     def softmax(self, x):
123         exp_x = np.exp(x - np.max(x, axis=1, keepdims=True)) # for numerical stability
124         return exp_x / np.sum(exp_x, axis=1, keepdims=True)
125
126     def forward(self,x,training = False):
127         for layer in self.layers:
128             x = layer.forward(x,training)
129             #print(f'x.shape{x.shape}')
130         return x
131
132     def cross_entropy_loss(self,y_true,y_pred):
133         y_pred = np.clip(y_pred,1e-7,1-1e-7) #avoid zero
134         return -np.sum(y_true*np.log(y_pred))/y_true.shape[0]

```

```

135
136 def one_hot_encoding(self,y):
137     y = y.reshape(-1)
138     one_hot = np.zeros((y.shape[0],self.output_size))
139     one_hot[np.arange(y.shape[0]),y] = 1
140     return one_hot
141
142 def compute_loss_gradient(self, y_pred, y_true):
143     # Compute the gradient of the loss with respect to the predicted probabilities
144     return (y_pred - y_true)
145
146 def backward(self, loss_grad, learning_rate):
147     for layer in reversed(self.layers):
148         loss_grad = layer.backward(loss_grad, learning_rate)
149
150 def process_batch(self, x_batch, y_batch, learning_rate):
151
152     # Forward pass
153
154     logits = self.forward(x_batch, training=True)
155     y_pred = self.softmax(logits)
156     print(f'prediction:{np.argmax(y_pred,axis=1)}')
157     print(f'true:{np.argmax(y_batch,axis=1)}')
158     # Compute loss for this batch
159     loss = self.cross_entropy_loss(y_batch, y_pred)
160     # Backward pass
161     loss_grad = self.compute_loss_gradient(y_pred, y_batch)
162     self.backward(loss_grad, learning_rate)
163
164     return loss
165
166
167 def evaluateLoss(self,y_val, y_pred):
168     """
169     Calculate the loss on the validation set.
170     Pass in validation data and predicted values and evaluate loss.
171     """
172     return self.cross_entropy_loss(y_val,y_pred) #note this returns non-normalized, so
        average over data
173
174 def calcAccuracy(self,y,y_pred):
175     """
176     Calculate model accuracy. Take in y validation data, and the prediction
177     from running self.forward(x).
178     """
179     y_pred = np.argmax(y_pred,axis=1)
180     y = np.argmax(y,axis=1)
181
182     return np.mean(y_pred == y)
183
184 def sgd(self,x_train,y_train,x_val,y_val, learning_rate=0.001,epochs = 10, batch_size =
    32):
185     """
186     Run the actual training. This will evaluate loss and accuracy each epoch.
187     """
188     num_samples = x_train.shape[0]
189     # One hot encode the true labels
190     y_train = self.one_hot_encoding(y_train)
191     y_val = self.one_hot_encoding(y_val)
192     # training loop; loop epochs and batches
193     for t in range(epochs):
194         batchIdx = np.random.choice(num_samples,size = num_samples,replace=False)
195         print(f'##### Epoch ={t} #####')
196         num_batches = num_samples // batch_size #number of full batches
197         remaining_samples = num_samples % batch_size #remaining samples
198         epochTrainingLoss = 0
199         for b in range(num_batches):
200             #slice batches

```

```

201         thisBatch = batchIdx[b*batch_size:(b+1)*batch_size]
202         x_batch = x_train[thisBatch,:]
203         y_batch = y_train[thisBatch,:]
204         loss = self.process_batch(x_batch,y_batch,learning_rate)
205         epochTrainingLoss += loss
206
207
208         if remaining_samples > 0:
209             x_batch = x_train[num_batches*batch_size:]
210             y_batch = y_train[num_batches*batch_size:]
211             loss = self.process_batch(x_batch,y_batch,learning_rate)
212             epochTrainingLoss += loss
213         self.epochTrainingLoss.append(epochTrainingLoss/num_samples)
214         #print(f'epochTrainingLoss{epochTrainingLoss/num_samples}')
215         y_predTrain = self.softmax(self.forward(x_train,training=False))
216
217         self.epochTrainingAccuracy.append(self.calcAccuracy(y_train,y_predTrain))
218         # evaluate performance
219         y_predVal = self.softmax(self.forward(x_val,training=False))
220         self.epochValidationLoss.append(self.cross_entropy_loss(y_val,y_predVal)/y_val.shape
221                                     [0])
222         self.epochValidationAccuracy.append(self.calcAccuracy(y_val,y_predVal))
223
224 # Make MLP
225 ## PART A
226
227 inputLayer = Layer(input_size = 784,output_size = 256,activation = 'relu')
228 outputLayer = Layer(input_size = 256,output_size = 10,activation = 'none')
229 mlpA = MLP(input_size = 784,output_size = 10)
230 mlpA.add_layer(inputLayer)
231 mlpA.add_layer(outputLayer)
232 #run sgd
233 mlpA.sgd(x_train, y_train,x_test,y_test, learning_rate, 10, 256)
234
235 ## PART B
236 # 512 hidden
237 mlpB512 = MLP(input_size = 784,output_size = 10)
238 inputLayer512 = Layer(input_size = 784,output_size = 512,activation = 'relu')
239 outputLayer512 = Layer(input_size = 512,output_size = 10,activation = 'none')
240 mlpB512.add_layer(inputLayer512)
241 mlpB512.add_layer(outputLayer512)
242 mlpB512.sgd(x_train, y_train,x_test,y_test, learning_rate, num_epochs, 256)
243
244 # 128 hidden
245 inputLayer128 = Layer(input_size = 784,output_size = 128,activation = 'relu')
246 outputLayer128 = Layer(input_size = 128,output_size = 10,activation = 'none')
247 mlpB128 = MLP(input_size = 784,output_size = 10)
248 mlpB128.add_layer(inputLayer128)
249 mlpB128.add_layer(outputLayer128)
250 mlpB128.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, 256)
251
252 # Part C with each size
253 mlpC256 = MLP(input_size = 784,output_size = 10)
254 inputLayer = Layer(input_size = 784,output_size = 256,activation = 'relu')
255 hiddenLayer = Layer(input_size = 256,output_size = 256,activation = 'relu')
256 outputLayer = Layer(input_size = 256,output_size = 10,activation = 'none')
257 mlpC256.add_layer(inputLayer)
258 mlpC256.add_layer(hiddenLayer)
259 mlpC256.add_layer(outputLayer)
260 mlpC256.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, 256)
261
262 mlpC512 = MLP(input_size = 784,output_size = 10)
263 inputLayer = Layer(input_size = 784,output_size = 512,activation = 'relu')
264 hiddenLayer = Layer(input_size = 512,output_size = 512,activation = 'relu')
265 outputLayer = Layer(input_size = 512,output_size = 10,activation = 'none')
266 mlpC512.add_layer(inputLayer)
267 mlpC512.add_layer(hiddenLayer)

```

```

268 mlpC512.add_layer(outputLayer)
269 mlpC512.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, 256)
270
271 mlpC128 = MLP(input_size = 784, output_size = 10)
272 inputLayer = Layer(input_size = 784, output_size = 128, activation = 'relu')
273 hiddenLayer = Layer(input_size = 128, output_size = 128, activation = 'relu')
274 outputLayer = Layer(input_size = 128, output_size = 10, activation = 'none')
275 mlpC128.add_layer(inputLayer)
276 mlpC128.add_layer(hiddenLayer)
277 mlpC128.add_layer(outputLayer)
278 mlpC128.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, 256)
279
280 ##PART D
281 # add dropout between in and out
282 inputLayer256 = Layer(input_size = 784, output_size = 256, activation = 'relu')
283 dropoutLayer256 = Layer(input_size = 256, output_size = 256, activation = 'dropout',
284                          dropout_prob = 0.5)
285 hiddenLayer256 = Layer(input_size = 256, output_size = 256, activation = 'relu')
286 outputLayer256 = Layer(input_size = 256, output_size = 10, activation = 'none')
287 mlp256_Dropout = MLP(input_size = 784, output_size = 10)
288 mlp256_Dropout.add_layer(inputLayer256)
289 mlp256_Dropout.add_layer(dropoutLayer256)
290 mlp256_Dropout.add_layer(outputLayer256)
291
292 mlp256_Dropout.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, B)
293
294 inputLayer128 = Layer(input_size = 784, output_size = 128, activation = 'relu')
295 dropoutLayer128 = Layer(input_size = 128, output_size = 128, activation = 'dropout',
296                          dropout_prob = 0.5)
297 hiddenLayer128 = Layer(input_size = 128, output_size = 128, activation = 'relu')
298 outputLayer128 = Layer(input_size = 128, output_size = 10, activation = 'none')
299 mlp128_Dropout = MLP(input_size = 784, output_size = 10)
300 mlp128_Dropout.add_layer(inputLayer128)
301 mlp128_Dropout.add_layer(dropoutLayer128)
302 mlp128_Dropout.add_layer(outputLayer128)
303
304 mlp128_Dropout.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, B)
305
306 inputLayer512 = Layer(input_size=784, output_size=512, activation='relu')
307 dropoutLayer512 = Layer(input_size=512, output_size=512, activation='dropout', dropout_prob
308                          =0.5)
309 hiddenLayer512 = Layer(input_size=512, output_size=512, activation='relu')
310 outputLayer512 = Layer(input_size=512, output_size=10, activation='none')
311 mlp512_Dropout = MLP(input_size=784, output_size=10)
312 mlp512_Dropout.add_layer(inputLayer512)
313 mlp512_Dropout.add_layer(dropoutLayer512)
314 mlp512_Dropout.add_layer(outputLayer512)
315
316 mlp512_Dropout.sgd(x_train, y_train, x_test, y_test, learning_rate, num_epochs, B)
317
318 fig,axs = plt.subplots(1,3,figsize = (15,5))
319 axs[0].plot(mlpA.epochTrainingAccuracy,label = '256 Training Accuracy', linestyle = '-')
320 axs[0].plot(mlpA.epochValidationAccuracy,label = '256 Validation Accuracy',linestyle = '-')
321 axs[0].plot(mlpB128.epochTrainingAccuracy,label = '128 Training Accuracy',linestyle = "--")
322 axs[0].plot(mlpB128.epochValidationAccuracy,label = '128 Validation Accuracy',linestyle = '--')
323
324 axs[0].plot(mlpB512.epochTrainingAccuracy,label = '512 Training Accuracy',linestyle = ':')
325 axs[0].plot(mlpB512.epochValidationAccuracy,label = '512 Validation Accuracy',linestyle = ':')
326
327 axs[0].set_title(f'Part A & B: Accuracy Vs. Epoch: eta = {learning_rate}')
328 axs[0].set_xlabel('Epoch')
329 axs[0].set_ylabel('Accuracy')
330 axs[0].legend()
331
332 axs[1].plot(mlpC256.epochTrainingAccuracy,label = '256 Training Accuracy', linestyle = '-')
333 axs[1].plot(mlpC256.epochValidationAccuracy,label = '256 Validation Accuracy',linestyle = '-')

```

```

    ')
331 axs[1].plot(mlpC128.epochTrainingAccuracy,label = '128 Training Accuracy',linestyle = "--")
332 axs[1].plot(mlpC128.epochValidationAccuracy,label = '128 Validation Accuracy',linestyle = '
--')
333 axs[1].plot(mlpC512.epochTrainingAccuracy,label = '512 Training Accuracy',linestyle = ':')
334 axs[1].plot(mlpC512.epochValidationAccuracy,label = '512 Validation Accuracy',linestyle = ':
')
335 axs[1].set_title(f'Part C: Accuracy Vs. Epoch: eta = {learning_rate}')
336 axs[1].set_xlabel('Epoch')
337 axs[1].set_ylabel('Accuracy')
338 axs[1].legend()
339
340 axs[2].plot(mlp256_Dropout.epochTrainingAccuracy,label = '256 Training Accuracy', linestyle
= '-')
341 axs[2].plot(mlp256_Dropout.epochValidationAccuracy,label = '256 Validation Accuracy',
linestyle = '-')
342 axs[2].plot(mlp128_Dropout.epochTrainingAccuracy,label = '128 Training Accuracy',linestyle =
"--")
343 axs[2].plot(mlp128_Dropout.epochValidationAccuracy,label = '128 Validation Accuracy',
linestyle = '--')
344 axs[2].plot(mlp512_Dropout.epochTrainingAccuracy,label = '512 Training Accuracy',linestyle =
':')
345 axs[2].plot(mlp512_Dropout.epochValidationAccuracy,label = '512 Validation Accuracy',
linestyle = ':')
346 axs[2].set_title(f'Part D: Accuracy Vs. Epoch: eta = {learning_rate}')
347 axs[2].set_xlabel('Epoch')
348 axs[2].set_ylabel('Accuracy')
349 axs[2].legend()
350 plt.show()
351 #plt.legend()
352 #plt.show()
353 #plt.plot(mlpB512.epochLoss,label = '512 hidden')
354 #plt.plot(mlpB128.epochLoss,label = '128 hidden')
355 #plt.xlabel('Epoch')
356 #plt.ylabel('Loss')
357 #plt.title(f'Loss vs Epoch; eta = {learning_rate}')
358 #plt.legend()
359 #plt.show()

```

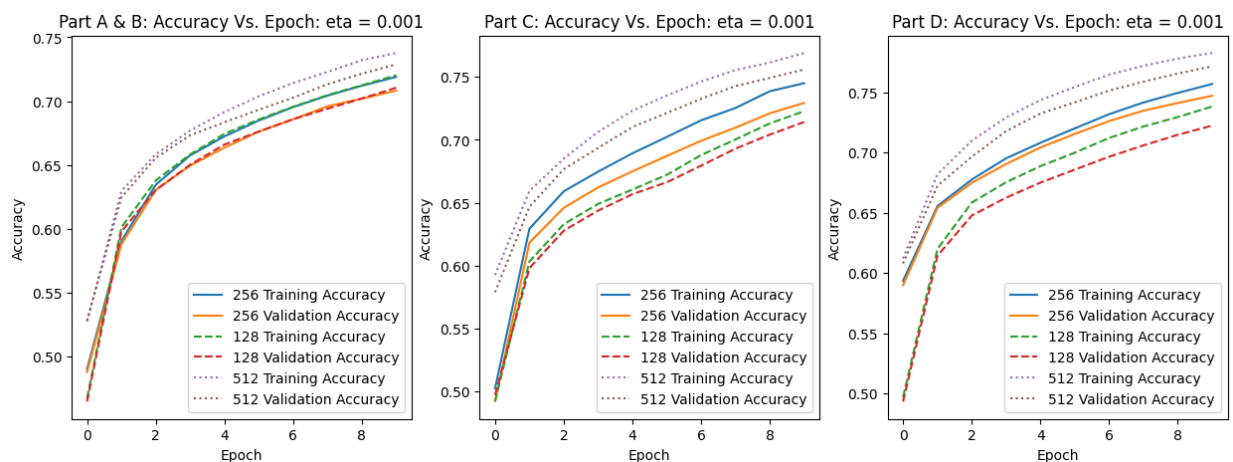


Figure 2: Accuracy Plots