# 1 Introduction

The purpose of this homework assignment is to understand step size optimization in deep learning. The assignment dives into and compares the efficiency of the Stochastic Gradient Descent (SGD), SGD+, and the Adam optimizers.

# 2 Theoretical Background

## 2.1 Gradient Descent

Gradient descent is an iterative approach where the objective is to take a step towards the descent (negative gradient of the hyperplane) and converge at the global minimum point of the loss function and identify the parameters that give the minimum loss. The equation for gradient descent is as follows:

$$p_{t+1} = p_t - \alpha \times g_{t+1} \tag{1}$$

Where $p_t$ is the learnable parameter at parameter $t$, $\alpha$ is the learning rate, and $g_t$ is the gradient of the cost function at parameter $t$

## 2.2 Stochastic Gradient Descent

Stochastic Gradient Descent (SGD) uses the same principle, eq. (1), as the normal gradient descent approach. The only thing it differs in is that unlike gradient descent, SGD doesn't use all the training data at once at each iteration. SGD updates its learnable parameters by using only small batches of randomly drawn training samples. This is because if the current solution point in the parameter hyperplane is at a local minimum in the cost-function surface corresponding to the current batch, it would be highly unlikely that the same point would be a local minimum for the cost-function surfaces corresponding to the future randomly drawn batches of samples.

## 2.3 Stochastic Gradient Descent Plus

Stochastic Gradient Descent Plus (SGD+) is an extension of SGD where we now introduce a factor called momentum, $\mu$. The purpose of momentum is to dampen the oscillation around the minimum by retaining a fraction of the previous gradient. The equation for SGD+ is therefore

$$v_{t+1} = (\mu \times v_t) + g_{t+1} \tag{2}$$

$$p_{t+1} = p_t - \alpha \times v_{t+1} \tag{3}$$

where $v_t$ is the step size at parameter $t$.

## 2.4   Adam Optimizer

The Adam optimizer keeps a running average of both the first and second moment of gradients, and takes both these moments into consideration for calculating the step size, thus adapting the learning rate and converging at the minimum quicker. The equations for Adam is as shown below

$$m_{t+1} = \beta_1 \times m_t + (1 - \beta_1) \times g_{t+1} \tag{4}$$

$$v_{t+1} = \beta_2 \times v_t + (1 - \beta_2) \times g_{t+1}^2 \tag{5}$$

$$p_{t+1} = p_t - \alpha \times \frac{\hat{m}}{\sqrt{\hat{v} + \epsilon}} \tag{6}$$

Where $m_t$ and $v_t$ are moments at parameter $t$ and $\hat{m}$ is

$$\hat{m}_k = \frac{m_k}{\sqrt{1 - \beta_1^k}} \tag{7}$$

and $\hat{v}$ is

$$\hat{v}_k = \frac{v_k}{\sqrt{1 - \beta_2^k}} \tag{8}$$

where $k$ is the iteration $k$ and $v_k$ and $m_k$ is the moment at parameter $t$ at iteration $k$.

# 3 Methodology

Modifications are made to Professor Kak's professor to implement SGD+ and the Adam optimizers.

1. We create two new sub classes that inherit from the **ComputationalGraphPrimer** class for the SGD+ and Adam approach

2. The **run_training_loop_one_neuron_model** and **run_training_loop_multi_neuron_model** methods are created again in their respective subclasses so that the method from the parent class is overridden. These methods are modified to initialize new variables for updating the learnable parameters

   (a) For SGD+, we introduce a list storing all the step sizes, momentum, and a new bias factor that updates the bias based on the momentum. The momentum is set to 0.85

   (b) For Adam, we introduce $\beta_1$, $\beta_2$, $\epsilon$, a list of weights for moment $m$ and moment $v$, and their respective biases. $\beta_1$ is set to 0.9, $\beta_2$ is set to 0.999, and $\epsilon$ is set to $1e-6$

3. The **backprop_and_update_params_one_neuron_model** and **backprop_and_update_params_multi_neuron_model** are created again in their respective subclasses so that the method from the parent class is overridden. These methods are modified to include Equations (2) and (3) for SGD+ and Equations (4)-(8) for Adam to update the step size, learnable parameters, and bias.

4. The losses for SGD, SGD+, and Adam are displayed in a graph for two different learning rates as shown in the results section

## 3.1 One Neuron Classifier

```
1  # Import Libraries
2  import random
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import operator
6  from ComputationalGraphPrimer import *
7
8  # Constants
9  SEED = 512
10 random.seed(SEED)
11 np.random.seed(SEED)
12
13
14 class ComputationalGraphPrimerSGDPlus(ComputationalGraphPrimer):
15     def __init__(self, *args, **kwargs):
16         super().__init__(*args, **kwargs) # Inheriting from the parent
   class
17
18     # Modifying and Overriding the run_training_loop_one_neuron_model to
   implement SGD+
```

```python
 # mu is between [0,1]
 def run_training_loop_one_neuron_model(self, training_data, mu=0.5):
     ####################################Copied from the original
function###################################
     self.vals_for_learnable_params = {param: random.uniform(0,1) for
param in self.learnable_params} # initializing learnable parameters
with random numbers from a uniform distribution over the interval (0,1)

     self.bias = random.uniform(0,1)                      ## Adding the
bias improves class discrimination.
                                                          ##    We
initialize it to a random number.

     class DataLoader:
         """
         To understand the logic of the dataloader, it would help if
you first understand how
         the training dataset is created.  Search for the following
function in this file:

                         gen_training_data(self)

         As you will see in the implementation code for this method,
the training dataset
         consists of a Python dict with two keys, 0 and 1, the former
points to a list of
         all Class 0 samples and the latter to a list of all Class 1
samples.  In each list,
         the data samples are drawn from a multi-dimensional Gaussian
distribution.  The two
         classes have different means and variances.  The
dimensionality of each data sample
         is set by the number of nodes in the input layer of the neural
 network.

         The data loader's job is to construct a batch of samples drawn
 randomly from the two
         lists mentioned above.  And it mush also associate the class
label with each sample
         separately.
         """
         def __init__(self, training_data, batch_size):
             self.training_data = training_data
             self.batch_size = batch_size
             self.class_0_samples = [(item, 0) for item in self.
training_data[0]]   ## Associate label 0 with each sample
             self.class_1_samples = [(item, 1) for item in self.
training_data[1]]   ## Associate label 1 with each sample

         def __len__(self):
             return len(self.training_data[0]) + len(self.training_data
[1])

         def _getitem(self):
```

```python
                    cointoss = random.choice([0,1])
       ## When a batch is created by getbatch(), we want the
       ##   samples to be chosen randomly from the two lists
                    if cointoss == 0:
                        return random.choice(self.class_0_samples)
                    else:
                        return random.choice(self.class_1_samples)

            def getbatch(self):
                batch_data,batch_labels = [],[]
       ## First list for samples, the second for labels
                maxval = 0.0
       ## For approximate batch data normalization
                for _ in range(self.batch_size):
                    item = self._getitem()
                    if np.max(item[0]) > maxval:
                        maxval = np.max(item[0])
                    batch_data.append(item[0])
                    batch_labels.append(item[1])
                batch_data = [item/maxval for item in batch_data]
       ## Normalize batch data
                batch = [batch_data, batch_labels]
                return batch

       #
       ################################################################################

       # Modified part of the function
       self.mu = mu
       self.bias_factor = 0 # Update the bias, the factor depends on the
       current mu
       self.step_sizes = [0 for i in range(len(self.learnable_params) +
       1)]

       ######################################## Copied from the original
       function#################################
       data_loader = DataLoader(training_data, batch_size=self.batch_size
       )
       loss_running_record = []
       i = 0
       avg_loss_over_iterations = 0.0
    ##  Average the loss over iterations for printing out
     ##    every N iterations during the training loop.
       for i in range(self.training_iterations):
           data = data_loader.getbatch()
           data_tuples = data[0]
           class_labels = data[1]
           y_preds, deriv_sigmoids =  self.forward_prop_one_neuron_model(
       data_tuples)                ##  FORWARD PROP of data
           loss = sum([(abs(class_labels[i] - y_preds[i]))**2 for i in
       range(len(class_labels))])  ##  Find loss
           loss_avg = loss / float(len(class_labels))
```

5

```python
                                ##  Average the loss over batch
            avg_loss_over_iterations += loss_avg
            if i%(self.display_loss_how_often) == 0:
                avg_loss_over_iterations /= self.display_loss_how_often
                loss_running_record.append(avg_loss_over_iterations)
                print("[iter=%d]  loss = %.4f" %  (i+1,
    avg_loss_over_iterations))                    ## Display average loss
                avg_loss_over_iterations = 0.0
                                ## Re-initialize avg loss
            y_errors = list(map(operator.sub, class_labels, y_preds))
            y_error_avg = sum(y_errors) / float(len(class_labels))
            deriv_sigmoid_avg = sum(deriv_sigmoids) / float(len(
    class_labels))
            data_tuple_avg = [sum(x) for x in zip(*data_tuples)]
            data_tuple_avg = list(map(operator.truediv, data_tuple_avg,
                                [float(len(class_labels))] * len(
    class_labels) ))
            self.backprop_and_update_params_one_neuron_model(y_error_avg,
    data_tuple_avg, deriv_sigmoid_avg)     ## BACKPROP loss
        # plt.figure()
        # plt.plot(loss_running_record)
        # plt.show()
        return loss_running_record


        #
    ########################################################################


     # Modify backpropagation function for one_neuron_model -
    backpropagating the loss and updating the values of the learnable
    parameters.
     def backprop_and_update_params_one_neuron_model(self, y_error,
    vals_for_input_vars, deriv_sigmoid):
        """
        As should be evident from the syntax used in the following call to
     backprop function,

            self.backprop_and_update_params_one_neuron_model( y_error_avg,
    data_tuple_avg, deriv_sigmoid_avg)
                                                              ^^^
                ^^^                   ^^^
        the values fed to the backprop function for its three arguments
    are averaged over the training
        samples in the batch.  This in keeping with the spirit of SGD that
     calls for averaging the
        information retained in the forward propagation over the samples
    in a batch.

        See Slide 59 of my Week 3 slides for the math of back propagation
    for the One-Neuron network.
        """
        input_vars = self.independent_vars
        vals_for_input_vars_dict = dict(zip(input_vars, list(
    vals_for_input_vars)))
```

```
129         vals_for_learnable_params = self.vals_for_learnable_params
130         for i,param in enumerate(self.vals_for_learnable_params):
131             ## Calculate the next step in the parameter hyperplane
132             ###############################################Modified
    ############################################
133             self.step_sizes[i + 1] = (self.mu * self.step_sizes[i]) + (
    self.learning_rate * y_error * vals_for_input_vars_dict[input_vars[i]]
    * deriv_sigmoid)
134             self.step_sizes[i] = self.step_sizes[i + 1] # Save the new
    current step size in the previous iteration position
135
136             ## Update the learnable parameters
137             # self.vals_for_learnable_params[param] += -self.learning_rate
     * self.step_sizes[i + 1]
138             self.vals_for_learnable_params[param] += self.step_sizes[i +
    1]
139
140         self.bias_factor = (self.mu * self.bias_factor) + (self.
    learning_rate * y_error * deriv_sigmoid) ## Update the bias
141         self.bias += self.bias_factor
142      #
    ###############################################################################

143
144 class ComputationalGraphPrimerAdam(ComputationalGraphPrimer):
145     def __init__(self, *args, **kwargs):
146         super().__init__(*args, **kwargs)
147
148     # Modifying and Overriding the run_training_loop_one_neuron_model to
    implement SGD+
149     # Beta1 and Beta2 are close to 1
150     def run_training_loop_one_neuron_model(self, training_data, beta1=0.9,
     beta2=0.999, e=1e-6):
151         ####################################Copied from the original
    function#################################
152         self.vals_for_learnable_params = {param: random.uniform(0,1) for
    param in self.learnable_params} # initializing learnable parameters
    with random numbers from a uniform distribution over the interval (0,1)
153
154         self.bias = random.uniform(0,1)                      ## Adding the
    bias improves class discrimination.
155                                                              ##   We
    initialize it to a random number.
156
157         class DataLoader:
158             """
159             To understand the logic of the dataloader, it would help if
    you first understand how
160             the training dataset is created.  Search for the following
    function in this file:
161
162                             gen_training_data(self)
163
164             As you will see in the implementation code for this method,
```

```python
        the training dataset
        consists of a Python dict with two keys, 0 and 1, the former
points to a list of
        all Class 0 samples and the latter to a list of all Class 1
samples.  In each list,
        the data samples are drawn from a multi-dimensional Gaussian
distribution.  The two
        classes have different means and variances.  The
dimensionality of each data sample
        is set by the number of nodes in the input layer of the neural
 network.

        The data loader's job is to construct a batch of samples drawn
 randomly from the two
        lists mentioned above.  And it mush also associate the class
label with each sample
        separately.
        """
        def __init__(self, training_data, batch_size):
            self.training_data = training_data
            self.batch_size = batch_size
            self.class_0_samples = [(item, 0) for item in self.
training_data[0]]    ## Associate label 0 with each sample
            self.class_1_samples = [(item, 1) for item in self.
training_data[1]]    ## Associate label 1 with each sample

        def __len__(self):
            return len(self.training_data[0]) + len(self.training_data
[1])

        def _getitem(self):
            cointoss = random.choice([0,1])
 ## When a batch is created by getbatch(), we want the

 ##   samples to be chosen randomly from the two lists
            if cointoss == 0:
                return random.choice(self.class_0_samples)
            else:
                return random.choice(self.class_1_samples)

        def getbatch(self):
            batch_data,batch_labels = [],[]
 ## First list for samples, the second for labels
            maxval = 0.0
 ## For approximate batch data normalization
            for _ in range(self.batch_size):
                item = self._getitem()
                if np.max(item[0]) > maxval:
                    maxval = np.max(item[0])
                batch_data.append(item[0])
                batch_labels.append(item[1])
            batch_data = [item/maxval for item in batch_data]
 ## Normalize batch data
            batch = [batch_data, batch_labels]
```

```
203                    return batch
204
205        #
    ##################################################################

206        # Modified part of the function
207        self.beta1, self.beta2 = beta1, beta2
208        self.e = e
209        self.m_db, self.v_db = 0, 0
210        self.m_dw = [0 for i in range(len(self.learnable_params) + 1)] # m
211        self.v_dw = [0 for i in range(len(self.learnable_params) + 1)] # v
212
213
214        ####################################Copied from the original
    function#################################
215        data_loader = DataLoader(training_data, batch_size=self.batch_size
    )
216        loss_running_record = []
217        i = 0
218        avg_loss_over_iterations = 0.0
    ##   Average the loss over iterations for printing out

219
     ##     every N iterations during the training loop.
220        for i in range(self.training_iterations):
221            data = data_loader.getbatch()
222            data_tuples = data[0]
223            class_labels = data[1]
224            y_preds, deriv_sigmoids =  self.forward_prop_one_neuron_model(
    data_tuples)                 ##  FORWARD PROP of data
225            loss = sum([(abs(class_labels[i] - y_preds[i]))**2 for i in
    range(len(class_labels))])  ##  Find loss
226            loss_avg = loss / float(len(class_labels))
                             ##   Average the loss over batch
227            avg_loss_over_iterations += loss_avg
228            if i%(self.display_loss_how_often) == 0:
229                avg_loss_over_iterations /= self.display_loss_how_often
230                loss_running_record.append(avg_loss_over_iterations)
231                print("[iter=%d]   loss = %.4f" %  (i+1,
    avg_loss_over_iterations))                     ## Display average loss
232                avg_loss_over_iterations = 0.0
                             ## Re-initialize avg loss
233            y_errors = list(map(operator.sub, class_labels, y_preds))
234            y_error_avg = sum(y_errors) / float(len(class_labels))
235            deriv_sigmoid_avg = sum(deriv_sigmoids) / float(len(
    class_labels))
236            data_tuple_avg = [sum(x) for x in zip(*data_tuples)]
237            data_tuple_avg = list(map(operator.truediv, data_tuple_avg,
238                                [float(len(class_labels))] * len(
    class_labels) ))
239            self.backprop_and_update_params_one_neuron_model(y_error_avg,
    data_tuple_avg, deriv_sigmoid_avg, i + 1)      ## BACKPROP loss
240        # plt.figure()
241        # plt.plot(loss_running_record)
242        # plt.show()
```

9

```
243
244        return loss_running_record
245
246      #
    ############################################################################

247
248    # Modify backpropagation function for one_neuron_model -
    backpropagating the loss and updating the values of the learnable
    parameters.
249    def backprop_and_update_params_one_neuron_model(self, y_error,
    vals_for_input_vars, deriv_sigmoid, k):
250        """
251        As should be evident from the syntax used in the following call to
    backprop function,
252
253            self.backprop_and_update_params_one_neuron_model( y_error_avg,
    data_tuple_avg, deriv_sigmoid_avg)
254                                                                    ^^^
                    ^^^                        ^^^
255        the values fed to the backprop function for its three arguments
    are averaged over the training
256        samples in the batch.  This in keeping with the spirit of SGD that
    calls for averaging the
257        information retained in the forward propagation over the samples
    in a batch.
258
259        See Slide 59 of my Week 3 slides for the math of back propagation
    for the One-Neuron network.
260        """
261        input_vars = self.independent_vars
262        vals_for_input_vars_dict = dict(zip(input_vars, list(
    vals_for_input_vars)))
263        vals_for_learnable_params = self.vals_for_learnable_params
264        for i,param in enumerate(self.vals_for_learnable_params):
265            ## Calculate the next step in the parameter hyperplane
266            ################################Modified
    ################################
267            self.m_dw[i + 1] = (self.beta1 * self.m_dw[i]) + ((1 - self.
    beta1) * (self.learning_rate * y_error * vals_for_input_vars_dict[
    input_vars[i]] * deriv_sigmoid))
268            self.m_dw[i] = self.m_dw[i + 1] # Save the new current moment
    in the previous iteration position
269
270            self.v_dw[i + 1] = (self.beta2 * self.v_dw[i]) + ((1 - self.
    beta2) * (self.learning_rate * y_error * vals_for_input_vars_dict[
    input_vars[i]] * deriv_sigmoid)**2)
271            self.v_dw[i] = self.v_dw[i + 1] # Save the new current moment
    in the previous iteration position
272
273            ## Update the learnable parameters
274            mk_hat = self.m_dw[i + 1] / (1 - self.beta1 ** k)
275            vk_hat = self.v_dw[i + 1] / (1 - self.beta2 ** k)
276
```

```
277             self.vals_for_learnable_params[param] += mk_hat / np.sqrt(
        vk_hat + self.e)

278
279         # Inspired by: https://towardsdatascience.com/how-to-implement-an-
        adam-optimizer-from-scratch-76e7b217f1cc
280         # Inspired by: https://www.youtube.com/watch?v=JXQT_vxqwIs&
        ab_channel=DeepLearningAI
281         self.m_db = (self.beta1 * self.m_db) + (1 - self.beta1) * (self.
        learning_rate * y_error * deriv_sigmoid)
282         self.v_db = (self.beta2 * self.v_db) + (1 - self.beta2) * (self.
        learning_rate * y_error * deriv_sigmoid) ** 2

283
284         m_db_hat = self.m_db / (1 - self.beta1 ** k)
285         v_db_hat = self.v_db / (1 - self.beta1 ** k)

286
287         self.bias += m_db_hat / np.sqrt(v_db_hat + self.e) ## Update the
        bias
288      #
        ####################################################################################

289
290 def sgd_plus(lr=1e-3, mu=0.9):
291     cgp = ComputationalGraphPrimerSGDPlus(
292                 one_neuron_model = True,
293                 expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'],
294                 output_vars = ['xw'],
295                 dataset_size = 5000,
296                 learning_rate = lr,
297                 training_iterations = 40000,
298                 batch_size = 8,
299                 display_loss_how_often = 100,
300                 debug = True,
301         )

302
303     cgp.parse_expressions()
304     # cgp.display_one_neuron_network()

305
306     training_data = cgp.gen_training_data()
307     loss_per_iteration = cgp.run_training_loop_one_neuron_model(
        training_data, mu=mu)

308
309     return loss_per_iteration

310
311 def adam(lr=1e-3):
312     cgp = ComputationalGraphPrimerAdam(
313                 one_neuron_model = True,
314                 expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'],
315                 output_vars = ['xw'],
316                 dataset_size = 5000,
317                 learning_rate = lr,
318                 training_iterations = 40000,
319                 batch_size = 8,
320                 display_loss_how_often = 100,
321                 debug = True,
```

```python
322          )
323
324      cgp.parse_expressions()
325      # cgp.display_one_neuron_network()
326
327      training_data = cgp.gen_training_data()
328      loss_per_iteration = cgp.run_training_loop_one_neuron_model(
      training_data )
329
330      return loss_per_iteration
331
332  def sgd(lr=1e-3):
333      cgp = ComputationalGraphPrimer(
334                  one_neuron_model = True,
335                  expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'],
336                  output_vars = ['xw'],
337                  dataset_size = 5000,
338                  learning_rate = lr,
339                  training_iterations = 40000,
340                  batch_size = 8,
341                  display_loss_how_often = 100,
342                  debug = True,
343      )
344
345      cgp.parse_expressions()
346      # cgp.display_one_neuron_network()
347
348      training_data = cgp.gen_training_data()
349      loss_per_iteration = cgp.run_training_loop_one_neuron_model(
      training_data )
350
351      return loss_per_iteration
352
353  def plot_losses(sgd, sgd_plus, adam, lr):
354      number_of_iterations = len(adam)
355      plt.plot(range(number_of_iterations), sgd, label="SGD Loss")
356      plt.plot(range(number_of_iterations), sgd_plus, label="SGD+ Loss")
357      plt.plot(range(number_of_iterations), adam, label="Adam Loss")
358
359      plt.title(f"Loss per Iteration for Different Optimizers for One Neuron
      Model for Learning Rate: {lr}")
360      plt.xlabel("Iteration Number")
361      plt.ylabel("Loss")
362      plt.legend(loc="upper left")
363
364      plt.show(); quit()
365      plt.savefig(r"/Users/nikitaravi/Documents/Academics/ECE 60146/HW3/
      one_neuron_" + str(lr) + "_learning_rate.png", dpi=200)
366
367  if __name__ == "__main__":
368      lr = 1e-3
369      sgd_loss = sgd(lr)
370      sgd_plus_loss = sgd_plus(lr)
371      adam_loss = adam(lr)
```

```
372
373    plot_losses(sgd_loss, sgd_plus_loss, adam_loss, lr)
```

<div align="center">Listing 1: SGD+ and Adam for One Neuron Classifier</div>

## 3.2  Multi-Neuron Classifier

```
1  # Import Libraries
2  import random
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import operator
6  from ComputationalGraphPrimer import *
7
8  # Constants
9  SEED = 1234
10 random.seed(SEED)
11 np.random.seed(SEED)
12
13 class ComputationalGraphPrimerSGDPlus(ComputationalGraphPrimer):
14     def __init__(self, *args, **kwargs):
15         super().__init__(*args, **kwargs) # Inheriting from the parent
   class
16
17     # Modifying and Overriding the run_training_loop_one_neuron_model to
   implement SGD+
18     # mu is between [0,1]
19
20     def run_training_loop_multi_neuron_model(self, training_data, mu=0.5):
21         #####################################Copied from the original
   function###################################
22         class DataLoader:
23             """
24             To understand the logic of the dataloader, it would help if
   you first understand how
25             the training dataset is created.  Search for the following
   function in this file:
26
27                            gen_training_data(self)
28
29             As you will see in the implementation code for this method,
   the training dataset
30             consists of a Python dict with two keys, 0 and 1, the former
   points to a list of
31             all Class 0 samples and the latter to a list of all Class 1
   samples.  In each list,
32             the data samples are drawn from a multi-dimensional Gaussian
   distribution.  The two
33             classes have different means and variances.  The
   dimensionality of each data sample
34             is set by the number of nodes in the input layer of the neural
    network.
35
```

```
36              The data loader's job is to construct a batch of samples drawn
       randomly from the two
37              lists mentioned above.  And it mush also associate the class
       label with each sample
38              separately.
39              """
40          def __init__(self, training_data, batch_size):
41              self.training_data = training_data
42              self.batch_size = batch_size
43              self.class_0_samples = [(item, 0) for item in self.
       training_data[0]]    ## Associate label 0 with each sample
44              self.class_1_samples = [(item, 1) for item in self.
       training_data[1]]    ## Associate label 1 with each sample
45
46          def __len__(self):
47              return len(self.training_data[0]) + len(self.training_data
       [1])
48
49          def _getitem(self):
50              cointoss = random.choice([0,1])
        ## When a batch is created by getbatch(), we want the

51
        ##    samples to be chosen randomly from the two lists
52              if cointoss == 0:
53                  return random.choice(self.class_0_samples)
54              else:
55                  return random.choice(self.class_1_samples)
56
57          def getbatch(self):
58              batch_data,batch_labels = [],[]
        ## First list for samples, the second for labels
59              maxval = 0.0
        ## For approximate batch data normalization
60              for _ in range(self.batch_size):
61                  item = self._getitem()
62                  if np.max(item[0]) > maxval:
63                      maxval = np.max(item[0])
64                  batch_data.append(item[0])
65                  batch_labels.append(item[1])
66              batch_data = [item/maxval for item in batch_data]
        ## Normalize batch data
67              batch = [batch_data, batch_labels]
68              return batch
69      #
    #######################################################################

70      # Modified part of the function
71      self.mu = mu
72      self.step_sizes = [0 for i in range(len(self.learnable_params) +
    1)]
73      self.bias_factor = 0

74
75      #########################################Copied from the original
    function#################################
```

```
76          """
77          The training loop must first initialize the learnable parameters.
     Remember , these are the
78          symbolic names in your input expressions for the neural layer that
     do not begin with the
79          letter 'x'.  In this case , we are initializing with random numbers
     from a uniform distribution
80          over the interval (0 ,1).
81          """
82          self.vals_for_learnable_params = {param: random.uniform (0 ,1) for
     param in self.learnable_params}
83
84          self.bias = [random.uniform (0 ,1) for _ in range (self.num_layers -1)
     ]       ## Adding the bias to each layer improves
85
                ##   class discrimination. We initialize it
86
                ##   to a random number.
87
88          data_loader = DataLoader (training_data , batch_size=self.batch_size
     )
89          loss_running_record = []
90          i = 0
91          avg_loss_over_iterations = 0.0
              ##  Average the loss over iterations for printing out
92
                ##     every N iterations during the training loop.
93          for i in range (self.training_iterations):
94              data = data_loader.getbatch ()
95              data_tuples = data [0]
96              class_labels = data [1]
97              self.forward_prop_multi_neuron_model (data_tuples)
                        ## FORW PROP works by side -effect
98              predicted_labels_for_batch = self.forw_prop_vals_at_layers [
     self.num_layers -1]      ## Predictions from FORW PROP
99              y_preds = [item for sublist in  predicted_labels_for_batch
     for item in sublist]  ## Get numeric vals for predictions
100             loss = sum ([( abs (class_labels [i] - y_preds [i]))**2 for i in
     range (len (class_labels))])  ## Calculate loss for batch
101             loss_avg = loss / float (len (class_labels))
                        ## Average the loss over batch
102             avg_loss_over_iterations += loss_avg
                        ## Add to Average loss over iterations
103             if i%(self.display_loss_how_often) == 0:
104                 avg_loss_over_iterations /= self.display_loss_how_often
105                 loss_running_record.append (avg_loss_over_iterations)
106                 print ("[iter=%d]  loss = %.4f" % (i+1,
     avg_loss_over_iterations))              ## Display avg loss
107                 avg_loss_over_iterations = 0.0
                        ## Re-initialize avg -over -iterations loss
108             y_errors = list (map (operator.sub, class_labels , y_preds))
109             y_error_avg = sum (y_errors) / float (len (class_labels))
110             self.backprop_and_update_params_multi_neuron_model (y_error_avg
     , class_labels)     ## BACKPROP loss
```

15

```
111         # plt.figure()
112         # plt.plot(loss_running_record)
113         # plt.show()
114
115         return loss_running_record
116
117         #
     ###########################################################################

118
119      # Modify backpropagation function for one_neuron_model -
     backpropagating the loss and updating the values of the learnable
     parameters.
120      def backprop_and_update_params_multi_neuron_model(self, y_error,
     class_labels):
121         """
122         First note that loop index variable 'back_layer_index' starts with
      the index of
123         the last layer.  For the 3-layer example shown for 'forward',
     back_layer_index
124         starts with a value of 2, its next value is 1, and that's it.
125
126         Stochastic Gradient Gradient calls for the backpropagated loss to
     be averaged over
127         the samples in a batch.  To explain how this averaging is carried
     out by the
128         backprop function, consider the last node on the example shown in
     the forward()
129         function above.  Standing at the node, we look at the 'input'
     values stored in the
130         variable "input_vals".  Assuming a batch size of 8, this will be
     list of
131         lists. Each of the inner lists will have two values for the two
     nodes in the
132         hidden layer. And there will be 8 of these for the 8 elements of
     the batch.  We average
133         these values 'input vals' and store those in the variable "
     input_vals_avg".  Next we
134         must carry out the same batch-based averaging for the partial
     derivatives stored in the
135         variable "deriv_sigmoid".
136
137         Pay attention to the variable 'vars_in_layer'.  These store the
     node variables in
138         the current layer during backpropagation.  Since back_layer_index
     starts with a
139         value of 2, the variable 'vars_in_layer' will have just the single
      node for the
140         example shown for forward(). With respect to what is stored in
     vars_in_layer', the
141         variables stored in 'input_vars_to_layer' correspond to the input
     layer with
142         respect to the current layer.
143         """
```

```
144        # backproped prediction error:
145        pred_err_backproped_at_layers = {i : [] for i in range(1,self.
   num_layers-1)}
146        pred_err_backproped_at_layers[self.num_layers-1] = [y_error]
147        for back_layer_index in reversed(range(1,self.num_layers)):
148            input_vals = self.forw_prop_vals_at_layers[back_layer_index
   -1]
149            input_vals_avg = [sum(x) for x in zip(*input_vals)]
150            input_vals_avg = list(map(operator.truediv, input_vals_avg, [
   float(len(class_labels))] * len(class_labels)))
151            deriv_sigmoid =  self.gradient_vals_for_layers[
   back_layer_index]
152            deriv_sigmoid_avg = [sum(x) for x in zip(*deriv_sigmoid)]
153            deriv_sigmoid_avg = list(map(operator.truediv,
   deriv_sigmoid_avg,
154                                                          [float(len(
   class_labels))] * len(class_labels)))
155            vars_in_layer  =  self.layer_vars[back_layer_index]
         ## a list like ['xo']
156            vars_in_next_layer_back  =  self.layer_vars[back_layer_index -
    1]    ## a list like ['xw', 'xz']
157
158            layer_params = self.layer_params[back_layer_index]
159            ## note that layer_params are stored in a dict like
160                 ##     {1: [['ap', 'aq', 'ar', 'as'], ['bp', 'bq', 'br', '
   bs']], 2: [['cp', 'cq']]}
161            ## "layer_params[idx]" is a list of lists for the link weights
    in layer whose output nodes are in layer "idx"
162            transposed_layer_params = list(zip(*layer_params))        ##
   creating a transpose of the link matrix
163
164            backproped_error = [None] * len(vars_in_next_layer_back)
165            for k,varr in enumerate(vars_in_next_layer_back):
166                for j,var2 in enumerate(vars_in_layer):
167                    backproped_error[k] = sum([self.
   vals_for_learnable_params[transposed_layer_params[k][i]] *
168
   pred_err_backproped_at_layers[back_layer_index][i]
169                                                        for i in range(len(
   vars_in_layer))])
170 #                                            deriv_sigmoid_avg[i] for i
    in range(len(vars_in_layer))])
171            pred_err_backproped_at_layers[back_layer_index - 1]  =
   backproped_error
172            input_vars_to_layer = self.layer_vars[back_layer_index-1]
173            for j,var in enumerate(vars_in_layer):
174                layer_params = self.layer_params[back_layer_index][j]
175                ##  Regarding the parameter update loop that follows, see
   the Slides 74 through 77 of my Week 3
176                ##  lecture slides for how the parameters are updated
   using the partial derivatives stored away
177                ##  during forward propagation of data. The theory
   underlying these calculations is presented
178                ##  in Slides 68 through 71.
```

```python
                for i,param in enumerate(layer_params):
                    gradient_of_loss_for_param = input_vals_avg[i] *
    pred_err_backproped_at_layers[back_layer_index][j]
                        ###############################################Modified
    #################################
                    self.step_sizes[i + 1] = (self.mu * self.step_sizes[i
    ]) + (self.learning_rate * gradient_of_loss_for_param *
    deriv_sigmoid_avg[j])
                    self.step_sizes[i] = self.step_sizes[i + 1]
                    self.vals_for_learnable_params[param] += self.
    step_sizes[i + 1]

            self.bias_factor = (self.mu * self.bias_factor) + (self.
    learning_rate * sum(pred_err_backproped_at_layers[back_layer_index]) \

     * sum(deriv_sigmoid_avg)/len(deriv_sigmoid_avg))
            self.bias[back_layer_index -1] += self.bias_factor
     #
    ###############################################################################


class ComputationalGraphPrimerAdam(ComputationalGraphPrimer):
    def __init__(self, *args, **kwargs):
        super().__init__(*args, **kwargs) # Inheriting from the parent
    class

    # Modifying and Overriding the run_training_loop_one_neuron_model to
    implement SGD+
    # mu is between [0,1]

    def run_training_loop_multi_neuron_model(self, training_data, beta1
    =0.9, beta2=0.999, e=1e-6):
        #####################################Copied from the original
    function#################################
        class DataLoader:
            """
            To understand the logic of the dataloader, it would help if
    you first understand how
            the training dataset is created.  Search for the following
    function in this file:

                            gen_training_data(self)

            As you will see in the implementation code for this method,
    the training dataset
            consists of a Python dict with two keys, 0 and 1, the former
    points to a list of
            all Class 0 samples and the latter to a list of all Class 1
    samples.  In each list,
            the data samples are drawn from a multi-dimensional Gaussian
    distribution.  The two
            classes have different means and variances.  The
    dimensionality of each data sample
            is set by the number of nodes in the input layer of the neural
```

```
         network.

         The data loader's job is to construct a batch of samples drawn
    randomly from the two
         lists mentioned above.  And it mush also associate the class
    label with each sample
         separately.
         """
         def __init__(self, training_data, batch_size):
             self.training_data = training_data
             self.batch_size = batch_size
             self.class_0_samples = [(item, 0) for item in self.
    training_data[0]]    ## Associate label 0 with each sample
             self.class_1_samples = [(item, 1) for item in self.
    training_data[1]]    ## Associate label 1 with each sample

         def __len__(self):
             return len(self.training_data[0]) + len(self.training_data
    [1])

         def _getitem(self):
             cointoss = random.choice([0,1])
     ## When a batch is created by getbatch(), we want the

     ##    samples to be chosen randomly from the two lists
             if cointoss == 0:
                 return random.choice(self.class_0_samples)
             else:
                 return random.choice(self.class_1_samples)

         def getbatch(self):
             batch_data,batch_labels = [],[]
     ## First list for samples, the second for labels
             maxval = 0.0
     ## For approximate batch data normalization
             for _ in range(self.batch_size):
                 item = self._getitem()
                 if np.max(item[0]) > maxval:
                     maxval = np.max(item[0])
                 batch_data.append(item[0])
                 batch_labels.append(item[1])
             batch_data = [item/maxval for item in batch_data]
     ## Normalize batch data
             batch = [batch_data, batch_labels]
             return batch
        #
    ##################################################################################

        # Modified part of the function
        self.beta1, self.beta2 = beta1, beta2
        self.e = e
        self.m_db, self.v_db = 0.0, 0.0
        self.m_dw = [0.0 for i in range(len(self.learnable_params) + 1)] #
    m
```

```python
253             self.v_dw = [0.0 for i in range(len(self.learnable_params) + 1)] #
      v


256             ########################################Copied from the original
      function#################################
257             """
258             The training loop must first initialize the learnable parameters.
      Remember, these are the
259             symbolic names in your input expressions for the neural layer that
      do not begin with the
260             letter 'x'.  In this case, we are initializing with random numbers
      from a uniform distribution
261             over the interval (0,1).
262             """
263             self.vals_for_learnable_params = {param: random.uniform(0,1) for
      param in self.learnable_params}

265             self.bias = [random.uniform(0,1) for _ in range(self.num_layers-1)
      ]        ## Adding the bias to each layer improves

266                 ##    class discrimination. We initialize it

267                 ##    to a random number.

269             data_loader = DataLoader(training_data, batch_size=self.batch_size
      )
270             loss_running_record = []
271             i = 0
272             avg_loss_over_iterations = 0.0
                 ##  Average the loss over iterations for printing out

273                 ##     every N iterations during the training loop.
274             for i in range(self.training_iterations):
275                 data = data_loader.getbatch()
276                 data_tuples = data[0]
277                 class_labels = data[1]
278                 self.forward_prop_multi_neuron_model(data_tuples)
                            ## FORW PROP works by side-effect
279                 predicted_labels_for_batch = self.forw_prop_vals_at_layers[
      self.num_layers-1]      ## Predictions from FORW PROP
280                 y_preds =  [item for sublist in  predicted_labels_for_batch
      for item in sublist]  ## Get numeric vals for predictions
281                 loss = sum([(abs(class_labels[i] - y_preds[i]))**2 for i in
      range(len(class_labels))])  ## Calculate loss for batch
282                 loss_avg = loss / float(len(class_labels))
                            ## Average the loss over batch
283                 avg_loss_over_iterations += loss_avg
                            ## Add to Average loss over iterations
284                 if i%(self.display_loss_how_often) == 0:
285                     avg_loss_over_iterations /= self.display_loss_how_often
286                     loss_running_record.append(avg_loss_over_iterations)
287                     print("[iter=%d]  loss = %.4f" %  (i+1,
      avg_loss_over_iterations))                ## Display avg loss
```

```
288              avg_loss_over_iterations = 0.0
                         ## Re-initialize avg-over-iterations loss
289          y_errors = list(map(operator.sub, class_labels, y_preds))
290          y_error_avg = sum(y_errors) / float(len(class_labels))
291          self.backprop_and_update_params_multi_neuron_model(y_error_avg
    , class_labels, i+1)      ## BACKPROP loss
292      # plt.figure()
293      # plt.plot(loss_running_record)
294      # plt.show()
295
296      return loss_running_record
297
298      #
    ##############################################################################

299
300   # Modify backpropagation function for one_neuron_model -
    backpropagating the loss and updating the values of the learnable
    parameters.
301   def backprop_and_update_params_multi_neuron_model(self, y_error,
    class_labels, iteration):
302      """
303      First note that loop index variable 'back_layer_index' starts with
    the index of
304      the last layer.  For the 3-layer example shown for 'forward',
    back_layer_index
305      starts with a value of 2, its next value is 1, and that's it.
306
307      Stochastic Gradient Gradient calls for the backpropagated loss to
    be averaged over
308      the samples in a batch.  To explain how this averaging is carried
    out by the
309      backprop function, consider the last node on the example shown in
    the forward()
310      function above.  Standing at the node, we look at the 'input'
    values stored in the
311      variable "input_vals".  Assuming a batch size of 8, this will be
    list of
312      lists. Each of the inner lists will have two values for the two
    nodes in the
313      hidden layer. And there will be 8 of these for the 8 elements of
    the batch.  We average
314      these values 'input vals' and store those in the variable "
    input_vals_avg".  Next we
315      must carry out the same batch-based averaging for the partial
    derivatives stored in the
316      variable "deriv_sigmoid".
317
318      Pay attention to the variable 'vars_in_layer'.  These store the
    node variables in
319      the current layer during backpropagation.  Since back_layer_index
    starts with a
320      value of 2, the variable 'vars_in_layer' will have just the single
     node for the
```

```
321        example shown for forward(). With respect to what is stored in
       vars_in_layer', the
322        variables stored in 'input_vars_to_layer' correspond to the input
       layer with
323        respect to the current layer.
324        """
325        # backproped prediction error:
326        pred_err_backproped_at_layers = {i : [] for i in range(1,self.
       num_layers-1)}
327        pred_err_backproped_at_layers[self.num_layers-1] = [y_error]
328        for back_layer_index in reversed(range(1,self.num_layers)):
329            input_vals = self.forw_prop_vals_at_layers[back_layer_index
       -1]
330            input_vals_avg = [sum(x) for x in zip(*input_vals)]
331            input_vals_avg = list(map(operator.truediv, input_vals_avg, [
       float(len(class_labels))] * len(class_labels)))
332            deriv_sigmoid =  self.gradient_vals_for_layers[
       back_layer_index]
333            deriv_sigmoid_avg = [sum(x) for x in zip(*deriv_sigmoid)]
334            deriv_sigmoid_avg = list(map(operator.truediv,
       deriv_sigmoid_avg,
335                                            [float(len(
       class_labels))] * len(class_labels)))
336            vars_in_layer  =  self.layer_vars[back_layer_index]
       ## a list like ['xo']
337            vars_in_next_layer_back  =  self.layer_vars[back_layer_index -
        1]    ## a list like ['xw', 'xz']
338
339            layer_params = self.layer_params[back_layer_index]
340            ## note that layer_params are stored in a dict like
341                ##    {1: [['ap', 'aq', 'ar', 'as'], ['bp', 'bq', 'br', '
       bs']], 2: [['cp', 'cq']]}
342            ## "layer_params[idx]" is a list of lists for the link weights
        in layer whose output nodes are in layer "idx"
343            transposed_layer_params = list(zip(*layer_params))        ##
       creating a transpose of the link matrix
344
345            backproped_error = [None] * len(vars_in_next_layer_back)
346            for k,varr in enumerate(vars_in_next_layer_back):
347                for j,var2 in enumerate(vars_in_layer):
348                    backproped_error[k] = sum([self.
       vals_for_learnable_params[transposed_layer_params[k][i]] *
349
       pred_err_backproped_at_layers[back_layer_index][i]
350                                            for i in range(len(
       vars_in_layer))])
351 #                                          deriv_sigmoid_avg[i] for i
        in range(len(vars_in_layer))])
352            pred_err_backproped_at_layers[back_layer_index - 1]  =
       backproped_error
353            input_vars_to_layer = self.layer_vars[back_layer_index-1]
354            for j,var in enumerate(vars_in_layer):
355                layer_params = self.layer_params[back_layer_index][j]
356                ##  Regarding the parameter update loop that follows, see
```

```python
                the Slides 74 through 77 of my Week 3
                ##  lecture slides for how the parameters are updated
    using the partial derivatives stored away
                ##  during forward propagation of data. The theory
    underlying these calculations is presented
                ##  in Slides 68 through 71.
                for i,param in enumerate(layer_params):
                    gradient_of_loss_for_param = input_vals_avg[i] *
    pred_err_backproped_at_layers[back_layer_index][j]
                    ############################################Modified
    ###################################
                    self.m_dw[i + 1] = (self.beta1 * self.m_dw[i]) + (1 -
    self.beta1) * (gradient_of_loss_for_param * deriv_sigmoid_avg[j])
                    self.m_dw[i] = self.m_dw[i + 1] # Save the new current
     moment in the previous iteration position

                    self.v_dw[i + 1] = (self.beta2 * self.v_dw[i]) + (1 -
    self.beta2) * (gradient_of_loss_for_param * deriv_sigmoid_avg[j]) ** 2
                    self.v_dw[i] = self.v_dw[i + 1] # Save the new current
     moment in the previous iteration position

                    ## Update the learnable parameters
                    mk_hat = self.m_dw[i + 1] / (1 - self.beta1 **
    iteration)
                    vk_hat = self.v_dw[i + 1] / (1 - self.beta2 **
    iteration)

                    self.vals_for_learnable_params[param] += self.
    learning_rate * mk_hat / np.sqrt(vk_hat + self.e)

            # Inspired by: https://towardsdatascience.com/how-to-implement
    -an-adam-optimizer-from-scratch-76e7b217f1cc
            # Inspired by: https://www.youtube.com/watch?v=JXQT_vxqwIs&
    ab_channel=DeepLearningAI
            self.m_db = (self.beta1 * self.m_db) + (1 - self.beta1) * (sum
    (pred_err_backproped_at_layers[back_layer_index]) \

     * sum(deriv_sigmoid_avg)/len(deriv_sigmoid_avg))
            self.v_db = (self.beta2 * self.v_db) + (1 - self.beta2) * (sum
    (pred_err_backproped_at_layers[back_layer_index]) \

     * sum(deriv_sigmoid_avg)/len(deriv_sigmoid_avg)) ** 2

            m_db_hat = self.m_db / (1 - self.beta1 ** iteration)
            v_db_hat = self.v_db / (1 - self.beta2 ** iteration)

            self.bias[back_layer_index-1] += self.learning_rate * m_db_hat
     / np.sqrt(v_db_hat + self.e) ## Update the bias

     #
    ####################################################################################


def sgd_plus(lr=1e-3, mu=0.9):
```

```python
      cgp = ComputationalGraphPrimerSGDPlus(
                num_layers = 3,
                layers_config = [4,2,1],                          # num of
    nodes in each layer
                expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
                                'xz=bp*xp+bq*xq+br*xr+bs*xs',
                                'xo=cp*xw+cq*xz'],
                output_vars = ['xo'],
                dataset_size = 5000,
                learning_rate = lr,
                training_iterations = 40000,
                batch_size = 8,
                display_loss_how_often = 100,
                debug = True,
        )

    cgp.parse_multi_layer_expressions()
    # cgp.display_multi_neuron_network()

    training_data = cgp.gen_training_data()
    loss_per_iteration = cgp.run_training_loop_multi_neuron_model(
    training_data, mu)

    return loss_per_iteration

def adam(lr=1e-3):
    cgp = ComputationalGraphPrimerAdam(
                num_layers = 3,
                layers_config = [4,2,1],                          # num of
    nodes in each layer
                expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
                                'xz=bp*xp+bq*xq+br*xr+bs*xs',
                                'xo=cp*xw+cq*xz'],
                output_vars = ['xo'],
                dataset_size = 5000,
                learning_rate = lr,
                training_iterations = 40000,
                batch_size = 8,
                display_loss_how_often = 100,
                debug = True,
        )

    cgp.parse_multi_layer_expressions()
    # cgp.display_multi_neuron_network()

    training_data = cgp.gen_training_data()
    loss_per_iteration = cgp.run_training_loop_multi_neuron_model(
    training_data )

    return loss_per_iteration

def sgd(lr=1e-3):
    cgp = ComputationalGraphPrimer(
                num_layers = 3,
```

```
440                    layers_config = [4,2,1],                                    # num of
      nodes in each layer
441                    expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
442                                   'xz=bp*xp+bq*xq+br*xr+bs*xs',
443                                   'xo=cp*xw+cq*xz'],
444                    output_vars = ['xo'],
445                    dataset_size = 5000,
446                    learning_rate = lr,
447                    training_iterations = 40000,
448                    batch_size = 8,
449                    display_loss_how_often = 100,
450                    debug = True,
451        )
452
453     cgp.parse_multi_layer_expressions()
454     training_data = cgp.gen_training_data()
455
456     loss_per_iteration = cgp.run_training_loop_multi_neuron_model(
      training_data )
457     return loss_per_iteration
458
459 def plot_losses(sgd, sgd_plus, adam, lr):
460     number_of_iterations = len(adam)
461     plt.plot(range(number_of_iterations), sgd, label="SGD Loss")
462     plt.plot(range(number_of_iterations), sgd_plus, label="SGD+ Loss")
463     plt.plot(range(number_of_iterations), adam, label="Adam Loss")
464
465     plt.title(f"Loss per Iteration for Different Optimizers for Multi-
      Neuron Model for Learning Rate: {lr}")
466     plt.xlabel("Iteration Number")
467     plt.ylabel("Loss")
468     plt.legend(loc="lower left")
469
470     plt.show(); quit()
471     plt.savefig(r"/Users/nikitaravi/Documents/Academics/ECE 60146/HW3/
      multi_neuron_" + str(lr) + "_learning_rate.png", dpi=200)
472
473 if __name__ == "__main__":
474     lr = 1e-3
475     sgd_loss = sgd(lr)
476     sgd_plus_loss = sgd_plus(lr)
477     adam_loss = adam(lr)
478
479     plot_losses(sgd_loss, sgd_plus_loss, adam_loss, lr)
```

Listing 2: SGD+ and Adam for Multi-Neuron Classifier

# 4   Results

The following figure illustrates the loss from SGD, SGD+, and Adam optimizers at each iteration for a learning rate of 0.001 and 0.005 when using either one-neuron classifier or multi-neuron classifier.
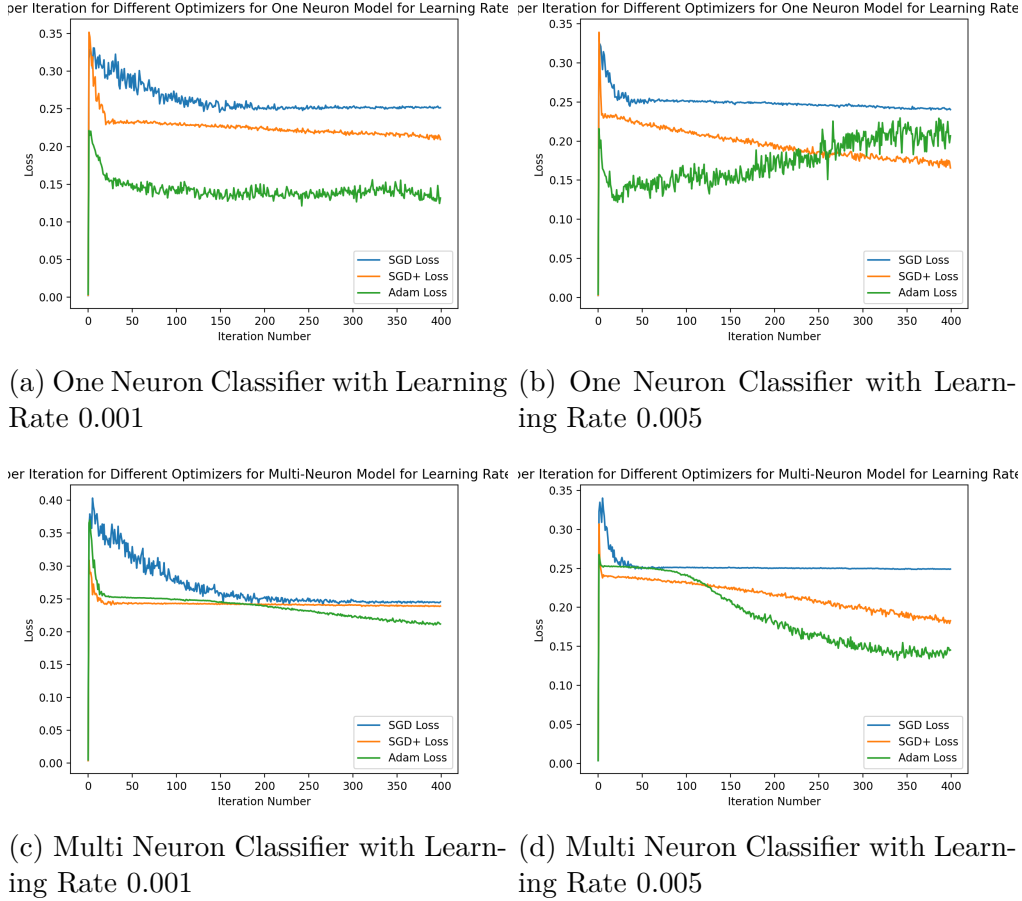


(a) One Neuron Classifier with Learning Rate 0.001

(b) One Neuron Classifier with Learning Rate 0.005

(c) Multi Neuron Classifier with Learning Rate 0.001

(d) Multi Neuron Classifier with Learning Rate 0.005

Figure 1: Losses per Iteration

# 5   Evaluation

In the one neuron classifier, the Adam optimizer outperformed SGD+ and SGD for a learning rate of 0.001 because it achieved less loss than the other two. However, the Adam optimizer oscillates a lot more as it attempts to converge at the global minimum unlike SGD+ which seems to have converged a lot quicker than SGD. On the other hand, with a learning rate of 0.005, SGD and SGD+ have similar performances but the Adam optimizer get significantly worse at each iteration.

In the multi-neuron classifier, the Adam optimizer significantly outperforms the SGD+ and SGD for both learning rates as it gets lower loss at each iteration. For the lower learning rate, both SGD+ and Adam are close to converging to a minimum whereas with the learning

rate of 0.005, both SGD+ and Adam keep oscillating and decreasing at each iteration, which suggests its taking longer to reach the minimum of the cost function.