# 1- Exploring the Primer

Figure 1 shows the loss-vs-iterations plot produced by multi_neuron_classifer.py and Figure 2 shows the loss-vs-iterations plot produced by verify_with_torchnn.py in the multi neuron configuration. From inspecting these plots, we can see that both models are learning since the loss is continually decreasing. However, when we take a closer look at the values of the loss. The torch.nn model performs worse in as its final loss value hovered around 6, while the manually created multi neuron model's loss hovered around 0.10. We also notice that in Figure 1 the loss didn't improve after about 60 iterations suggesting that learning is complete. However, in Figure 2 we can see that the loss is continually improving suggesting that more iterations of learning can take place. We can also notice that torch.nn is more smooth in its loss than the manual method.
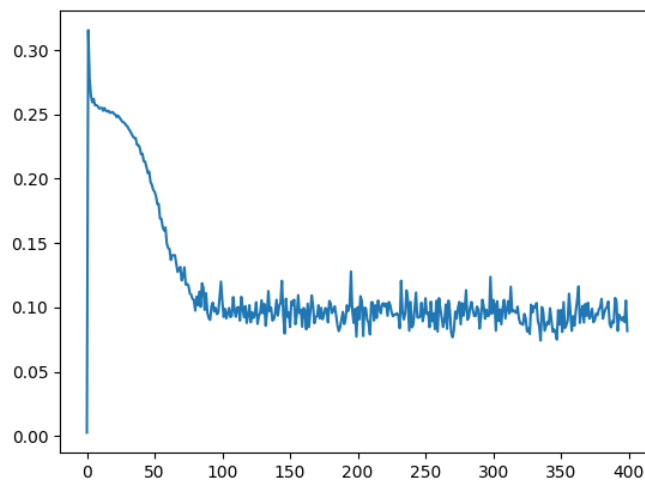


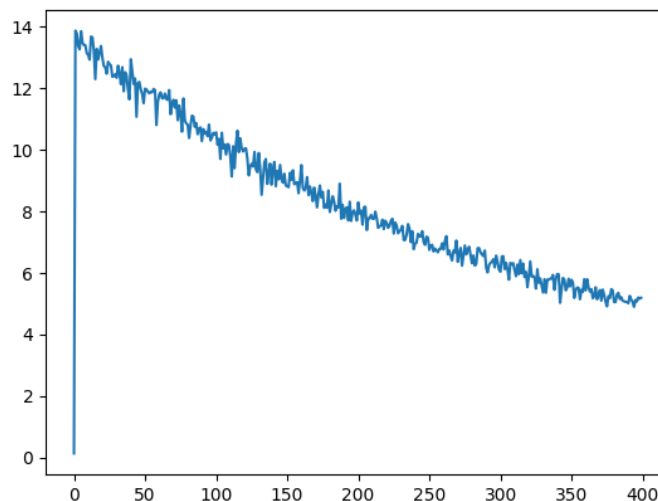Figure 1: loss-vs-iteration plot of multi_neuron_classifer.py



Figure 2: loss-vs-iteration plot of verify_with_torchnn.py (multi-neuron config)

Figure 3 and Figure 4 show the same plots as Figure 1 and Figure 2 but for the one-neuron scenario. A similarity in Figures 3 and 4 is that both have decreasing loss suggesting that both models are indeed learning. However, it's the rate at which they are learning is the difference. For instance, Figure 3's loss is more linear whereas Figure 4's loss is more exponential decay. However, Figure 4's model performs much worse as its loss hovers around 1 at the end of training, but the loss for Figure 3's model hovers around 0.20 at the end of training. Nevertheless, using torch.nn is still better to use as it includes optimizations to better update the parameters.
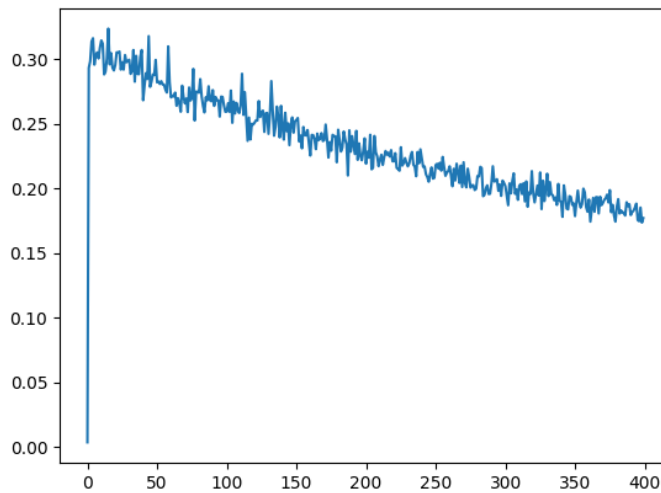


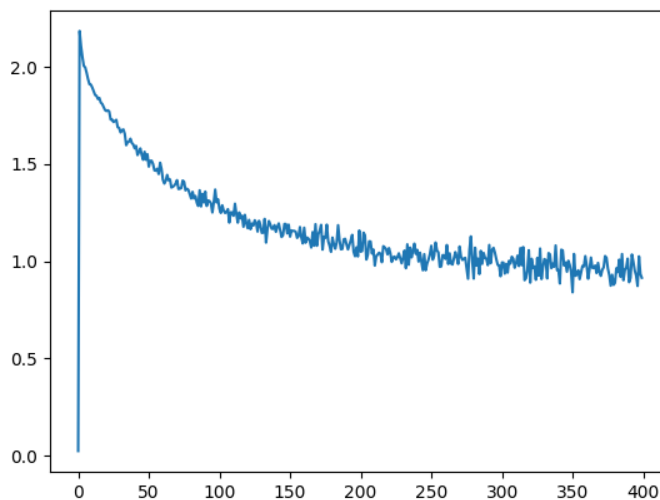Figure 3: loss-vs-iteration plot of one_neuron_classifer.py



Figure 4: loss-vs-iteration plot of verify_with_torchnn.py (one-neuron config)

## 2- Description of SGD+ and Adam
### a. SGD+

We all know that Stochastic Gradient Decent (SGD) is a popular method to optimize a loss function in any machine learning model. SGD+ is an extension of SGD that uses the concept of momentum to speed up convergence. It does so by increasing the values of the gradients of the parameters. In SGD's journey of finding the minimum of the loss's hyperplane, if we realize that we have been correctly descending to the bottom of the valley, then we can make a greater jump to get to the bottom faster. This faster jump or the idea of taking a larger leap is called momentum.

The typical update rule for SGD is as follows, $p_{t+1} = p_t - lr \cdot g_{t+1}$. Here $p_{t+1}$ is the update to be made, the current update $p_t$, the learning rate $lr$, and the next gradient $g_{t+1}$.

To convert the mentioned equation to SGD+, we need to consider the previous update (the step that was taken last iteration). We still want to make sure we consider the current update as well, since if the current gradient is the opposite direction of the previous update, then we know we went passed the minimum and we need to self-correct. Therefore, we take a fraction of the previous update. The fractional amount of the previous update to consider is determined by the variable called momentum. The concepts can be mathematically represented by the following equations.

$$v_{t+1} = \mu \cdot v_t + g_{t+1}$$
$$p_{t+1} = p_t - lr \cdot v_{t+1}$$

In these equations, the momentum is $\mu$, the previous update is $v_t$, and $v_{t+1}$ can be seen as a temporary variable.

When comparing the equations of SGD and SGD+, we can see that there is an additional step before we modify the update in SGD+. That additional step is essentially, $\mu \cdot v_t$. The value of $\mu$ can range from 0 to 1. A value of 1 means we consider the entire amount of the previous update. Alternatively, a value of 0 means we don't consider the previous update at all. Notice that setting $\mu = 0$ simplifies the equations to be the same as SGD.

## b. ADAM

Adaptive Moment Estimation, or ADAM for short, is an improved optimizer and is known to work very well. Adam combines the concepts of momentum and adaptability. We already know that momentum helps reach convergence faster. It does so by checking the previous update and if it points at the same direction as the current update, then it will take a larger jump. The concept of adaptability addresses the issue of sparse gradients. It is often the case that gradients are zeros, meaning there is not much to learn at a particular iteration in training. However, in the rare instance that there is a non-zero gradient, then we know this particular gradient contains vital information, and as a result we should give it more attention.

The following equations show mathematically what is mentioned in the previous paragraph.

$$m_{t+1} = \beta_1 \cdot m_t + (1 - \beta_1) \cdot g_{t+1}$$

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

$$p_{t+1} = p_t - lr \cdot \frac{\widehat{m_{t+1}}}{\sqrt{\widehat{v_{t+1}}} + \epsilon}$$

The first equation incorporates momentum by taking a running weighted average of the gradients. As we can see, the weight of the previous gradients is determined by $\beta_1$. Setting $\beta_1 = 0$ results in not using momentum at all. The second equation incorporates adaptability by taking a running weighted average of the gradient's variance. Like the case with $\beta_1$, setting $\beta_2 = 0$ will mean not considering the previous variances. The third equation is our parameter update equation. The difference here from the original SGD update equation is that the gradient term is replaced with $\frac{\widehat{m_{t+1}}}{\sqrt{\widehat{v_{t+1}}}+\epsilon}$. The numerator of this fraction is essentially the gradient with momentum. The denominator of the fraction handles the inverse relationship with variance. The logic behind this is that when high variance is observed Adam needs to slow down its learning. Additionally, when low variance is observed Adam needs to speed up learning because the gradients don't contain new knowledge to gain. Lastly, the $\epsilon$ is added to avoid a divide by zero error: typically, $\epsilon = 1e - 8$.

While this is the essence of Adam, an issue occurs in practice that requires us to make one small change. In the initial iterations of training, $m_t = v_t = 0$. These sway the averages of the mean and variance to 0 in the initial iterations. As a result, the following equations are used as corrected values for $m_t$ and $v_t$.

$$m_{\text{corrected}} = \frac{m}{1-\beta_1^k} \qquad\qquad v_{\text{corrected}} = \frac{v}{1-\beta_2^k}$$

Here k indicates what iteration the training loop is in. The logic here is the following. In the beginning stages of training (small k), the denominator will be a small value resulting in large adjustments to $m_t$ and $v_t$. While we approach the later stages of training (large k), the denominator will approach 1 and not correct our $m_t$ and $v_t$. This is perfect because it's in the beginning where we need correction to avoid the 0 bias caused by setting $m_t = v_t = 0$.

## 3- Implementation
### a. One Neuron

Iterations vs Training Loss lr = 0.0001

Iterations vs Training Loss lr = 0.001

Figures 5: One Neuron Loss Plot (lr = 0.0001)

Figures 6: One Neuron Loss Plot (lr = 0.001)

Iterations vs Training Loss lr = 0.005

Figures 7: One Neuron Loss Plot (lr = 0.005)

Figure 5 – 7 show the iterations vs loss plot of the one neuron model at different learning rates with 3 different optimizers: SGD, SGD+, and Adam. In regard to smoothness, it is very evident that Adam has the smoothest loss plots. We can see this since the green line in all these plots are very thin meaning there is less variation. SGD+ is also considerably smooth but not as smooth as the Adam. However, the least smooth is SGD. Figure 5 shows that when the learning rate is low, the SGD loss is very volatile and not smooth at all. Regarding convergence, we can

see that with a small learning rate (0.0001) all the optimizers take longer to learn. While Adam stabilizes at a loss slightly above 0.20 very quickly, SGD+ can continually decrease its loss ever so slowly: shown in Figure 5. With a medium learning rate (0.001), SGD performs much better than how it did with a slower learning rate (Figure 6). Adam can stabilize its loss much quicker but doesn't decrease the loss at all (Figure 6). However, SGD+ outperforms both optimizers by getting to a loss below 0.15 (Figure 6). Lastly, with a high learning rate (0.005) we in fact see divergence in SGD. This is likely due to SGD overshooting the minimum and never being able to return back to the bottom of the valley (Figure 7). This is where we see the beauty of momentum as SGD+ is able to remember the past gradients and continue its path to convergence (Figure 7). Adam again stabilizes very quickly but is not able to improve its loss (Figure 7).

## b. Multi Neuron



Figures 8: Multi Neuron Loss Plot (lr = 0.0001)
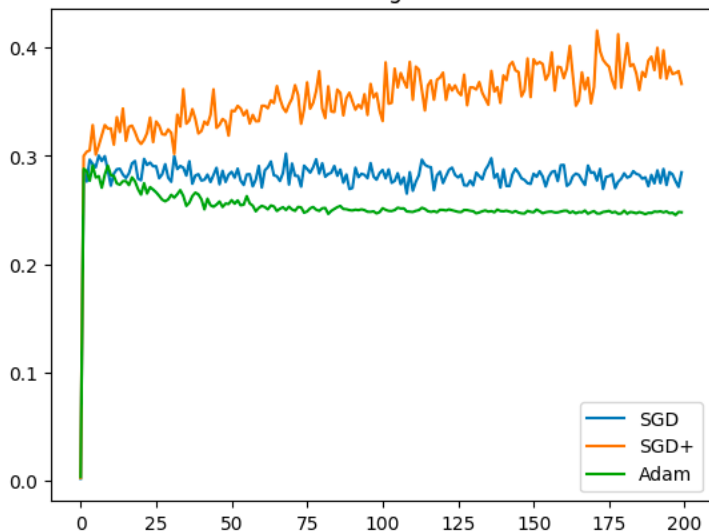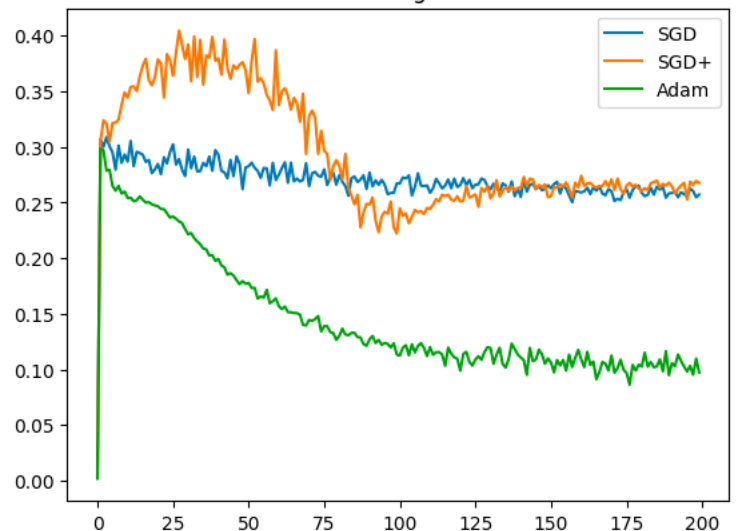
Figures 9: Multi Neuron Loss Plot (lr = 0.01)

Figures 10: Multi Neuron Loss Plot (lr = 0.01)

Figure 8 – 10 show the iterations vs loss plot of the multi neuron model at different learning rates with 3 different optimizers: SGD, SGD+, and Adam. In regard to smoothness, a high learning rate yielded the smoothest loss for all 3 optimizers (Figure 10). In general, Adam had smooth loss curves except in the case of Figure 9 where Adam had larger variances towards the end of training. SGD has overall the least smoothness. SGD did however smooth out at the end of training in Figure 9. Regarding convergence, different learning rates didn't change if a model converged or not, except for the case in SGD+ with a small learning rate (Figure 8). This observation is peculiar and could be a result of bad initialization weights. It is clear that a learning rate of 0.001 with Adam yielded the lowest loss (Figure 9). Additionally, SGD and SGD+ performed the best with a higher learning rate of 0.01 (Figure 10). In fact, SGD's loss starts to decrease toward the end of training suggesting that the loss can be further improved with more iterations (Figure 10).

## c. Findings

A lot can be learned from the conducted experiments. The overall take from this is that Adam tends to produce smoother loss curves which can be attributed to the consideration of the second moment. Additionally, Adam may not converge to the lowest loss, but it will stabilize the fastest. This is likely due to the combination of momentum and adaptability. If we were to rank the optimizer, SGD+ will be the runner up as it too can converge quickly but not as fast as Adam. This again is due to the consideration of momentum. The additional benefit of momentum is that SGD+ is able to avoid divergence where ordinary SGD cannot. As seen in the case of Figure 7, SGD overshot its estimation but since SGD+ is able to remember its previous gradient it avoided that issue. This brings SGD to last place where it experiences high variance and generally takes longer to converge.

It is important to mention that the different learning rates play a curial role in the performance of these optimizers. Given enough effort, one can easily find the right learning

rate that will sway one optimizer to perform better than the other. This only shows that hyperparameter tuning is critical when training deep learning models.

## d. Adam Optimizer Performance

Figure 11 shows the results of testing Adam optimizer with various beta 1 and beta 2 values. From the results the duration of training with different beta values plays no significant role. All the models finished training in about 8.8 to 9.1 seconds. The worst performing models occurred in trials 8 and 9. This was when the beta 1 value was 0.99 and beta 2 values were 0.9 and 0.95. The best performing model was trial 1 where beta 1 = 0.8 and beta 2 = 0.89. The second best model was trial 4 with beta 1 = 0.95 and beta 2 = 0.89. It's important to make clear that I determine "best model" by looking at the lowest final loss. Based on these trends, it seems that Adam performs better when both beta values are a little less than the normal 0.9 and 0.99 values that are traditionally used. This means Adam performs better when given slightly more importance to the current first and second moments. That is essentially the effect that the beta values play in Adam. The lower the beta values, the more weightage is given to the current moments when averaging. Inversely, the higher the beta values, the less weightage is given to the current moments when averaging.

| Trial | B1 | B2 | Time | Final-L | Min-L |
|-------|------|------|------|---------|--------|
| 1 | 0.8 | 0.89 | 8.9 | 0.0952 | 0.0027 |
| 2 | 0.8 | 0.9 | 8.9 | 0.1974 | 0.0011 |
| 3 | 0.8 | 0.95 | 9.0 | 0.13 | 0.0026 |
| 4 | 0.95 | 0.89 | 9.0 | 0.1006 | 0.0018 |
| 5 | 0.95 | 0.9 | 8.8 | 0.1887 | 0.0019 |
| 6 | 0.95 | 0.95 | 8.9 | 0.1908 | 0.0037 |
| 7 | 0.99 | 0.89 | 8.9 | 0.1996 | 0.0048 |
| 8 | 0.99 | 0.9 | 8.9 | 0.2403 | 0.0033 |
| 9 | 0.99 | 0.95 | 9.1 | 0.2005 | 0.0038 |

Figures 11: Table of Tested B1 and B2 Values for Adam

## 4- Source Code

```
5- # %%
6- import sys
7- sys.path.append("/Users/nikkhil/Documents/ECE 60146/hw3/ComputationalGraphPrimer-
   1.1.4/ComputationalGraphPrimer")
8-
9- import random
10-import numpy as np
```

```python
import matplotlib.pyplot as plt
import operator
from ComputationalGraphPrimer import *
import time

seed = 0
random.seed(seed)
np.random.seed(seed)

# %% [markdown]
# Create Subclass that returns the loss

# %%
# NOTE: This subclass is created so the loss is returned so the values can be used
# for plotting
class ComputationalGraphPrimer_ReturnLoss(ComputationalGraphPrimer):
    def __init__(self, *args, **kwargs):
        super(ComputationalGraphPrimer_ReturnLoss, self).__init__(*args, **kwargs)

    # NOTE: Code is taken from
  https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
  1.1.4_CodeOnly.html
    #        The only modification made is the loss_running_record list is returned
  and the generated plots are commented out
    def run_training_loop_one_neuron_model(self, training_data):
        self.vals_for_learnable_params = {param: random.uniform(0,1) for param in
  self.learnable_params}
        self.bias = random.uniform(0,1)

        class DataLoader:
            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)
```

```
50-                  else:
51-                      return random.choice(self.class_1_samples)
52-
53-          def getbatch(self):
54-                  batch_data,batch_labels = [],[]                          ## First
    list for samples, the second for labels
55-                  maxval = 0.0                                            ## For
    approximate batch data normalization
56-                  for _ in range(self.batch_size):
57-                      item = self._getitem()
58-                      if np.max(item[0]) > maxval:
59-                          maxval = np.max(item[0])
60-                      batch_data.append(item[0])
61-                      batch_labels.append(item[1])
62-                  batch_data = [item/maxval for item in batch_data]        ##
    Normalize batch data
63-                  batch = [batch_data, batch_labels]
64-                  return batch
65-
66-          data_loader = DataLoader(training_data, batch_size=self.batch_size)
67-          loss_running_record = []
68-          i = 0
69-          avg_loss_over_iterations = 0.0                                  ##
    Average the loss over iterations for printing out
70-                                                                          ##
    every N iterations during the training loop.
71-          for i in range(self.training_iterations):
72-              data = data_loader.getbatch()
73-              data_tuples_in_batch = data[0]
74-              class_labels_in_batch = data[1]
75-              y_preds, deriv_sigmoids =
    self.forward_prop_one_neuron_model(data_tuples_in_batch)     ##  FORWARD PROP of
    data
76-              loss = sum([(abs(class_labels_in_batch[i] - y_preds[i]))**2 for i in
    range(len(class_labels_in_batch))])  ##  Find loss
77-              avg_loss_over_iterations += loss / float(len(class_labels_in_batch))
78-              if i%(self.display_loss_how_often) == 0:
79-                  avg_loss_over_iterations /= self.display_loss_how_often
80-                  loss_running_record.append(avg_loss_over_iterations)
81-                  print("[iter=%d]  loss = %.4f" % (i+1, avg_loss_over_iterations))
    ## Display average loss
82-                  avg_loss_over_iterations = 0.0
    ## Re-initialize avg loss
83-              y_errors_in_batch = list(map(operator.sub, class_labels_in_batch,
    y_preds))
84-              self.backprop_and_update_params_one_neuron_model(data_tuples_in_batch,
    y_preds, y_errors_in_batch, deriv_sigmoids)  ## BACKPROP loss
85-          #plt.figure()
```

```
86-            #plt.plot(loss_running_record)
87-            #plt.show()
88-            return loss_running_record
89-
90-     # NOTE: Code is taken from
   https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
   1.1.4_CodeOnly.html
91-     #        The only modification made is the loss_running_record list is returned
   and the generated plots are commented out
92-     def run_training_loop_multi_neuron_model(self, training_data):
93-
94-         class DataLoader:
95-             def __init__(self, training_data, batch_size):
96-                 self.training_data = training_data
97-                 self.batch_size = batch_size
98-                 self.class_0_samples = [(item, 0) for item in
   self.training_data[0]]    ## Associate label 0 with each sample
99-                 self.class_1_samples = [(item, 1) for item in
   self.training_data[1]]    ## Associate label 1 with each sample
100-
101-             def __len__(self):
102-                 return len(self.training_data[0]) + len(self.training_data[1])
103-
104-             def _getitem(self):
105-                 cointoss = random.choice([0,1])                          ##
   When a batch is created by getbatch(), we want the
106-                                                                         ##
   samples to be chosen randomly from the two lists
107-                 if cointoss == 0:
108-                     return random.choice(self.class_0_samples)
109-                 else:
110-                     return random.choice(self.class_1_samples)
111-
112-             def getbatch(self):
113-                 batch_data,batch_labels = [],[]                          ##
   First list for samples, the second for labels
114-                 maxval = 0.0                                            ##
   For approximate batch data normalization
115-                 for _ in range(self.batch_size):
116-                     item = self._getitem()
117-                     if np.max(item[0]) > maxval:
118-                         maxval = np.max(item[0])
119-                     batch_data.append(item[0])
120-                     batch_labels.append(item[1])
121-                 batch_data = [item/maxval for item in batch_data]         ##
   Normalize batch data
122-                 batch = [batch_data, batch_labels]
123-                 return batch
```

```python
124-
125-            self.vals_for_learnable_params = {param: random.uniform(0,1) for param
    in self.learnable_params}
126-            self.bias =   {i : [random.uniform(0,1) for j in range(
    self.layers_config[i] ) ]   for i in range(1, self.num_layers)}
127-            data_loader = DataLoader(training_data, batch_size=self.batch_size)
128-            loss_running_record = []
129-            i = 0
130-            avg_loss_over_iterations = 0.0
    ## Average the loss over iterations for printing out
131-
    ##    every N iterations during the training loop.
132-            for i in range(self.training_iterations):
133-                data = data_loader.getbatch()
134-                data_tuples = data[0]
135-                class_labels = data[1]
136-                self.forward_prop_multi_neuron_model(data_tuples)
    ## FORW PROP works by side-effect
137-                predicted_labels_for_batch =
    self.forw_prop_vals_at_layers[self.num_layers-1]         ## Predictions from FORW
    PROP
138-                y_preds =  [item for sublist in  predicted_labels_for_batch  for
    item in sublist]      ## Get numeric vals for predictions
139-                loss = sum([(abs(class_labels[i] - y_preds[i]))**2 for i in
    range(len(class_labels))])  ## Calculate loss for batch
140-                loss_avg = loss / float(len(class_labels))
    ## Average the loss over batch
141-                avg_loss_over_iterations += loss_avg
    ## Add to Average loss over iterations
142-                if i%(self.display_loss_how_often) == 0:
143-                    avg_loss_over_iterations /= self.display_loss_how_often
144-                    loss_running_record.append(avg_loss_over_iterations)
145-                    print("[iter=%d]  loss = %.4f" %  (i+1,
    avg_loss_over_iterations))            ## Display avg loss
146-                    avg_loss_over_iterations = 0.0
    ## Re-initialize avg-over-iterations loss
147-                y_errors_in_batch = list(map(operator.sub, class_labels, y_preds))
148-                self.backprop_and_update_params_multi_neuron_model(y_preds,
    y_errors_in_batch)
149-            #plt.figure()
150-            #plt.plot(loss_running_record)
151-            #plt.show()
152-            return loss_running_record
153-
154-    # %% [markdown]
155-    # Create SGD+ (one & multi neuron)
156-
157-    # %%
```

```python
158-    # NOTE: Create a subclass of ComputationalGraphPrimer - SGD+
159-    #       In this class the backprop functions are modified to use SGD+
160-    class
   ComputationalGraphPrimer_ReturnLoss_SGD_Plus(ComputationalGraphPrimer_ReturnLoss):
161-        def __init__(self, *args, **kwargs):
162-            super(ComputationalGraphPrimer_ReturnLoss_SGD_Plus,
   self).__init__(*args, **kwargs)
163-
164-        # NOTE: Function created to set up all variables needed for carrying out
   momentum
165-        def initializeMomentumStuff_OneNeuron(self, momentum):
166-            self.mu = momentum
167-            # Previous updates for weights - set everything to 0 at first
168-            self.param_updates = {key: 0.0 for key in self.learnable_params}
169-            # Previous update for bias - set to 0 at first
170-            self.bias_update = 0.0
171-
172-        # NOTE: Function created to set up all variables needed for carrying out
   momentum
173-        def initializeMomentumStuff_MultiNeuron(self, momentum):
174-            self.muMulti = momentum
175-            # Previous updates for weights - set everything to 0 at first
176-            self.param_updates = {key: 0.0 for key in self.learnable_params}
177-            # Previous update for bias - set to 0 at first
178-            self.bias_updates = {i: [0.0 for _ in range(self.layers_config[i])] for
   i in range(1, self.num_layers)}
179-
180-        # NOTE: Code is taken from
   https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
   1.1.4_CodeOnly.html
181-        #        Most of this function is the same as Prof Avinash Kak's.
182-        #        Modifications made are indicated by comments.
183-        def backprop_and_update_params_one_neuron_model(self, data_tuples_in_batch,
   predictions, y_errors_in_batch, deriv_sigmoids):
184-            input_vars = self.independent_vars
185-            input_vars_to_param_map = self.var_to_var_param[self.output_vars[0]]
   ## These two statements align the
186-            param_to_vars_map = {param : var for var, param in
   input_vars_to_param_map.items()}   ##   the input vars
187-            vals_for_learnable_params = self.vals_for_learnable_params
188-            for i,param in enumerate(self.vals_for_learnable_params):
189-                ## For each param, sum the partials from every training data sample
   in batch
190-                partial_of_loss_wrt_param = 0.0
191-                for j in range(self.batch_size):
192-                    vals_for_input_vars_dict =  dict(zip(input_vars,
   list(data_tuples_in_batch[j])))
```

```python
193-                    partial_of_loss_wrt_param   +=   -  y_errors_in_batch[j] *
    vals_for_input_vars_dict[param_to_vars_map[param]] * deriv_sigmoids[j]
194-                partial_of_loss_wrt_param /=  float(self.batch_size)
195-                # Modification ---- Start
196-                # Here I am using the past parameter update and momentum to slighly
    modify the gradient
197-                # This is then used to update self.vals_for_learnable_params[param]
198-                # Equation set 2 from the hw description file is being implemented
    here
199-                self.param_updates[param] = self.mu * self.param_updates[param] -
    self.learning_rate * partial_of_loss_wrt_param
200-                self.vals_for_learnable_params[param] += self.param_updates[param]
201-                # Modification ---- End
202-
203-            y_error_avg = sum(y_errors_in_batch) / float(self.batch_size)
204-            deriv_sigmoid_avg = sum(deriv_sigmoids) / float(self.batch_size)
205-
206-            # Modification ---- Start
207-            # Like the weights the bias is also being updated in the same manner
208-            # Again equation set 2 from the hw description file is being
    implemented here
209-            self.bias_update = self.mu * self.bias_update + self.learning_rate *
    y_error_avg * deriv_sigmoid_avg
210-            self.bias += self.bias_update
211-            # Modification ---- End
212-
213-        # NOTE: Code is taken from
    https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
    1.1.4_CodeOnly.html
214-        #        Most of this function is the same as Prof Avinash Kak's.
215-        #        Modifications made are indicated by comments.
216-        def backprop_and_update_params_multi_neuron_model(self, predictions,
    y_errors):
217-            ## Eq. (24) on Slide 73 of my Week 3 lecture says we need to store
    backproped errors in each layer leading up to the last:
218-            pred_err_backproped_at_layers =   [ {i : [None for j in range(
    self.layers_config[i] ) ]
219-                                                          for i in
    range(self.num_layers)} for _ in range(self.batch_size) ]
220-            ## This will store "\delta L / \delta w" you see at the LHS of the
    equations on Slide 73:
221-            partial_of_loss_wrt_params = {param : 0.0 for param in self.all_params}
222-            ## For estimating the changes to the bias to be made on the basis of
    the derivatives of the Sigmoids:
223-            bias_changes =   {i : [0.0 for j in range( self.layers_config[i] ) ]
    for i in range(1, self.num_layers)}
224-            for b in range(self.batch_size):
```

```
225-                    pred_err_backproped_at_layers[b][self.num_layers - 1] = [
    y_errors[b] ]
226-                    for back_layer_index in reversed(range(1,self.num_layers)):
    ## For the 3-layer network, the first val for back_layer_index is 2 for the 3rd
    layer
227-                        input_vals = self.forw_prop_vals_at_layers[back_layer_index -1]
    ## This is a list of 8 two-element lists  --- since we have two nodes in the 2nd
    layer
228-                        deriv_sigmoids =
    self.gradient_vals_for_layers[back_layer_index]   ## This is a list eight one-
    element lists, one for each batch element
229-                        vars_in_layer  =  self.layer_vars[back_layer_index]
    ## A list like ['xo']
230-                        vars_in_next_layer_back  =  self.layer_vars[back_layer_index -
    1]   ## A list like ['xw', 'xz']
231-                        vals_for_input_vars_dict = dict(zip(vars_in_next_layer_back,
    self.forw_prop_vals_at_layers[back_layer_index - 1][b]))
232-                        ## For the next statement, note that layer_params are stored in
    a dict like
233-                        ##      {1: [['ap', 'aq', 'ar', 'as'], ['bp', 'bq', 'br',
    'bs']], 2: [['cp', 'cq']]}
234-                        ## "layer_params[idx]" is a list of lists for the link weights
    in layer whose output nodes are in layer "idx"
235-                        layer_params = self.layer_params[back_layer_index]
236-                        transposed_layer_params = list(zip(*layer_params))
    ## Creating a transpose of the link matrix, See Eq. 30 on Slide 77
237-                        for k,var1 in enumerate(vars_in_next_layer_back):
238-                            for j,var2 in enumerate(vars_in_layer):
239-                                pred_err_backproped_at_layers[b][back_layer_index -
    1][k] = sum([self.vals_for_learnable_params[transposed_layer_params[k][i]]
240-
    * pred_err_backproped_at_layers[b][back_layer_index][i]
241-
    for i in range(len(vars_in_layer)))])
242-                            for j,var in enumerate(vars_in_layer):
243-                            layer_params = self.layer_params[back_layer_index][j]
    ##  ['cp', 'cq']   for the end layer
244-                            input_vars_to_param_map = self.var_to_var_param[var]
    ## These two statements align the    {'xw': 'cp', 'xz': 'cq'}
245-                            param_to_vars_map = {param : var for var, param in
    input_vars_to_param_map.items()}   ##   and the input vars   {'cp': 'xw', 'cq':
    'xz'}
246-
247-                                ##  Update the partials of Loss wrt to the learnable
    parameters between the current layer
248-                                ##  and the previous layer. You are accumulating these
    partials over the different training
```

```
249-                          ##  data samples in the batch being processed.  For each
   training data sample, the formula
250-                          ##  being used is shown in Eq. (29) on Slide 77 of my Week
   3 slides:
251-                          for i,param in enumerate(layer_params):
252-                              partial_of_loss_wrt_params[param]    +=
   pred_err_backproped_at_layers[b][back_layer_index][j] * \
253-
   vals_for_input_vars_dict[param_to_vars_map[param]] * deriv_sigmoids[b][j]
254-                      ##  We will now estimate the change in the bias that needs to
   be made at each node in the previous layer
255-                      ##  from the derivatives the sigmoid at the nodes in the
   current layer and the prediction error as
256-                      ##  backproped to the previous layer nodes:
257-                      for k,var1 in enumerate(vars_in_next_layer_back):
258-                          for j,var2 in enumerate(vars_in_layer):
259-                              if back_layer_index-1 > 0:
260-                                  bias_changes[back_layer_index-1][k] +=
   pred_err_backproped_at_layers[b][back_layer_index - 1][k] * deriv_sigmoids[b][j]
261-
262-          ## Now update the learnable parameters.  The loop shown below carries
   out SGD mandated averaging
263-          for param in partial_of_loss_wrt_params:
264-              partial_of_loss_wrt_param = partial_of_loss_wrt_params[param] /
   float(self.batch_size)
265-              # Modification ---- Start
266-              # Here I am using the past parameter update and momentum to slighly
   modify the gradient
267-              # This is then used to update self.vals_for_learnable_params[param]
268-              # Equation set 2 from the hw description file is being implemented
   here
269-              self.param_updates[param] = self.muMulti *
   self.param_updates[param] - self.learning_rate * partial_of_loss_wrt_param
270-              self.vals_for_learnable_params[param] += self.param_updates[param]
271-              # Modification ---- End
272-
273-          ##  Finally we update the biases at all the nodes that aggregate data:
274-          for layer_index in range(1,self.num_layers):
275-              for k in range(self.layers_config[layer_index]):
276-                  # Modification ---- Start
277-                  # Like the weights the bias is also being updated in the same
   manner
278-                  # Again equation set 2 from the hw description file is being
   implemented here
279-                  # A temporary variable 'temp_bias_change' is created to avoid
   having 1 long line of code
280-                  temp_bias_change = bias_changes[layer_index][k] /
   float(self.batch_size)
```

```
281-                        self.bias_updates[layer_index][k] = self.muMulti *
    self.bias_updates[layer_index][k] + self.learning_rate * temp_bias_change
282-                        self.bias[layer_index][k]  +=
    self.bias_updates[layer_index][k]
283-                    # Modification ---- End
284-
285-    # %% [markdown]
286-    # Create ADAM (one & multi neuron)
287-
288-    # %%
289-    # NOTE: Create a subclass of ComputationalGraphPrimer - ADAM
290-    class
    ComputationalGraphPrimer_ReturnLoss_ADAM(ComputationalGraphPrimer_ReturnLoss):
291-        def __init__(self, *args, **kwargs):
292-            super(ComputationalGraphPrimer_ReturnLoss_ADAM, self).__init__(*args,
    **kwargs)
293-
294-        # NOTE: Function created to set up all variables needed for carrying out
    adam
295-        # Epsilon is hard coded to be 1e-8 as it is the standard mentioned in
    PyTorch documentation
296-        def initializeADAMStuff_OneNeuron(self, beta1, beta2, epsilon=1e-8):
297-            self.b1 = beta1
298-            self.b2 = beta2
299-            self.ep = epsilon
300-            self.timestep = 0
301-            self.m = {param: 0 for param in self.learnable_params}
302-            self.v = {param: 0 for param in self.learnable_params}
303-
304-        # NOTE: Function created to set up all variables needed for carrying out
    adam
305-        # Epsilon is hard coded to be 1e-8 as it is the standard mentioned in
    PyTorch documentation
306-        def initializeADAMStuff_MultiNeuron(self, beta1, beta2, epsilon=1e-8):
307-            self.b1 = beta1
308-            self.b2 = beta2
309-            self.ep = epsilon
310-            self.timestep = 0
311-            # For weights
312-            self.m = {param: 0 for param in self.learnable_params}
313-            self.v = {param: 0 for param in self.learnable_params}
314-            # For bias
315-            self.m_bias = {layer: [0 for _ in range(self.layers_config[layer])] for
    layer in range(1, self.num_layers)}
316-            self.v_bias = {layer: [0 for _ in range(self.layers_config[layer])] for
    layer in range(1, self.num_layers)}
317-
```

```
318-        # NOTE: Code is taken from
   https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
   1.1.4_CodeOnly.html
319-        #       Most of this function is the same as Prof Avinash Kak's.
320-        #       Modifications made are indicated by comments.
321-        def backprop_and_update_params_one_neuron_model(self, data_tuples_in_batch,
   predictions, y_errors_in_batch, deriv_sigmoids):
322-
323-            # Modification ---- Start
324-            self.timestep += 1
325-            # Modification ---- End
326-
327-            input_vars = self.independent_vars
328-            input_vars_to_param_map = self.var_to_var_param[self.output_vars[0]]
   ## These two statements align the
329-            param_to_vars_map = {param : var for var, param in
   input_vars_to_param_map.items()}   ##   the input vars
330-            vals_for_learnable_params = self.vals_for_learnable_params
331-            for i,param in enumerate(self.vals_for_learnable_params):
332-                ## For each param, sum the partials from every training data sample
   in batch
333-                partial_of_loss_wrt_param = 0.0
334-                for j in range(self.batch_size):
335-                    vals_for_input_vars_dict =  dict(zip(input_vars,
   list(data_tuples_in_batch[j])))
336-                    partial_of_loss_wrt_param   +=   - y_errors_in_batch[j] *
   vals_for_input_vars_dict[param_to_vars_map[param]] * deriv_sigmoids[j]
337-                partial_of_loss_wrt_param /=  float(self.batch_size)
338-
339-                # Modification ---- Start
340-                # Here I am using Equation set 3 from the hw description file
341-                # The m_corrected and v_corrected are from Week 3's presentaion
   slide 118
342-                # These values are then used to update
   self.vals_for_learnable_params[param]
343-                # The bias is updated in teh same manner
344-                self.m[param] = self.b1 * self.m[param] + (1-self.b1) *
   partial_of_loss_wrt_param
345-                self.v[param] = self.b2 * self.v[param] + (1-self.b2) *
   (partial_of_loss_wrt_param ** 2)
346-
347-                m_corr = self.m[param] / (1 - (self.b1**self.timestep))
348-                v_corr = self.v[param] / (1 - (self.b2**self.timestep))
349-
350-                self.vals_for_learnable_params[param] -= self.learning_rate *
   m_corr / ((v_corr + self.ep) ** 0.5)
351-                self.bias -= self.learning_rate * m_corr / ((v_corr + self.ep) **
   0.5)
```

```
352-                  # Modification ---- End
353-
354-          # NOTE: Code is taken from
    https://engineering.purdue.edu/kak/distCGP/ComputationalGraphPrimer-
    1.1.4_CodeOnly.html
355-          #        Most of this function is the same as Prof Avinash Kak's.
356-          #        Modifications made are indicated by comments.
357-          def backprop_and_update_params_multi_neuron_model(self, predictions,
    y_errors):
358-
359-              # Modification ---- Start
360-              self.timestep += 1
361-              # Modification ---- End
362-
363-              ## Eq. (24) on Slide 73 of my Week 3 lecture says we need to store
    backproped errors in each layer leading up to the last:
364-              pred_err_backproped_at_layers =   [ {i : [None for j in range(
    self.layers_config[i] ) ]
365-                                                          for i in
    range(self.num_layers)} for _ in range(self.batch_size) ]
366-              ## This will store "\delta L / \delta w" you see at the LHS of the
    equations on Slide 73:
367-              partial_of_loss_wrt_params = {param : 0.0 for param in self.all_params}
368-              ## For estimating the changes to the bias to be made on the basis of
    the derivatives of the Sigmoids:
369-              bias_changes =   {i : [0.0 for j in range( self.layers_config[i] ) ]
    for i in range(1, self.num_layers)}
370-              for b in range(self.batch_size):
371-                  pred_err_backproped_at_layers[b][self.num_layers - 1] = [
    y_errors[b] ]
372-                  for back_layer_index in reversed(range(1,self.num_layers)):
    ## For the 3-layer network, the first val for back_layer_index is 2 for the 3rd
    layer
373-                      input_vals = self.forw_prop_vals_at_layers[back_layer_index -1]
    ## This is a list of 8 two-element lists  --- since we have two nodes in the 2nd
    layer
374-                      deriv_sigmoids =
    self.gradient_vals_for_layers[back_layer_index]   ## This is a list eight one-
    element lists, one for each batch element
375-                      vars_in_layer  = self.layer_vars[back_layer_index]
    ## A list like ['xo']
376-                      vars_in_next_layer_back  =  self.layer_vars[back_layer_index -
    1]   ## A list like ['xw', 'xz']
377-                      vals_for_input_vars_dict = dict(zip(vars_in_next_layer_back,
    self.forw_prop_vals_at_layers[back_layer_index - 1][b]))
378-                      ## For the next statement, note that layer_params are stored in
    a dict like
```

```
379-                    ##        {1: [['ap', 'aq', 'ar', 'as'], ['bp', 'bq', 'br',
   'bs']], 2: [['cp', 'cq']]}
380-                    ## "layer_params[idx]" is a list of lists for the link weights
   in layer whose output nodes are in layer "idx"
381-                    layer_params = self.layer_params[back_layer_index]
382-                    transposed_layer_params = list(zip(*layer_params))
   ## Creating a transpose of the link matrix, See Eq. 30 on Slide 77
383-                    for k,var1 in enumerate(vars_in_next_layer_back):
384-                        for j,var2 in enumerate(vars_in_layer):
385-                            pred_err_backproped_at_layers[b][back_layer_index -
   1][k] = sum([self.vals_for_learnable_params[transposed_layer_params[k][i]]
386-
   * pred_err_backproped_at_layers[b][back_layer_index][i]
387-
   for i in range(len(vars_in_layer))])
388-                        for j,var in enumerate(vars_in_layer):
389-                            layer_params = self.layer_params[back_layer_index][j]
   ##  ['cp', 'cq']   for the end layer
390-                            input_vars_to_param_map = self.var_to_var_param[var]
   ## These two statements align the    {'xw': 'cp', 'xz': 'cq'}
391-                            param_to_vars_map = {param : var for var, param in
   input_vars_to_param_map.items()}   ##   and the input vars   {'cp': 'xw', 'cq':
   'xz'}
392-
393-                            ##  Update the partials of Loss wrt to the learnable
   parameters between the current layer
394-                            ##  and the previous layer. You are accumulating these
   partials over the different training
395-                            ##  data samples in the batch being processed.  For each
   training data sample, the formula
396-                            ##  being used is shown in Eq. (29) on Slide 77 of my Week
   3 slides:
397-                            for i,param in enumerate(layer_params):
398-                                partial_of_loss_wrt_params[param]   +=
   pred_err_backproped_at_layers[b][back_layer_index][j] * \
399-
   vals_for_input_vars_dict[param_to_vars_map[param]] * deriv_sigmoids[b][j]
400-                            ##  We will now estimate the change in the bias that needs to
   be made at each node in the previous layer
401-                            ##  from the derivatives the sigmoid at the nodes in the
   current layer and the prediction error as
402-                            ##  backproped to the previous layer nodes:
403-                            for k,var1 in enumerate(vars_in_next_layer_back):
404-                                for j,var2 in enumerate(vars_in_layer):
405-                                    if back_layer_index-1 > 0:
406-                                        bias_changes[back_layer_index-1][k] +=
   pred_err_backproped_at_layers[b][back_layer_index - 1][k] * deriv_sigmoids[b][j]
407-
```

```
408-            ## Now update the learnable parameters.  The loop shown below carries
   out Adam
409-            for param in partial_of_loss_wrt_params:
410-                partial_of_loss_wrt_param = partial_of_loss_wrt_params[param] /
   float(self.batch_size)
411-                # Modification ---- Start
412-                # Here I am using Equation set 3 from the hw description file
413-                # The m_corrected and v_corrected are from Week 3's presentaion
   slide 118
414-                # These values are then used to update
   self.vals_for_learnable_params[param]
415-                self.m[param] = self.b1 * self.m[param] + (1-self.b1) *
   partial_of_loss_wrt_param
416-                self.v[param] = self.b2 * self.v[param] + (1-self.b2) *
   (partial_of_loss_wrt_param ** 2)
417-
418-                m_corr = self.m[param] / (1 - (self.b1 ** self.timestep))
419-                v_corr = self.v[param] / (1 - (self.b2 ** self.timestep))
420-
421-                self.vals_for_learnable_params[param] += self.learning_rate *
   m_corr / ((v_corr + self.ep) ** 0.5)
422-                # Modification ---- End
423-
424-            ##  Finally we update the biases at all the nodes that aggregate data:
425-            for layer_index in range(1,self.num_layers):
426-                for k in range(self.layers_config[layer_index]):
427-                    # Modification ---- Start
428-                    # The bias is updated in the same manner as the weights were
   updated
429-                    # Again temp variable temp_bias_change for easy code
   readability
430-                    temp_bias_change = (bias_changes[layer_index][k] /
   float(self.batch_size))
431-
432-                    self.m_bias[layer_index][k] = self.b1 *
   self.m_bias[layer_index][k] + (1 - self.b1) * temp_bias_change
433-                    self.v_bias[layer_index][k] = self.b2 *
   self.v_bias[layer_index][k] + (1 - self.b2) * (temp_bias_change ** 2)
434-
435-                    m_bias_corr = self.m_bias[layer_index][k] / (1 - (self.b1 **
   self.timestep))
436-                    v_bias_corr = self.v_bias[layer_index][k] / (1 - (self.b2 **
   self.timestep))
437-
438-                    self.bias[layer_index][k] += self.learning_rate * m_bias_corr /
   ((v_bias_corr + self.ep) ** 0.5)
439-                    # Modification ---- End
440-
```

```python
441-
442-    # %% [markdown]
443-    # Test One Neuron (SGD, SGD+, ADAM)
444-
445-    # %%
446-    lr = 5e-3
447-
448-    # %%
449-    # Train One Neuron - SGD
450-    oneN_sgd = ComputationalGraphPrimer_ReturnLoss(one_neuron_model = True,
    expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'], output_vars = ['xw'], dataset_size =
    5000,
451-                                                  learning_rate = lr,
452-                                                  training_iterations = 40000,
    batch_size = 8, display_loss_how_often = 100, debug = False,)
453-
454-    oneN_sgd.parse_expressions()
455-    training_data = oneN_sgd.gen_training_data()
456-    loss_oneN_sgd = oneN_sgd.run_training_loop_one_neuron_model(training_data)
457-
458-    # %%
459-    # Train One Neuron - SGD+
460-    oneN_sgdP = ComputationalGraphPrimer_ReturnLoss_SGD_Plus(one_neuron_model =
    True, expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'], output_vars = ['xw'],
    dataset_size = 5000,
461-                                                  learning_rate = lr,
462-                                                  training_iterations =
    40000, batch_size = 8, display_loss_how_often = 100,debug = False,)
463-
464-    oneN_sgdP.parse_expressions()
465-    training_data = oneN_sgdP.gen_training_data()
466-    oneN_sgdP.initializeMomentumStuff_OneNeuron(momentum=0.9)
467-    loss_oneN_sgdP = oneN_sgdP.run_training_loop_one_neuron_model(training_data)
468-
469-    # %%
470-    # Train One Neuron - ADAM
471-    oneN_adam = ComputationalGraphPrimer_ReturnLoss_ADAM(one_neuron_model = True,
    expressions = ['xw=ab*xa+bc*xb+cd*xc+ac*xd'], output_vars = ['xw'], dataset_size =
    5000,
472-                                                  learning_rate = lr,
473-                                                  training_iterations =
    40000, batch_size = 8, display_loss_how_often = 100, debug = False,)
474-
475-    oneN_adam.parse_expressions()
476-    training_data = oneN_adam.gen_training_data()
477-    oneN_adam.initializeADAMStuff_OneNeuron(beta1=0.9, beta2=0.99)
478-    loss_oneN_adam = oneN_adam.run_training_loop_one_neuron_model(training_data)
479-
```

```
480-    # %%
481-    # Pot all the losses
482-    plt.figure()
483-    plt.plot(loss_oneN_sgd, label='SGD')
484-    plt.plot(loss_oneN_sgdP, label='SGD+')
485-    plt.plot(loss_oneN_adam, label='ADAM')
486-    plt.legend()
487-    plt.title(f"Iterations vs Training Loss lr = {lr}")
488-    plt.show()
489-
490-    # %% [markdown]
491-    # Test Multi Neuron (SGD, SGD+, ADAM)
492-
493-    # %%
494-    #lr = 9e-2
495-    lr = 1e-2
496-    #lr = 5e-2
497-    #lr = .01
498-
499-    # %%
500-    multiN_sgd = ComputationalGraphPrimer_ReturnLoss(num_layers = 3, layers_config
    = [4,2,1], expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
501-
    'xz=bp*xp+bq*xq+br*xr+bs*xs',
502-
    'xo=cp*xw+cq*xz'], output_vars = ['xo'], dataset_size = 5000,
503-                                                    learning_rate = lr,
504-                                                    training_iterations = 20000,
    batch_size = 8, display_loss_how_often = 100, debug = False,)
505-
506-    multiN_sgd.parse_multi_layer_expressions()
507-    training_data = multiN_sgd.gen_training_data()
508-    loss_multiN_sgd =
    multiN_sgd.run_training_loop_multi_neuron_model(training_data)
509-
510-    # %%
511-    multiN_sgdP = ComputationalGraphPrimer_ReturnLoss_SGD_Plus(num_layers =
    3,layers_config = [4,2,1], expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
512-
    'xz=bp*xp+bq*xq+br*xr+bs*xs',
513-
    'xo=cp*xw+cq*xz'], output_vars = ['xo'], dataset_size = 5000,
514-                                                        learning_rate = lr,
515-                                                        training_iterations
    = 20000, batch_size = 8, display_loss_how_often = 100, debug = False,)
516-
517-    multiN_sgdP.parse_multi_layer_expressions()
518-    training_data = multiN_sgdP.gen_training_data()
```

```python
519-    multiN_sgdP.initializeMomentumStuff_MultiNeuron(momentum=0.9)
520-    loss_multi_SGD =
    multiN_sgdP.run_training_loop_multi_neuron_model(training_data)
521-
522-    # %%
523-    multiN_adam = ComputationalGraphPrimer_ReturnLoss_ADAM(num_layers =
    3,layers_config = [4,2,1], expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
524-
    'xz=bp*xp+bq*xq+br*xr+bs*xs',
525-
    'xo=cp*xw+cq*xz'], output_vars = ['xo'], dataset_size = 5000,
526-                                                       learning_rate = lr,
527-                                                       training_iterations =
    20000, batch_size = 8, display_loss_how_often = 100, debug = False,)
528-
529-    multiN_adam.parse_multi_layer_expressions()
530-    training_data = multiN_adam.gen_training_data()
531-    multiN_adam.initializeADAMStuff_MultiNeuron(beta1=0.9, beta2=0.99)
532-    loss_multi_Adam =
    multiN_adam.run_training_loop_multi_neuron_model(training_data)
533-
534-    # %%
535-    # Pot all the losses
536-    plt.figure()
537-    plt.plot(loss_multiN_sgd, label='SGD')
538-    plt.plot(loss_multi_SGD, label='SGD+')
539-    plt.plot(loss_multi_Adam, label='Adam')
540-    plt.title(f"Iterations vs Training Loss lr = {lr}")
541-    plt.legend()
542-    plt.show()
543-
544-    # %% [markdown]
545-    # ADAM — Hyperparamter Tuning — B1 and B2
546-
547-    # %%
548-    # NOTE: The following code loops through all beta1 and beta2 values.
549-    #       The code then creates a new cgp and trains with the beta values
550-    #       The results are then saved to the results list
551-
552-    B1 = [0.8 , 0.95, 0.99]
553-    B2 = [0.89, 0.9 , 0.95]
554-
555-    i = 1
556-    results = []
557-    for b1 in B1:
558-        for b2 in B2:
559-
```

```
560-             testAdam = ComputationalGraphPrimer_ReturnLoss_ADAM(num_layers =
    3,layers_config = [4,2,1], expressions = ['xw=ap*xp+aq*xq+ar*xr+as*xs',
561-
    'xz=bp*xp+bq*xq+br*xr+bs*xs',
562-
    'xo=cp*xw+cq*xz'], output_vars = ['xo'], dataset_size = 5000,
563-                                               learning_rate = lr,
564-                                               training_iterations =
    20000, batch_size = 8, display_loss_how_often = 100, debug = False,)
565-             testAdam.parse_multi_layer_expressions()
566-             training_data = testAdam.gen_training_data()
567-             testAdam.initializeADAMStuff_MultiNeuron(beta1=b1, beta2=b2)
568-
569-             # Start time measurement
570-             start_time = time.time()
571-
572-             loss_testAdam =
    testAdam.run_training_loop_multi_neuron_model(training_data)
573-
574-             # End time measurement
575-             end_time = time.time()
576-             duration = end_time - start_time
577-
578-             # Save metrics results list
579-             finalLoss = loss_testAdam[-1]
580-             minLoss = min(loss_testAdam)
581-
582-             result = f"{i}\t{b1}\t{b2}\t{round(duration, 1)}\t{round(finalLoss,
    4)}\t{round(minLoss, 4)}"
583-             results.append(result)
584-             i+=1
585-
586-   # %%
587-   # Results are printed
588-   print(f"Trial\tB1\tB2\tTime\tFinal-L\tMin-L")
589-   for r in results:
590-       print(r)
591-
```