

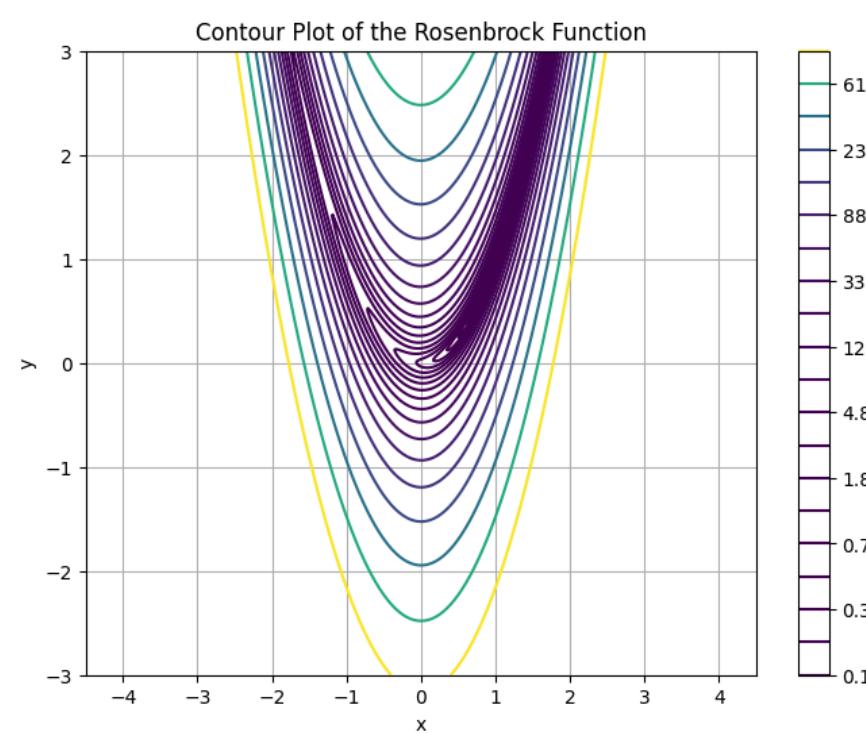
Souleymane Doumbia

## Homework 7

### Problem 1: Comparing Optimization Algorithms

#### 1(a): Contour Plot of the Objective Function

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Function and gradient definition based on the Rosenbrock function:
5 # f(x, y) = (1 - x)^2 + 100*(y - x^2)^2
6
7 def f(X):
8     x, y = X[0], X[1]
9     return (1 - x)**2 + 100 * (y - x**2)**2
10
11 x_vals = np.linspace(-4.5, 4.5, 400)
12 y_vals = np.linspace(-3.0, 3.0, 400)
13 X, Y = np.meshgrid(x_vals, y_vals)
14 Z = np.array([[f([x, y]) for x in x_vals] for y in y_vals])
15
16 plt.figure(figsize=(8, 6))
17 cp = plt.contour(X, Y, Z, levels=np.logspace(-1, 3, 20), cmap='viridis')
18 plt.colorbar(cp)
19 plt.title('Contour Plot of the Rosenbrock Function')
20 plt.xlabel('x')
21 plt.ylabel('y')
22 plt.grid(True)
23 plt.show()
```



#### 1(b): Optimization Algorithms: Setup for Gradient Descent Comparison

```

1 # Run each algorithm starting from (-4, -2), eta = 1e-4, and n_epochs = 100000
2 # Plot their trajectories and log10 of distance to minimum at (1, 1)
3
4 # Gradient of the function
5 def grad(X):
6     x, y = X[0], X[1]
7     dfdx = -2 * (1 - x) - 400 * x * (y - x**2)
8     dfdy = 200 * (y - x**2)
9     return np.array([dfdx, dfdy])
10
11 def gd(grad, init, n_epochs=1000, eta=10**-4):
12     params = np.array(init)
13     param_traj = np.zeros([n_epochs+1, 2])
14     param_traj[0,] = init
15     v = 0;
16     for j in range(n_epochs):
17         v = eta * (np.array(grad(params)))
18         params = params - v
19         param_traj[j+1,] = params
20     return param_traj
21
22 # --- 1(b) Trajectory Visualization and Log Error Plot ---
23

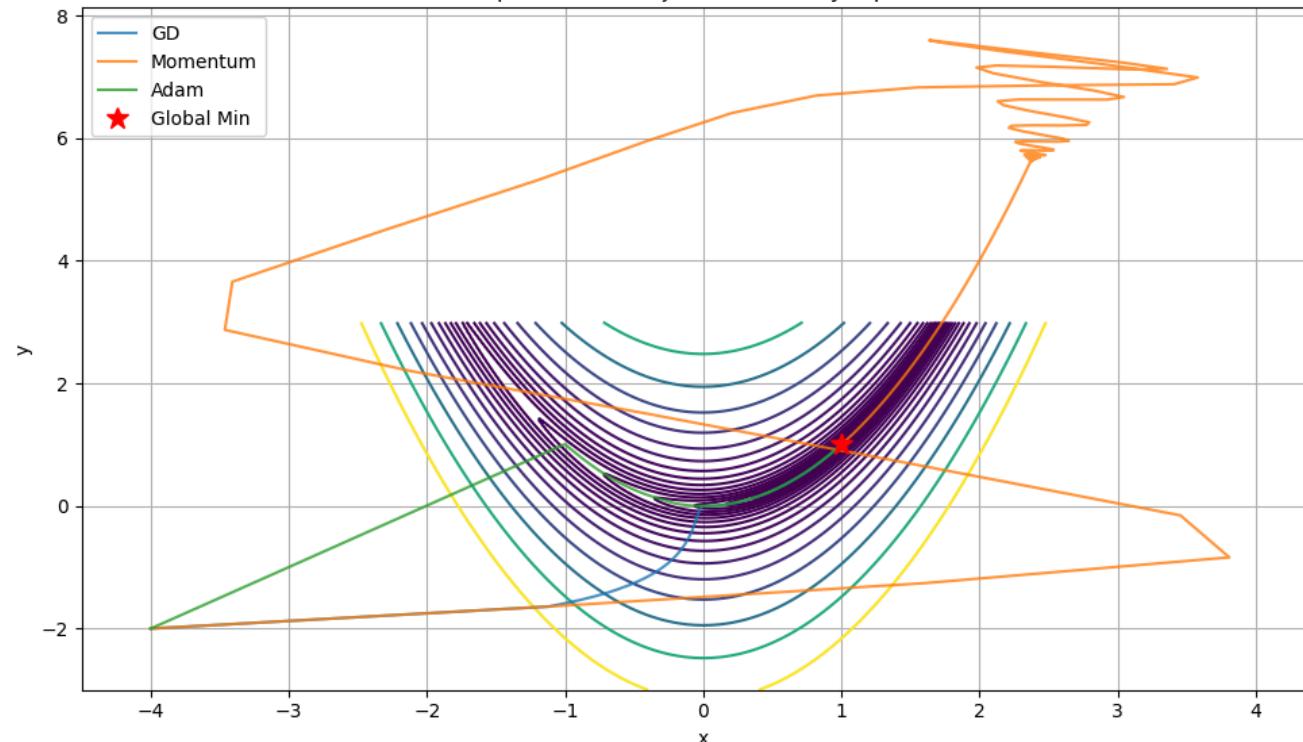
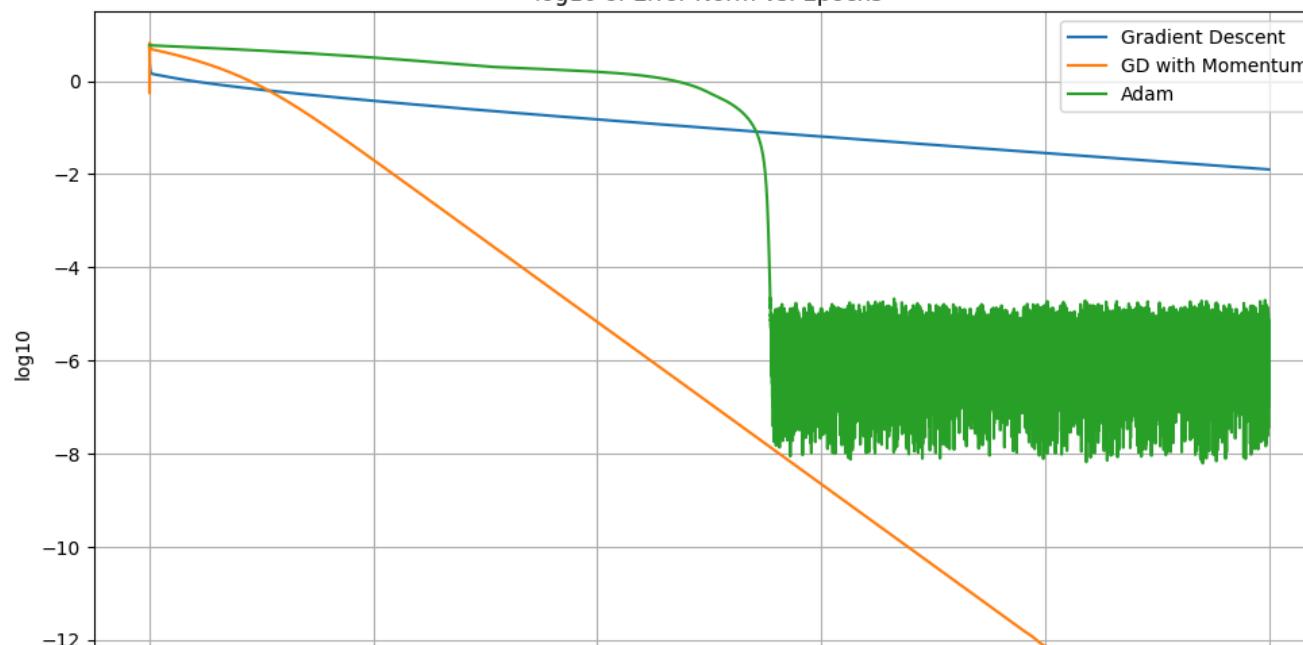
```

```

24 def gd_with_mom(grad, init, n_epochs=5000, eta=10**-4, beta=0.9, gamma=0.9):
25     params=np.array(init) # Start with initial condition
26     param_traj=np.zeros([n_epochs+1,2]) # Save the entire trajecotry
27     param_traj[0,]=init # Also save the initial condition to the trajectory
28
29 v=0 # Starting with 0 momentum
30
31 # Epochs is borrowing term from machine learning
32 # Here it means timestep
33
34 for j in range(n_epochs):
35     v=gamma*v+(np.array(grad(params))) # Compute v
36     params=params-eta*v # Update the location
37     param_traj[j+1,]=params # Save the trajectory
38 return param_traj
39
40
41 def adams(grad, init, n_epochs=5000, eta=10**-4, gamma=0.9, beta=0.99, epsilon=10**-8):
42     params = np.array(init)
43     param_traj = np.zeros([n_epochs+1, 2])
44     param_traj[0,] = init
45     v = 0
46     grad_sq = 0
47     for j in range(n_epochs):
48         g = np.array(grad(params))
49         v = gamma * v + (1 - gamma) * g
50         grad_sq = beta * grad_sq + (1 - beta) * g * g
51         v_hat = v / (1 - gamma**(j + 1))
52         grad_sq_hat = grad_sq / (1 - beta**(j + 1))
53         params = params - eta * np.divide(v_hat, np.sqrt(grad_sq_hat + epsilon))
54         param_traj[j+1,] = params
55     return param_traj
56
57 # Plot the trajectories of each algorithm and the base 10 of the error rate...
58 init_point = [-4, -2]
59 n_epochs = 100000
60 eta = 1e-4
61
62 traj_gd = gd(grad, init_point, n_epochs=n_epochs, eta=eta)
63 traj_mom = gd_with_mom(grad, init_point, n_epochs=n_epochs, eta=eta)
64 traj_adam = adams(grad, init_point, n_epochs=n_epochs, eta=eta)
65
66 # Plot optimization trajectories in the (x, y) plane
67 plt.figure(figsize=(10, 6))
68 plt.contour(X, Y, Z, levels=np.logspace(-1, 3, 20), cmap='viridis')
69 plt.plot(traj_gd[:, 0], traj_gd[:, 1], label='GD', alpha=0.8)
70 plt.plot(traj_mom[:, 0], traj_mom[:, 1], label='Momentum', alpha=0.8)
71 plt.plot(traj_adam[:, 0], traj_adam[:, 1], label='Adam', alpha=0.8)
72 plt.plot(1, 1, 'r*', markersize=12, label='Global Min')
73 plt.title('Optimization Trajectories in (x, y) Space')
74 plt.xlabel('x')
75 plt.ylabel('y')
76 plt.legend()
77 plt.grid(True)
78 plt.tight_layout()
79 plt.show()
80

```

```
81 ## Compute log10 error over time relative to the global minimum (1, 1)
82 optimum = np.array([1.0, 1.0])
83 log_err_gd = np.log10(np.linalg.norm(traj_gd - optimum, axis=1))
84 log_err_mom = np.log10(np.linalg.norm(traj_mom - optimum, axis=1))
85 log_err_adam = np.log10(np.linalg.norm(traj_adam - optimum, axis=1))
86
87 ## Plot the log10 error decay
88 plt.figure(figsize=(10, 6))
89 plt.plot(log_err_gd, label='Gradient Descent')
90 plt.plot(log_err_mom, label='GD with Momentum')
91 plt.plot(log_err_adam, label='Adam')
92 plt.title('log10 of Error Norm vs. Epochs')
93 plt.xlabel('Epoch')
94 plt.ylabel('log10')
95 plt.legend()
96 plt.grid(True)
97 plt.tight_layout()
98 plt.show()
99
100 # --- Report final point and function value (minimum) found by each method ---
101 print("Final point (Gradient Descent):", traj_gd[-1])
102 print("Final function value:", f(traj_gd[-1]))
103
104 print("Final point (GD with Momentum):", traj_mom[-1])
105 print("Final function value:", f(traj_mom[-1]))
106
107 print("Final point (Adam):", traj_adam[-1])
108 print("Final function value:", f(traj_adam[-1]))
```

Optimization Trajectories in  $(x, y)$  Spacelog<sub>10</sub> of Error Norm vs. Epochs



```

Final point (Gradient Descent): [0.99433318 0.98867571]
Final function value: 3.216465854863436e-05
Final point (GD with Momentum): [1. 1.]
Final function value: 2.0850579801122868e-26
Final point (Adam): [0.99999479 1.00000519]
Final function value: 2.4411637529060323e-08

```

**Discussion:** What do you notice about the performance of the different algorithms, both in terms of convergence speed and ultimate accuracy?

#### Gradient Descent (GD)

- **Convergence Speed:** Very slow and linear. Error decreases steadily but slowly over 100,000 iterations.
- **Final Accuracy:** Final point: [0.99433318, 0.98867571], function value: 3.216465854863436e-05 .

#### Gradient Descent with Momentum

- **Convergence Speed:** Very fast. Achieves rapid convergence due to momentum accumulation.
- **Final Accuracy:** Final point: [1. 1.], function value: 2.0850579801122868e-26 (close to machine precision).

#### Adam

- **Convergence Speed:** Fast, especially in the middle phase. However, it shows high-frequency fluctuations near the end.
- **Final Accuracy:** Final point: [0.99999479, 1.00000519], function value: 2.4411637529060323e-08 .

#### Overall Observations

- GD is robust but slow.
- Momentum significantly accelerates convergence and achieves highest precision.
- Adam is efficient and smooth at first but slightly unstable near the minimum due to adaptive learning rate dynamics.

### ✓ 1(c) Comparison with Higher Learning Rate: ADAM and GD with Momentum using eta = 1e-3

```

1 # Compare ADAM and GD with Momentum using eta = 1e-3
2 eta_c = 1e-3
3
4 traj_mom_c = gd_with_mom(grad, init_point, n_epochs=n_epochs, eta=eta_c)
5 traj_adam_c = adams(grad, init_point, n_epochs=n_epochs, eta=eta_c)
6
7 # Compute log10 error over time relative to (1,1)
8 log_err_mom_c = np.log10(np.linalg.norm(traj_mom_c - optimum, axis=1))
9 log_err_adam_c = np.log10(np.linalg.norm(traj_adam_c - optimum, axis=1))
10
11 # Plot the log10 error decay
12 plt.figure(figsize=(10, 6))
13 plt.plot(log_err_mom_c, label='Momentum (eta=1e-3)')
14 plt.plot(log_err_adam_c, label='Adam (eta=1e-3)')
15 plt.title('log10 Error vs Epochs for eta = 1e-3')
16 plt.xlabel('Epoch')

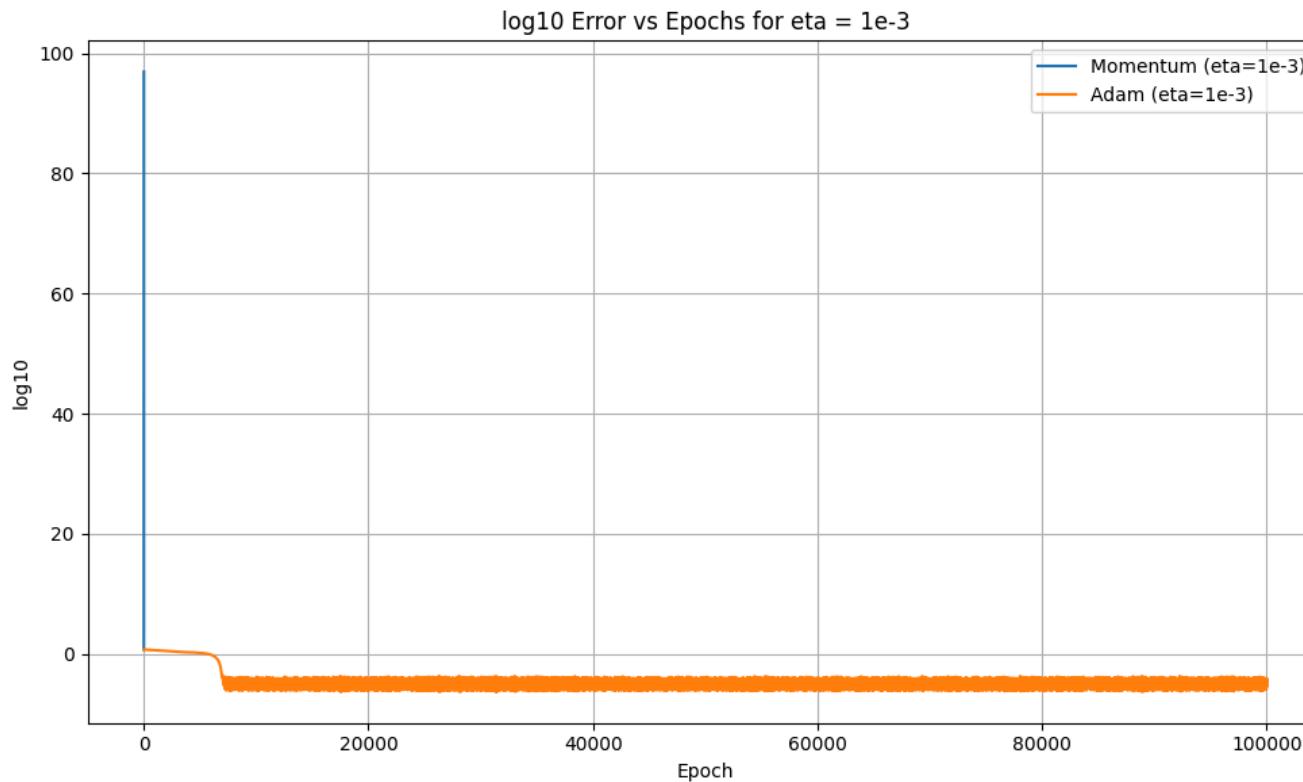
```

```

17 plt.ylabel('log10')
18 plt.legend()
19 plt.grid(True)
20 plt.tight_layout()
21 plt.show()
22
23 # Print final results
24 print("Final point (GD with Momentum, eta=1e-3):", traj_mom_c[-1])
25 print("Final function value:", f(traj_mom_c[-1]))
26
27 print("Final point (Adam, eta=1e-3):", traj_adam_c[-1])
28 print("Final function value:", f(traj_adam_c[-1]))

```

↳ <ipython-input-2-082c0dc1750b>:7: RuntimeWarning: overflow encountered in scalar power  
 $\text{dfdx} = -2 * (1 - x) - 400 * x * (y - x**2)$   
<ipython-input-2-082c0dc1750b>:8: RuntimeWarning: overflow encountered in scalar power  
 $\text{dfdy} = 200 * (y - x**2)$   
<ipython-input-2-082c0dc1750b>:7: RuntimeWarning: invalid value encountered in scalar subtract  
 $\text{dfdx} = -2 * (1 - x) - 400 * x * (y - x**2)$   
<ipython-input-2-082c0dc1750b>:8: RuntimeWarning: invalid value encountered in scalar subtract  
 $\text{dfdy} = 200 * (y - x**2)$   
/usr/local/lib/python3.11/dist-packages/numpy/linalg/\_linalg.py:2772: RuntimeWarning: overflow encountered in multiply  
 $s = (x.\text{conj}() * x).\text{real}$



Final point (GD with Momentum, eta=1e-3): [nan nan]  
Final function value: nan  
Final point (Adam, eta=1e-3): [1.00001807 0.99998026]  
Final function value: 3.1256313541136063e-07

### Discussion: Performance at High Learning Rate ( $\eta = 1e-3$ )

Gradient Descent with Momentum:

- Failed to converge. Encountered numerical overflow and returned NaN (Not a Number) for both coordinates and function value.
- This is visible in the early spike in the error plot and invalid gradient operations.

Adam:

- Successfully converged despite the high learning rate.
- Final point: [1.00001807, 0.99998026], function value: 3.1256313541136063e-07
- Slightly less precise than with smaller eta, but remained stable.

Observation:

- Adam handled the high learning rate robustly, while momentum-based gradient descent became unstable and unusable.
- Thus, **Adam worked better with  $\eta = 1e-3$** .

### ✓ 1(d) Trade-Off Comparison: ADAM ( $\eta=1e-2$ ) vs Momentum ( $\eta=1e-4$ )

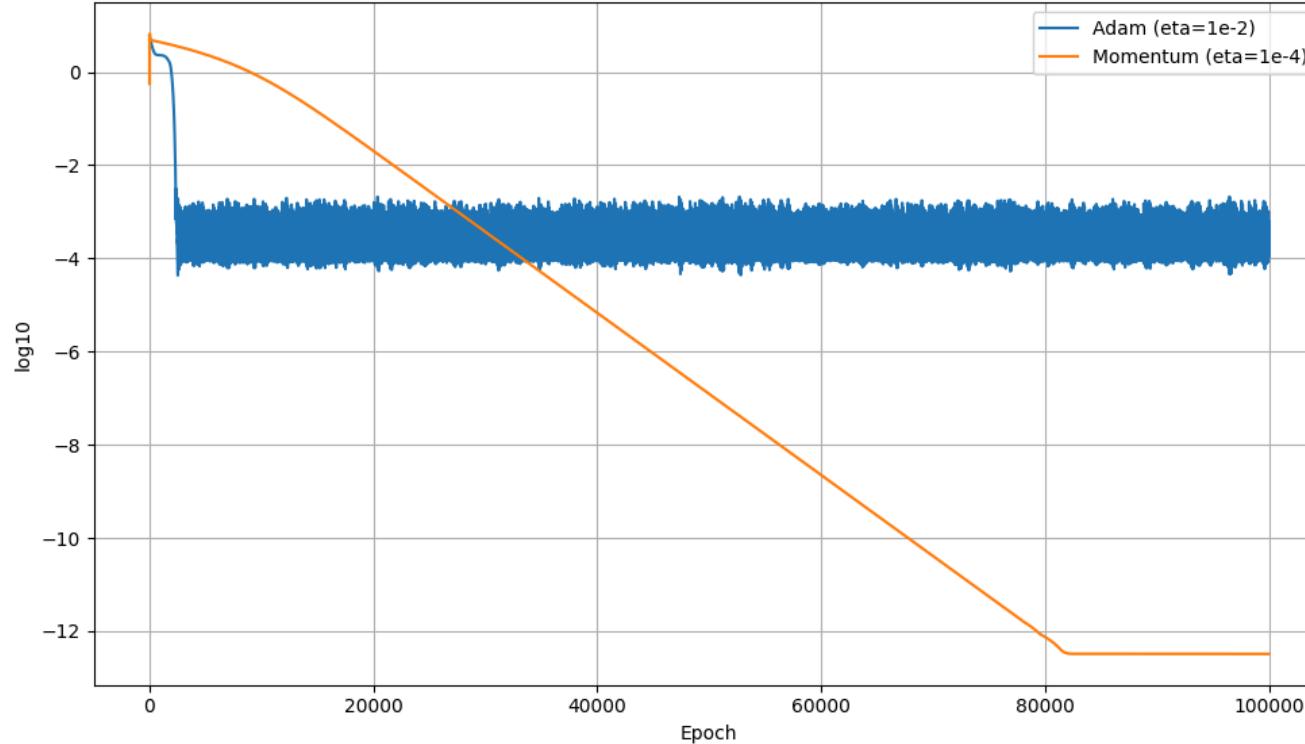
```

1 # Setup for different learning rates
2 eta_adam_d = 1e-2
3 eta_mom_d = 1e-4
4
5 # Run both optimizers
6 traj_adam_d = adams(grad, init_point, n_epochs=n_epochs, eta=eta_adam_d)
7 traj_mom_d = gd_with_mom(grad, init_point, n_epochs=n_epochs, eta=eta_mom_d)
8
9 # Compute error over time
10 log_err_adam_d = np.log10(np.linalg.norm(traj_adam_d - optimum, axis=1))
11 log_err_mom_d = np.log10(np.linalg.norm(traj_mom_d - optimum, axis=1))
12
13 # Plot
14 plt.figure(figsize=(10, 6))
15 plt.plot(log_err_adam_d, label='Adam (\eta=1e-2)')
16 plt.plot(log_err_mom_d, label='Momentum (\eta=1e-4)')
17 plt.title('log10 Error vs Epochs')
18 plt.xlabel('Epoch')
19 plt.ylabel('log10')
20 plt.legend()
21 plt.grid(True)
22 plt.tight_layout()
23 plt.show()
24
25 # Final point & function values
26 print("Final point (Adam, \eta=1e-2):", traj_adam_d[-1])
27 print("Final function value:", f(traj_adam_d[-1]))
28
29 print("Final point (GD with Momentum, \eta=1e-4):", traj_mom_d[-1])
30 print("Final function value:", f(traj_mom_d[-1]))

```



log10 Error vs Epochs



Final point (Adam, eta=1e-2): [1.00005344 0.99980312]

Final function value: 9.230096545900223e-06

Final point (GD with Momentum, eta=1e-4): [1. 1.]

Final function value: 2.0850579801122868e-26

### Discussion: Trade-Offs Between Adam (eta=1e-2) and Momentum (eta=1e-4)

Comparing ADAM with a high learning rate (eta=1e-2) against gradient descent with momentum using a smaller eta=1e-4

ADAM (eta=1e-2):

- Converges quickly at first but oscillates heavily near the minimum.
- Final point: [1.00005344, 0.99980312], function value: 9.230096545900223e-06
- Less accurate and noisy due to aggressive step size.

GD with Momentum (eta=1e-4):

- Converges more slowly but steadily, with minimal noise.
- Final point: [1.0, 1.0], function value: 2.0850579801122868e-26 (machine precision)

Trade-Off Summary:

- **Adam** offers speed and tolerance to noise at the cost of final precision.
- **Momentum** delivers higher accuracy but requires more iterations and a carefully tuned smaller learning rate.

✓ **Problem 2:** Shallow Nets and MNIST

✓ **2(a)** Validates Default Training and Explores Learning Rate Effects

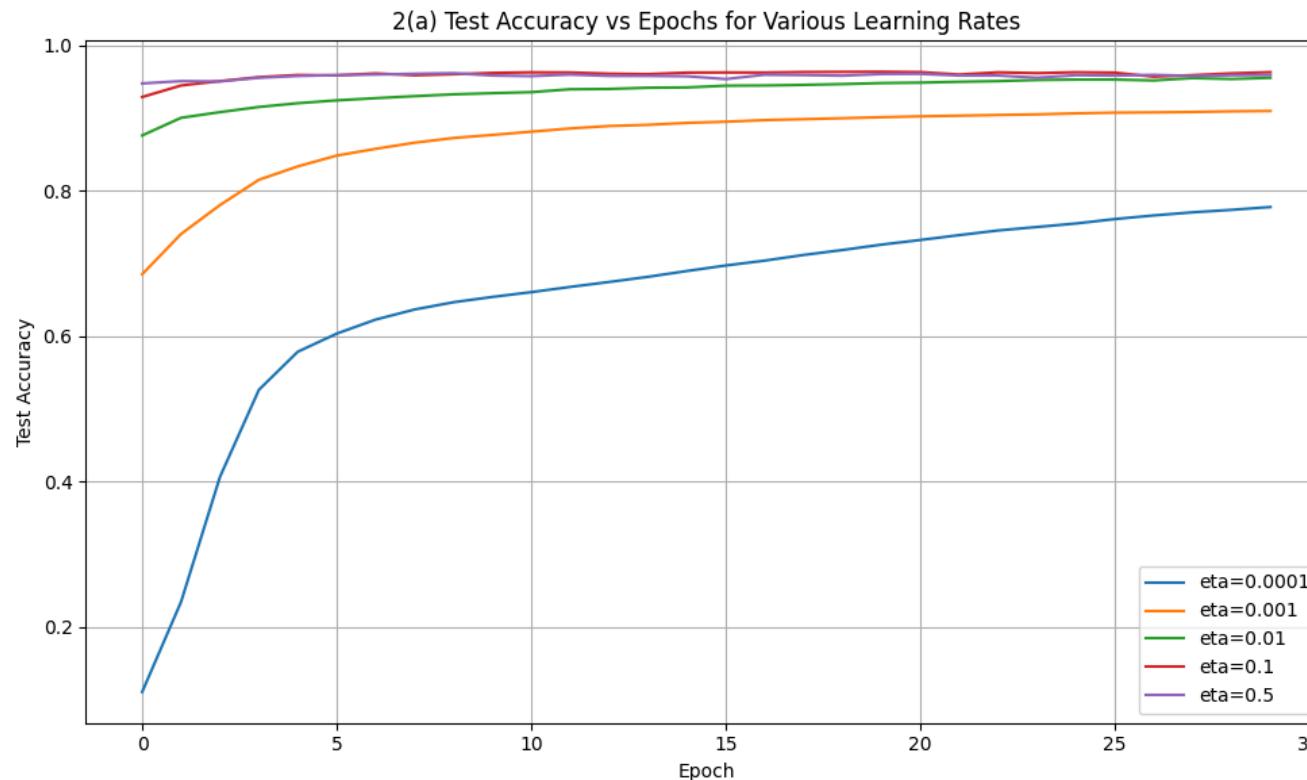
```

1 # Setup: Dataset and Network Initialization
2
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7 from torch.utils.data import DataLoader
8 from torchvision import datasets, transforms
9
10 # Define transformation: convert to tensor and flatten 28x28 -> 784
11 transform = transforms.Compose([
12     transforms.ToTensor(),
13     transforms.Lambda(lambda x: x.view(-1)) # flatten
14 ])
15
16 # Load MNIST dataset
17 train_dataset = datasets.MNIST('data/', train=True, download=True, transform=transform)
18 test_dataset = datasets.MNIST('data/', train=False, transform=transform)
19
20 train_loader = DataLoader(train_dataset, batch_size=10, shuffle=True)
21 test_loader = DataLoader(test_dataset, batch_size=10, shuffle=False)
22
23 # Network Class Definition (From Helper Notebook)
24
25 class Network(nn.Module):
26     def __init__(self, sizes):
27         super(Network, self).__init__()
28         self.sizes = sizes
29         self.num_layers = len(sizes)
30         self.layers = nn.ModuleList()
31         for i in range(self.num_layers - 1):
32             layer = nn.Linear(sizes[i], sizes[i+1])
33             nn.init.xavier_normal_(layer.weight) # suited for shallow/sigmoid nets
34
35             nn.init.zeros_(layer.bias) # initialize the bias to 0
36             self.layers.append(layer)
37
38     # Forward is the method that calculates the value of the neural network. Basically we recursively apply the activations in each layer
39     def forward(self, x):
40         for layer in self.layers[:-1]:
41             x = F.sigmoid(layer(x)) # default is sigmoid for shallow nets
42         x = self.layers[-1](x) # last layer: logits
43         return x
44
45 # Train Function
46
47 def train(network, train_loader, epochs, eta, test_data):
48     optimizer = optim.SGD(network.parameters(), lr=eta)
49     loss_fn = nn.CrossEntropyLoss()

```

```
50     test_accuracies = []
51
52     for epoch in range(epochs):
53         network.train()
54         for images, labels in train_loader:
55             output = network(images)
56             loss = loss_fn(output, labels)
57             optimizer.zero_grad()
58             loss.backward()
59             optimizer.step()
60
61     ## Evaluate on test set
62     network.eval()
63     correct = 0
64     total = 0
65
66     ## We have the "with torch.no_grad()" line for efficiency purposes
67     with torch.no_grad():
68         for images, labels in test_data:
69             output = network(images)
70             preds = output.argmax(dim=1)
71             correct += (preds == labels).sum().item()
72             total += labels.size(0)
73     accuracy = correct / total
74     test_accuracies.append(accuracy)
75
76 return test_accuracies
77
78 # Sweep Over Different Learning Rates
79 etas = [0.0001, 0.001, 0.01, 0.1, 0.5]
80 final_accuracies = []
81 all_curves = {}
82
83 for eta in etas:
84     net = Network([784, 30, 10])
85     acc_curve = train(net, train_loader, epochs=30, eta=eta, test_data=test_loader)
86     all_curves[eta] = acc_curve
87     final_accuracies.append((eta, acc_curve[-1]))
88
89 # Plot test accuracy curves
90 import matplotlib.pyplot as plt
91 plt.figure(figsize=(10, 6))
92 for eta, curve in all_curves.items():
93     plt.plot(curve, label=f'eta={eta}')
94 plt.xlabel('Epoch')
95 plt.ylabel('Test Accuracy')
96 plt.title('2(a) Test Accuracy vs Epochs for Various Learning Rates')
97 plt.legend()
98 plt.grid(True)
99 plt.tight_layout()
100 plt.show()
101
102 # Print best learning rate and its final accuracy
103 best_eta = max(final_accuracies, key=lambda x: x[1])
104 print("Best learning rate:", best_eta[0])
105 print("Final test accuracy:", best_eta[1])
```

```
100%|██████████| 9.91M/9.91M [00:01<00:00, 5.10MB/s]
100%|██████████| 28.9k/28.9k [00:00<00:00, 135kB/s]
100%|██████████| 1.65M/1.65M [00:06<00:00, 247kB/s]
100%|██████████| 4.54K/4.54K [00:00<00:00, 8.50MB/s]
```



Best learning rate: 0.1  
Final test accuracy: 0.9635

### Discussion: Evaluation of Learning Rate Sweep

We trained a shallow neural network ( $784 \rightarrow 30 \rightarrow 10$ ) using stochastic gradient descent (SGD) across a range of learning rates (etas = [0.0001, 0.001, 0.01, 0.1, 0.5]). Each configuration was run for 30 epochs on the MNIST dataset, and test accuracy was tracked.

We Observe:

- **Low learning rates** ( $\eta = 0.0001$  and  $0.001$ ) resulted in slow learning and lower accuracy.
- **Moderate learning rates** ( $\eta = 0.01$  and  $0.1$ ) achieved rapid convergence and high final accuracy.
- **High learning rate** ( $\eta = 0.5$ ) showed signs of instability (oscillations), although still high in accuracy.

Best Result:

- **Best learning rate:**  $\eta = 0.1$
- **Final test accuracy:** 0.9635

This suggests that for sigmoid-activated shallow networks on MNIST using SGD, a learning rate around 0.1 provides the best balance between convergence speed and stability.

## ✓ 2(b) Comparison with ADAM Optimizer

```

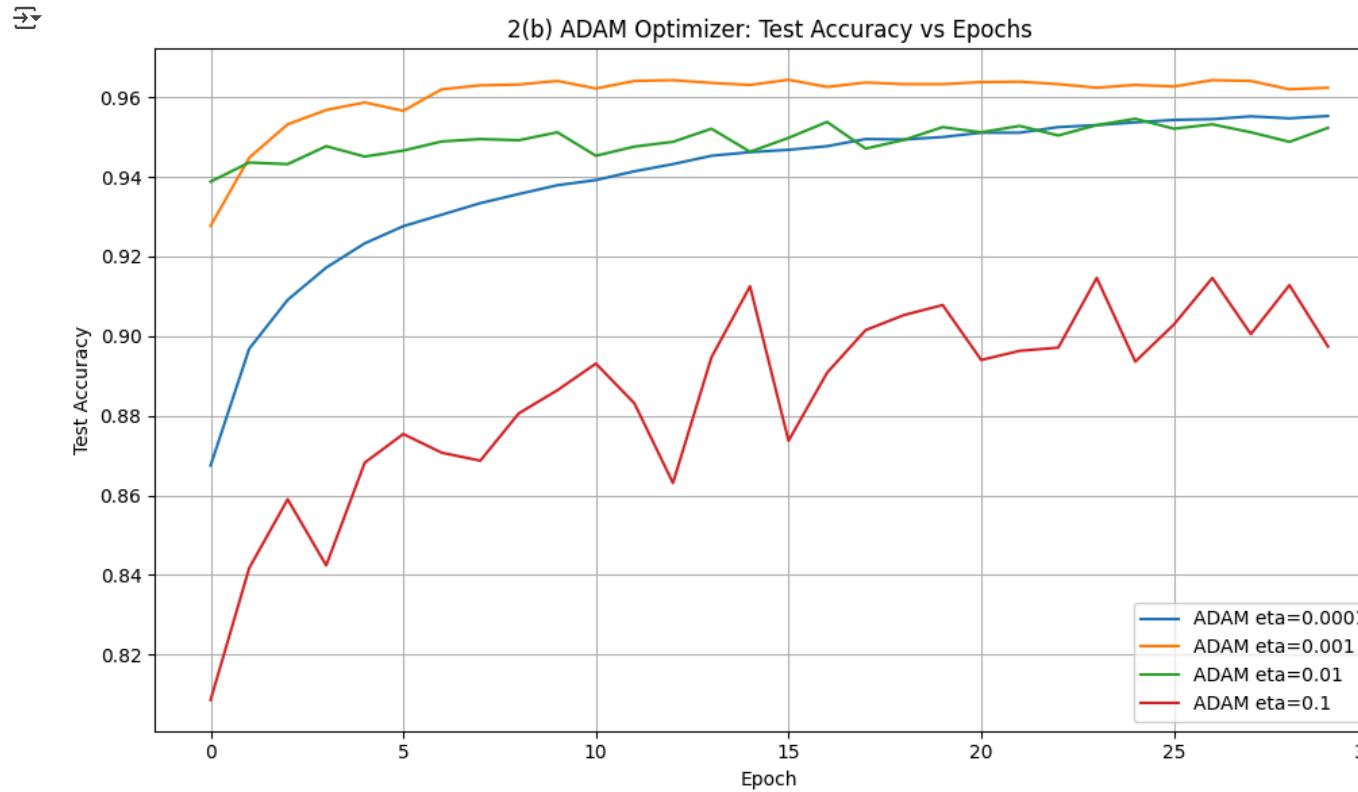
1 # This block tests ADAM optimizer using similar network and training setup
2
3 adam_etas = [0.0001, 0.001, 0.01, 0.1]
4 adam_final_accuracies = []
5 adam_curves = {}
6
7 def train_adam(network, train_loader, epochs, eta, test_data):
8     optimizer = optim.Adam(network.parameters(), lr=eta)
9     loss_fn = nn.CrossEntropyLoss()
10    test_accuracies = []
11
12    for epoch in range(epochs):
13        network.train()
14        for images, labels in train_loader:
15            output = network(images)
16            loss = loss_fn(output, labels)
17            optimizer.zero_grad()
18            loss.backward()
19            optimizer.step()
20
21        network.eval()
22        correct = 0
23        total = 0
24        with torch.no_grad():
25            for images, labels in test_data:
26                output = network(images)
27                preds = output.argmax(dim=1)
28                correct += (preds == labels).sum().item()
29                total += labels.size(0)
30        accuracy = correct / total
31        test_accuracies.append(accuracy)
32
33    return test_accuracies
34
35 # Run experiments for different ADAM learning rates
36 for eta in adam_etas:
37     net = Network([784, 30, 10])
38     acc_curve = train_adam(net, train_loader, epochs=30, eta=eta, test_data=test_loader)
39     adam_curves[eta] = acc_curve
40     adam_final_accuracies.append((eta, acc_curve[-1]))
41
42 # Plot ADAM test accuracy curves
43 plt.figure(figsize=(10, 6))
44 for eta, curve in adam_curves.items():
45     plt.plot(curve, label=f'ADAM eta={eta}')
46 plt.xlabel('Epoch')
47 plt.ylabel('Test Accuracy')
48 plt.title('2(b) ADAM Optimizer: Test Accuracy vs Epochs')
49 plt.legend()
50 plt.grid(True)

```

```

51 plt.tight_layout()
52 plt.show()
53
54 # Print best learning rate and accuracy for ADAM
55 best_adam_eta = max(adam_final_accuracies, key=lambda x: x[1])
56 print("Best ADAM learning rate:", best_adam_eta[0])
57 print("Final test accuracy:", best_adam_eta[1])

```



### Discussion: ADAM Optimizer Evaluation and Comparison with Stochastic Gradient Descent (SGD) in 2(a)

We evaluated the ADAM optimizer across several learning rates using the same shallow neural network ( $784 \rightarrow 30 \rightarrow 10$ ) on the MNIST dataset and compared it to the results obtained in Part 2(a) using SGD.

Observations with ADAM:

- $\eta = 0.001$  provided the best performance, reaching **0.9624** test accuracy.
- $\eta = 0.0001$  converged more slowly but stably.
- $\eta = 0.01$  achieved high accuracy but slightly less consistent than 0.001.
- $\eta = 0.1$  caused instability and lower accuracy due to overly aggressive updates.

Comparison with SGD (from 2a):

Optimizer	Best $\eta$	Final Accuracy
SGD	0.1	0.9635
ADAM	0.001	0.9624

- SGD with  $\eta=0.1$  reached the highest final test accuracy overall (0.9635).
- ADAM with  $\eta=0.001$  was more stable across training and easier to tune.
- ADAM showed better initial convergence at lower  $\eta$  but plateaued slightly earlier.

Conclusion:

While SGD slightly outperformed ADAM in final accuracy, ADAM showed strong early learning performance and required less tuning effort.

## ✓ 2(c) Hyperparameter and Architecture Tuning

In this part we explore how to improve performance. We vary hidden layer size, optimizer settings, and regularization.

```

1 # We vary hidden layer size, optimizer settings, and regularization.
2 custom_configs = [
3     {"hidden": 50, "optimizer": "adam", "eta": 0.001, "weight_decay": 1e-5},
4     {"hidden": 100, "optimizer": "adam", "eta": 0.001, "weight_decay": 1e-4},
5     {"hidden": 50, "optimizer": "sgd", "eta": 0.1, "momentum": 0.8},
6     {"hidden": 100, "optimizer": "sgd", "eta": 0.1, "momentum": 0.9}
7 ]
8
9 results_2c = []
10
11 def train_custom(network, optimizer, train_loader, epochs, test_data):
12     loss_fn = nn.CrossEntropyLoss()
13     test_accuracies = []
14     for epoch in range(epochs):
15         network.train()
16         for images, labels in train_loader:
17             output = network(images)
18             loss = loss_fn(output, labels)
19             optimizer.zero_grad()
20             loss.backward()
21             optimizer.step()
22
23         network.eval()
24         correct = 0
25         total = 0
26         with torch.no_grad():
27             for images, labels in test_data:
28                 output = network(images)
29                 preds = output.argmax(dim=1)
30                 correct += (preds == labels).sum().item()
31                 total += labels.size(0)
32         accuracy = correct / total
33         test_accuracies.append(accuracy)
34
35     return test_accuracies
36
37 for config in custom_configs:
38     net = Network([784, config["hidden"], 10])
39     if config["optimizer"] == "adam":

```

```

40     optimizer = optim.Adam(net.parameters(), lr=config["eta"], weight_decay=config["weight_decay"])
41 else:
42     optimizer = optim.SGD(net.parameters(), lr=config["eta"], momentum=config["momentum"])
43
44 acc_curve = train_custom(net, optimizer, train_loader, epochs=30, test_data=test_loader)
45 results_2c.append((config, acc_curve[-1]))
46
47 # Report top performing configuration
48 best_config = max(results_2c, key=lambda x: x[1])
49 print("Best custom configuration (2c):")
50 print("Hidden units:", best_config[0]["hidden"])
51 print("Optimizer:", best_config[0]["optimizer"])
52 print("Learning rate:", best_config[0]["eta"])
53 print("Final accuracy:", best_config[1])

```

→ Best custom configuration (2c):  
 Hidden units: 100  
 Optimizer: sgd  
 Learning rate: 0.1  
 Final accuracy: 0.9801

## Discussion: Custom Network & Hyperparameter Tuning

In this one, we explored various combinations of hidden layer size, optimizer choice, learning rate, and regularization to improve performance.

Configurations Tested:

- **Hidden units:** 50 and 100
- **Optimizers:** SGD with momentum, and ADAM with weight decay
- **Learning rates:** 0.1 (SGD), 0.001 (ADAM)

Best Configuration:

- **Hidden units:** 100
- **Optimizer:** sgd
- **Learning rate:** 0.1
- **Final test accuracy:** 0.9801

Interpretation:

- Increasing the number of hidden units to 100 provided better representational power.
- A learning rate of 0.1 with SGD led to rapid convergence and high generalization.
- Among all configurations tested, this combination yielded the best accuracy (0.9801).

Conclusion:

With careful tuning of both the **network architecture** and **training hyperparameters**, we achieved a final test accuracy of **0.9801**, the best observed in this set of experiments.

## ✓ Problem 3: Deep Nets: Overcoming Gradients

### ✓ 3(a) Setup: Deep Network and Gradient Ratio Measurement

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import torch.optim as optim
5 from torch.utils.data import DataLoader
6 from torchvision import datasets, transforms
7
8 # Load MNIST (same transform as Problem 2)
9 transform = transforms.Compose([
10     transforms.ToTensor(),
11     transforms.Lambda(lambda x: x.view(-1))
12 ])
13
14 train_dataset = datasets.MNIST('data/', train=True, download=True, transform=transform)
15 test_dataset = datasets.MNIST('data/', train=False, transform=transform)
16
17 train_loader = DataLoader(train_dataset, batch_size=10, shuffle=True)
18
19 # Deep Network with Gradient Tracking (sigmoid + xavier)
20
21 class DeepNetwork(nn.Module):
22     def __init__(self, sizes):
23         super(DeepNetwork, self).__init__()
24         self.sizes = sizes
25         self.num_layers = len(sizes)
26         self.layers = nn.ModuleList()
27         for i in range(self.num_layers - 1):
28             layer = nn.Linear(sizes[i], sizes[i+1])
29             nn.init.xavier_normal_(layer.weight)
30             nn.init.zeros_(layer.bias)
31             self.layers.append(layer)
32
33     def forward(self, x):
34         self.activations = []
35         for i, layer in enumerate(self.layers[:-1]):
36             x = F.sigmoid(layer(x))
37             self.activations.append(x)
38         x = self.layers[-1](x) # no softmax needed for CrossEntropyLoss
39         return x
40
41 # Function to compute and return gradient ratio
42
43 def compute_gradient_ratio(model):
44     grad_norm_input = model.layers[0].weight.grad.norm().item()
45     grad_norm_output = model.layers[-1].weight.grad.norm().item()
46     ratio = grad_norm_input / grad_norm_output if grad_norm_output != 0 else float('inf')
47     return grad_norm_input, grad_norm_output, ratio
48
49 # Initialize deep net and train for 1 epoch to observe gradient vanishing
50 net = DeepNetwork([784, 30, 30, 30, 30, 30, 30, 30, 30, 10])
51 optimizer = optim.SGD(net.parameters(), lr=0.01)
52 loss_fn = nn.CrossEntropyLoss()
53
54 for images, labels in train_loader:
55     output = net(images)
56     loss = loss_fn(output, labels)
```

```

57     optimizer.zero_grad()
58     loss.backward()
59     optimizer.step()
60     break # Just one batch for gradient observation
61
62 # Gradient stats
63 grad_in, grad_out, ratio = compute_gradient_ratio(net)
64 print("3(a) Gradient Norm of Input Layer:", grad_in)
65 print("3(a) Gradient Norm of Output Layer:", grad_out)
66 print("3(a) Gradient Ratio (Input/Output):", ratio)

→ 3(a) Gradient Norm of Input Layer: 4.663690197048709e-05
3(a) Gradient Norm of Output Layer: 1.0840625762939453
3(a) Gradient Ratio (Input/Output): 4.3020488844770725e-05

```

## Observing Vanishing Gradients in Deep Sigmoid Network

In this part, we constructed a deep feedforward neural network with 9 hidden layers, all using **sigmoid activation** functions and **Xavier initialization**. We trained it on a single batch of MNIST data and calculated the gradient norms at the input and output layers.

Observed Values:

- **Input Layer Gradient Norm:** 4.66e-05
- **Output Layer Gradient Norm:** 1.08
- **Gradient Ratio (Input/Output):** 4.30e-05

Interpretation:

- The input layer's gradients are **orders of magnitude smaller** than the output layer's.
- This confirms the presence of **vanishing gradients**, which makes it difficult for early layers to learn meaningful features during backpropagation.

## 3(b) ReLU Activation + Kaiming Initialization

Replacing sigmoid with ReLU and xavier with kaiming initialization.

```

1 # Replace sigmoid with ReLU and xavier with kaiming initialization.
2
3 class DeepNetworkReLU(nn.Module):
4     def __init__(self, sizes):
5         super(DeepNetworkReLU, self).__init__()
6         self.sizes = sizes
7         self.num_layers = len(sizes)
8         self.layers = nn.ModuleList()
9         for i in range(self.num_layers - 1):
10             layer = nn.Linear(sizes[i], sizes[i+1])
11             nn.init.kaiming_uniform_(layer.weight, nonlinearity='relu')
12             nn.init.zeros_(layer.bias)
13             self.layers.append(layer)
14
15     def forward(self, x):
16         self.activations = []
17         for i, layer in enumerate(self.layers[:-1]):

```

```

18         x = F.relu(layer(x))
19         self.activations.append(x)
20     x = self.layers[-1](x)
21     return x
22
23 # Initialize deep ReLU network
24 net_relu = DeepNetworkReLU([784, 30, 30, 30, 30, 30, 30, 30, 30, 10])
25 optimizer_relu = optim.SGD(net_relu.parameters(), lr=0.01)
26
27 # Train on one batch to observe gradient behavior
28 for images, labels in train_loader:
29     output = net_relu(images)
30     loss = loss_fn(output, labels)
31     optimizer_relu.zero_grad()
32     loss.backward()
33     optimizer_relu.step()
34     break
35
36 # Compute and print gradient stats for ReLU
37 grad_in_relu, grad_out_relu, ratio_relu = compute_gradient_ratio(net_relu)
38 print("3(b) Gradient Norm of Input Layer (ReLU):", grad_in_relu)
39 print("3(b) Gradient Norm of Output Layer (ReLU):", grad_out_relu)
40 print("3(b) Gradient Ratio (Input/Output, ReLU):", ratio_relu)

```

3(b) Gradient Norm of Input Layer (ReLU): 4.0457539558410645  
 3(b) Gradient Norm of Output Layer (ReLU): 0.8337193131446838  
 3(b) Gradient Ratio (Input/Output, ReLU): 4.852657113796479

### Discussion: ReLU Activation and Kaiming Initialization

To resolve the vanishing gradient problem observed in Part 3(a), we replaced sigmoid activations with **ReLU**, and used **Kaiming uniform initialization**, which is well-suited for ReLU-based deep networks.

Observed Gradient Norms:

- **Input Layer Gradient Norm:** 4.05
- **Output Layer Gradient Norm:** 0.83
- **Gradient Ratio (Input/Output):** 4.85

Interpretation (Gradient Flow):

- The input layer gradient is significantly stronger than the output layer's, indicating **no vanishing** and healthy gradient propagation.
- This is a major improvement over 3(a), where the input gradient was nearly 5 orders of magnitude smaller than the output.

Conclusion:

Switching to **ReLU + Kaiming** resolved vanishing gradients and unlocked the potential for training deeper networks.

### 3(c) Deep Network Improvement Experiment

Experimenting with deeper networks or decreasing hidden layer sizes layer-by-layer

```

1 class DeepCustomNet(nn.Module):
2     def __init__(self, sizes):
3         super(DeepCustomNet, self).__init__()

```

```
super(DeepCustomNet, self).__init__()
4     self.sizes = sizes
5     self.num_layers = len(sizes)
6     self.layers = nn.ModuleList()
7     for i in range(self.num_layers - 1):
8         layer = nn.Linear(sizes[i], sizes[i+1])
9         nn.init.kaiming_uniform_(layer.weight, nonlinearity='relu')
10        nn.init.zeros_(layer.bias)
11        self.layers.append(layer)
12
13    def forward(self, x):
14        for i, layer in enumerate(self.layers[:-1]):
15            x = F.relu(layer(x))
16        return self.layers[-1](x)
17
18 # Try an architecture that tapers from 784 → 512 → 256 → 128 → 64 → 10
19 net_custom = DeepCustomNet([784, 512, 256, 128, 64, 10])
20 optimizer_custom = optim.Adam(net_custom.parameters(), lr=0.001, weight_decay=1e-4)
21 loss_fn = nn.CrossEntropyLoss()
22
23 def evaluate_accuracy(model):
24     model.eval()
25     correct = 0
26     total = 0
27     with torch.no_grad():
28         for images, labels in DataLoader(test_dataset, batch_size=64):
29             output = model(images)
30             preds = output.argmax(dim=1)
31             correct += (preds == labels).sum().item()
32             total += labels.size(0)
33     return correct / total
34
35 # Train for a few epochs
36 epochs = 5
37 for epoch in range(epochs):
38     net_custom.train()
39     for images, labels in train_loader:
40         output = net_custom(images)
41         loss = loss_fn(output, labels)
42         optimizer_custom.zero_grad()
43         loss.backward()
44         optimizer_custom.step()
45
```