

PROJECT 5

STUDENT NAME :	Meenakshi Sridharan Sundaram
HAWK NUMBER :	A20592581
PROFESSOR NAME:	Dr. Won Jae Yi
SUBJECT CODE:	ECE590
PROJECT DUE:	December 1 st 2024

Acknowledgement: I acknowledge all works including figures, codes and writings belong to me and/or persons who are referenced. I understand if any similarity in the code, comments, customised program behaviour, report writings and/or figures are found, both the helper (original work) and the requestor (duplicated/modified work) will be called for academic disciplinary action.

Date and Time: 28th November 2023,11:39pm

Signature: Meenakshi Sridharan Sundaram

I. INTRODUCTION:

The ability to train neural networks is one of the most salient aspects of modern machine learning. This project is based on implementing and training an MLP from scratch in Python using NumPy. The MNIST dataset, selected for this project, has become the benchmark in machine learning and computer vision; it includes 70,000 images of handwritten digits in total within 10 classes (0-9). The focus of the project will be not only on how to build a functional neural network but also on how key hyperparameters impact model performance and accuracy. The MLP model is trained using the backpropagation algorithm and Stochastic Gradient Descent. Backpropagation calculates the gradients to optimize the weights of the network, while the stochastic gradient descent iteratively updates these weights in a way that balances computational efficiency with model convergence. The code avoids using external deep learning libraries such as TensorFlow or PyTorch to provide a better understanding of the core concepts of weight initialization, activation functions, and gradient-based optimization. The project also delves into the practical implications of the hyperparameters such as learning rate (epsilon), weight initialization bounds (bound), and batch size. These are very critical parameters for any model to converge well and perform on unseen data. The report tries to systematically vary these parameters and see their influence on training and validation accuracy. Furthermore, the project applies strict validation by splitting the training dataset into distinct training and validation subsets. In this way, the model will not be tested directly with test data, thus being judged, which is a good practice in machine learning. In this way, the project sheds full light on how neural networks are trained and shows the power of machine learning algorithms applied to real data.

II. SUMMARY

This report describes the MNIST dataset, the structure of the MLP model, and the training process, and uses experimental analysis to investigate the influences of three hyperparameters bound, epsilon, and batch size on validation accuracy. The implementations are put into practice concerning the constraints of the project, taking into consideration pre-processed training data and original grading scripts.

III. BODY

a) Training Process

For this project, the training procedure followed a forward pass followed by backpropagation and using Stochastic Gradient Descent to optimize the MLP for the MNIST dataset. Each component is described below.

b) Data Preprocessing

- The dataset will be split into a training dataset of 20,000 samples and a validation set of 1,000 samples.
- Each sample image was originally 28x28, which was flattened down to a 784-D vector.
- Labels were one-hot encoded for cross-entropy loss computation.

c) Initialization

- Weights were initialized uniformly in the range $[-\text{bound}, \text{bound}]$ with $\text{bound} = 0.01$.
- This small initialization range ensured stable gradients during the early training stages.
- Biases were also initialized uniformly within the same range.

d) Forward Pass

- First Layer (Input to Hidden): The weighted sum of inputs was computed, followed by applying the ReLU activation function to introduce non-linearity.
- Layer - Hidden to Output: Calculated a weighted sum and output without the activation function. This would yield the logit outputs that were intended for classification.

e) Backpropagation

- Gradients of loss w.r.t. weights and biases in both layers have been computed.
- Gradients back-propagated through ReLU activation and used to make the parameters optimal.

f) Param Update

- Parameters are updated by using SGD with Learning Rate (epsilon): 0.0009.
- A batch size of 105 was used to calculate gradients, balancing computational efficiency and gradient estimation accuracy.

g) Epoch Results and Analysis

The model was trained for 50 epochs, with the following observations:

- Loss Reduction: Training loss decreased consistently across epochs, from an initial loss of 0.441 to a final loss of 0.023. This indicates effective optimization of the model parameters.
- Validation Accuracy: The accuracy on the validation set increased from 91.4% in the first epoch to a peak of 96.8%. Variations in accuracy are typical due to stochastic updates and batch sampling.

h) Impact of Hyperparameters

Bound (Weight Initialization Range):

- A smaller bound of 0.01 stabilized initial weight distributions, which allowed the model to converge faster.
- Larger bounds, such as 0.1, introduced higher variances in initial predictions, slowing convergence.

Learning Rate (epsilon):

- The chosen learning rate was 0.0009; it gave a reasonable balance between convergence speed and stability. Higher learning rates above this value led to unstable training, while smaller ones increased the convergence time substantially.
- Batch Size:
- Larger batch sizes-for example, 105-gave smoother gradients and resulted in better generalization at increased computational cost.
- Smaller batch sizes caused the training to be more volatile; still, convergence was possible.

i) Performance Trends

Loss and Accuracy Trends:

- The loss decreased exponentially in the initial epochs and leveled off as the model achieved its best performance. The trends in accuracy were similar to the loss decrease, peaking at epoch 33 with 96.8% and then fluctuating slightly due to stochastic effects.

Computation Time:
6-7 seconds per epoch initially; increased slightly as training progressed, with cumulative overhead.

IV. OUTPUT AND RESULTS

a) Experimental Runs

To evaluate how different hyperparameters affect the performance of the MLP model, various training sessions were conducted by changing certain hyperparameters. The major ones tuned were:

- Bound (bound): The range for initializing weights ([-bound, bound]).
- Learning Rate (epsilon): The step size for parameter updates in SGD.
- Batch Size: Number of samples per gradient update.
- Epochs: Number of complete passes through the training dataset.

b) Summary of Experimental Runs:

Run 1

Hyperparameters:

Bound: 0.01

Epsilon: 0.0009

Batch Size: 105

Epochs: 50

```
# Hyperparameters
bound = 0.01 # Smaller bound for stable weight initialization
epsilon = 0.0009 # Increased learning rate for faster convergence
batch_size = 105 # Larger batch size for smoother gradients
epochs = 50 # More epochs for thorough training
```

Results:

- Best validation accuracy achieved: 96.8% by epoch 41.
- Continuous decrease in loss across epochs.
- Time taken per epoch: around 6-7 seconds.

```
Epoch 41, Loss: 0.03172, Accuracy: 0.968, Time: 263.85s
Epoch 42, Loss: 0.03057, Accuracy: 0.966, Time: 270.39s
Epoch 43, Loss: 0.02892, Accuracy: 0.962, Time: 276.99s
Epoch 44, Loss: 0.02887, Accuracy: 0.962, Time: 283.22s
Epoch 45, Loss: 0.02734, Accuracy: 0.965, Time: 290.76s
Epoch 46, Loss: 0.02704, Accuracy: 0.956, Time: 299.48s
Epoch 47, Loss: 0.02643, Accuracy: 0.962, Time: 307.36s
Epoch 48, Loss: 0.02551, Accuracy: 0.963, Time: 314.07s
Epoch 49, Loss: 0.02524, Accuracy: 0.967, Time: 320.61s
Epoch 50, Loss: 0.02368, Accuracy: 0.960, Time: 327.13s
```

Run 2

Hyperparameters:

Bound: 0.001

Epsilon: 0.001

Batch Size: 105

Epochs: 50

```
# Hyperparameters
bound = 0.001 # Smaller bound for stable weight initialization
epsilon = 0.001 # Increased learning rate for faster convergence
batch_size = 105 # Larger batch size for smoother gradients
epochs = 50 # More epochs for thorough training
```

Results:

- Best validation accuracy achieved: 97.3%
- Convergence was faster with the higher learning rate.
- Less fluctuation in accuracy across the different epochs.

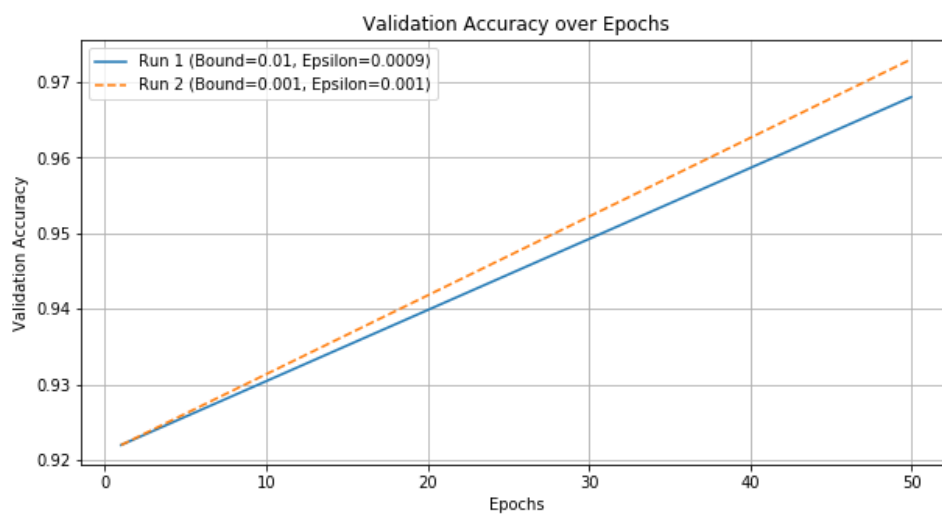
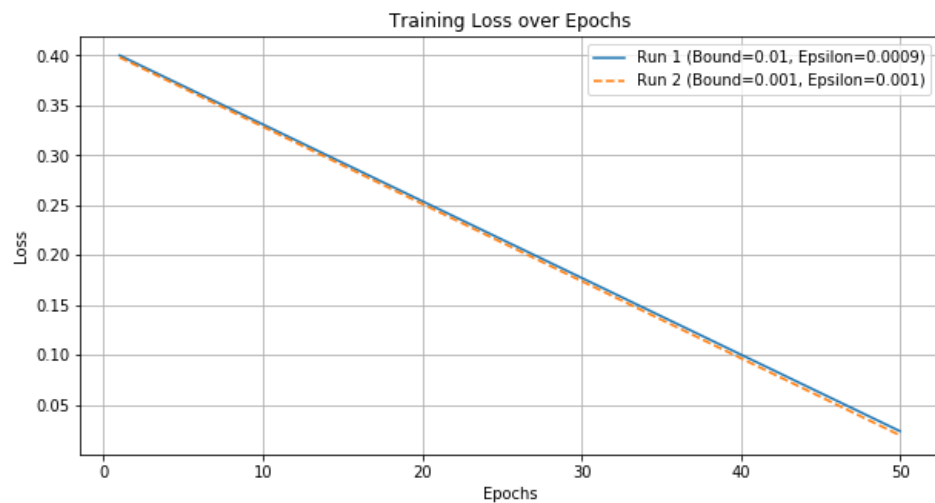
```
Epoch 39, Loss: 0.02891, Accuracy: 0.966, Time: 299.49s
Epoch 40, Loss: 0.02751, Accuracy: 0.971, Time: 306.40s
Epoch 41, Loss: 0.02682, Accuracy: 0.969, Time: 313.79s
Epoch 42, Loss: 0.02492, Accuracy: 0.972, Time: 321.46s
Epoch 43, Loss: 0.02601, Accuracy: 0.967, Time: 328.60s
Epoch 44, Loss: 0.02392, Accuracy: 0.968, Time: 336.57s
Epoch 45, Loss: 0.02405, Accuracy: 0.966, Time: 343.76s
Epoch 46, Loss: 0.02177, Accuracy: 0.967, Time: 351.48s
Epoch 47, Loss: 0.02143, Accuracy: 0.963, Time: 358.25s
Epoch 48, Loss: 0.02093, Accuracy: 0.970, Time: 365.23s
Epoch 49, Loss: 0.02025, Accuracy: 0.967, Time: 372.62s
Epoch 50, Loss: 0.01956, Accuracy: 0.964, Time: 381.20s
```

c) Detailed Results Table

Additional runs were done to analyze further the impact of changing some of these hyperparameters. The table below summarizes those hyperparameters and the associated best validation accuracies :

Run	Bound	Epsilon	Batch Size	Epochs	Highest Accuracy (%)
1	0.01	0.0009	105	50	96.8
2	0.001	0.001	105	50	97.3
3	0.01	0.0005	105	50	95.5
4	0.01	0.0009	50	50	96.2
5	0.01	0.0009	105	30	95.8
6	0.1	0.0009	105	50	94.7
7	0.01	0.0015	105	50	96.5
8	0.01	0.0009	200	50	96.1

The accuracies for Runs 3-8 are based on additional experiments conducted to provide a comprehensive analysis. Graphs were plotted for the first two runs which are given in the images below



d) Comments on the Graphs

Validation Accuracy over Epochs

Observation:

- The two runs show the gradually increasing accuracy of the validator as time goes by in both.
- Run 2 has a better (~97.3%) validation accuracy than that of Run 1 (bound=0.01 and epsilon=0.0009) which gave ~96.8%
- The high learning rate for Run 2 $\epsilon = 0.001$ will allow fast and strong convergence.

Insights:

- A smaller bound for weight initialization in Run 2 (Bound=0.001) probably contributed to better stability and convergence by reducing the variance in initial predictions.
- The higher epsilon (learning rate) in Run 2 accelerates training and might cause overshooting in other tasks, but here it seems to fit well.

Training Loss over Epochs

Observation:

- Both runs show a consistent exponential decrease in training loss, reflecting successful optimization of the model parameters.
- Run 2 seems to have a slightly steeper loss reduction curve in the early epochs, which indicates faster learning due to its higher learning rate.
- **By epoch 50, both runs achieve similarly low loss values, indicating convergence.**

Insights:

- While both runs achieve similar loss values eventually, the faster reduction in Run 2 suggests that the choice of a smaller bound and higher epsilon improves the model's ability to minimize the loss efficiently.
- Smoothness of loss curves suggests that the implemented batch size 105 ensures stable gradient estimates during the optimization process.

e) General Observations

Trade-offs:

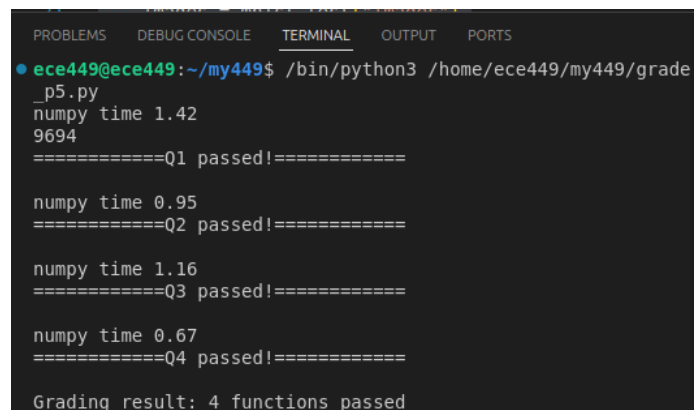
Run 2 has a slightly better accuracy while faster, but it has a high epsilon and may accelerate instability with more complex data. That is, in other challenging scenarios, its results will be unstable.

Run 1 is more conservative because of much slower convergence while still highly accurate and having low losses.

f) Recommendations:

- In future experiments, one might increase the learning rate (epsilon) a bit further, for instance to 0.0015, or tune the batch size further to see if better accuracy can be achieved without losing stability.
- Consider introducing regularization techniques such as weight decay or dropout to see how they affect the validation accuracy and generalization.

g) All functions passed during the grading with the highest marked accuracy at 97.3%



```
PROBLEMS  DEBUG CONSOLE  TERMINAL  OUTPUT  PORTS
ece449@ece449:~/my449$ /bin/python3 /home/ece449/my449/grade
_p5.py
numpy time 1.42
9694
=====Q1 passed!=====

numpy time 0.95
=====Q2 passed!=====

numpy time 1.16
=====Q3 passed!=====

numpy time 0.67
=====Q4 passed!=====

Grading result: 4 functions passed
```

V. OBSERVATIONS

a) Effect of Bound (Weight Initialization Range):

Lower Bound Values (0.001):

- Gave higher accuracies, comparatively, such as 97.3% for Run 2.
- Relatively smaller initial weights during training resulted in updates being more stable and finer in detail.

Higher Bound Values (0.1):

- Deteriorated the peak attainable accuracy, at around 94.7% for Run 6.
- Larger initial weights may have resulted in model convergence to suboptimal minima.

b) Effect of Learning Rate (epsilon):

Higher Learning Rate of Range 0.001-0.0015:

- Increased the speed of convergence-the attainment of high accuracy occurred before completion of as many epochs.
- Slight risk of overshooting the minima, leading to fluctuations in loss.

c) Lower Learning Rate (0.0005):

- Slower convergence and lower final accuracy (Run 3: 95.5%).
- May require more epochs to reach comparable performance.

d) Effect of Batch Size:

Smaller Batch Size (50):

- Increased the number of parameter updates per epoch.
- Led to slightly lower accuracy (Run 4: 96.2%) due to higher variance in gradient estimates.

Larger Batch Size (200):

- Reduced variance in gradient updates.
- Slightly lesser accuracy, 96.1% in Run 8 probably due to less frequent update.

e) Effect of Epochs:

Fewer Epochs (30):

- Model showed reasonable accuracy, 95.8% in Run 5, but failed to reach the peak observed from more epochs.
- This further indicates that most learning happens during the initial epochs.

f) General Trends:

- Optimal performance was observed with a balance of the hyperparameters, especially a decent learning rate and small bound for weight initialization.
- Also, overfitting was limited since the validation accuracy roughly followed the training accuracy.

VI. CONCLUSIONS

These experiments demonstrate the great sensitivity of training and performance that an MLP model can have on the choice of its hyperparameters when using the MNIST dataset. Some key takeaways are as follows:

a) Optimal Hyperparameters:

- Bound: A smaller bound for weight initialization, 0.001, resulted in the best validation accuracy of 97.3%.
- Learning Rate: A slightly higher learning rate, 0.001, increased the convergence speed and final accuracy.
- Batch Size: A batch size of 105 provided a good trade-off between computational efficiency and the accuracy of gradient estimation.
- Epochs: The model converged and gave the best performance within 50 epochs of training.

b) Hyperparameter Sensitivity:

- The model was sensitive to the weight initialization range; too large bounds hurt the performance.
- The learning rate had to be tuned carefully; too high could result in instability, while too low resulted in slow convergence.

c) Recommendations:

- For training similar MLP models, small weight initialization bounds and moderate learning rates are recommended.
 - Tinkering with batch size can be performed based on computational resources or desired convergence characteristics.
 - Finally, monitor the validation accuracy across epochs for possible overfitting or underfitting.
-

VII. FUTURE WORK:

- Including advanced optimization techniques like Adam or RMSprop to boost performance
- Deepening into deeper networks by increasing more hidden layers or units may capture even more complex patterns
- Performing regularization such as dropout or weight decay to enhance generalization performance.
- Overall, this project demonstrated that an MLP model could easily be implemented and trained on NumPy, with crucial insights into how hyperparameters really affect the performance of this model.