

Perceptron Training for Handwritten Digit Classification using MNIST Dataset

KHADKA PILOT^{1*} AND POUDEL KSHITIZ^{2*}

¹Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: pilot.076bct025@tcioe.edu.np)

²Institute of Engineering, Thapathali Campus, Thapathali, Nepal (e-mail: kshitizpoudel18@gmail.com)

Corresponding author: Khadka Pilot (e-mail: pilot.076bct025@tcioe.edu.np).

* Authors contributed equally

ABSTRACT Artificial Neural Network are computational system inspired from human brain to learn human like abilities. By harnessing the power of the concepts we can design any learning systems advanced weight initialization techniques like Xavier random and He initialization, we elevate recognition precision. Moreover, we systematically investigate diverse neural network architectures, manipulating neuron quantities and hidden layers with the MNIST dataset. This empirical exploration highlights the significant influence of architecture on accuracy and efficiency. Additionally, our study delves into the impact of various activation functions on model performance, aiming to optimize results.

INDEX TERMS ANN, MNIST dataset

I. INTRODUCTION

AN Artificial Neural Network (ANN) is a computational model inspired by the structure and functioning of the human brain's interconnected neurons. It's widely used in machine learning for tasks like pattern recognition, classification, regression, and more. The basic building block of an ANN is a neuron, which takes in inputs, performs computations, and produces an output. This project focuses on understanding artificial neural networks, which are important tools in computing and learning. These networks can find complex patterns in data. The main goal is to grasp how they work internally, including their forward and backward processes. These processes help improve the networks over time.

In today's machine learning world, neural networks are crucial. They're used in tasks like recognizing images and understanding language. This project aims to show how neural networks work, especially how they process data and get better with practice. We'll change settings, use different functions, and start from different points to see how they affect network performance. Knowing about neural networks is valuable as our world relies on data. Businesses and schools use data to learn things, and understanding neural networks gives us an important skill. This project wants to uncover how neural networks turn raw data into smart predictions. By learning theory and doing experiments, we can add to what people already know about these networks. In this experiment we will build a neural network from scratch by mathemati-

cally formulating the process and then experiment with different hyperparameters to observe their effect on final result. It will also help us to understand the mechanism of ANNs in more detail. The main goal is to make neural networks easier to understand and get ideas to make better models and choices.

II. METHODOLOGY

A. THE NEURAL SYSTEM: RECEPTORS, NEURONS, AND EFFECTORS

The human nervous system can be conceptualized as a triad of components: Receptors, Neuronal Network, and Effectors. The journey begins with Receptors, which transmute both external and internal stimuli into electrical impulses, signaling towards the central hub of the nervous system, the brain. Here, the Neural Network comes into play, adeptly processing and decoding the received information, subsequently formulating an appropriate response. Finally, the Effectors translates the brain's electrical instructions into tangible output responses.

Central to this system are neurons, the building blocks of the brain. These neurons interact through specialized units known as synapses, which act as bridges between them. Notably, Chemical synapses, the most prevalent type, convert electrical signals into chemical messengers and then reconver them back into electrical impulses.

The neurons themselves encode their outputs in the form of voltage pulses, famously termed "action potentials," which traverse the neuron's anatomy with constant velocity and

amplitude.

B. ANATOMY OF A NEURON

At the core of neuronal operation lie three fundamental components:

- **Synaptic Assemblies:** These intricate ensembles represents an array of synapses, where each synapse exhibits a quantifiable weight or synaptic strength. The input signal x_i , arriving at synapse j and connected to neuron k transforms via multiplication by the corresponding weight w_{kj} .
- **Summative Integration:** The adder serves as an essential component responsible for the aggregation of input signals, each multiplicatively weighted by the corresponding input signals.
- **Activation Function:** It provides gating mechanism in regulating the neuron's output. The neuron's behavioral output is governed by an activation function, a nonlinear operator that limits the magnitude of the response. By introducing a thresholding effect on the cumulative input, it controls response to varying levels of input stimulation.

C. FORWARD PROPAGATION

Forward propagation involves calculating the output for given inputs. For the first pass, the weights within the network can be zero, but it is not a very good idea. Therefore, there are many initialization techniques, the most basic being random initialization where every weight from one neuron in layer L to another neuron in layer $L + 1$ is randomly chosen from between -1 to 1 or another range of values.

After the weights are initialized, we need to compute the weighted sum coming to each neuron in the first hidden layer. Each neuron in the hidden layer has a weight associated with every neuron in the input layer. The forward propagation process is done in batches, where each batch consists of a number of training/test instances. Therefore, it is easier and faster to compute the weighted sum in matrix form as:

$$Z^{(1)} = W^{(1)T} \cdot A^{(0)} + B^{(1)} \quad (1)$$

Where,

$Z^{(n)}$ represents the weighted sum going to the n th layer

$A^{(n)}$ represents the activation of the n th layer the 0th layer corresponds to the input layer

$W^{(n)T}$ represents the weight matrix, where the j th column represents the weights for the j th neuron in the hidden layer

$B^{(n)}$ represents the bias of the n th layer

The weighted sum of the activation of the previous layer and the weights is then given as input to an activation function:

$$A^{(1)} = \text{activation}\{Z^{(1)}\} \quad (2)$$

The forward propagation continues until we obtain the activations for the output layer. Those activations represent the probabilities of the given instances belonging to each of

the 10 classes. In other words, if n hidden layers are used, then the output vector is given by:

$$\text{Output vector} = \text{Activation}(Z^{(n)}) = \text{softmax}(W^{(n)T} \cdot A^{(n-1)} + B^{(n)}) \quad (3)$$

D. BACKWARD PROPAGATION

After the forward pass is made, we obtain predictions for the classes. Now, the weights need to be updated for the learning process to occur. The goal of learning is to minimize the difference between the predicted output and the actual class label. Thus, a loss function is required to measure this difference.

Mean Squared Error (MSE) loss works well when there are multiple classes of output labels. Our goal is to find the combination of parameters (weights + bias) that minimizes the MSE loss. The MSE loss is given as:

$$MSE = \frac{1}{2n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Where, Y_i is the ground truth and \hat{Y}_i is the predicted value, and n is the total number of training examples.

Now we need to calculate the derivative of the loss function with respect to every weight. This is known as the gradient. Moving opposite to the direction of the gradient will result in a local minimum or global minimum. A local minimum is the point at which the loss is the lowest. It defines the combination of parameters for which the loss is the lowest.

$$\frac{\partial E}{\partial w_{ij}^{(k)}} = \frac{\partial E}{\partial a_j^{(k)}} \cdot \frac{\partial a_j^{(k)}}{\partial w_{ij}^{(k)}} \quad (4)$$

$$\frac{\partial a_j^{(k)}}{\partial w_{ij}^{(k)}} = \frac{\partial}{\partial w_{ij}^{(k)}} \left(\sum_{i'=0}^{r_{k-1}-1} w_{ij}^{(k)} o_{i'}^{(k-1)} \right) = o_i^{(k-1)} \quad (5)$$

Therefore,

$$\frac{\partial E}{\partial w_{ij}^{(k)}} = \frac{\partial E}{\partial a_j^{(k)}} o_i^{(k-1)} \quad (6)$$

III. ALGORITHM

Algorithm 1 Perceptron Training using Training Data

Require: Training Data: $(\mathbf{x}_i, y_i); \forall i \in \{1, 2, \dots, N\}$, Learning Rate: η

Ensure: Separating Hyper-plane coefficients: \mathbf{w}^*

- 1: Initialize $\mathbf{w} \leftarrow \mathbf{0}$
- 2: **repeat**
- 3: Get example (\mathbf{x}_i, y_i)
- 4: $\hat{y}_i \leftarrow \mathbf{w}^T \mathbf{x}_i$
- 5: **if** $\hat{y}_i y_i \leq 0$ **then**
- 6: $\mathbf{w} \leftarrow \mathbf{w} + \eta y_i \mathbf{x}_i$
- 7: **end if**
- 8: **until** convergence

Algorithm 2 Gradient Descent for Training a Linear Unit

Require: Training Data: $(\mathbf{x}_i, y_i); \forall i \in \{1, 2, \dots, N\}$, Learning Rate: η

Ensure: Optimal Hyper-plane coefficients based on squared loss: \mathbf{w}^*

- 1: Initialize $\mathbf{w} \leftarrow$ random weights
- 2: **repeat**
- 3: Calculate ∇E
- 4: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E$
- 5: **until** convergence

A. SYSTEM BLOCK DIAGRAM

Appendix

B. INSTRUMENTATION

Here are some key NumPy functions utilized in our neural network implementation:

- **np.random:** This function is employed to generate random values sampled from a uniform or normal distribution. In our context, it's utilized for initializing weights and biases with random values.
- **np.dot:** The `np.dot` function performs matrix multiplication between the weight matrix and input matrix. It's a fundamental operation in neural network feedforward computations.
- **np.sum:** In various instances, `np.sum` is used to compute the sum over derivatives or values within specific layers of the network. It plays a role in backpropagation and loss calculation.
- **np.exp:** The `np.exp` function is employed for softmax calculations, which are crucial for obtaining probability distributions over multiple classes.
- **np.arange:** `np.arange` is used for creating a one-hot encoded representation of the labels. It helps in converting categorical labels into a binary format suitable for neural network training.
- **np.mean:** This function is used to calculate the mean value, often across all samples in a dataset. It might be employed, for instance, to calculate the mean loss over a batch of training examples.
- **np.argmax:** The `np.argmax` function helps in obtaining the index of the maximum value in a column. It's frequently used when making predictions or evaluating model outputs.
- **np.sum:** Another use of `np.sum` is to calculate the sum of values that match a specific target. This operation is valuable when computing the count of correct predictions in classification tasks.
- **plt.plot:** The `plt.plot` function from the `matplotlib.pyplot` library is employed for visualizing the accuracy and losses of the neural network during training. It aids in understanding the learning progress and potential issues.

IV. WORKING PRINCIPLE**A. DATASET OVERVIEW**

The MNIST dataset, which stands for the Modified National Institute of Standards and Technology dataset, is a widely used benchmark in machine learning and computer vision. It consists of handwritten digit images along with corresponding labels, making it suitable for tasks like image classification and digit recognition. The MNIST dataset has played a significant role in advancing these fields and bench-marking various algorithms and models.

The MNIST dataset comprises a total of 70,000 gray-scale images depicting handwritten digits. Typically, this dataset is segregated into a training set, encompassing 60,000 images, and a distinct testing set, comprising 10,000 images. However, for our specific utilization, we've opted to work with a subset of the MNIST dataset, containing 42,000 instances.

B. DATASET CONTENTS**1) Images as csv**

The images in the MNIST dataset are single-channel (grayscale) images. Each pixel's intensity value ranges from 0 to 255, representing the darkness of the pixel. The images contain handwritten digits from 0 to 9.

Each image is associated with a label corresponding to the digit it represents. The labels are integers ranging from 0 to 9, indicating the actual value of the digit in the image.

C. DATASET PREPARATION

Preparing the MNIST dataset for machine learning involves several key steps:

1) Data Loading Process

The data loading process involved utilizing the `'read_csv'` function to import the CSV file. Within this dataset, there are a total of 785 columns. The initial column serves as the label, while the subsequent columns correspond to each individual pixel within a 28x28 image. Notably, frameworks such as TensorFlow and PyTorch commonly offer convenient tools for both acquiring and loading such datasets.

2) Data Preprocessing

Normalize the raw pixel values in the images to a range between 0 and 1 for faster convergence during training. This is achieved by dividing all pixel values by 255. Additionally, reshape the 28x28 images into a flattened vector of length 784 to serve as input to the neural network. In order to send it into a neural network, each image is transformed into a 1-D vector, which can be fed as input into the network. Moreover, for convenience, the vectorized image instances are transposed, effectively arranging each instance as a column in the dataset.

3) Division into Training and Testing Sets

To ensure the model's ability to generalize to unseen data, the dataset is partitioned into training and testing subsets. This segregation enables the model to learn and then evaluate its

performance on previously unseen examples. For this purpose, we opt for a ratio of 80% for the training set and 20% for the validation set, effectively maintaining an 8:2 distribution.

4) One-Hot Encoding

In tasks such as digit recognition using the MNIST dataset, the labels represent categorical values (ranging from 0 to 9), signifying the respective class of the image (i.e., the digit it portrays).

During the training of a classification model, a loss function quantifies the dissimilarity between predicted and actual labels. The use of one-hot encoded labels allows the calculation of these losses. This convenience arises from the fact that each class is distinctly represented by an individual binary element within the encoded vector. This process involves the conversion of a single label (e.g., 3) into a binary vector (e.g., [0, 0, 0, 1, 0, 0, 0, 0, 0, 0]), clearly indicating the class membership.

D. MODEL ARCHITECTURE

The basic architecture of model consists of input layer few hidden layers and 10 output layers. Since there are 784 pixels, where each pixel acts as a feature, the input layer always has 784 neurons and 1 extra bias term. The number of hidden layers and number of neurons in each layer is a hyper-parameter which can be experimented. Finally the output layer contains 10 neurons one for each class.

Besides that each neuron must contain an activation function to introduce non linearity. If there were no any activation then the entire network becomes a linear system and is just equivalent to a single neuron no matter how deep the network is. By adding activation function are trying to separate the classes with non linear decision boundaries. Some choices of activation functions we experimented with are: A) Sigmoid activation B) Relu Activation C) Tanh activation

V. RESULT ANALYSIS

For initializing initial values for the neural network's weights and biases, our approach involved drawing values from a uniform distribution spanning between -0.5 and 0.5.

A. WEIGHT INITIALIZATION

TABLE 1. Performance for Different weight initialization

Activation	Initialization method	Accuracy	Loss
Relu	Random uniform	0.7977	0.8490
Relu	He	0.8512	0.6239
Relu	Xavier	0.8516	0.6234
Tanh	Random uniform	0.7461	0.8548
Tanh	He	0.8594	0.5792
Tanh	Xavier	0.8572	0.5879
Sigmoid	Random uniform	0.6290	1.1680
Sigmoid	He	0.7400	1.6020
Sigmoid	Xavier	0.7287	0.6668

During the initial phases of training, the random weight initialization yielded higher loss values in the initial iterations,

accompanied by diminished accuracy. However, the accuracy exhibited an upward trend as training progressed, and the loss gradually diminished. Interestingly, it was observed that the randomly initialized model achieved a high final accuracy and lower loss.

Comparing the outcomes of random initialization to those of Xavier and He initialization revealed notable differences. Both the Xavier and He initialization methods showcased substantially lower loss values during the initial iterations, contributing to an accelerated convergence rate of the model. In particular, for a model featuring a single hidden layer comprising 256 nodes, the loss for random initialization was approximately 0.84, whereas Xavier and He initialization yielded losses around 0.62.

Moreover, the weight initialization techniques were applied across various activation functions, including ReLU, tanh, and sigmoid. Remarkably, this effect persisted consistently across all three activation functions.

This underlines the crucial role that weight initialization plays in the model's convergence and overall performance, with Xavier and He initialization providing a beneficial head start in terms of convergence speed and loss reduction.

B. DETERMINING THE NUMBER OF NEURONS IN A HIDDEN LAYER

Determining the number of nodes in each layer of the neural network involves considering various factors. The input layer's size should correspond to the input dimensions, which are 784 in this case. Additionally, the output layer needs a node for every potential output label. With 10 labels (0-9 numbers) in this scenario, the output layer will encompass 10 neurons.

Nonetheless, the quantity of neurons within the hidden layers remains a matter of design. The network's complexity relies on considerations like:

- **Data Complexity:** Complex datasets often require larger networks. If the data exhibits intricate patterns or relationships, a larger network may be necessary.
- **Overfitting:** Too many neurons can lead to overfitting, where the model becomes overly specialized in training data and doesn't generalize to new data. Regularization techniques can help prevent this.
- **Computational Resources:** Larger networks demand more computational power and memory, influencing the network's size.
- **Empirical Experimentation:** Experimentation is crucial. Starting with a small network, evaluating its performance, and incrementally increasing its size helps identify optimal points between diminishing returns and overfitting.
- **Preceding Architectures:** Existing architectures or state-of-the-art models in your domain can guide the range of neurons used effectively.

To determine a suitable count, we initiated a neural network with xavier weight initialization. First, we implemented a network with:

- 784 input neurons
- variable hidden layer 1 neurons
- 128 hidden layers 2 neurons
- 10 output neurons

For that we calculated the accuracy and loss, and increased the number of hidden neurons by increments of 5. in our observation we saw that the hidden layer containing nodes less than 50 had significantly lower accuracy and higher losses than the accuracy and loss with higher neuron.

Interestingly, our investigation revealed a point of saturation in performance enhancement at approximately 125 neurons within the hidden layer. As seen in the figure (5), this plateau persisted until reaching the 300-neuron layer. The rationale behind this phenomenon can be attributed to the fact that around 125 neurons, the network possesses adequate capacity to capture the intrinsic patterns within the MNIST dataset. Beyond this threshold, further escalation in neuron count fails to yield corresponding accuracy improvements. The phenomenon points to the presence of diminishing returns, where excessive neuron augmentation does not proportionally enhance the model's predictive capabilities.

C. SINGLE VS MULTIPLE HIDDEN LAYERS

We conducted a performance analysis by comparing two distinct neural network architectures using the MNIST dataset. The first architecture featured a single hidden layer containing 128 neurons, while the second configuration comprised multiple hidden layers with 512, 224, and 128 neurons. Both models underwent training for 100 epochs, employing Xavier weight initialization to ensure a balanced starting point.

TABLE 2. Performance for Different Network Architectures

Architecture	Epoch	Accuracy	Loss
[784, 512, 224, 128, 10]	100	0.8968	0.3833
[784, 128, 10]	100	0.8949	0.3922

Surprisingly, our results revealed minimal disparity in performance between the two architectures. The observed similarity in performance prompts us to consider potential reasons behind this outcome. One plausible explanation is rooted in the nature of the MNIST dataset itself. Given the dataset's relative simplicity and the prominence of its inherent patterns, it is conceivable that these patterns were adequately captured by the single hidden layer network we employed initially. While we cannot definitively conclude this to be the sole reason, it remains a valid hypothesis.

It's worth highlighting that the complexity of the dataset plays a pivotal role in determining the efficacy of different network architectures. While our results may not starkly differentiate the two architectures on the MNIST dataset, this should not overshadow the potential advantages of multi-layer networks in more intricate datasets. These deep architectures have a unique capability to discern intricate hierarchical features that might elude a shallower, single-layer network.

However, it's essential to exercise caution in interpreting these findings. While we anticipate discernible differences on more complex datasets, the interplay of various factors—such as activation functions, optimization techniques, and overall architecture capacity—can confound straightforward expectations. Therefore, future work should encompass a broader spectrum of datasets and rigorous cross-validation to comprehensively evaluate how different architectures perform across diverse scenarios.

D. ACTIVATION FUNCTIONS

We tested different activation functions in the hidden layer, departing from our previous use of the ReLU activation function. To expedite convergence, we employed Xavier weight initialization without dropout for all three experiments, resulting in notable speed enhancements in learning.

VI. NETWORK ARCHITECTURE

TABLE 3. Performance for Different Network Architectures and Activation Functions

Architecture	Activation Function	Accuracy	Loss
[784, 128, 10]	Relu	0.87	0.50
[784, 128, 10]	Tanh	0.87	0.51
[784, 128, 10]	Sigmoid	0.73	1.45

With the ReLU activation function, we observed faster convergence. By the 10th epoch, accuracy had already reached 65.18%, climbing to 87.25% by the 80th epoch, when our training concluded.

Switching to the tanh activation function further accelerated convergence, surpassing the progress achieved by ReLU. Accuracy stood at 70.77% by the 10th epoch and peaked at 86.86% by the end of the 80th epoch.

However, the sigmoid activation function demonstrated a slower trajectory among the three. Its accuracy was a mere 34.01% in the 10th epoch, gradually increasing to 73.97% by the 80th epoch.

The variation in performance across activation functions (ReLU, tanh, sigmoid) can be attributed to their distinct traits and their fit with the dataset and learning process.

ReLU: ReLU's simplicity and ability to tackle the vanishing gradient issue make it effective. It allows faster learning for positive inputs, which aligns well with the image data in MNIST. The dataset's clear pixel value transitions suit ReLU's activation patterns, aiding in capturing these patterns effectively. In the context of the MNIST dataset, where images represent handwritten digits, these transitions could represent edges, boundaries, or distinct features within the digit images.'

tanh: tanh's range from -1 to 1 suits datasets with zero-centered data like normalized MNIST pixels. Despite its vanishing gradient issue for extreme inputs, tanh can model both positive and negative data relationships due to its symmetry around zero.

Sigmoid: Sigmoid's (0, 1) range and sensitivity around zero could explain its performance. Yet, its rapid gradi-

ent saturation for strong inputs hampers learning, especially in deeper networks. For MNIST's pattern-centric data, sigmoid's saturation for intense inputs can limit its pattern recognition ability.

Weight Initialization: Xavier initialization's role in gradient normalization helps counter the vanishing gradient problem, benefiting ReLU and tanh activations.

Activation function performance depends on their characteristics, compatibility with data distribution, and interaction with weight initialization. The discrepancies observed in accuracy and convergence stem from these factors, highlighting how certain activations align better with MNIST's intrinsic attributes.

A. DROPOUT

Dropout is a regularization technique used in neural networks to prevent overfitting and improve generalization performance. Overfitting occurs when a neural network learns to perform exceptionally well on the training data but struggles to perform well on unseen data (validation or test data). Dropout is a simple yet effective way to combat overfitting by introducing randomness during training.

We used neuron deactivation for our dropout. To experiment with dropout, we masked a fraction of neurons in the layers randomly by passing a parameter `dropout_prob`. This creates a dropout mask so that during the training phase, the output of masked neurons is set to zero.

First we trained the models for 50 epoch, the table below shows the result. But we realized the network could further be trained. So, for the second run, we a model 512 hidden

TABLE 4. Effect of Dropout Probability on Network Performance

Network	epoch	Dropout Probability	Accuracy	Loss
[784, 224, 10]	50	0	0.85	0.63
[784, 224, 10]	50	0.1	0.83	0.67
[784, 224, 10]	50	0.2	0.81	0.75
[784, 224, 10]	50	0.5	0.70	1.03
[784, 224, 10]	50	0.8	0.44	1.67

neurons for 150 epochs. We trained three single hidden layer

TABLE 5. Effect of Dropout Probability on Network Performance

Network	Epoch	Dropout Probability	Accuracy	Loss
[784, 512, 10]	150	0	0.89	0.3778
[784, 512, 10]	150	0.1	0.89	0.3978
[784, 512, 10]	150	0.2	0.88	0.4243

neural networks with 0%, 10%, and 20% of the neurons deactivated during training. we observed a training accuracy of 89% for the models with no dropout and 10% dropout. For the model with 20% dropout, we observed the accuracy to be 88%. Moving to test accuracy and loss, as seen in the figure 11, the accuracy was almost the same in all three models. However, the loss slightly increased as the dropout increased.

It's important to consider a few potential reasons for this result:

- **Dataset Size and Complexity:** The MNIST dataset is relatively small and contains well-defined patterns. In such cases, the benefits of dropout regularization might not be as pronounced as in larger and more complex datasets. Dropout is particularly effective when dealing with larger and more intricate datasets, where overfitting is a more significant concern.
- **Model Architecture:** The architecture of neural network can also influence the effectiveness of dropout. If your network is already relatively simple or shallow, overfitting might not be a major issue even without dropout. Dropout tends to provide more pronounced benefits in deeper networks.

VII. CONCLUSION

This project aimed to enhance the performance of a Multi-Layer Perceptron (MLP) on the MNIST dataset by systematically exploring various factors that influence its training process and final accuracy. Through the meticulous design of experiments, we investigated the impact of different weight initialization techniques, activation functions, and dropout regularization on the model's performance.

Our findings revealed that weight initialization plays a crucial role in the convergence speed and overall performance of the MLP. Certain techniques, such as Xavier/Glorot initialization, demonstrated a clear advantage in facilitating faster convergence and higher accuracy. Moreover, the choice of activation functions showcased distinct effects on the model's ability to capture complex patterns within the dataset. We observed that ReLU and its variants tend to outperform others by mitigating the vanishing gradient problem and enabling more effective learning.

Incorporating dropout regularization during training can show positive effect on the model's generalization ability. This technique effectively reduced overfitting and improved the MLP's performance on unseen data, leading to a more robust and reliable model.

As machine learning practitioners, understanding the nuances of these factors is pivotal in designing effective neural network architectures. This project underscores the importance of a systematic approach to hyperparameter tuning and experimentation. By optimizing these key components—weight initialization, activation functions, and dropout regularization—we not only enhanced the accuracy of our MLP on the MNIST dataset but also gained valuable insights into the inner workings of neural networks.

In the ever-evolving field of deep learning, this project serves as a testament to the significance of thorough experimentation and analysis. As we continue to push the boundaries of AI research, the knowledge gained from this project will undoubtedly contribute to the development of more efficient and accurate neural network models for a wide range of applications, ultimately advancing the field as a whole.

REFERENCES

- [1] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *"Deep Learning."* MIT Press.
- [2] Daniel Jurafsky & James H. Martin. *Speech and Language Processing*.
- [3] M. Narasimha Murty, V. Susheela Devi. *Pattern Recognition: An Algorithmic Approach*.



PILOT KHADKA is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. During his time at Thapathali Campus, Pilot has actively engaged in various academic and extracurricular activities. He has participated in coding competitions, collaborated on software development projects, and demonstrated a keen interest in exploring new technologies.



KSHITIZ POUDEL is a student at Institute of Engineering, Thapathali Campus. He is expected to graduate in Bachelor of Computer Engineering in 2024. His relentless pursuit of knowledge and dedication to uplifting and empowering others make him an exceptional contributor and an invaluable asset to the academic community.

VIII. APPENDIX

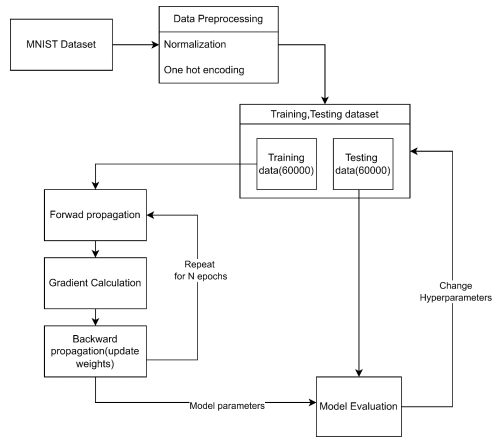


FIGURE 1. System Block Diagram

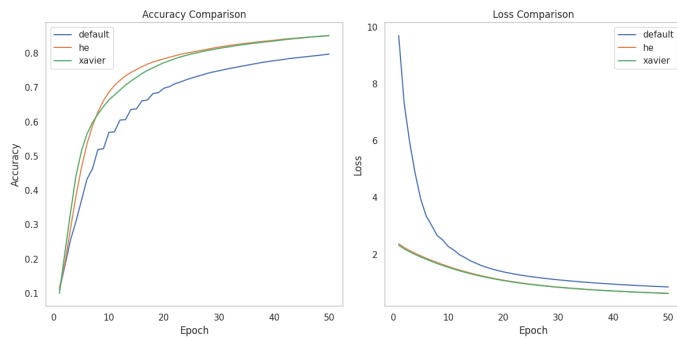


FIGURE 2. Weight initialization for Relu activation

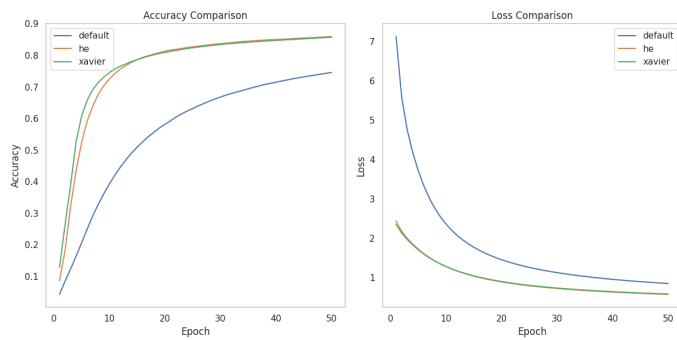


FIGURE 3. Weight initialization for Tanh activation

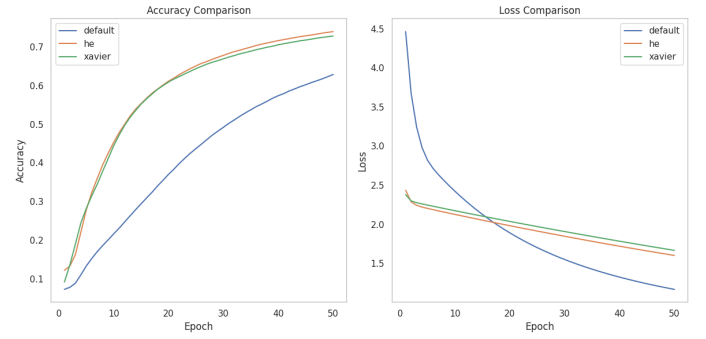


FIGURE 4. Weight initialization for sigmoid activation

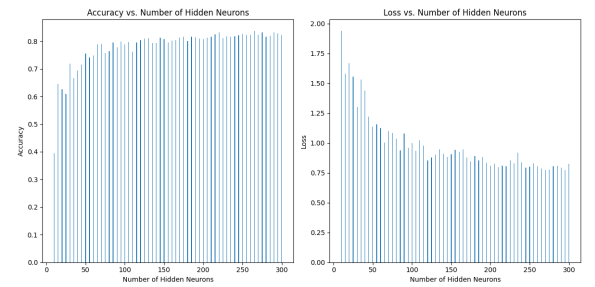


FIGURE 5. System Block Diagram

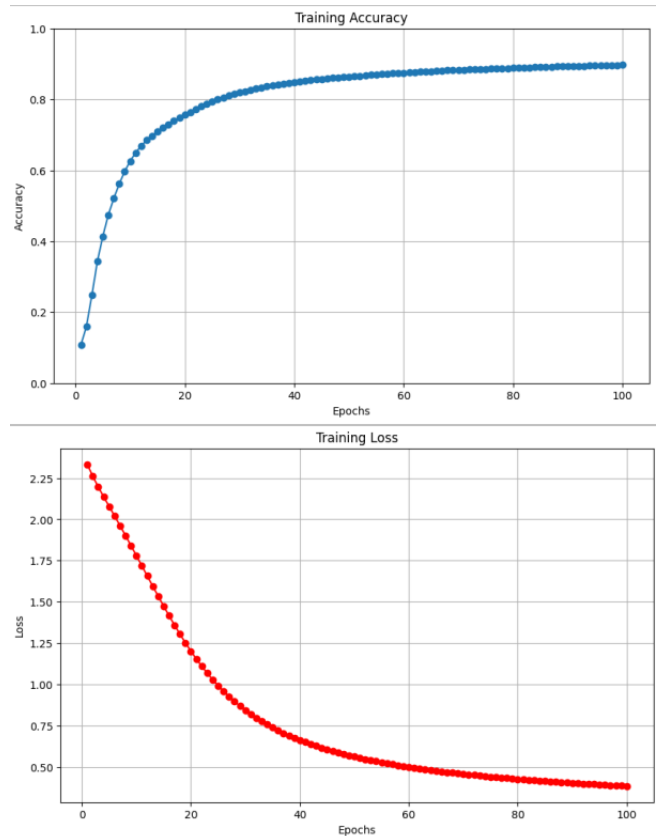


FIGURE 6. Accuracy/Loss curve of Deeper Network

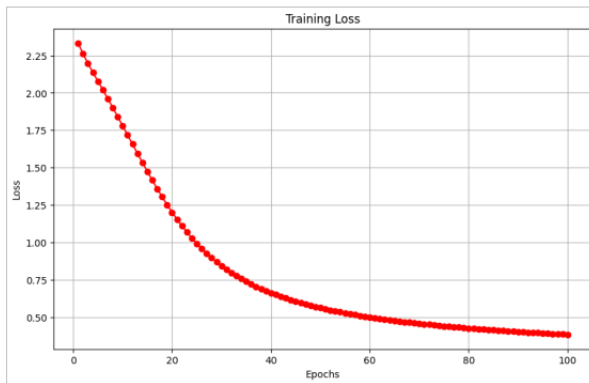
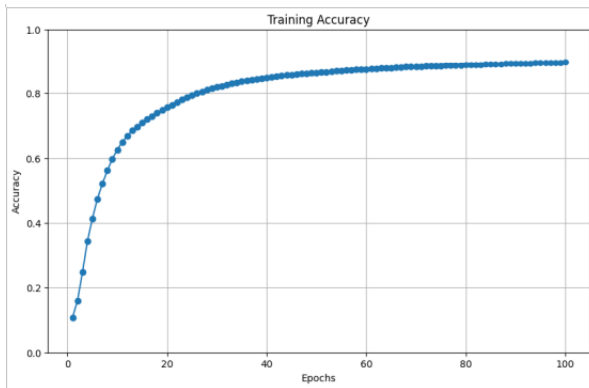


FIGURE 7. Accuracy/Loss curve of Shallow Network

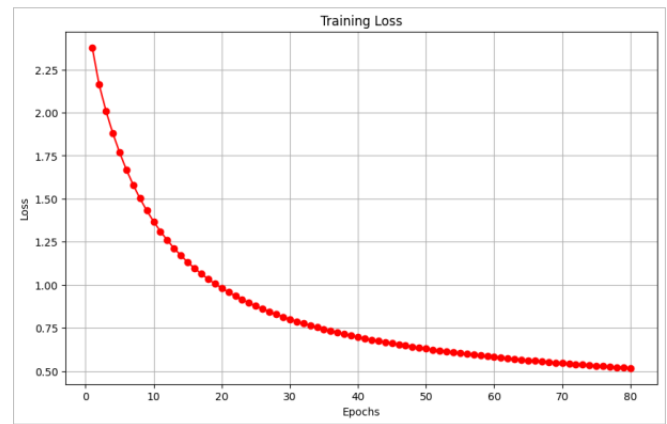
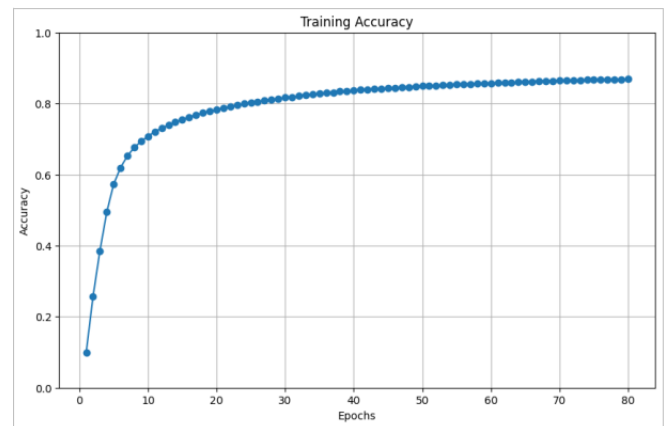


FIGURE 9. Accuracy/Loss curve of network with tanh

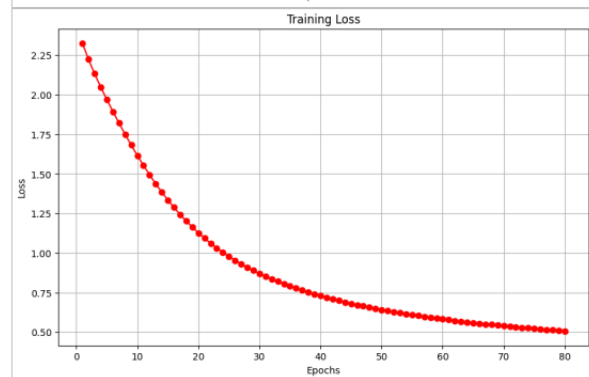
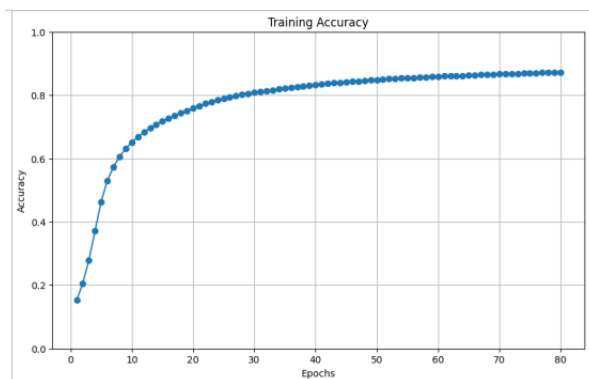


FIGURE 8. Accuracy/Loss curve of Network with relu

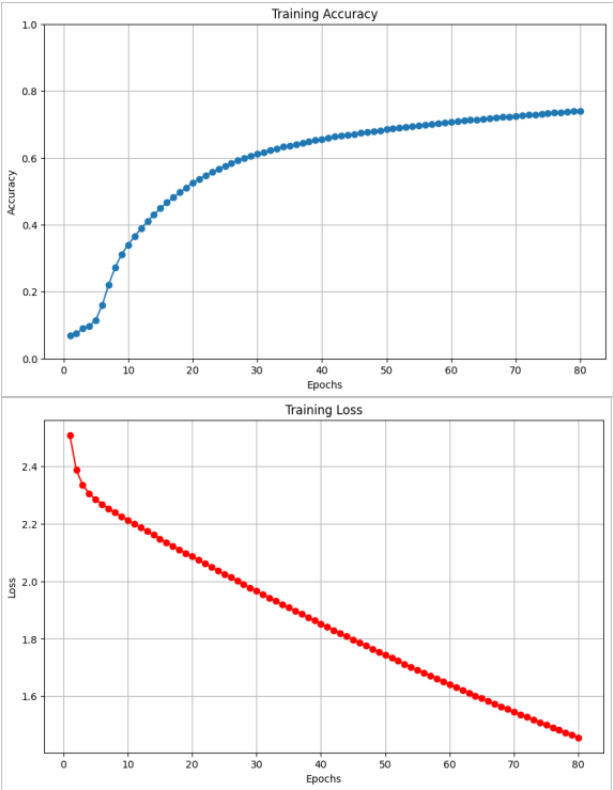


FIGURE 10. Accuracy/Loss curve of network with sigmoid

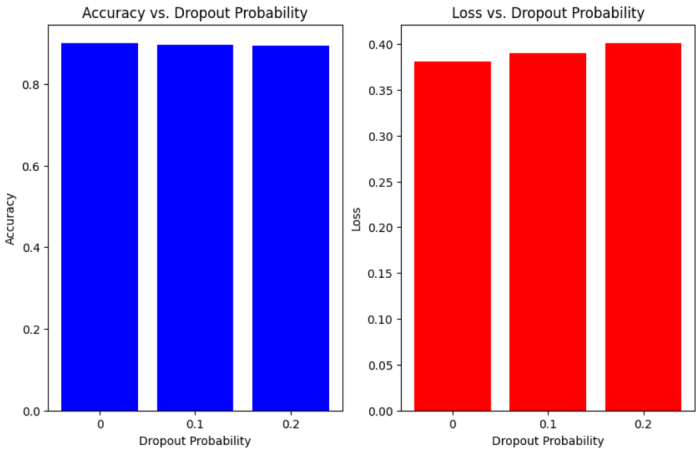


FIGURE 12. Accuracy/loss comparison with dropouts 0,0.1,0.2

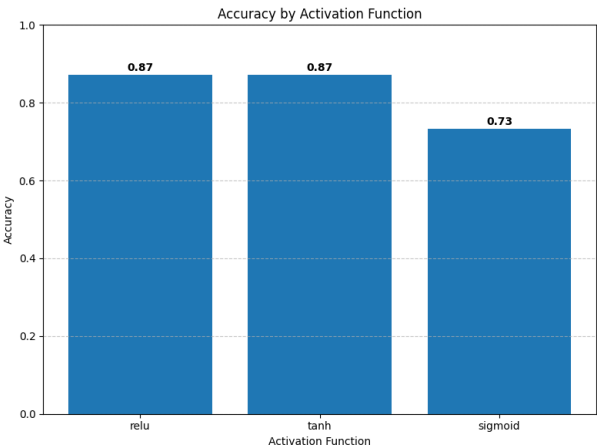


FIGURE 11. accuracy comparison of three activations

IX. CODE

```

1
2
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 def init_params(layers_dims, initialization="
7     default"):
8     """
9     Initializes the parameters (weights and biases
10    ) for a neural network.
11
12    Parameters:
13    layers_dims (list): List of integers
14    representing the dimensions of each layer in
15    the neural network.
16    initialization (str): Initialization method to
17    use for parameter initialization.
18        - "default": Random initialization using a
19        uniform distribution between -0.5 and 0.5.
20        - "xavier": Xavier/Glorot initialization
21        for better training convergence.
22        - "he": He initialization for training
23        deeper networks.
24
25    Returns:
26    params (dict): A dictionary containing the
27    initialized parameters for each layer.
28    """
29    params = {}
30
31    for layer in range(1, len(layers_dims)):
32        input_dim = layers_dims[layer - 1]
33        output_dim = layers_dims[layer]
34
35        # Xavier weight initialization
36        if initialization == "xavier":
37            params["W" + str(layer)] = np.random.
38            uniform(
39                -np.sqrt(6 / (input_dim +
40                output_dim)),
41                np.sqrt(6 / (input_dim +
42                output_dim)),
43                (output_dim, input_dim),
44            )
45            params["b" + str(layer)] = np.random.
46            uniform(
47                -np.sqrt(6 / (input_dim +
48                output_dim)),
49                np.sqrt(6 / (input_dim +
50                output_dim)),
51                (output_dim, 1),
52            )
53
54        # He weight initialization
55        elif initialization == "he":
56            params["W" + str(layer)] = np.random.
57            normal(
58                0, np.sqrt(2 / input_dim), (
59                output_dim, input_dim)
60            )
61            params["b" + str(layer)] = np.random.
62            normal(
63                0, np.sqrt(2 / input_dim), (
64                output_dim, 1)
65            )
66
67        # Default initialization uniform
68        distribution
69        else:
70            params["W" + str(layer)] = np.random.
71            uniform(
72                -0.5, 0.5, (output_dim, input_dim)
73            )
74
75            params["b" + str(layer)] = np.random.
76            uniform(-0.5, 0.5, (output_dim, 1))
77            return params
78
79 def forward_prop(X, params, activation="relu",
80 dropout_prob=0):
81     """
82     Perform forward propagation through the neural
83     network to compute activations.
84
85     Parameters:
86     X (numpy.ndarray): Input data of shape (
87     input_size, m), where input_size is the number
88     of features and m is the number of examples.
89     params (dict): Dictionary containing
90     network parameters including weights (W) and
91     biases (b) for each layer.
92     dropout_prob (float): Dropout probability
93     for applying dropout regularization to hidden
94     layers (default 0).
95
96     Returns:
97     activations (dict): Dictionary containing
98     computed activations for each layer.
99     dropout_masks (dict): Dictionary
100    containing dropout masks applied to hidden
101    layers.
102    """
103    # Get the number of layers directly from the
104    length of W parameters
105    L = len(params) // 2
106    activations = {}
107    activations["A0"] = X # Input activations
108    dropout_masks = {} # Dictionary to store
109    dropout masks
110
111    for l in range(1, L):
112        # Calculate Z and A for intermediate
113        layers using ReLU activation
114        activations["Z" + str(l)] = np.dot(params[
115        "W" + str(l)], activations["A" + str(l - 1)])
116        + params["b" + str(l)]
117
118        if activation == 'tanh':
119            activations["A" + str(l)] = _tanh(
120            activations["Z" + str(l)]) # Tanh activation
121        elif activation == 'sigmoid':
122            activations["A" + str(l)] = 1 / (1 +
123            np.exp(-activations["Z" + str(l)])) # Sigmoid
124            activation
125        else :
126            activations["A" + str(l)] = _relu(
127            activations["Z" + str(l)]) # ReLU activation
128
129        if l < L:
130            dropout_mask = np.random.rand(*
131            activations["A" + str(l)].shape) < (1 -
132            dropout_prob) # Inverted dropout mask
133            activations["A" + str(l)] *=
134            dropout_mask
135            activations["A" + str(l)] /= (1 -
136            dropout_prob) # Scale to maintain expected
137            value
138            dropout_masks["D" + str(l)] =
139            dropout_mask
140
141    # Calculate Z and A for the output layer using
142    softmax activation
143    activations["Z" + str(L)] = np.dot(params["W"
144    + str(L)], activations["A" + str(L - 1)]) +
145    params["b" + str(L)]
146    exp_scores = np.exp(activations["Z" + str(L)])
147    activations["A" + str(L)] = exp_scores / np.
148    sum(exp_scores, axis=0, keepdims=True) #

```

```

Softmax activation
95
96 return activations, dropout_masks
97
98
99 def back_prop(activations, params, Y, dropout_masks
100 , activation="relu"):
101     """
102     Perform backpropagation to compute gradients
103     of the loss with respect to parameters.
104
105     Parameters:
106         activations (dict): Dictionary containing
107         computed activations for each layer during
108         forward propagation.
109         params (dict): Dictionary containing
110         network parameters including weights (W) and
111         biases (b) for each layer.
112         Y (numpy.ndarray): True labels (ground
113         truth) of shape (1, m), where m is the number
114         of examples.
115         dropout_masks (dict): Dictionary
116         containing dropout masks applied to hidden
117         layers.
118
119     Returns:
120         grads (dict): Dictionary containing
121         computed gradients of the loss with respect to
122         parameters.
123     """
124     L = len(params) // 2
125     one_hot_Y = one_hot_encode(Y)
126     m = one_hot_Y.shape[1]
127
128     derivatives = {}
129     grads = {}
130
131     # for layer L
132     derivatives["dZ" + str(L)] = activations["A" +
133 str(L)] - one_hot_Y
134     grads["dW" + str(L)] = (
135         1 / m * np.dot(derivatives["dZ" + str(L)],
136         activations["A" + str(L - 1)].T)
137     )
138     grads["db" + str(L)] = 1 / m * np.sum(
139     derivatives["dZ" + str(L)])
140
141     # for layers L-1 to 1
142     for l in reversed(range(1, L)):
143         if activation == 'relu':
144             _activation_derivative =
145             _relu_derivative
146         elif activation == 'tanh':
147             _activation_derivative =
148             _tanh_derivative
149         elif activation == 'sigmoid':
150             _activation_derivative =
151             _sigmoid_derivative
152
153         derivatives["dZ" + str(l)] = np.dot(
154             params["W" + str(l + 1)].T,
155             derivatives["dZ" + str(l + 1)]
156         ) * _activation_derivative(activations["Z"
157 + str(l)])
158
159         # apply dropout mask
160         if l < L:
161             derivatives["dZ" + str(l)] *=
162             dropout_masks["D" + str(l)]
163
164         grads["dW" + str(l)] = (
165             1 / m * np.dot(derivatives["dZ" + str(
166 l)], activations["A" + str(l - 1)].T)
167         )
168
169         grads["db" + str(l)] = (
170             1 / m * np.sum(derivatives["dZ" + str(
171 l)], axis=1, keepdims=True)
172         )
173
174     return grads
175
176 def update_params(params, grads, alpha):
177     """
178     Update network parameters using gradient
179     descent optimization.
180
181     Parameters:
182         params (dict): Dictionary containing
183         network parameters including weights (W) and
184         biases (b) for each layer.
185         grads (dict): Dictionary containing
186         gradients of the loss with respect to
187         parameters.
188         alpha (float): Learning rate, controlling
189         the step size of parameter updates.
190
191     Returns:
192         params_updated (dict): Dictionary
193         containing updated network parameters.
194     """
195     # number of layers
196     L = len(params) // 2
197
198     params_updated = {}
199     for l in range(1, L + 1):
200         params_updated["W" + str(l)] = (
201             params["W" + str(l)] - alpha * grads["
202 dW" + str(l)]
203         )
204         params_updated["b" + str(l)] = (
205             params["b" + str(l)] - alpha * grads["
206 db" + str(l)]
207         )
208
209     return params_updated
210
211 def train(X, Y, params, max_iter=10, learning_rate
212 =0.1, dropout_prob=0, activation="relu"):
213     """
214     Trains a neural network model using gradient
215     descent optimization.
216
217     Parameters:
218         X (numpy.ndarray): Input data of shape (
219         input_size, num_samples).
220         Y (numpy.ndarray): Ground truth labels of
221         shape (1, num_samples).
222         params (dict): Dictionary containing the
223         parameters of the neural network.
224             Keys are "W1", "b1", ..., "WL",
225             "bL" where L is the number of layers.
226         max_iter (int, optional): Maximum number of
227         training iterations. Default is 10.
228         learning_rate (float, optional): Learning rate
229         for gradient descent. Default is 0.1.
230         dropout_prob (float, optional): Dropout
231         probability for regularization. Default is 0.
232         activation (str, optional): Activation
233         function to use in hidden layers. Default is "
234 relu".
235
236     Returns:
237         tuple: A tuple containing updated parameters,
238         list of accuracies per iteration, and list of
239         losses per iteration.
240     """
241     # Initialize parameters W1, b1 for layers 1

```

```

197     =1,...,L
198     L = len(params) // 2
199     accuracies = []
200     losses = []
201
202     for iteration in range(1, max_iter + 1):
203         # Forward propagation
204         activations, dropout_mask = forward_prop(X,
205           params, activation, dropout_prob)
206
207         # Make predictions
208         Y_hat = get_predictions(activations["A" +
209           str(L)])
210
211         # Compute accuracy
212         accuracy = get_accuracy(Y_hat, Y)
213         accuracies.append(accuracy)
214
215         # Compute loss (cross-entropy)
216         loss = cross_entropy(one_hot_encode(Y),
217           activations["A" + str(L)])
218         losses.append(loss)
219
220         # Backpropagation
221         gradients = back_prop(activations, params,
222           Y, dropout_mask, activation)
223
224         # Update parameters
225         params = update_params(params, gradients,
226           learning_rate)
227
228         # Print progress
229         if iteration == 1 or (iteration % 5) == 0:
230             print("Iteration {}: Accuracy = {},
231               Loss = {}".format(iteration, accuracy, loss))
232
233     return params, accuracies, losses
234
235 def test(X_test, Y_test, params, activation):
236     """
237     Evaluates the trained neural network model on
238     test data.
239
240     Parameters:
241     X_test (numpy.ndarray): Test input data of
242       shape (input_size, num_samples).
243     Y_test (numpy.ndarray): Ground truth labels
244       for test data of shape (1, num_samples).
245     params (dict): Dictionary containing the
246       parameters of the neural network.
247       Keys are "W1", "b1", ..., "WL",
248       "bL" where L is the number of layers.
249     activation (str): Activation function used in
250       the hidden layers during forward propagation.
251
252     Returns:
253     tuple: A tuple containing test accuracy and
254       loss.
255     """
256     L = len(params) // 2
257
258     # Forward propagation
259     activations, _ = forward_prop(X_test, params,
260       activation, 0)
261
262     # Make predictions
263     Y_hat = get_predictions(activations["A" + str(
264       L)])
265
266     # Compute test accuracy
267     test_accuracy = get_accuracy(Y_hat, Y_test)
268     loss = cross_entropy(one_hot_encode(Y_test),
269       activations["A" + str(L)])
270
271     return test_accuracy, loss
272
273 def one_hot_encode(Y):
274     Y_one_hot = np.zeros((Y.shape[0], Y.max() + 1)
275       )
276     # set to 1 the correct indices
277     Y_one_hot[np.arange(Y.shape[0]), Y] = 1
278     # transpose
279     Y_one_hot = Y_one_hot.T
280     return Y_one_hot
281
282 def cross_entropy(Y_one_hot, Y_hat, epsilon=1e-10):
283     # clip predictions to avoid values of 0 and 1
284     Y_hat = np.clip(Y_hat, epsilon, 1.0 - epsilon)
285     # sum on the columns of Y_hat * np.log(Y),
286     # then take the mean
287     # between the m samples
288     cross_entropy = -np.mean(np.sum(Y_one_hot * np
289       .log(Y_hat), axis=0))
290     return cross_entropy
291
292 def shuffle_rows(data):
293     # Convert input dataframe to ndarray
294     data = np.array(data)
295     np.random.shuffle(data)
296     return data
297
298 def normalize_pixels(data):
299     return data / 255.0
300
301 def _relu(Z):
302     return np.maximum(Z, 0)
303
304 def _softmax(Z):
305     A = np.exp(Z) / sum(np.exp(Z))
306     return A
307
308 def _relu_derivative(Z):
309     return Z > 0
310
311 def _softmax_derivative(Z):
312     dZ = np.exp(Z) / sum(np.exp(Z)) * (1.0 - np
313       .exp(Z) / sum(np.exp(Z)))
314     return dZ
315
316 def _tanh(x):
317     return (np.exp(x) - np.exp(-x)) / (np.exp(x) +
318       np.exp(-x))
319
320 def _tanh_derivative(x):
321     return 1 - np.tanh(x) ** 2
322
323 def _sigmoid_derivative(x):
324     sigmoid = 1 / (1 + np.exp(-x))
325     return sigmoid * (1 - sigmoid)
326
327 def get_predictions(AL):
328     # get the max index by the columns
329     return np.argmax(AL, axis=0)
330
331 def get_accuracy(Y_hat, Y):
332     return np.sum(Y_hat == Y) / Y.size
333
334 def plot_accuracy_and_loss(accuracies, losses,
335   max_iter):
336     # Plot training accuracy
337     plt.figure(figsize=(10, 6))

```



```

321 plt.plot(range(1, max_iter + 1), accuracies,
322          marker='o')
323 plt.title("Training Accuracy")
324 plt.xlabel("Epochs")
325 plt.ylabel("Accuracy")
326 plt.grid()
327 plt.show()
328
329 # Plot training loss
330 plt.figure(figsize=(10, 6))
331 plt.plot(range(1, max_iter + 1), losses,
332          marker='o', color='r')
333 plt.title("Training Loss")
334 plt.xlabel("Epochs")
335 plt.ylabel("Loss")
336 plt.grid()
337 plt.show()
338
339 # load training data
340 import pandas as pd
341 df_train = pd.read_csv("mnist_train.csv")
342
343 # shuffle the data
344 df_train = shuffle_rows(df_train)
345
346 # split train and validation set
347 train_val_split = 0.8
348 train_size = round(df_train.shape[0] *
349                    train_val_split)
350 data_train = df_train[:train_size, :].T
351 data_val = df_train[train_size:, :].T
352
353 # divide input features and target feature
354 X_train = data_train[1:]
355 y_train = data_train[0]
356 X_test = data_val[1:]
357 y_test = data_val[0]
358
359 # normalize training and val sets
360 X_train = normalize_pixels(X_train)
361 X_test = normalize_pixels(X_test)
362
363 print(X_train.shape)
364 print(y_train.shape)
365
366 # set network and optimizer parameters
367 layers_dims = [784, 128, 10]
368 # layers_dims = [784, 10, 10]
369 max_iter = 80
370 alpha = 0.1
371 dropout_prob = 0
372 accuracies = []
373 losses = []
374
375 activation = ["relu", "tanh", "sigmoid"]
376 for a in activation:
377     params = init_params(layers_dims, 'xavier')
378
379     # train the network
380     trained_params, train_acc, train_loss = train(
381         X_train, y_train, params, max_iter, alpha,
382         dropout_prob, a
383     )
384     plot_accuracy_and_loss(train_acc, train_loss,
385                           max_iter)
386     accuracy, loss = test(X_test, y_test,
387                          trained_params, a)
388     accuracies.append(accuracy)
389     losses.append(loss)
390
391 import seaborn as sns
392
393 plt.figure(figsize=(8, 6)) # Adjust the figure
394 size
395 # Create the bar plot
396 plt.bar(activation, accuracies)
397
398 # Adding title and labels
399 plt.title('Accuracy by Activation Function')
400 plt.xlabel('Activation Function')
401 plt.ylabel('Accuracy')
402
403 # Adding data labels on top of the bars
404 for i, acc in enumerate(accuracies):
405     plt.text(i, acc + 0.01, f'{acc:.2f}', ha=
406              'center', color='black', fontweight='bold')
407
408 plt.ylim(0, 1) # Set y-axis limits
409 plt.grid(axis='y', linestyle='--', alpha=0.7) #
410 Add horizontal grid lines
411
412 plt.tight_layout() # Adjust spacing
413
414 plt.show() # Display the plot
415
416 test_accuracy, test_loss = test(X_test, y_test,
417                                trained_params)
418 print("Test Accuracy: {:.2f}, Test loss: {:.2f}".
419       format(test_accuracy * 100, test_loss))
420
421 hidden_neurons_range = range(10, 201, 5)
422 accuracies = []
423 losses = []
424
425 for hidden_neurons in hidden_neurons_range:
426     layers_dims = [784, hidden_neurons, 128, 10]
427     # Update hidden layer neurons
428     params = init_params(layers_dims, 'xavier') #
429     Reinitialize parameters
430
431     # Train the network
432     trained_params, _, _ = train(X_train, y_train,
433                                 params, max_iter, alpha, dropout_prob)
434
435     accuracy, loss = test(X_test, y_test,
436                          trained_params)
437     accuracies.append(accuracy)
438     losses.append(loss)
439
440 plt.figure(figsize=(12, 6))
441
442 plt.subplot(1, 2, 1)
443 plt.bar(hidden_neurons_range, accuracies)
444 plt.xlabel('Number of Hidden Neurons')
445 plt.ylabel('Accuracy')
446 plt.title('Accuracy vs. Number of Hidden Neurons')
447
448 plt.subplot(1, 2, 2)
449 plt.bar(hidden_neurons_range, losses)
450 plt.xlabel('Number of Hidden Neurons')
451 plt.ylabel('Loss')
452 plt.title('Loss vs. Number of Hidden Neurons')
453
454 plt.tight_layout()
455 plt.show()
456
457 # Deeper networks
458 m_layers_dims = [784, 512, 224, 128, 10] # Update
459 hidden layer neurons
460 m_params = init_params(layers_dims, 'xavier') #
461 Reinitialize parameters
462
463 max_iter=100
464 # Train the network
465 m_trained_params, m_train_accuracy, m_train_loss =
466 train(X_train, y_train, params, max_iter,

```

```

alpha, dropout_prob)
451 m_accuracy,m_loss = test(X_test,y_test,
452     trained_params)
453
454 plot_accuracy_and_loss(m_train_accuracy,
455     m_train_loss, max_iter)
456
457 s_layers_dims = [784,128, 10] # Update hidden
458     layer neurons
459 s_params = init_params(layers_dims,'xavier') #
460     Reinitialize parameters
461 max_iter=100
462 # Train the network
463 s_trained_params, s_train_accuracy,s_train_loss =
464     train(X_train, y_train, params, max_iter,
465     alpha, dropout_prob)
466
467 s_accuracy,s_loss = test(X_test,y_test,
468     trained_params)
469
470 plot_accuracy_and_loss(s_train_accuracy,
471     s_train_loss,max_iter)
472
473 #Compare initialization methods
474
475 import matplotlib.pyplot as plt
476 # set network and optimizer parameters
477 layers_dims = [784, 256, 10]
478 # layers_dims = [784, 10, 10]
479 max_iter = 50
480 alpha = 0.1
481 dropout_prob= 00
482
483 # Define the initialization methods you want to
484     compare
485 initialization_methods = ['default', 'he', 'xavier
486     ']
487 activation = ['default','relu','sigmoid']
488 # Initialize an empty dictionary to store
489     accuracies and losses for each method
490 results = {}
491
492 # Loop through each initialization method
493 for a in activation:
494     for method in initialization_methods:
495         # Initialize parameters using the current
496         method
497         params = init_params(layers_dims, method)
498
499         # Train the network and collect accuracies
500         and losses
501         updated_params,accuracies, losses = train(
502             X_train, y_train, params, max_iter, alpha,
503             dropout_prob,a)
504
505         # Store the results for the current method
506         results[method] = {'accuracies':
507             accuracies, 'losses': losses}
508
509 # Plotting the results
510 plt.figure(figsize=(12, 6))
511
512 # Plot accuracy for each method
513 plt.subplot(1, 2, 1)
514 for method, data in results.items():
515     plt.plot(range(1, max_iter + 1), data['
516         accuracies'], label=method)
517 plt.xlabel('Epoch')
518 plt.ylabel('Accuracy')
519 plt.title('Accuracy Comparison')
520 plt.grid()
521 plt.legend()
522
523 # Plot loss for each method
524 plt.subplot(1, 2, 2)
525 for method, data in results.items():
526     plt.plot(range(1, max_iter + 1), data['
527         losses'], label=method)
528 plt.xlabel('Epoch')
529 plt.ylabel('Loss')
530 plt.title('Loss Comparison')
531 plt.grid()
532 plt.legend()
533
534 plt.tight_layout()
535 plt.show()
536
537 dropout_prob = [num * 0.1 for num in range(0, int
538     (1 / 0.1) + 1)]
539 for alpha in dropout_prob:
540     print(alpha)
541
542 # Dropout during training
543 # set network and optimizer parameters
544 layers_dims = [784,224, 10]
545 max_iter = 50
546 alpha = 0.1
547
548 dropout_prob =[0,0.1,0.2,0.5,0.8]
549 accuracies =[]
550 losses =[]
551
552 for d in dropout_prob:
553     params = init_params(layers_dims,'xavier')
554     # train the network
555     trained_params,train_acc,train_loss = train(
556         X_train, y_train,params, max_iter, alpha,d
557         ,"relu"
558     )
559     accuracy,loss = test(X_test,y_test,
560         trained_params,a)
561     plot_accuracy_and_loss(train_acc,train_loss,
562         max_iter)
563     accuracies.append(accuracy)
564     losses.append(loss)
565
566 # Dropout during training
567 # set network and optimizer parameters
568 layers_dims = [784,512, 10]
569 max_iter = 150
570 alpha = 0.1
571
572 dropout_prob =[0,0.1,0.2]
573 accuracies =[]
574 losses =[]
575
576 for d in dropout_prob:
577     params = init_params(layers_dims,'xavier')
578     # train the network
579     trained_params,train_acc,train_loss = train(
580         X_train, y_train,params, max_iter, alpha,d
581         ,"relu"
582     )
583     accuracy,loss = test(X_test,y_test,
584         trained_params,"relu")
585     plot_accuracy_and_loss(train_acc,train_loss,
586         max_iter)
587     accuracies.append(accuracy)
588     losses.append(loss)
589
590 print(accuracies)
591
592 print(losses)
593
594 print(dropout_prob)
595
596 plt.figure(figsize=(10, 6))

```

```

573
574 # Bar plot for accuracies
575 plt.subplot(1, 2, 1)
576 plt.bar(['0', '0.1', '0.2'], accuracies, color='blue')
577 plt.title('Accuracy vs. Dropout Probability')
578 plt.xlabel('Dropout Probability')
579 plt.ylabel('Accuracy')
580
581 # Bar plot for losses
582 plt.subplot(1, 2, 2)
583 plt.bar(['0', '0.1', '0.2'], losses, color='red')
584 plt.title('Loss vs. Dropout Probability')
585 plt.xlabel('Dropout Probability')
586 plt.ylabel('Loss')
587
588 plt.show()
589
590 # Define hyperparameter grid
591 layer_options = [
592     [784, 256, 10],
593     [784, 256, 256, 10],
594     [784, 256, 128, 64, 10],
595 ]
596 alpha=0.1
597 init_methods = ['default', 'xavier', 'he']
598 max_iter=50
599 best_accuracy = 0
600 best_params = None
601
602 for layers_dims in layer_options:
603     for initialization_method in init_methods:
604         # Initialize Neural network
605         params = init_params(layers_dims,
606                               initialization_method)
607
608         # Train the neural network
609         updated_param, accuracy, loss = train(
610             X_train, y_train, params, max_iter, alpha)
611
612         # Check if this set of hyperparameters is
613         # the best so far
614         if accuracy > best_accuracy:
615             best_accuracy = accuracy
616             best_params = {
617                 'layers_dims': layers_dims,
618                 'alpha': alpha,
619                 'initialization_method':
620                     initialization_method
621             }
622
623 print("Best hyperparameters:", best_params)
624 print("Best validation accuracy:", best_accuracy)
625
626 # Test best model on test data
627 test_accuracy = evaluate_model(model, X_test,
628                                y_test)
629 print("Test accuracy:", test_accuracy)

```

...