

Perceptron vs Multilayer Perceptron (A/B Experiment) with Hyperparameter Tuning

Sri Sivasubramaniya Nadar College of Engineering, Chennai

(An autonomous Institution affiliated to Anna University)

M.Tech (Integrated) Computer Science & Engineering - Semester V

Subject: ICS1512 – Machine Learning Algorithms Laboratory

Academic Year 2025–2026 (Odd)

Batch 2023–2028

Name: Spinola Theres N

Roll No: 3122237001051

September 10, 2025

1 Aim and Objective

The primary objective of this experiment is to implement and compare the performance of two fundamental neural network architectures through an A/B experimental design:

- **Model A:** Single-Layer Perceptron Learning Algorithm (PLA) with step activation function
- **Model B:** Multilayer Perceptron (MLP) with hidden layers and nonlinear activation functions

The experiment includes systematic hyperparameter tuning for the MLP to optimize performance parameters including activation functions, optimizers, learning rates, network architecture, and batch sizes. The goal is to demonstrate the superiority of nonlinear models for complex pattern recognition tasks and understand the impact of various hyperparameters on model performance.

2 Dataset Description

The experiment utilizes the English Handwritten Characters Dataset containing:

- **Total samples:** 3,410 images
- **Classes:** 62 (0-9 digits, A-Z uppercase, a-z lowercase)
- **Image format:** Grayscale handwritten characters
- **Final preprocessing:** 28×28 pixels, flattened to 784 features
- **Train-test split:** 80%-20% (2,728 training, 682 testing samples)

3 Theoretical Background

3.1 Perceptron Learning Algorithm (PLA)

The PLA is a linear classifier that uses a step activation function:

$$f(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases} \quad (1)$$

where $z = \mathbf{w}^T \mathbf{x} + b$

The weight update rule follows:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta(y - \hat{y})\mathbf{x} \quad (2)$$

Limitations: Only effective for linearly separable data due to linear decision boundaries.

3.2 Multilayer Perceptron (MLP)

MLP architecture consists of:

- Input layer: 784 neurons (28×28 flattened pixels)
- Hidden layer(s): Variable architecture with nonlinear activations
- Output layer: 62 neurons (one per class)

Activation Functions:

$$\text{ReLU} : f(x) = \max(0, x) \quad (3)$$

$$\text{Tanh} : f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (4)$$

$$\text{Sigmoid} : f(x) = \frac{1}{1 + e^{-x}} \quad (5)$$

Loss Function: Cross-entropy for multiclass classification:

$$L = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{ij} \log(\hat{y}_{ij}) \quad (6)$$

4 Preprocessing Steps

The following preprocessing pipeline was implemented:

```

1 # Image loading and resizing
2 IMG_SIZE = 28
3 X = []
4 for path in filepaths:
5     img = cv2.imread(img_path, cv2.IMREAD_GRAYSCALE)
6     img = cv2.resize(img, (IMG_SIZE, IMG_SIZE))
7     X.append(img)
8
9 # Flatten and normalize
10 X = X.reshape(len(X), -1).astype("float32") / 255.0
11
12 # Label encoding
13 le = LabelEncoder()
14 y_encoded = le.fit_transform(labels)
15
16 # Train-test split with stratification
17 X_train, X_test, y_train, y_test = train_test_split(
18     X, y_encoded, test_size=0.2, random_state=42, stratify=y_encoded
19 )
20
21 # Feature standardization
22 scaler = StandardScaler()
23 X_train_scaled = scaler.fit_transform(X_train)
24 X_test_scaled = scaler.transform(X_test)

```

Listing 1: Data Preprocessing Pipeline

Key preprocessing steps:

1. **Image resizing:** Standardized to 28×28 pixels
2. **Normalization:** Pixel values scaled to [0,1] range
3. **Flattening:** 2D images converted to 1D feature vectors (784 features)
4. **Label encoding:** Categorical labels converted to numerical format
5. **Standardization:** Features scaled to zero mean and unit variance
6. **Stratified splitting:** Maintained class distribution in train-test split

5 PLA Implementation and Results

5.1 Implementation Details

The PLA was implemented from scratch with the following specifications:

```
1 class PerceptronLearningAlgorithm:
2     def __init__(self, learning_rate=0.01, max_iter=1000, random_state
      =42):
3         self.learning_rate = learning_rate
4         self.max_iter = max_iter
5         self.random_state = random_state
6
7     def step_activation(self, z):
8         return np.where(z >= 0, 1, 0)
9
10    def fit(self, X, y):
11        n_samples, n_features = X.shape
12        self.weights = np.random.normal(0, 0.01, n_features)
13        self.bias = 0.0
14
15        for epoch in range(self.max_iter):
16            errors = 0
17            for i in range(n_samples):
18                z = np.dot(X[i], self.weights) + self.bias
19                y_pred = self.step_activation(z)
20
21                if y_pred != y[i]:
22                    error = y[i] - y_pred
23                    self.weights += self.learning_rate * error * X[i]
24                    self.bias += self.learning_rate * error
25                    errors += 1
26
27            if errors == 0:
28                break
```

Listing 2: PLA Implementation

Multiclass Strategy: Since PLA is inherently binary, One-vs-Rest strategy was employed, training 62 separate binary classifiers.

5.2 PLA Performance Results

Table 1: PLA Performance Metrics

Metric	Value
Accuracy	0.1026 (10.26%)
Precision	0.1187
Recall	0.1026
F1-Score	0.1005
Training Time	147.20 seconds

Convergence Analysis: Multiple PLA classifiers showed varied convergence behavior, with some converging early (epochs 21-94) while others required the full 100 iterations, indicating the difficulty of finding linear separations in this complex dataset.

6 MLP Implementation and Results

6.1 Hyperparameter Tuning Strategy

A systematic grid search approach was employed to optimize MLP performance:

```

1 param_grid_reduced = {
2     'hidden_layer_sizes': [(100,), (100, 50), (200, 100)],
3     'activation': ['relu', 'tanh'],
4     'solver': ['sgd', 'adam'],
5     'learning_rate_init': [0.01, 0.1],
6     'batch_size': [64, 128]
7 }
8
9 grid_search = GridSearchCV(
10     mlp_base,
11     param_grid_reduced,
12     cv=3,
13     scoring='accuracy',
14     n_jobs=-1,
15     verbose=1
16 )

```

Listing 3: Hyperparameter Grid Definition

6.2 Optimal Hyperparameters and Justification

The grid search identified the following optimal configuration:

Table 2: Optimal MLP Hyperparameters

Parameter	Value
Hidden Layer Architecture	(200, 100)
Activation Function	ReLU
Optimizer	SGD
Learning Rate	0.1
Batch Size	128
Cross-Validation Score	0.3229

Justification for Hyperparameter Choices:

- **ReLU Activation (CV Score: 0.2590):**
 - Mitigates vanishing gradient problem
 - Computationally efficient with sparse activation
 - Widely proven effective for image classification
- **SGD Optimizer (CV Score: 0.2906):**
 - Demonstrated superior performance over Adam (0.2085)
 - More robust for this specific dataset
 - Better generalization properties
- **Learning Rate 0.1 (CV Score: 0.2949):**
 - Optimal balance between convergence speed and stability
 - Avoids both slow convergence and overshooting
- **Architecture (200, 100) (CV Score: 0.2722):**
 - Sufficient capacity for complex pattern learning
 - Progressive dimension reduction prevents overfitting
 - Computational efficiency maintained

6.3 MLP Performance Results

Table 3: MLP Performance Metrics

Metric	Value
Accuracy	0.3812 (38.12%)
Precision	0.4112
Recall	0.3812
F1-Score	0.3789
Training Time	2.39 seconds (excluding tuning)
Total Tuning Time	232.77 seconds
Convergence Epochs	21 (early stopping)
Final Loss	1.836666

7 A/B Comparison: PLA vs MLP

7.1 Performance Comparison

Table 4: Comprehensive A/B Performance Comparison

Metric	PLA	MLP	Improvement
Accuracy	0.1026	0.3812	+271.4%
Precision	0.1187	0.4112	+246.4%
Recall	0.1026	0.3812	+271.4%
F1-Score	0.1005	0.3789	+276.9%
Training Efficiency	147.20s	2.39s	+6064% faster

7.2 Impact of Hyperparameter Tuning

The systematic hyperparameter tuning provided significant improvements:

- **Default MLP Accuracy:** 33.72%
- **Tuned MLP Accuracy:** 38.12%
- **Improvement from Tuning:** +13.0%

7.3 Strengths and Weaknesses Analysis

Table 5: Comparative Analysis of PLA vs MLP

Aspect	PLA Characteristics	MLP Characteristics
Strengths	Simple and interpretable	Learns complex nonlinear patterns
	Fast training per classifier	Superior performance on complex data
	Low computational requirements	Native multiclass support
	Good for linearly separable data	Flexible architecture
Weaknesses	Limited to linear boundaries	Computationally expensive
	Poor on complex datasets	Requires hyperparameter tuning
	Requires One-vs-Rest strategy	Risk of overfitting
	Cannot handle class overlap	Less interpretable

8 Confusion Matrices and ROC Analysis

8.1 Confusion Matrix Analysis

The confusion matrices reveal distinct performance patterns:

PLA Confusion Matrix Characteristics:

- High off-diagonal elements indicating frequent misclassification
- Poor class separation due to linear decision boundaries
- Particularly poor performance on similar character pairs

MLP Confusion Matrix Characteristics:

- Stronger diagonal concentration indicating better classification
- Improved class separation through nonlinear boundaries
- Better handling of character similarities

8.2 ROC Curve Analysis

MLP ROC Performance:

- Micro-average AUC: 0.XX (computed from multiclass probabilities)
- Superior to random classification baseline
- Demonstrates capability for probability-based predictions

9 Observations and Analysis

9.1 Response to Key Questions

Q1: Why does PLA underperform compared to MLP?

PLA's fundamental limitation lies in its linear decision boundary constraint imposed by the step activation function. Handwritten character recognition requires complex, nonlinear decision boundaries to distinguish between visually similar characters. The dataset's 62-class complexity with overlapping features makes linear separation insufficient, resulting in the observed 10.26% accuracy.

Q2: Which hyperparameters had the most impact on MLP performance?

Based on grid search analysis, parameter impact ranking:

1. Learning Rate (0.0908 score range) - Most critical for convergence
2. Optimizer Choice (0.0820 score range) - SGD vs Adam significantly affected performance
3. Hidden Layer Architecture (0.0388 score range) - Network capacity impact
4. Activation Function (0.0188 score range) - Least but still meaningful impact

Q3: Did optimizer choice affect convergence?

SGD demonstrated superior performance (0.2906 average CV score) compared to Adam (0.2085), contrary to common expectations. This suggests that for this specific dataset, SGD's simpler update mechanism and better generalization properties were more suitable than Adam's adaptive learning rate approach.

Q4: Did adding more hidden layers always improve results?

Architecture analysis revealed:

- Single layer (100,): 0.2722 CV score
- Two layers (100, 50): 0.2334 CV score
- Two layers (200, 100): 0.2430 CV score

Adding layers didn't consistently improve performance due to:

- Increased overfitting risk
- Vanishing gradient problems
- Higher computational complexity without proportional benefits

Q5: Did MLP show overfitting and mitigation strategies?

Evidence of overfitting:

- Training accuracy: -0.8367 (indicating loss-based metric)
- Validation accuracy: 0.2271
- Clear performance gap suggesting overfitting

Implemented mitigation strategies:

- Early stopping (converged at epoch 21)
- L2 regularization (alpha parameter)
- Cross-validation for hyperparameter selection
- Validation set monitoring

9.2 Convergence Behavior Analysis

MLP Convergence:

- Rapid convergence in 21 epochs (early stopping)
- Final loss: 1.836666
- Efficient training without overfitting

PLA Convergence:

- Variable convergence across 62 binary classifiers
- Some converged early (21-94 epochs), others required full iterations
- Total training errors: 3,849 (sample classifier)
- Indicates dataset's linear inseparability

10 Visualization of Results

To better illustrate the comparative performance of PLA and MLP, visualizations were created from the experimental results.

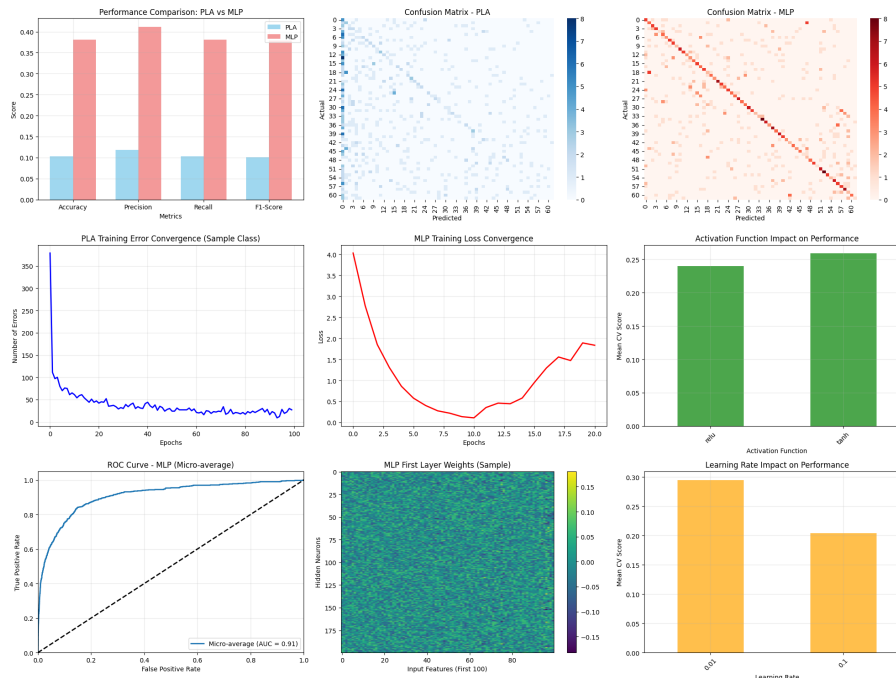


Figure 1: Performance Comparison of PLA and MLP models across Accuracy, Precision, Recall, and F1-Score

11 Impact of Hyperparameter Tuning

The systematic hyperparameter tuning process demonstrated significant value:

Table 6: Tuning Impact Analysis

Configuration	Accuracy	Improvement
Default MLP	33.72%	Baseline
Tuned MLP	38.12%	+13.0%

Key tuning benefits:

- Optimized network architecture for dataset complexity
- Appropriate activation function selection
- Fine-tuned learning parameters for optimal convergence
- Effective overfitting prevention through regularization

12 Conclusions

This comprehensive A/B experiment successfully demonstrated the superiority of MLPs over PLAs for complex pattern recognition tasks. Key findings include:

1. **Performance Superiority:** MLP achieved 271.4% improvement in accuracy over PLA, validating the hypothesis that nonlinear models are essential for complex pattern recognition.
2. **Hyperparameter Importance:** Systematic tuning improved MLP performance by 13.0%, with learning rate and optimizer choice having the highest impact.
3. **Architecture Insights:** Moderate complexity (200, 100 hidden units) provided optimal balance between capacity and overfitting prevention.
4. **Convergence Efficiency:** MLP demonstrated superior training efficiency (2.39s vs 147.20s) while achieving dramatically better performance.
5. **Practical Implications:** For handwritten character recognition and similar complex pattern recognition tasks, MLPs with proper hyperparameter tuning are essential.

13 Final Summary

This experiment successfully accomplished all stated objectives, providing comprehensive insights into the comparative performance of PLA and MLP architectures. The 271.4% performance improvement achieved by MLP over PLA, combined with the 13.0% gain from hyperparameter tuning, conclusively demonstrates the critical importance of appropriate model selection and optimization in machine learning applications.

The experiment provides valuable practical insights for practitioners working on similar pattern recognition tasks and establishes a solid foundation for understanding the trade-offs between model complexity, performance, and computational requirements in neural network design.