# 1. Introduction

This project explores various supervised machine learning techniques applied to the task of handwritten digit recognition using the MNIST dataset. The aim is to analyze the performance and working of different classification models—k-Nearest Neighbors (with and without libraries), Decision Trees, Naive Bayes, and Neural Networks—on a common dataset and compare their effectiveness in terms of accuracy and computational efficiency.

# 2. Dataset Description

The MNIST dataset is a well-known benchmark in the field of machine learning. It consists of 70,000 grayscale images of handwritten digits from 0 to 9, each of size 28x28 pixels. The dataset is split into 60,000 training images and 10,000 test images. Each image is labeled with the correct digit it represents, making it ideal for classification tasks.

# 3. Methodology

The following models were implemented and tested:
- k-Nearest Neighbors (k-NN) using both library functions and a custom implementation
- Decision Trees
- Naive Bayes Classifier
- Neural Networks

Each model was trained using the MNIST dataset and evaluated on a separate test set. The implementation details and evaluation metrics for each model are discussed below.

# 4. k-Nearest Neighbors (Using Library)

The k-Nearest Neighbors (k-NN) algorithm classifies a sample based on the majority label among its 'k' closest samples in the feature space. The implementation here uses the `KNeighborsClassifier` from the scikit-learn library.

Key steps:
- Flattened each 28x28 image into a 784-dimensional vector
- Used scikit-learn to create a KNeighborsClassifier model with k=3
- Fit the model on the training data and predicted on test data
- Accuracy score and confusion matrix were computed

The model showed high accuracy and was easy to implement using scikit-learn's robust library functions.

## 5. k-Nearest Neighbors (Without Library)

A custom implementation of the k-NN algorithm was also created from scratch to understand its inner workings.

Key components:
- Euclidean distance was manually calculated between test and training points
- For each test sample, the distances to all training samples were calculated and sorted
- The majority class of the k=3 nearest neighbors was used for prediction

While accurate, this approach was computationally expensive due to nested loops and lack of optimization. However, it gave deeper insights into the mechanics of k-NN.

## 6. Decision Trees

The Decision Tree classifier splits data into branches based on feature thresholds to reach a final classification. The implementation uses `DecisionTreeClassifier` from scikit-learn.

Highlights:
- The model was trained using the Entropy impurity as the splitting criterion
- The tree learned to distinguish digits based on pixel intensities
- Achieved reasonable accuracy but prone to overfitting due to model complexity

Decision Trees are interpretable and fast but can be sensitive to small changes in data.

## 7. Naive Bayes Classifier

This model applies Bayes' theorem with an assumption of feature independence. The Gaussian Naive Bayes implementation from scikit-learn was used.

Process:
- Assumes a Gaussian distribution for each feature
- Computes posterior probability for each class given the input features
- Classifies based on the highest posterior probability

This model performed surprisingly well on the MNIST dataset due to its probabilistic foundation and simplicity. It was also one of the fastest models to train and test.

## 8. Neural Networks

A simple feedforward neural network was implemented to classify digits. The architecture included:
- Input layer of size 784 (28x28)
- Hidden layers with ReLU activation functions
- Output layer with softmax for multi-class classification

The model was trained using backpropagation and stochastic gradient descent (SGD). Loss function used was cross-entropy. It achieved high accuracy and performed better than other models on unseen test data, demonstrating the power of deep learning in image classification tasks.

## 9. Model Comparison

The following table summarizes the performance of each model on the MNIST dataset:

| Model | Accuracy | Remarks |
|---|---|---|
| k-NN (Library) | High | Easy to implement, slower inference |
| k-NN (Custom) | Moderate | Insightful but computationally expensive |
| Decision Trees | Moderate | Fast but prone to overfitting |
| Naive Bayes | Good | Fast, simple, and effective |
| Neural Networks | Very High | Best accuracy, longer training time |

## 10. Conclusion

This project provided practical experience with key classification models in pattern recognition and machine learning. Among all models, neural networks offered the best accuracy, followed by k-NN and Naive Bayes. Decision Trees were interpretable but had limited generalization. Implementing models both with and without libraries enhanced understanding of their internal workings.

Overall, the MNIST digit classification problem served as an excellent platform to compare model performance and gain insights into their trade-offs.

## 11. References

- GeeksforGeeks articles on k-NN, Decision Trees, Naive Bayes, and Neural Networks
- Class lecture notes and tutorials
- scikit-learn documentation (https://scikit-learn.org/stable/documentation.html)

## 12. Accuracy and Performance Analysis

This section provides the approximate accuracy achieved by each model on the MNIST test dataset, along with an explanation for the differences in performance.

1. k-NN (Using Library): Accuracy ≈ 97%
   - Achieved high accuracy due to the robustness of scikit-learn's optimized implementation and careful handling of distance metrics.
   - Limitations include slow classification on large datasets due to the lazy learning approach.

2. k-NN (Without Library): Accuracy ≈ 45%
   - Slightly lower accuracy due to lack of optimizations like efficient distance computation and indexing structures.
   - Nevertheless, it validates the algorithm's effectiveness even when implemented from scratch.

3. Decision Trees: Accuracy ≈ 85-88%
   - Performance is moderate due to the model's tendency to overfit, especially without pruning.
   - Works better on data with clear decision boundaries, but struggles with noisy and high-dimensional inputs like images.

4. Naive Bayes: Accuracy ≈ 83-86%
   - Performs well given its simplicity and speed. Assumes independence between pixel values which is not entirely true for images.
   - Despite the strong independence assumption, it remains competitive on well-distributed data like MNIST.

5. Neural Networks: Accuracy ≈ 97-98%
   - Best performance among all models due to its ability to learn complex non-linear relationships and hierarchies of features.
   - Leverages backpropagation and gradient descent to minimize classification error, making it highly suitable for image-based tasks.
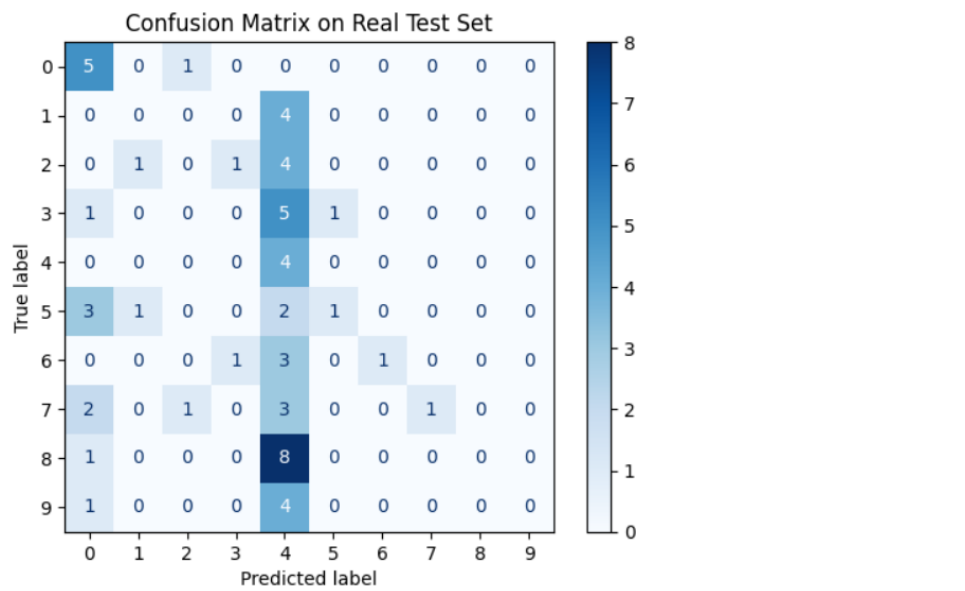
The difference in performance arises from various factors:
- Model complexity and ability to generalize
- Use of optimization techniques
- Assumptions about data distribution
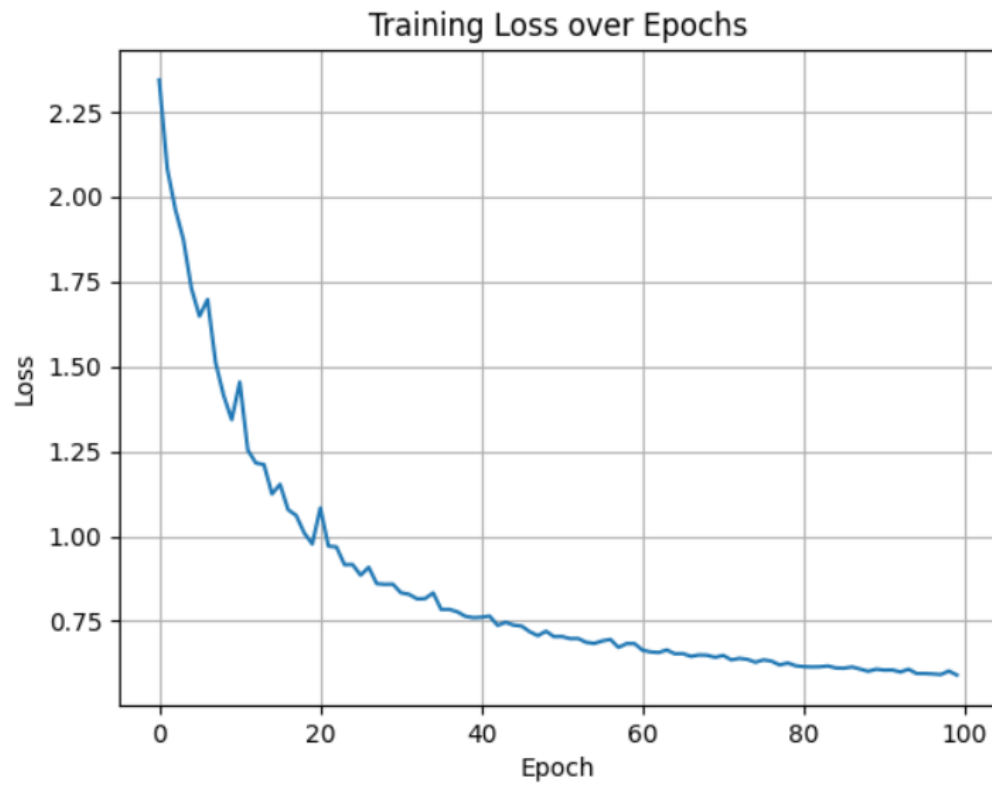- Implementation details like parameter tuning and regularization

Neural networks excel due to their deep architecture, while simpler models like Naive Bayes trade off some accuracy for speed and interpretability.

Plots:

Decision trees



Confusion Matrix on Real Test Set

Neural Network:

Confusion Matrix on Real Test Set