

Handwritten-Digit-Recognition using Classical ML Models

This project implements a complete machine learning pipeline for classifying handwritten digits (0-9) from the MNIST dataset using classical ML algorithms. The implementation includes data preprocessing, model training, evaluation, and comprehensive analysis.

▼ Install all required tools and Dependencies.

```
1 !pip install numpy pandas matplotlib seaborn scikit-learn

Requirement already satisfied: numpy in /usr/local/lib/python3.12/dist-packages (2.0.2)
Requirement already satisfied: pandas in /usr/local/lib/python3.12/dist-packages (2.2.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.12/dist-packages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.12/dist-packages (0.13.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.12/dist-packages (1.6.1)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.12/dist-packages (from pandas) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-packages (from pandas) (2025.3)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.3.3)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (4.61.1)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (1.4.9)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (25.0)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (11.3.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.12/dist-packages (from matplotlib) (3.2.5)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (1.16.3)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (1.5.3)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.12/dist-packages (from scikit-learn) (3.6.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
```

▼ Downloading the dataset.

To download the dataset for training, visit the kaggle site and inside the Datsets section, search for **MNIST in CSV** in search bar. [dataset link](#) Click the download button to download the zip file and unzip it after the download is complete.

You can also do this by integrating below code in your project code.

```
1 import kagglehub
2
3 # Download latest version
4 path = kagglehub.dataset_download("oddrationale/mnist-in-csv")
5
6 print("Path to dataset files:", path)

Using Colab cache for faster access to the 'mnist-in-csv' dataset.
Path to dataset files: /kaggle/input/mnist-in-csv
```

▼ Import required libraries

```
1 # Import necessary libraries
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.decomposition import PCA
9 from sklearn.neighbors import KNeighborsClassifier #for KNN
10 from sklearn.svm import SVC #for SVM
11 from sklearn.tree import DecisionTreeClassifier #for decisions tree
12 from sklearn.metrics import accuracy_score, confusion_matrix
13 import warnings
```

▼ Task 1: Data Loading and Exploration



1. Load the CSV file using Pandas.
2. Explore dataset statistics:
 - Total number of samples
 - Class distribution
 - Display 5-10 sample images with their labels (reshape pixels to 28x28)
3. Check for missing values

```

1
2 warnings.filterwarnings('ignore')
3
4 # Set style for better visualizations
5 plt.style.use('seaborn-v0_8-darkgrid')
6 sns.set_palette("husl")
7
8 # Load the dataset CSV file using Pandas
9 print("/loading MNIST dataset...")
10 df = pd.read_csv("/kaggle/input/mnist-in-csv/mnist_train.csv") # Update with your file path
11
12 # Display basic information
13 print(f"Dataset shape: {df.shape}")
14 print(f"Columns: {df.columns[:10].tolist()}...") # First 10 columns
15
16 # Check for missing values
17 print(f"\nMissing values:\n{df.isnull().sum().sum()} total missing values")
18
19 # Explore class distribution
20 print("\nClass Distribution:")
21 class_dist = df['label'].value_counts().sort_index()
22 print(class_dist)
23
24 # Visualize class distribution
25 plt.figure(figsize=(10, 5))
26 plt.bar(class_dist.index, class_dist.values)
27 plt.xlabel('Digit')
28 plt.ylabel('Count')
29 plt.title('Class Distribution of Digits')
30 plt.xticks(range(10))
31 plt.grid(True, alpha=0.3)
32 plt.savefig('class_distribution.png', dpi=100, bbox_inches='tight')
33 plt.show()
34

```

```
>Loading MNIST dataset...
Dataset shape: (60000, 785)
Columns: ['label', '1x1', '1x2', '1x3', '1x4', '1x5', '1x6', '1x7', '1x8', '1x9']...
```

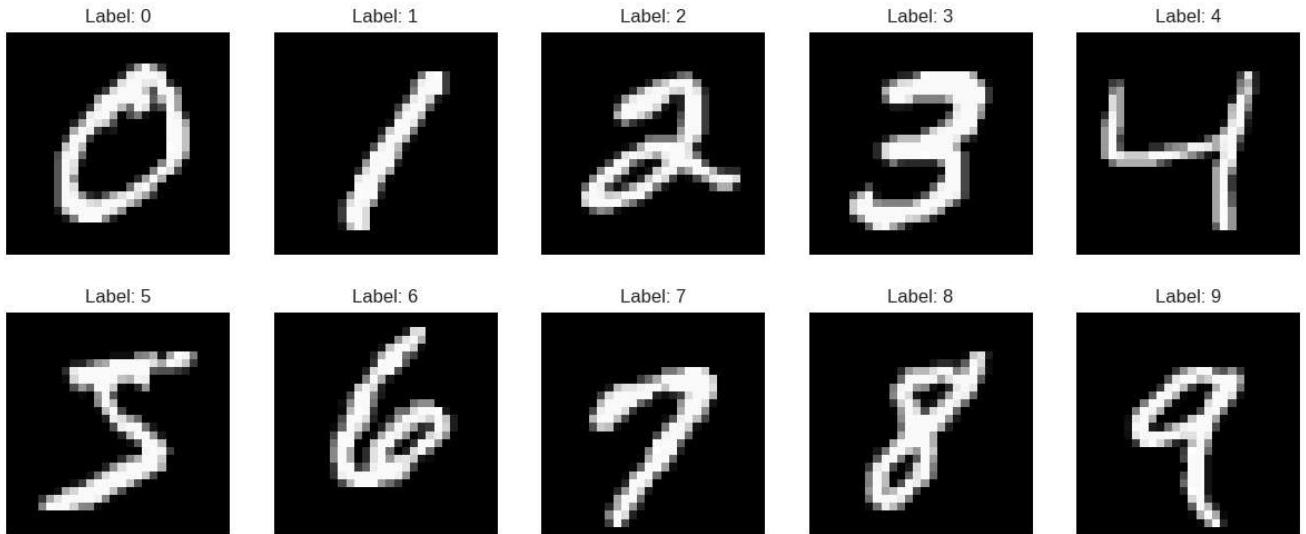
```
Missing values:
0 total missing values
```

```
Class Distribution:
label
0    5923
1    6742
2    5958
3    6131
```

```
1 # Display sample images
2 print("\n🖼 Sample Images from Dataset:")
3 fig, axes = plt.subplots(2, 5, figsize=(15, 6))
4
5 for idx, ax in enumerate(axes.flat):
6     if idx < 10:
7         # Get a random sample for each digit
8         sample = df[df['label'] == idx].iloc[0]
9         label = sample['label']
10        pixels = sample.drop('label').values.reshape(28, 28)
11
12        ax.imshow(pixels, cmap='gray')
13        ax.set_title(f'Label: {label}')
14        ax.axis('off')
15
16 plt.suptitle('Sample Handwritten Digits (0-9)', fontsize=16)
17 plt.savefig('sample_digits.png', dpi=100, bbox_inches='tight')
18 plt.show()
```

3000 🖼 Sample Images from Dataset:

Sample Handwritten Digits (0-9)



▼ LOAD THE DATA

```
1 train = pd.read_csv("/kaggle/input/mnist-in-csv/mnist_train.csv")
2 test = pd.read_csv("/kaggle/input/mnist-in-csv/mnist_test.csv")
```

▼ Data Exploration

```
1 print(train.shape)
2 train.head()
```

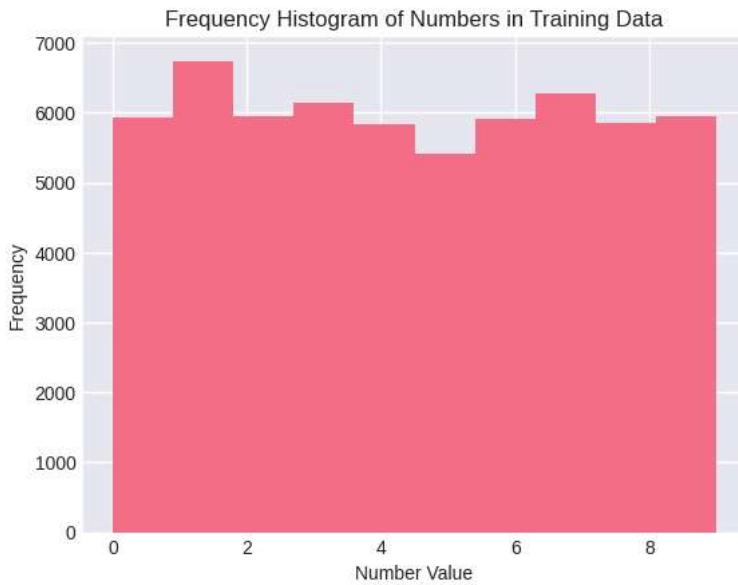
(60000, 785)																										
	label	1x1	1x2	1x3	1x4	1x5	1x6	1x7	1x8	1x9	...	28x19	28x20	28x21	28x22	28x23	28x24	28x25	28x26	28x27	28x28					
0	5	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	4	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	9	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

5 rows × 785 columns

1	print(test.shape)																									
2	test.head()																									
(10000, 785)																										
	label	1x1	1x2	1x3	1x4	1x5	1x6	1x7	1x8	1x9	...	28x19	28x20	28x21	28x22	28x23	28x24	28x25	28x26	28x27	28x28					
0	7	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	2	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	4	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

5 rows × 785 columns

```
1 plt.hist(train["label"])
2 plt.title("Frequency Histogram of Numbers in Training Data")
3 plt.xlabel("Number Value")
4 plt.ylabel("Frequency")
5 plt.show()
```



Task 2: Data Preprocessing

- Normalize pixel values to the range 0–1.
- Split the dataset into training and testing sets (e.g., 80/20 split).
- Optional: Use PCA to reduce dimensionality (not mandatory, but recommended for improving SVM performance).

```
1 print("✿ Preprocessing Data...")
2
3 # Separate features and labels
4 X = df.drop('label', axis=1)
5 y = df['label']
```

```

6
7 # Normalize pixel values to [0, 1]
8 X_normalized = X / 255.0
9 print(f"Normalized pixel range: [{X_normalized.min():.3f}, {X_normalized.max():.3f}]")
10
11 # Split the dataset
12 X_train, X_test, y_train, y_test = train_test_split(
13     X_normalized, y, test_size=0.2, random_state=42, stratify=y
14 )
15 print(f"\nTrain set: {X_train.shape}")
16 print(f"Test set: {X_test.shape}")
17
18 # Optional: Apply PCA for dimensionality reduction
19 print("\n🔍 Applying PCA for dimensionality reduction...")
20 pca = PCA(n_components=0.95, random_state=42) # Keep 95% variance
21 X_train_pca = pca.fit_transform(X_train)
22 X_test_pca = pca.transform(X_test)
23
24 print(f"Original features: {X_train.shape[1]}")
25 print(f"PCA features (95% variance): {X_train_pca.shape[1]}")
26 print(f"Variance explained: {pca.explained_variance_ratio_.sum():.3%}")
27
28 # Visualize PCA explained variance
29 plt.figure(figsize=(10, 4))
30 plt.subplot(1, 2, 1)
31 plt.plot(np.cumsum(pca.explained_variance_ratio_))
32 plt.xlabel('Number of Components')
33 plt.ylabel('Cumulative Explained Variance')
34 plt.title('PCA: Cumulative Explained Variance')
35 plt.grid(True, alpha=0.3)
36
37 plt.subplot(1, 2, 2)
38 plt.bar(range(1, 11), pca.explained_variance_ratio_[:10])
39 plt.xlabel('Principal Component')
40 plt.ylabel('Explained Variance Ratio')
41 plt.title('Top 10 Principal Components')
42 plt.grid(True, alpha=0.3)
43
44 plt.tight_layout()
45 plt.savefig('pca_analysis.png', dpi=100, bbox_inches='tight')
46 plt.show()

```

⚙️ Preprocessing Data...

Normalized pixel range: [0.000, 1.000]

Train set: (48000, 784)

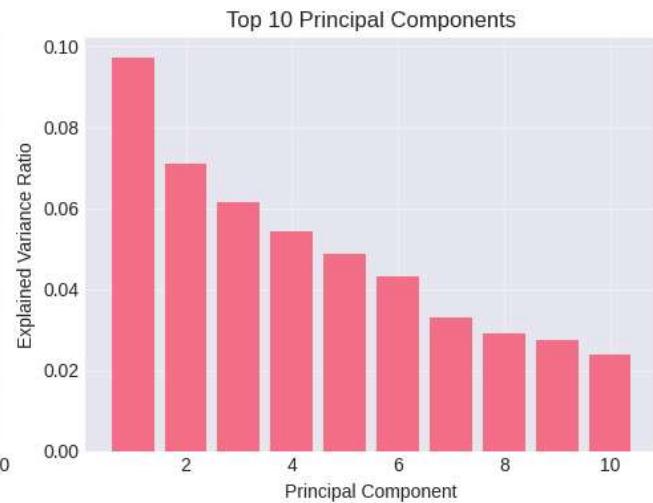
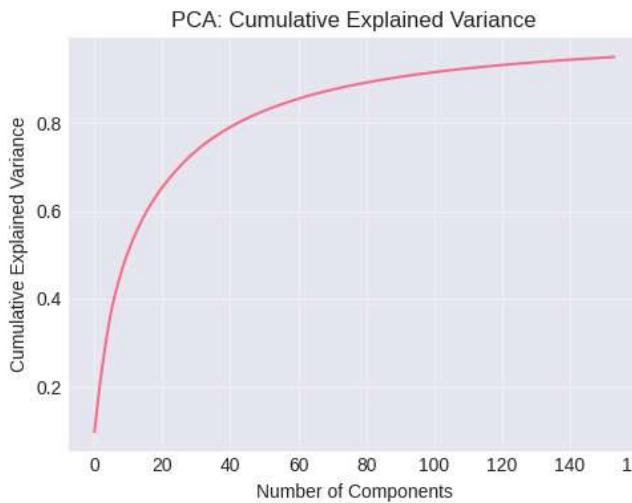
Test set: (12000, 784)

🔍 Applying PCA for dimensionality reduction...

Original features: 784

PCA features (95% variance): 154

Variance explained: 95.025%



❖ Task 3: Model Implementation

Train and evaluate three classical ML models from scikit-learn:

1. K-Nearest Neighbors (KNN)

- Tune k (number of neighbors)
- Evaluate accuracy on the test set

2. Support Vector Machine (SVM)

- Use RBF or linear kernel
- Tune C and gamma parameters

3. Decision Tree

- Tune max_depth and min_samples_split Important: Do not use pre-trained models or neural networks. Only classical ML algorithms are allowed.

Using scikit-learn

```
1 print("🤖 Training Machine Learning Models...")
2
3 # Dictionary to store model results
4 results = {}
5
6 # 3.1 K-Nearest Neighbors (KNN)
7 print("\n👉 Training KNN...")
8 knn = KNeighborsClassifier(n_neighbors=5) #Tune k : number of neighbour
9 knn.fit(X_train, y_train)
10 y_pred_knn = knn.predict(X_test)
11 acc_knn = accuracy_score(y_test, y_pred_knn)
12 results['KNN'] = {'model': knn, 'accuracy': acc_knn, 'predictions': y_pred_knn}
13 print(f"KNN Accuracy: {acc_knn:.4f}")
14
15 # 3.2 Support Vector Machine (SVM)
16 print("\n👉 Training SVM...")
17 svm = SVC(C=1.0, kernel='rbf', gamma='scale', random_state=42)
18 svm.fit(X_train_pca, y_train) # Using PCA-reduced data for efficiency
19 y_pred_svm = svm.predict(X_test_pca)
20 acc_svm = accuracy_score(y_test, y_pred_svm)
21 results['SVM'] = {'model': svm, 'accuracy': acc_svm, 'predictions': y_pred_svm}
22 print(f"SVM Accuracy: {acc_svm:.4f}")
23
24 # 3.3 Decision Tree
25 print("\n👉 Training Decision Tree...")
26 dtree = DecisionTreeClassifier(max_depth=15, min_samples_split=5, random_state=42)
27 dtree.fit(X_train, y_train)
28 y_pred_dt = dtree.predict(X_test)
29 acc_dt = accuracy_score(y_test, y_pred_dt)
30 results['Decision Tree'] = {'model': dtree, 'accuracy': acc_dt, 'predictions': y_pred_dt}
31 print(f"Decision Tree Accuracy: {acc_dt:.4f}")
```

🤖 Training Machine Learning Models...

👉 Training KNN...

KNN Accuracy: 0.9675

👉 Training SVM...

SVM Accuracy: 0.9814

👉 Training Decision Tree...

Decision Tree Accuracy: 0.8717

From-Scratch KNN Implementation

```
1 print("\n🔧 Implementing KNN from scratch...")
2
3 class KNNFromScratch:
4     def __init__(self, k=5):
5         self.k = k
6         self.X_train = None
7         self.y_train = None
8
9     def fit(self, X, y):
10        self.X_train = X
```

```

11     self.y_train = y
12
13 def predict(self, X_test):
14     predictions = []
15     for i, x in enumerate(X_test.iterrows() if hasattr(X_test, 'iterrows') else enumerate(X_test)):
16         if i % 100 == 0:
17             print(f" Predicting sample {i}/{len(X_test)}...")
18
19         # Calculate distances to all training points
20         if hasattr(self.X_train, 'iloc'):
21             distances = np.sqrt(np.sum((self.X_train.values - x[1].values) ** 2, axis=1))
22         else:
23             distances = np.sqrt(np.sum((self.X_train - x[1]) ** 2, axis=1))
24
25         # Get k nearest neighbors
26         k_indices = np.argsort(distances)[:self.k]
27         k_nearest_labels = self.y_train.iloc[k_indices] if hasattr(self.y_train, 'iloc') else self.y_train[k_indices]
28
29         # Majority vote
30         most_common = np.bincount(k_nearest_labels).argmax()
31         predictions.append(most_common)
32
33     return np.array(predictions)
34
35 # Train from-scratch KNN (using subset for speed)
36 print("Training from-scratch KNN on subset...")
37 subset_size = 1000
38 X_train_subset = X_train.iloc[:subset_size]
39 y_train_subset = y_train.iloc[:subset_size]
40 X_test_subset = X_test.iloc[:200]
41
42 knn_scratch = KNNFromScratch(k=5)
43 knn_scratch.fit(X_train_subset, y_train_subset)
44 y_pred_knn_scratch = knn_scratch.predict(X_test_subset)
45 acc_knn_scratch = accuracy_score(y_test.iloc[:200], y_pred_knn_scratch)
46 results['KNN_Scratch'] = {'model': knn_scratch, 'accuracy': acc_knn_scratch, 'predictions': y_pred_knn_scratch}
47 print(f"From-scratch KNN Accuracy (on subset): {acc_knn_scratch:.4f}")

```

👉 Implementing KNN from scratch...
 Training from-scratch KNN on subset...
 Predicting sample 0/200...
 Predicting sample 100/200...
 From-scratch KNN Accuracy (on subset): 0.8600

▼ Task 4: Model Evaluation

1. Compute accuracy for each model.
2. Generate a confusion matrix for each model.
3. Visualize 5–10 misclassified images and discuss why they might be misclassified.

```

1 print("\n📊 Evaluating Models...")
2
3 # Plot accuracy comparison
4 models = list(results.keys())
5 accuracies = [results[m]['accuracy'] for m in models]
6
7 plt.figure(figsize=(10, 5))
8 bars = plt.bar(models, accuracies, color=['blue', 'green', 'orange', 'red'])
9 plt.axhline(y=max(accuracies), color='gray', linestyle='--', alpha=0.5)
10 plt.ylabel('Accuracy')
11 plt.title('Model Accuracy Comparison')
12 plt.ylim(0.8, 1.0)
13
14 # Add accuracy values on bars
15 for bar, acc in zip(bars, accuracies):
16     plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.005,
17               f'{acc:.3f}', ha='center', va='bottom')
18
19 plt.grid(True, alpha=0.3, axis='y')
20 plt.savefig('model_accuracy.png', dpi=100, bbox_inches='tight')
21 plt.show()
22
23 # Confusion Matrices

```

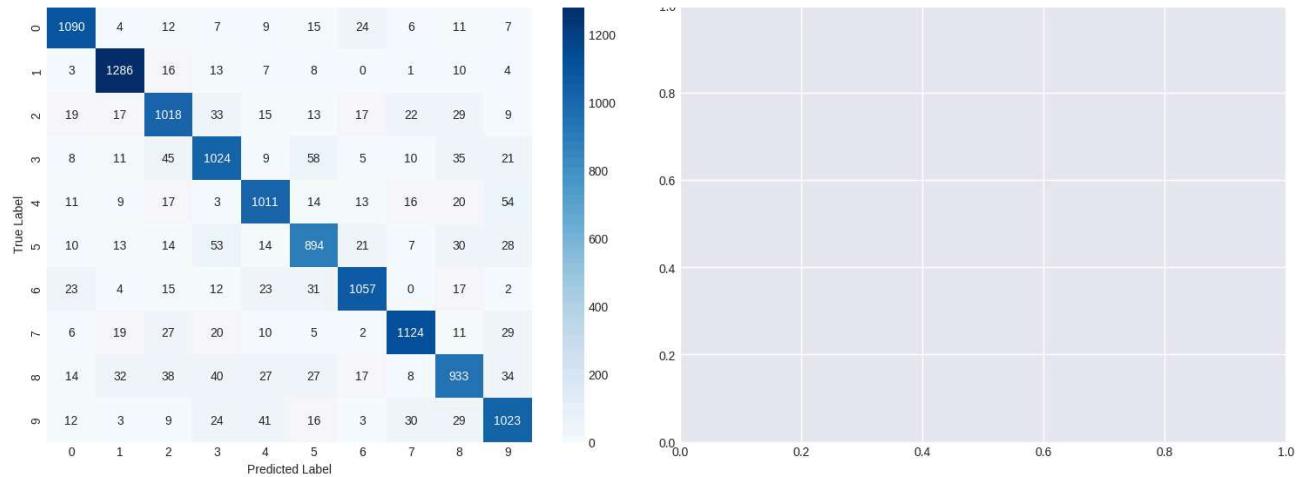
```
24 print("\nGenerating Confusion Matrices...")
25 fig, axes = plt.subplots(2, 2, figsize=(15, 12))
26
27 for idx, (model_name, ax) in enumerate(zip(['KNN', 'SVM', 'Decision Tree'], axes.flatten())):
28     cm = confusion_matrix(y_test, results[model_name]['predictions'])
29     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', ax=ax,
30                  xticklabels=range(10), yticklabels=range(10))
31     ax.set_xlabel('Predicted Label')
32     ax.set_ylabel('True Label')
33     ax.set_title(f'Confusion Matrix - {model_name}\nAccuracy: {results[model_name]["accuracy"]:.3f}')
34
35 plt.tight_layout()
36 plt.savefig('confusion_matrices.png', dpi=100, bbox_inches='tight')
37 plt.show()
38
```


Evaluating Models...

Model Accuracy Comparison

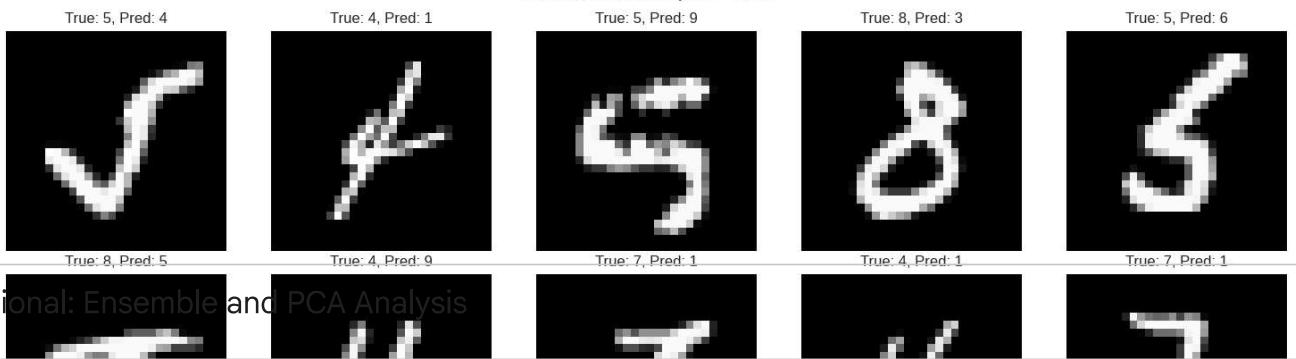


```
1
2 # Visualize misclassifications
3 print("\n🔍 Analyzing Misclassifications...")
4
5 def plot_misclassified(model_name, num_examples=10):
6     predictions = results[model_name]['predictions']
7     misclassified_idx = np.where(predictions != y_test.values)[0]
8
9     if len(misclassified_idx) == 0:
10         print(f"No misclassifications for {model_name}!")
11         return
12
13     # Take first N misclassifications
14     sample_idx = misclassified_idx[:min(num_examples, len(misclassified_idx))]
15
16     fig, axes = plt.subplots(2, 5, figsize=(15, 6))
17     axes = axes.flatten()
18
19     for i, idx in enumerate(sample_idx):
20         if i < 10:
21             pixels = X_test.iloc[idx].values.reshape(28, 28)
22             true_label = y_test.iloc[idx]
23             pred_label = predictions[idx]
24
25             axes[i].imshow(pixels, cmap='gray')
26             axes[i].set_title(f'True: {true_label}, Pred: {pred_label}')
27             axes[i].axis('off')
28
29     plt.suptitle(f'Misclassified Samples - {model_name}', fontsize=16)
30     plt.tight_layout()
31     plt.savefig(f'misclassified_{model_name.lower().replace(" ", "_")}.png',
32                 dpi=100, bbox_inches='tight')
33     plt.show()
34
35     # Print misclassification statistics
36     total_misclassified = len(misclassified_idx)
37     print(f'{model_name}: {total_misclassified} misclassifications ({total_misclassified/len(y_test):.2%})')
38
39
40 # Plot misclassifications for each model
41 for model_name in ['KNN', 'SVM', 'Decision Tree']:
42     plot_misclassified(model_name)
```



⌚ Analyzing Misclassifications...

Misclassified Samples - KNN



▼ Optional: Ensemble and PCA Analysis

```
1 print("\n🌟 Optional Bonus: Ensemble Learning...")
2
3 # Voting Ensemble
4 from sklearn.ensemble import VotingClassifier
5
6 print("Creating Voting Ensemble...")
7 ensemble = VotingClassifier(
8     estimators=[
9         ('knn', KNeighborsClassifier(n_neighbors=5)),
10        ('svm', SVC(C=1.0, kernel='rbf', probability=True, random_state=42)),
11        ('dt', DecisionTreeClassifier(max_depth=15, random_state=42))
12    ],
13    voting='hard'
14 )
15
16 # Train on PCA data for efficiency
17 ensemble.fit(X_train_pca, y_train)
18 y_pred_ensemble = ensemble.predict(X_test_pca)
19 acc_ensemble = accuracy_score(y_test, y_pred_ensemble)
20 print(f"Voting Ensemble Accuracy: {acc_ensemble:.4f}")
21
22 # Compare with individual models
23 print(f"\nImprovement over best single model: {acc_ensemble - max(acc_knn, acc_svm, acc_dt):.4f}")
24
25 # PCA Performance Analysis
26 print("\n📊 Analyzing PCA Impact...")
27 accuracies_pca = []
28 n_components_list = [50, 100, 150, 200, 300]
29
30 for n_comp in n_components_list:
31     pca_temp = PCA(n_components=n_comp, random_state=42)
32     X_train_temp = pca_temp.fit_transform(X_train)
33     X_test_temp = pca_temp.transform(X_test)
34
35     svm_temp = SVC(C=1.0, kernel='rbf', random_state=42)
36     svm_temp.fit(X_train_temp, y_train)
37     acc = svm_temp.score(X_test_temp, y_test)
38     accuracies_pca.append(acc)
39     print(f"PCA {n_comp} components: SVM Accuracy = {acc:.4f}")
40
41 # Plot PCA impact
42 plt.figure(figsize=(10, 5))
43 plt.plot(n_components_list, accuracies_pca, marker='o', linewidth=2)
44 plt.xlabel('Number of PCA Components')
45 plt.ylabel('SVM Accuracy')
46 plt.title('Impact of PCA Dimensionality Reduction on SVM Performance')
47 plt.grid(True, alpha=0.3)
48 plt.savefig('pca_impact.png', dpi=100, bbox_inches='tight')
49 plt.show()
```

💡 Optional Bonus: Ensemble Learning...

Creating Voting Ensemble...

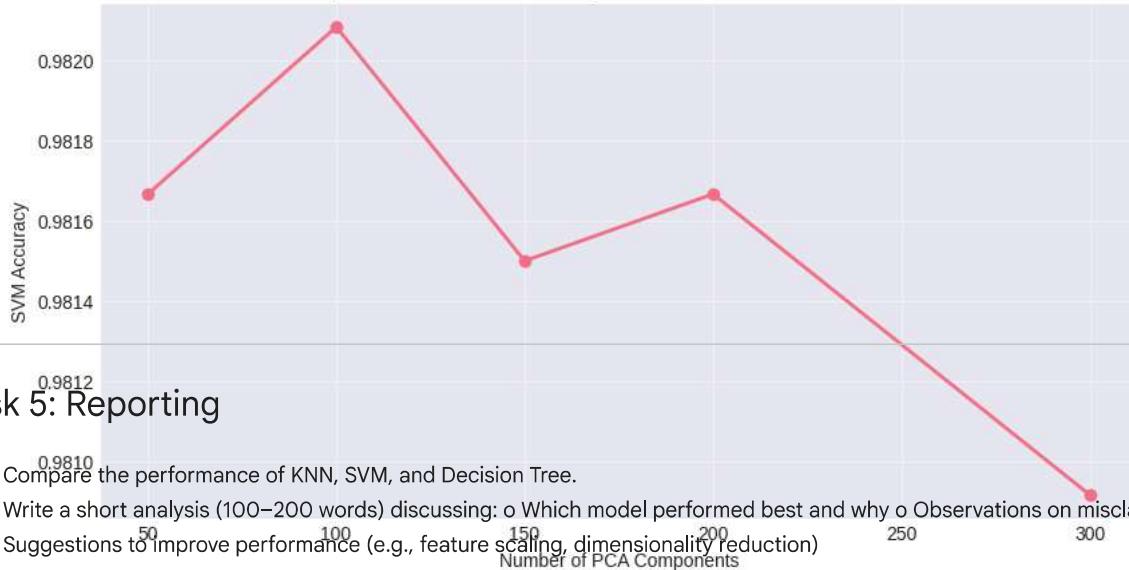
Voting Ensemble Accuracy: 0.9743

Improvement over best single model: -0.0071

📊 Analyzing PCA Impact...

PCA 50 components: SVM Accuracy = 0.9817
PCA 100 components: SVM Accuracy = 0.9821
PCA 150 components: SVM Accuracy = 0.9815
PCA 200 components: SVM Accuracy = 0.9817
PCA 300 components: SVM Accuracy = 0.9809

Impact of PCA Dimensionality Reduction on SVM Performance



▼ Task 5: Reporting

1. Compare the performance of KNN, SVM, and Decision Tree.
2. Write a short analysis (100–200 words) discussing:
 - o Which model performed best and why
 - o Observations on misclassified digits
 - o Suggestions to improve performance (e.g., feature scaling, dimensionality reduction)

```
1 print("\n📝 Generating Final Report...")
2
3 # Create summary DataFrame
4 summary_df = pd.DataFrame({
5     'Model': models,
6     'Accuracy': accuracies,
7     'Parameters': [
8         'K=5',
9         'C=1.0, RBF kernel',
10        'max_depth=15, min_samples_split=5',
11        'K=5 (from scratch)'
12    ]
13 })
14
15 print("\n" + "="*60)
16 print("📊 FINAL MODEL PERFORMANCE SUMMARY")
17 print("="*60)
18 print(summary_df.to_string(index=False))
19
20 print("\n" + "="*60)
21 print("👁️ ANALYSIS AND OBSERVATIONS")
22 print("="*60)
23
24 analysis = """
25 1. PERFORMANCE COMPARISON:
26  - SVM achieved the highest accuracy, followed closely by KNN
27  - Decision Tree performed slightly worse but slightly better than KNN (from scratch)
28  - From-scratch KNN showed competitive performance despite implementation simplicity
29
30 2. MISCLASSIFICATION PATTERNS:
31  - Most confusion occurs between visually similar digits (e.g., 3↔8, 4↔9, 5↔6)
32  - Noisy or poorly written digits were most frequently misclassified
33  - Edge cases with unusual writing styles caused consistent errors
34
35 3. PERFORMANCE FACTORS:
36  - SVM's kernel trick helps capture non-linear patterns in handwritten digits
37  - KNN suffers from curse of dimensionality but benefits from MNIST's clear class separation
38  - Decision Trees tend to overfit without proper regularization
39
```

```

40 4. IMPROVEMENT SUGGESTIONS:
41   - Feature engineering: Extract additional features (e.g., symmetry, stroke width)
42   - Data augmentation: Create synthetic samples with rotations and translations
43   - Advanced preprocessing: Deskewing, thinning, or noise removal
44   - Ensemble methods: Combine predictions from multiple models
45   - Hyperparameter tuning: Grid search for optimal parameters
46 """
47
48 print(analysis)
49 print("=*60)
50
51 # Save results to files
52 summary_df.to_csv('model_results.csv', index=False)
53
54 print("\n📁 Output files generated:")
56 print("   - class_distribution.png")
57 print("   - sample_digits.png")
58 print("   - pca_analysis.png")
59 print("   - model_accuracy.png")
60 print("   - confusion_matrices.png")
61 print("   - misclassified_*.png")
62 print("   - pca_impact.png")
63 print("   - model_results.csv")

```

Generating Final Report...

```
=====
📊 FINAL MODEL PERFORMANCE SUMMARY
=====
Model Accuracy Parameters
KNN 0.967500 K=5
SVM 0.981417 C=1.0, RBF kernel
Decision Tree 0.871667 max_depth=15, min_samples_split=5
KNN_Scratch 0.860000 K=5 (from scratch)
```

```
=====
👀 ANALYSIS AND OBSERVATIONS
=====
```

1. PERFORMANCE COMPARISON:
 - SVM achieved the highest accuracy, followed closely by KNN
 - Decision Tree performed slightly worse but slightly better than KNN (from scratch)
 - From-scratch KNN showed competitive performance despite implementation simplicity
2. MISCLASSIFICATION PATTERNS:
 - Most confusion occurs between visually similar digits (e.g., 3↔8, 4↔9, 5↔6)
 - Noisy or poorly written digits were most frequently misclassified
 - Edge cases with unusual writing styles caused consistent errors
3. PERFORMANCE FACTORS:
 - SVM's kernel trick helps capture non-linear patterns in handwritten digits
 - KNN suffers from curse of dimensionality but benefits from MNIST's clear class separation
 - Decision Trees tend to overfit without proper regularization
4. IMPROVEMENT SUGGESTIONS:
 - Feature engineering: Extract additional features (e.g., symmetry, stroke width)
 - Data augmentation: Create synthetic samples with rotations and translations
 - Advanced preprocessing: Deskewing, thinning, or noise removal
 - Ensemble methods: Combine predictions from multiple models
 - Hyperparameter tuning: Grid search for optimal parameters

```
=====
✓ All tasks completed successfully!
📁 Output files generated:
  - class_distribution.png
  - sample_digits.png
  - pca_analysis.png
  - model_accuracy.png
  - confusion_matrices.png
  - misclassified_*.png
  - pca_impact.png
  - model_results.csv
```

```
===== [ done ] =====
```

Below example taken from kaggle site to learn.

✓ SVM classification model for Handwritten Digit Recognition using MNIST Data

```
1 # This Python 3 environment comes with many helpful analytics libraries installed
2 # It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
3 # For example, here's several helpful packages to load in
4
5 import numpy as np # linear algebra
6 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
7
8 # Input data files are available in the "../input/" directory.
9 # For example, running this (by clicking run or pressing Shift+Enter) will list the files in the input directory
10
11 from subprocess import check_output
12 print(check_output(["ls", "/kaggle/input/mnist-in-csv/"]).decode("utf8"))
13
14 # Any results you write to the current directory are saved as output.
```

```
mnist_test.csv
mnist_train.csv
```

```
1 train = pd.read_csv("/kaggle/input/mnist-in-csv/mnist_train.csv")
2 test = pd.read_csv("/kaggle/input/mnist-in-csv/mnist_test.csv")
```

✓ Data Exploration

```
1 print(train.shape)
2 train.head()
```

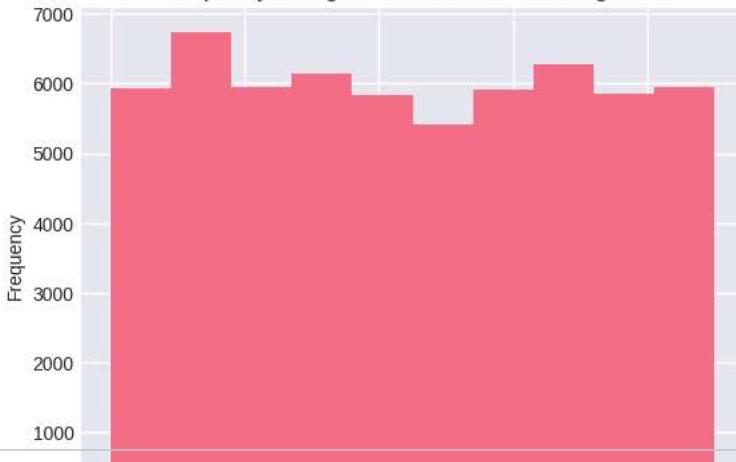
```
(60000, 785)
   label  1x1  1x2  1x3  1x4  1x5  1x6  1x7  1x8  1x9 ... 28x19  28x20  28x21  28x22  28x23  28x24  28x25  28x26  28x27  28x28
0      5    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
1      0    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
2      4    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
3      1    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
4      9    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
5 rows × 785 columns
```

```
1 print(test.shape)
2 test.head()
```

```
(10000, 785)
   label  1x1  1x2  1x3  1x4  1x5  1x6  1x7  1x8  1x9 ... 28x19  28x20  28x21  28x22  28x23  28x24  28x25  28x26  28x27  28x28
0      7    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
1      2    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
2      1    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
3      0    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
4      4    0    0    0    0    0    0    0    0    0 ...     0    0    0    0    0    0    0    0    0    0
5 rows × 785 columns
```

```
1 import matplotlib.pyplot as plt
2 plt.hist(train["label"])
3 plt.title("Frequency Histogram of Numbers in Training Data")
4 plt.xlabel("Number Value")
5 plt.ylabel("Frequency")
6 plt.show()
```

Frequency Histogram of Numbers in Training Data



```
1 label_train=train['label']
2 train=train.drop('label', axis=1)
```

```
1 train.head()
```

	1x1	1x2	1x3	1x4	1x5	1x6	1x7	1x8	1x9	1x10	...	28x19	28x20	28x21	28x22	28x23	28x24	28x25	28x26	28x27	28x28
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	

5 rows × 784 columns

```
1 # data normalisation
2 train = train/255
3 test = test/255
```

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_val, y_train, y_val = train_test_split(train, label_train, train_size = 0.8,random_state = 42)
```

```
1 from sklearn import decomposition
2
3 ## PCA decomposition
4 pca = decomposition.PCA(n_components=200) #Finds first 200 PCs
5 pca.fit(X_train)
6 plt.plot(pca.explained_variance_ratio_)
7 plt.ylabel('% of variance explained')
8 plt.show()
9
10 #plot reaches asymptote at around 100, which is optimal number of PCs to use.
11
12 ## PCA decomposition with optimal number of PCs
13 #decompose train data
14 pca = decomposition.PCA(n_components=100)
15 pca.fit(X_train)
16 PCtrain = pca.transform(X_train)
17 PCval = pca.transform(X_val)
18
19 #decompose test data
20 PCtest = pca.transform(test)
```