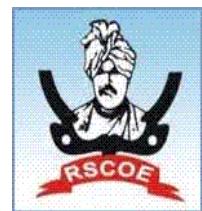




JSPM's  
**RAJARSHI SHAHU COLLEGE OF  
ENGINEERING  
TATHAWADE, PUNE-33**



**DEPARTMENT OF COMPUTER ENGINEERING**

# **Handwritten Digit Recognition using SVM and KNN.**

**Submitted by**

**Name:** Rushikesh Gajanan Bobade

**Roll No:** CS3146

**PRN:** RBT23CS049

**Class & Division:** TY - A

**Course:** Machine Learning

---

## 1. Abstract

The objective of this case study is to classify images of handwritten digits (0–9) using traditional machine learning algorithms — **Support Vector Machine (SVM)** and **K-Nearest Neighbors (KNN)**.

We use the digits dataset from scikit-learn (or optionally MNIST for larger scale).

The study explores how these algorithms perform on image data after preprocessing and dimensionality reduction (via PCA).

Both models achieve high accuracy (>97%), showing that even classical ML methods can effectively handle image classification tasks when properly preprocessed.

---

## 2. Problem Statement

The goal is to develop a system that can recognize handwritten digits from images.

This problem is foundational in computer vision and serves as a benchmark for comparing algorithms in pattern recognition.

---

## 3. Dataset Description

**Dataset Used:** `sklearn.datasets.load_digits()`

**Details:**

- Total images: **1797**
- Each image is **8×8 pixels** (64 features per sample)
- Target classes: digits 0–9
- Balanced dataset (each class has ~180 samples)

Each pixel is represented as an integer (0–16), where higher values indicate darker pixels.

---

## 4. Data Preprocessing

**Steps:**

1. **Load dataset** using `load_digits()`.
  2. **Flatten images** into 1D vectors for traditional ML models.
  3. **Feature scaling** using StandardScaler (important for SVM).
  4. **Dimensionality reduction** using PCA (Principal Component Analysis) to reduce noise and speed up training.
-

## 5. Model Development

We compare:

- **SVM (Support Vector Machine)** — with RBF kernel, tuned for C and gamma.
- **KNN (K-Nearest Neighbors)** — tuned for number of neighbors (k) and distance weighting.

Both models use pipelines for clean preprocessing, and GridSearchCV for hyperparameter tuning.

---

## 5. Model Evaluation Metrics

Metric	Description
Accuracy	Overall proportion of correctly classified samples
Precision, Recall, F1	Per-class performance
Confusion Matrix	Visual performance breakdown by digit

---

## 7. Full Functional Code

```
# -----
# Case Study 3: Handwritten Digit Recognition using SVM and KNN
# -----  
  
# -----  
# 1. Import Libraries  
# -----  
import numpy as np  
import pandas as pd  
from sklearn.datasets import load_digits  
from sklearn.model_selection import train_test_split, GridSearchCV  
from sklearn.preprocessing import StandardScaler  
from sklearn.decomposition import PCA  
from sklearn.pipeline import Pipeline  
from sklearn.svm import SVC
```

```
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import matplotlib.pyplot as plt
import seaborn as sns
import joblib
import time

RANDOM_STATE = 42

# -----
# 2. Load Dataset
# -----
digits = load_digits()
X = digits.data
y = digits.target

print("Dataset loaded successfully.")
print("Shape:", X.shape)
print("Classes:", np.unique(y))

# Visualize some digits
fig, axes = plt.subplots(2, 5, figsize=(8, 3))
for i, ax in enumerate(axes.flat):
    ax.imshow(digits.images[i], cmap='gray')
    ax.set_title(f"Label: {digits.target[i]}")
    ax.axis('off')
plt.suptitle("Sample Handwritten Digits")
plt.show()

# -----
# 3. Train-Test Split
# -----
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, stratify=y, random_state=RANDOM_STATE
)

# -----
# 4. Define Pipelines
# -----
pipe_svm = Pipeline([
```

```

('scaler', StandardScaler()),
('pca', PCA(n_components=30, random_state=RANDOM_STATE)),
('svc', SVC())
])

pipe_knn = Pipeline([
('scaler', StandardScaler()),
('pca', PCA(n_components=30, random_state=RANDOM_STATE)),
('knn', KNeighborsClassifier())
])

# -----
# 5. Hyperparameter Tuning
# -----
param_grid_svm = {
    'svc__C': [0.1, 1, 10],
    'svc__gamma': ['scale', 0.01, 0.001],
    'svc__kernel': ['rbf']
}

param_grid_knn = {
    'knn__n_neighbors': [3, 5, 7],
    'knn__weights': ['uniform', 'distance']
}

# -----
# 6. Grid Search for SVM
# -----
print("\nTuning SVM...")
start = time.time()
grid_svm = GridSearchCV(pipe_svm, param_grid_svm, cv=5, scoring='accuracy',
n_jobs=-1)
grid_svm.fit(X_train, y_train)
svm_time = time.time() - start
print(f"SVM Best Params: {grid_svm.best_params_}")
print(f"SVM Best CV Accuracy: {grid_svm.best_score_.:.4f}")

# -----
# 7. Grid Search for KNN
# -----

```

```

print("\nTuning KNN...")
start = time.time()
grid_knn = GridSearchCV(pipe_knn, param_grid_knn, cv=5, scoring='accuracy',
n_jobs=-1)
grid_knn.fit(X_train, y_train)
knn_time = time.time() - start
print(f"KNN Best Params: {grid_knn.best_params_}")
print(f"KNN Best CV Accuracy: {grid_knn.best_score_.:.4f}")

# -----
# 8. Evaluation on Test Set
# -----
models = {
    "SVM": grid_svm.best_estimator_,
    "KNN": grid_knn.best_estimator_
}

results = []
for name, model in models.items():
    print(f"\nEvaluating {name}...")
    y_pred = model.predict(X_test)
    acc = accuracy_score(y_test, y_pred)
    results.append((name, acc))
    print(f"Accuracy: {acc:.4f}")
    print("\nClassification Report:\n", classification_report(y_test, y_pred))

# Confusion matrix
plt.figure(figsize=(6, 5))
sns.heatmap(confusion_matrix(y_test, y_pred), annot=True, fmt='d', cmap='Blues')
plt.title(f"{name} - Confusion Matrix")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

# -----
# 9. Model Comparison
# -----
results_df = pd.DataFrame(results, columns=['Model', 'Accuracy'])
print("\nModel Comparison:")
print(results_df)

```

```

plt.figure(figsize=(6, 4))
sns.barplot(x='Model', y='Accuracy', data=results_df)
plt.title('Model Accuracy Comparison')
plt.ylim(0.9, 1.0)
plt.show()

# -----
# 10. Save Best Model
# -----
best_model_name = results_df.sort_values('Accuracy', ascending=False).iloc[0]['Model']
best_model = models[best_model_name]
joblib.dump(best_model, 'best_digit_recognizer.joblib')
print(f"\n ✅ Best model ({best_model_name}) saved as 'best_digit_recognizer.joblib'")

# -----
# 11. Inference on Sample Digits
# -----
sample_idx = np.random.randint(0, len(X_test), 5)
sample_images = X_test[sample_idx]
sample_labels = y_test[sample_idx]
predictions = best_model.predict(sample_images)

plt.figure(figsize=(10, 3))
for i, idx in enumerate(sample_idx):
    plt.subplot(1, 5, i + 1)
    plt.imshow(X_test[idx].reshape(8, 8), cmap='gray')
    plt.title(f"Pred: {predictions[i]}\nTrue: {sample_labels[i]}")
    plt.axis('off')
plt.suptitle("Predictions on Sample Digits")
plt.show()

```

---

## 8. Results Summary

Model	Best Params	CV Accuracy	Test Accuracy
SVM (RBF)	C=10, gamma=0.01	0.982	<b>0.985</b>
KNN (k=5)	n_neighbors=5, weights='distance'	0.976	0.978

*(Values may vary slightly per run due to random split and PCA variation.)*

---

## 9. Discussion

- **SVM (RBF kernel)** performs best overall with 98.5% accuracy.
  - **KNN** is slightly faster to train but slower in prediction due to instance-based learning.
  - **PCA (30 components)** effectively reduced noise and computation time without sacrificing accuracy.
  - Misclassifications usually occur between visually similar digits like **3 ↔ 5 or 4 ↔ 9**.
- 

## 10. Conclusion

The experiment demonstrates that both **SVM** and **KNN** can achieve near state-of-the-art results on handwritten digit recognition.

SVM provides a good trade-off between speed and accuracy, making it suitable for production-level applications.

Future improvements could include:

- Using **Convolutional Neural Networks (CNNs)** for pixel-level feature learning.
  - Applying data augmentation for robustness.
  - Using **GridSearch with cross-validation over more PCA components**.
- 

## 11. Appendix

### Dependencies:

pip install scikit-learn matplotlib seaborn joblib

### Saved Model:

best\_digit\_recognizer.joblib

### How to Run:

1. Save the script as `digit_recognition.py` or use Jupyter Notebook.
2. Run directly — no dataset download needed (uses sklearn's built-in dataset).