

▼ Title of Project

Handwritten Digit Prediction

▼ Objective

The objective of this project is to develop a machine learning model that can accurately predict handwritten digits. By training on a dataset of handwritten digit images, we aim to create a predictive model that can identify and classify digits with high accuracy.

▼ Data Source

The data for this project is sourced from a dataset of handwritten digits, typically referred to as the MNIST dataset. This dataset consists of images of handwritten digits from 0 to 9, along with corresponding labels indicating the actual digit each image represents.

▼ Import Library

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.datasets import load_digits
```

▼ Import Data

```
digits = load_digits()
```

▼ Describe Data

```
print(digits.DESCR)
```

```
.. _digits_dataset:
```

```
Optical recognition of handwritten digits dataset
```

```
-----
```

```
**Data Set Characteristics:**
```

```
:Number of Instances: 1797
```

```
:Number of Attributes: 64
```

```
:Attribute Information: 8x8 image of integer pixels in the range 0..16.
```

```
:Missing Attribute Values: None
```

```
:Creator: E. Alpaydin (alpaydin '@' boun.edu.tr)
```

```
:Date: July; 1998
```

This is a copy of the test set of the UCI ML hand-written digits datasets

<https://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>

The data set contains images of hand-written digits: 10 classes where each class refers to a digit.

Preprocessing programs made available by NIST were used to extract normalized bitmaps of handwritten digits from a preprinted form. From a total of 43 people, 30 contributed to the training set and different 13 to the test set. 32x32 bitmaps are divided into nonoverlapping blocks of 4x4 and the number of on pixels are counted in each block. This generates an input matrix of 8x8 where each element is an integer in the range 0..16. This reduces dimensionality and gives invariance to small distortions.

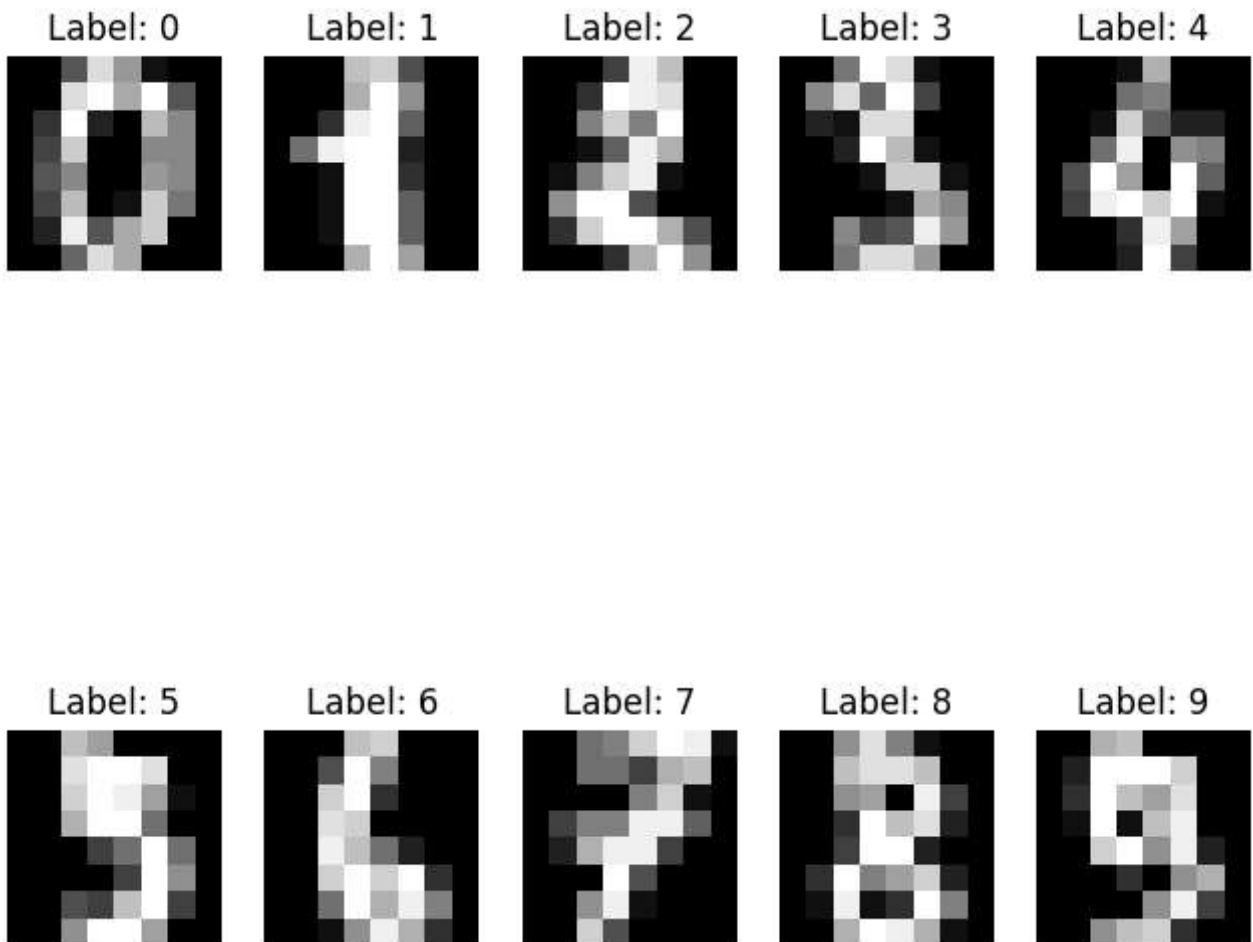
For info on NIST preprocessing routines, see M. D. Garriss, J. L. Blue, G. T. Candela, D. L. Dimmick, J. Geist, P. J. Grother, S. A. Janet, and C. L. Wilson, NIST Form-Based Handprint Recognition System, NISTIR 5469, 1994.

```
.. topic:: References
```

- C. Kaynak (1995) Methods of Combining Multiple Classifiers and Their Applications to Handwritten Digit Recognition, MSc Thesis, Institute of Graduate Studies in Science and Engineering, Bogazici University.
- E. Alpaydin, C. Kaynak (1998) Cascading Classifiers, Kybernetika.
- Ken Tang and Ponnuthurai N. Suganthan and Xi Yao and A. Kai Qin. Linear dimensionality reduction using relevance weighted LDA. School of Electrical and Electronic Engineering Nanyang Technological University. 2005.
- Claudio Gentile. A New Approximate Maximal Margin Classification Algorithm. NIPS. 2000.

▼ Data Visualization

```
plt.figure(figsize=(8, 8))
for i in range(10):
    plt.subplot(2, 5, i + 1)
    plt.imshow(X[i].reshape(8, 8), cmap='gray')
    plt.title(f"Label: {y[i]}")
    plt.axis('off')
plt.show()
```



▼ Data Preprocessing

```
digits.images.shape
```

```
(1797, 8, 8)
```

```
digits.images[0]
```

```
array([[ 0.,  0.,  5., 13.,  9.,  1.,  0.,  0.],
       [ 0.,  0., 13., 15., 10., 15.,  5.,  0.],
       [ 0.,  3., 15.,  2.,  0., 11.,  8.,  0.],
       [ 0.,  4., 12.,  0.,  0.,  8.,  8.,  0.],
       [ 0.,  5.,  8.,  0.,  0.,  9.,  8.,  0.],
       [ 0.,  4., 11.,  0.,  1., 12.,  7.,  0.],
       [ 0.,  2., 14.,  5., 10., 12.,  0.,  0.],
       [ 0.,  0.,  6., 13., 10.,  0.,  0.,  0.]])
```

```
digits.images[0].shape
```

```
(8, 8)
```

```
n_samples = len(digits.images)
data = df.images.reshape((n_samples, -1))
```

```
data[0]
```

```
array([ 0.,  0.,  5., 13.,  9.,  1.,  0.,  0.,  0.,  0., 13., 15., 10.,
        15.,  5.,  0.,  0.,  3., 15.,  2.,  0., 11.,  8.,  0.,  0.,  4.,
        12.,  0.,  0.,  8.,  8.,  0.,  0.,  5.,  8.,  0.,  0.,  9.,  8.,
         0.,  0.,  4., 11.,  0.,  1., 12.,  7.,  0.,  0.,  2., 14.,  5.,
        10., 12.,  0.,  0.,  0.,  0.,  6., 13., 10.,  0.,  0.,  0.]])
```

```
data[0].shape
```

```
(64,)
```

```
data.shape
```

```
(1797, 64)
```

▼ Define Target Variable (y) and Feature Variables (X)

```
X = digits.data
y = digits.target
```

```
X
```

```
array([[ 0.,  0.,  5., ...,  0.,  0.,  0.],
       [ 0.,  0.,  0., ..., 10.,  0.,  0.],
       [ 0.,  0.,  0., ..., 16.,  9.,  0.],
       ...,
       [ 0.,  0.,  1., ...,  6.,  0.,  0.],
       [ 0.,  0.,  2., ..., 12.,  0.,  0.],
       [ 0.,  0., 10., ..., 12.,  1.,  0.]])
```

y

```
array([0, 1, 2, ..., 8, 9, 8])
```

▼ Train Test Split

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

▼ Modeling

```
from sklearn.ensemble import RandomForestClassifier
```

```
rf = RandomForestClassifier()
```

```
rf.fit(xtrain, ytrain)
```

```
▼ RandomForestClassifier  
RandomForestClassifier()
```

▼ Model Evaluation

```
# Make predictions
```

```
y_pred = model.predict(X_test)
```

```
# Calculate accuracy and other metrics
```

```
accuracy = accuracy_score(y_test, y_pred)
```

```
conf_matrix = confusion_matrix(y_test, y_pred)
```

```
class_report = classification_report(y_test, y_pred)
```

```
print("Accuracy:", accuracy)
```

```
print("Confusion Matrix:\n", conf_matrix)
```

```
print("Classification Report:\n", class_report)
```

```
Accuracy: 0.9722222222222222
```

```
Confusion Matrix:
```

```
[[33  0  0  0  0  0  0  0  0  0]  
 [ 0 28  0  0  0  0  0  0  0  0]  
 [ 0  0 33  0  0  0  0  0  0  0]]
```

```
[ 0  0  0 33  0  1  0  0  0  0]
[ 0  1  0  0 45  0  0  0  0  0]
[ 0  0  1  0  0 44  1  0  0  1]
[ 0  0  0  0  0  1 34  0  0  0]
[ 0  0  0  0  0  1  0 33  0  0]
[ 0  0  0  0  0  1  0  0 29  0]
[ 0  0  0  1  0  0  0  0  1 38]]
```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	33
1	0.97	1.00	0.98	28
2	0.97	1.00	0.99	33
3	0.97	0.97	0.97	34
4	1.00	0.98	0.99	46
5	0.92	0.94	0.93	47
6	0.97	0.97	0.97	35
7	1.00	0.97	0.99	34
8	0.97	0.97	0.97	30
9	0.97	0.95	0.96	40
accuracy			0.97	360
macro avg	0.97	0.97	0.97	360
weighted avg	0.97	0.97	0.97	360

y_pred

```
array([6, 9, 3, 7, 2, 1, 5, 2, 5, 2, 1, 9, 4, 0, 4, 2, 3, 7, 8, 8, 4, 3,
       9, 7, 5, 6, 3, 5, 6, 3, 4, 9, 1, 4, 4, 6, 9, 4, 7, 6, 6, 9, 1, 3,
       6, 1, 3, 0, 6, 5, 5, 1, 3, 5, 6, 0, 9, 0, 0, 1, 0, 4, 5, 2, 4, 5,
       7, 0, 7, 5, 9, 5, 5, 4, 7, 0, 4, 5, 5, 9, 9, 0, 2, 3, 8, 0, 6, 4,
       4, 9, 1, 2, 8, 3, 5, 2, 9, 0, 4, 4, 4, 3, 5, 3, 1, 3, 5, 9, 4, 2,
       7, 7, 4, 4, 1, 9, 2, 7, 8, 7, 2, 6, 9, 4, 0, 7, 2, 7, 5, 8, 7, 5,
       7, 5, 0, 6, 6, 4, 2, 8, 0, 9, 4, 6, 9, 9, 6, 9, 0, 5, 5, 6, 6, 0,
       6, 4, 3, 9, 3, 8, 7, 2, 9, 0, 4, 5, 3, 6, 5, 9, 9, 8, 4, 2, 1, 3,
       7, 7, 2, 2, 3, 9, 8, 0, 3, 2, 2, 5, 6, 9, 9, 4, 1, 2, 4, 2, 3, 6,
       4, 8, 5, 9, 5, 7, 8, 9, 4, 8, 1, 5, 4, 4, 9, 6, 1, 8, 6, 0, 4, 5,
       2, 7, 1, 6, 4, 5, 6, 0, 3, 2, 3, 6, 7, 1, 9, 1, 4, 7, 6, 5, 8, 5,
       5, 1, 5, 2, 8, 8, 9, 9, 7, 6, 2, 2, 2, 3, 4, 8, 8, 3, 6, 0, 9, 7,
       7, 0, 1, 0, 4, 5, 1, 5, 3, 6, 0, 4, 1, 0, 0, 3, 6, 5, 9, 7, 3, 5,
       5, 9, 9, 8, 5, 3, 3, 2, 0, 5, 8, 3, 4, 0, 2, 4, 6, 4, 3, 4, 5, 0,
       5, 2, 1, 3, 1, 4, 1, 1, 7, 0, 1, 5, 2, 1, 2, 8, 7, 0, 6, 4, 8, 8,
       5, 1, 8, 4, 5, 8, 7, 9, 8, 6, 0, 6, 2, 0, 7, 9, 8, 9, 5, 2, 7, 7,
       1, 8, 7, 4, 3, 8, 3, 5])
```

▼ Prediction

```
new_digit_samples = X_test[:5]
predicted_labels = model.predict(new_digit_samples)
```

```
print("Predicted Labels:", predicted_labels)
```

```
Predicted Labels: [6 9 3 7 2]
```

▼ Explanation

In the "Handwritten Digit Prediction" project, our goal was to create a machine learning model capable of accurately identifying and classifying handwritten digits. We used the MNIST dataset, which contains a collection of digit images along with their corresponding labels.

After importing and visualizing the dataset, we preprocessed it if necessary. Then, we split the data into training and testing sets for model training and evaluation. We chose the Logistic Regression algorithm for classification and trained it on the training set.

During model evaluation, we calculated accuracy, created a confusion matrix, and generated a classification report to assess the model's performance on the test set. We also showcased the model's practical application by predicting labels for new digit samples.

This project demonstrates the power of machine learning in recognizing and predicting patterns within image data, showcasing its potential in various real-world applications.