

MNIST Mini Project

Mnist database : The MNIST handwritten digit classification problem is a standard dataset used in computer vision and deep learning. Although the dataset is effectively solved, it can be used as the basis for learning and practicing how to develop, evaluate, and use convolutional deep learning neural networks for image classification from scratch. This includes how to develop a robust test harness for estimating the performance of the model, how to explore improvements to the model, and how to save the model and later load it to make predictions on new data. In this tutorial, you will discover how to develop a convolutional neural network for handwritten digit classification from scratch

Sklearn : Scikit-learn (Sklearn) is the most useful and robust library for machine learning in Python. It provides a selection of efficient tools for machine learning and statistical modeling including classification, regression, clustering and dimensionality reduction via a consistence interface in Python. This library, which is largely written in Python, is built upon NumPy, SciPy and Matplotlib.

Numpy : NumPy is the fundamental package for scientific computing in Python

Pandas : Pandas is an open source Python package that is most widely used for data science/data analysis and machine learning tasks. It is built on top of another package named Numpy, which provides support for multi-dimensional arrays. As one of the most popular data wrangling packages, Pandas works well with many other data science modules inside the Python ecosystem, and is typically included in every Python distribution, from those that come with your operating system to commercial vendor distributions like ActiveState's ActivePython.

seaborn : Seaborn is a data visualization library built on top of matplotlib and closely integrated with pandas data structures in Python.

matplotlib : Matplotlib is a plotting library available for the Python programming language as a component of NumPy, a big data numerical handling resource. Matplotlib uses an object oriented API to embed plots in Python applications.

Importing libraries

In [1]:

```
import numpy as np
import pandas as pd
import seaborn as sn
import matplotlib.pyplot as plt
%matplotlib inline
```

Loading data

In [2]:

```
from sklearn.datasets import fetch_openml
mnist = fetch_openml('mnist_784')
```

creating dataframe of mnist data

In [3]:

```
dir(mnist)
```

Out[3]:

```
['DESCR',
 'categories',
 'data',
 'details',
 'feature_names']
```

In [4]:

printing first value of data

In [5]:

Out[5]:

[illegible]

In [10]:

```
labels.head()
```

Out[10]:

class	
0	5
1	0
2	4
3	1
4	9

In [11]:

```
pixels.head()
```

Out[11]:

	pixel1	pixel2	pixel3	pixel4	pixel5	pixel6	pixel7	pixel8	pixel9	pixel10	...	pixel775	pixel776	pixel777	pixel778	pixel779
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	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	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	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	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0

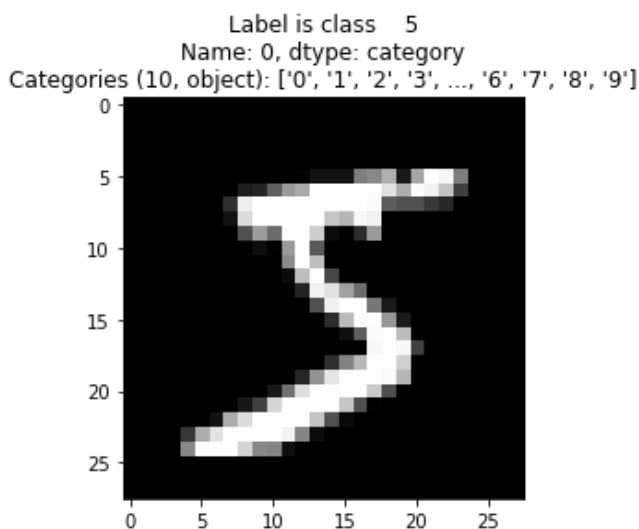
5 rows × 784 columns



printing data and target value

In [12]:

```
label = labels.loc[0]
pixel = pixels.loc[0]
pixel = np.array(pixel, dtype='uint8')
pixel = pixel.reshape((28,28))
plt.title('Label is {label}'.format(label=label))
plt.imshow(pixel, cmap='gray')
plt.show()
```



splitting the data into training data and testing data

In [13]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(mnist.data, mnist.target, test_size=
1/7.0)
```

Using RandomForestClassifier to train the model

In [14]:

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
model.fit(X_train, y_train)
```

Out[14]:

```
▼ RandomForestClassifier
RandomForestClassifier()
```

measuring the score and accuracy of model

In [15]:

```
model.score(X_test, y_test)
```

Out[15]:

0.9688

In [18]:

```
model.predict([pixels.loc[3].values]) #predicting the value of data at index 3
```

```
C:\Users\acer\anaconda3\lib\site-packages\sklearn\base.py:450: UserWarning: X does not ha
ve valid feature names, but RandomForestClassifier was fitted with feature names
  warnings.warn(
```

Out[18]:

```
array(['1'], dtype=object)
```

calculating cofusion matrix

In [20]:

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

In [21]:

```
from sklearn.metrics import accuracy_score
accuracy_score(y_pred, y_test)
```

Out[21]:

0.9688

In [22]:

```
from sklearn.metrics import confusion_matrix
```

In [23]:

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

Out[23]:

```
array([[55,  2,  0,  0,  0,  0,  0,  0,  0,  0],
       [ 2, 55,  0,  0,  0,  0,  0,  0,  0,  0],
       [ 0,  0, 55,  0,  0,  0,  0,  0,  0,  0],
       [ 0,  0,  0, 55,  0,  0,  0,  0,  0,  0],
       [ 0,  0,  0,  0, 55,  0,  0,  0,  0,  0],
       [ 0,  0,  0,  0,  0, 55,  0,  0,  0,  0],
       [ 0,  0,  0,  0,  0,  0, 55,  0,  0,  0],
       [ 0,  0,  0,  0,  0,  0,  0, 55,  0,  0],
       [ 0,  0,  0,  0,  0,  0,  0,  0, 55,  0],
       [ 0,  0,  0,  0,  0,  0,  0,  0,  0, 55]])
```

```
array([[ 985,    0,    0,    2,    0,    3,    4,    0,    5,    2],
       [   0, 1098,    6,    2,    1,    0,    0,    3,    2,    3],
       [  10,    2,  992,    2,    3,    0,    1,    6,    6,    1],
       [   1,    0,   12,  994,    0,   14,    1,    8,    8,    4],
       [   0,    1,    3,    0,  971,    0,    1,    2,    2,   17],
       [   4,    1,    1,   16,    1,  833,    9,    0,    8,    2],
       [   3,    4,    0,    0,    0,    4,  970,    0,    3,    0],
       [   1,    7,   14,    0,    3,    0,    0, 1011,    2,   10],
       [   5,    2,    6,    6,    2,    8,    2,    1,   936,    3],
       [   6,    3,    1,   11,   11,    2,    1,    5,    6,   898]],
      dtype=int64)
```

In [24]:

```
plt.figure(figsize=(10,7))
sn.heatmap(cm,annot=True, fmt="d")
plt.xlabel('Predicted')
plt.ylabel('Truth')
```

Out[24]:

Text(69.0, 0.5, 'Truth')

