

1.导入数据

In [4]:

```
import os
```

In [5]:

```
import pandas as pd
```

In [6]:

```
dir_root = os.path.join(os.path.curdir, 'digits')
```

In [7]:

```
dir_training = os.path.join(dir_root, 'trainingDigits') # 训练集
```

In [8]:

```
dir_test = os.path.join(dir_root, 'testDigits') # 测试集
```

In [9]:

```
def df_exact(dir_name):
    for root, dirs, files in os.walk(dir_name):
        df_digit = pd.DataFrame(columns=['X', 'Y'])
        for f in files:
            list_f = f.split('_')
            y = int(list_f[0])
            df_data = pd.read_csv(os.path.join(dir_name, f), header=None)
            list_x = df_data.iloc[:, 0]
            s = ''
            x = []
            for d in list_x:
                s = s + d
            for i in s:
                x.append(int(i))
            new = pd.DataFrame({'X':x, 'Y':y})
            df_digit = df_digit.append(new)
    print(df_digit)
    return df_digit
```

In [10]:

```
ds_training = df_exact(dir_training)
```

```

      X Y
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, ... 0
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, ... 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, 1, 1, ... 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, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, ... 9
```

[1934 rows x 2 columns]

In [11]:

```
ds_test = df_exact(dir_test)
```

```

      X Y
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 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, 0, 0, 0, 1, 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, 0, 0, 0, 0, 0, 0, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, ... 9
0  [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... 9
```

[946 rows x 2 columns]

In [12]:

```
ds_training = ds_training.reset_index()
ds_test = ds_test.reset_index()
```

2、采用sklearn完成knn算法

In [21]:

```
import numpy as np
from sklearn.neighbors import KNeighborsClassifier
```

In [25]:

```
array_x = np.array([i for i in ds_training.X])
array_y = np.array([i for i in ds_training.Y])
arr_testx = np.array([i for i in ds_test.X])
arr_testy = np.array([i for i in ds_test.Y])
```

In [26]:

```
neigh = KNeighborsClassifier(n_neighbors=15)
```

In [27]:

```
neigh.fit(array_x, array_y)
```

Out[27]:

```
KNeighborsClassifier(n_neighbors=15)
```

In [28]:

```
y_p = neigh.predict(arr_testx)
```

y_p [illegible]

In [32]:

```
def distEclud(vecA, vecB):
    return np.linalg.norm(vecA - vecB)
```

In [36]:

```
def knn(ds_training, ds_test, k = 10):
    Y_p = []
    for p in ds_test.X:
        ds = ds_training.copy(deep=True)
        D = []
        for q in ds_training.X:
            vecA = np.array(p)
            vecB = np.array(q)
            d = distEclud(vecA, vecB)
            D.append(d)
        ds['D'] = np.array(D)
        ds = ds.sort_values(by='D')
        ds = ds.reset_index()
        Y = ds.loc[0:k-1, ['Y']]
        # 打印Y
        rank_y = Y['Y'].value_counts()
        # 打印rank_y
        print(rank_y.index[0], end=',')
        y_p = rank_y.index[0]
        Y_p.append(y_p)
    return Y_p
```

In [37]:

```
Y_p = knn(ds_training, ds_test, k = 15)
```

[illegible]

In []:

