

课后作业：数据预处理、降维、特征提取及聚类

作者：欧新宇 (Xinyu OU)

本文档所展示的测试结果，均运行于：Intel Core i7-7700K CPU 4.2GHz

【作业提交】

将分类结果保存到文本文档进行提交(写上每一题的题号和题目，然后再贴答案)，同时提交源代码。

1. 测试结果命名为: ex09-结果-你的学号-你的姓名.txt
2. 输出图片命名为: ex09-性能对比图-你的学号-你的姓名.png (.jpg)
3. 源代码命名为:
 - ex09-01Ori-你的学号-你的姓名.py
 - ex09-02Preprocessing-你的学号-你的姓名.py
 - ex09-03PCA-你的学号-你的姓名.py

结果文件，要求每小标题注题号，两题之间要求空一行

要求在 "MNIST人脸识别" 数据集上完成以下任务，要求如下：

1. 要求使用原始数据集的10%完成以下习题。（已给出载入数据集和数据预处理部分的代码）
2. 使用MLP模型进行训练和测试，基本参数设置为：

```
1 solver='lbfgs', hidden_layer_sizes=[100, 100], activation='relu', alpha=1e-5,
  random_state=62,
```

3. 使用原始数据进行预测 (ex09-01Ori)
4. 测试六种预处理方向对性能的影响 (ex09-02Preprocessing)

```
1 methods = ['StandardScaler', 'MinMaxScaler', 'MaxAbsScaler', 'RobustScaler',
  'Normalizer', 'Binarizer']
```

在完成了六种预处理方法的性能输出后，对比预处理结果和原始数据的结果，从7种结果中选出性能最好的一种方法完成后续的实验。

5. 使用PCA进行降维，并测试性能，要求测试PCA的参数范围为 $n_components = [0.4:0.99]$ ，可以使用以下代码设置pca的参数范围 (ex09-03PCA, ex09-性能对比图)

```
1 num=20
2 scores = np.zeros([3,num])
3 scores[0,:] = np.linspace(0.40, 0.99, num)
```

6. 可视化出PCA的性能曲线 (ex09-03PCA, ex09-性能对比图)
7. (选做) 可视化出PCA曲线中最大值的点，以及未进行PCA降维的性能值点 (ex09-03PCA, ex09-性能对比图)

1. 载入MNIST数据集

```

1 import sys
2 import os
3 sys.path.append(os.path.join(os.getcwd(), '..', 'Modules'))
4 import load_MNIST
5
6 import time
7 start = time.time()
8
9 train_images = load_MNIST.load_train_images()
10 train_labels = load_MNIST.load_train_labels()
11 test_images = load_MNIST.load_test_images()
12 test_labels = load_MNIST.load_test_labels()
13
14 print("载入数据集共耗时: {:.3f}s".format(time.time() - start))

```

```

1 开始载入MNIST手写数字数据集:
2 训练集图片大小: 28*28, 已载入60000/60000.
3 训练集标签数量: 60000...已完成。
4 测试集图片大小: 28*28, 已载入10000/10000.
5 测试集标签数量: 10000...已完成。
6 载入数据集共耗时: 2.788s

```

2. 数据预处理 I

1. 将图像数据转换成二维矩阵，并归一化到 0 – 1 之间

```

1 # 标准调整形态的方法
2 # X_train = train_images.reshape(train_images.shape[0],
3   train_images.shape[1]*train_images.shape[2])/255
4 # 此处，因为我们已经知道的样本的形态，所以可以直接书写值
5
6 X_train = train_images.reshape(60000, 28*28)/255
7 y_train = train_labels
8 X_test = test_images.reshape(10000, 28*28)/255
9 y_test = test_labels

```

2. 为了加速运算，将样本缩减到10%的比例进行处理

```

1 # 为了提高训练速度，我们只提取10%的样本进行演示
2 X_train_lite = X_train[0:5999,:]
3 y_train_lite = y_train[0:5999]
4 X_test_lite = X_test[0:999,:]
5 y_test_lite = y_test[0:999]

```

3. 使用原始数据进行预测

```

1 # 导入多层感知机MLP神经网络
2 from sklearn.neural_network import MLPClassifier
3 import time
4
5 start = time.time()
6
7 mlp = MLPClassifier(solver='lbfgs', hidden_layer_sizes=[100, 100],
8   activation='relu', alpha=1e-5, random_state=62, verbose=2)

```

```

8 mlp.fit(X_train_lite, y_train_lite)
9
10 score_ori_train = mlp.score(X_train_lite, y_train_lite)
11 score_ori_test = mlp.score(X_test_lite, y_test_lite)
12
13 print('训练结束, 用时{:.2f}s.'.format(time.time() - start))
14 print('训练集得分: {:.4f}, 测试集得分: {:.4f}'.format(mlp.score(X_train_lite,
    y_train_lite), mlp.score(X_test_lite, y_test_lite)))

```

```

1 训练结束, 用时5.72s.
2 训练集得分: 1.0000, 测试集得分: 0.9239

```

4. 测试预处理对性能的影响

```

1 # 导入preprocessing预处理器
2 from sklearn import preprocessing
3
4 methods = ['StandardScaler', 'MinMaxScaler', 'MaxAbsScaler',
5           'RobustScaler', 'Normalizer', 'Binarizer']
6
7 for str in methods:
8     scaler = eval('preprocessing.' + str + '().fit(X_train_lite)')
9     X_train_scaled = scaler.transform(X_train_lite)
10    X_test_scaled = scaler.transform(X_test_lite)
11    mlp.fit(X_train_scaled, y_train_lite)
12    print('预处理方法: {}, 测试集得分: {:.4f}'
13          .format(str, mlp.score(X_test_scaled, y_test_lite)))

```

```

1 预处理方法: StandardScaler, 测试集得分: 0.9199
2 预处理方法: MinMaxScaler, 测试集得分: 0.9219
3 预处理方法: MaxAbsScaler, 测试集得分: 0.9219
4 预处理方法: RobustScaler, 测试集得分: 0.9059
5 预处理方法: Normalizer, 测试集得分: 0.9239
6 预处理方法: Binarizer, 测试集得分: 0.9139

```

【小结】

从以上结果来看，六种预处理方法都不如未经过预处理的结果。因此，后续的实验我们将继续采用未经过预处理数据完成。

5. 使用PCA进行降维，并测试性能

```

1 from sklearn.decomposition import PCA# 导入多层感知机MLP神经网络
2
3 #设置主成分数量为2以便我们进行可视化
4 pca = PCA(n_components=0.9)
5 pca.fit(X_train_lite)
6 X_train_pca = pca.transform(X_train_lite)
7 X_test_pca = pca.transform(X_test_lite)
8 print(X_train_lite.shape, X_test_lite.shape)
9 print(X_train_pca.shape, X_test_pca.shape)

```

```
1 | (5999, 784) (999, 784)
2 | (5999, 84) (999, 84)
```

```
1 | from sklearn.neural_network import MLPClassifier
2 | import time
3 |
4 | start = time.time()
5 |
6 | mlp = MLPClassifier(solver='lbfgs', hidden_layer_sizes=[100, 100],
7 | activation='relu', alpha=1e-5, random_state=62, verbose=2)
8 | mlp.fit(X_train_pca, y_train_lite)
9 |
10 | score_train = mlp.score(X_train_pca, y_train_lite)
11 | score_test = mlp.score(X_test_pca, y_test_lite)
12 |
13 | print('训练结束, 用时{:.2f}s.'.format(time.time() - start))
14 | print('训练集得分: {:.4f}, 测试集得分: {:.4f}'.format(score_train, score_test))
```

```
1 | 训练结束, 用时1.39s.
2 | 训练集得分: 1.0000, 测试集得分: 0.9309
```

6. 按要求计算范围内的PCA降维后的性能

```
1 | from sklearn.neural_network import MLPClassifier
2 | import numpy as np
3 | import time
4 |
5 | num=20
6 | scores = np.zeros([3,num])
7 | scores[0,:] = np.linspace(0.40, 0.99, num)
8 |
9 | start = time.time()
10 |
11 | # 基于信息量百分比
12 | n = 0
13 | for i in scores[0,:]:
14 |     # TODO: 3. 进行PCA降维
15 |     pca = PCA(n_components = i)
16 |     pca.fit(X_train_lite)
17 |     X_train_pca = pca.transform(X_train_lite)
18 |     X_test_pca = pca.transform(X_test_lite)
19 |
20 |     # TODO: 4. 训练MLP模型
21 |     mlp = MLPClassifier(solver='lbfgs', hidden_layer_sizes=[100, 100],
22 | activation='relu', alpha=1e-5, random_state=62, verbose=2)
23 |     mlp.fit(X_train_pca, y_train_lite)
24 |
25 |     score_train = mlp.score(X_train_pca, y_train_lite)
26 |     score_test = mlp.score(X_test_pca, y_test_lite)
27 |     scores[1, n] = score_train
28 |     scores[2, n] = score_test
29 |
30 |     t = time.time() - start
```

```

30     print("n_components={0:.2f}: 训练集得分 {1:.4f}, 测试集得分 {2:.4f}, t =
      {3:.2f}".format(i, score_train, score_test, t))
31
32     n = n + 1
33

```

```

1  n_components=0.40: 训练集得分 0.9948, 测试集得分 0.8368, t = 7.12
2  n_components=0.43: 训练集得分 1.0000, 测试集得分 0.8519, t = 13.99
3  n_components=0.46: 训练集得分 1.0000, 测试集得分 0.8689, t = 20.57
4  n_components=0.49: 训练集得分 1.0000, 测试集得分 0.8899, t = 26.59
5  n_components=0.52: 训练集得分 1.0000, 测试集得分 0.8989, t = 30.50
6  n_components=0.56: 训练集得分 1.0000, 测试集得分 0.9119, t = 33.75
7  n_components=0.59: 训练集得分 1.0000, 测试集得分 0.9109, t = 36.63
8  n_components=0.62: 训练集得分 1.0000, 测试集得分 0.9219, t = 39.38
9  n_components=0.65: 训练集得分 1.0000, 测试集得分 0.9269, t = 41.85
10 n_components=0.68: 训练集得分 1.0000, 测试集得分 0.9229, t = 44.80
11 n_components=0.71: 训练集得分 1.0000, 测试集得分 0.9389, t = 47.59
12 n_components=0.74: 训练集得分 1.0000, 测试集得分 0.9269, t = 49.92
13 n_components=0.77: 训练集得分 1.0000, 测试集得分 0.9359, t = 52.23
14 n_components=0.80: 训练集得分 1.0000, 测试集得分 0.9289, t = 54.58
15 n_components=0.83: 训练集得分 1.0000, 测试集得分 0.9309, t = 56.64
16 n_components=0.87: 训练集得分 1.0000, 测试集得分 0.9259, t = 58.68
17 n_components=0.90: 训练集得分 1.0000, 测试集得分 0.9269, t = 60.74
18 n_components=0.93: 训练集得分 1.0000, 测试集得分 0.9259, t = 62.83
19 n_components=0.96: 训练集得分 1.0000, 测试集得分 0.9279, t = 65.00
20 n_components=0.99: 训练集得分 1.0000, 测试集得分 0.9209, t = 67.34

```

7.可视化性能

```

1  import matplotlib.pyplot as plt
2
3  plt.rcParams['font.sans-serif'] = [u'Microsoft YaHei']
4
5  plt.plot(scores[0,:], scores[1,:], 'g--^', label='训练集 (信息量)')
6  plt.plot(scores[0,:], scores[2,:], 'b--o', label='测试集 (信息量)')
7  plt.legend(['训练集', '测试集'], loc='best')
8  plt.title('基于信息量百分比的性能对比图')
9
10 max_y = np.round(np.max(scores[2,:]), 4)
11 max_index = np.where(scores[2,:]==np.max(scores[2,:]))
12 max_x = scores[0, max_index]
13
14 score_ori_test_y = np.round(score_ori_test, 4)
15 score_ori_text_x = 1
16
17 plt.plot(max_x, max_y, 'ks')
18 plt.annotate(max_y,
19              xy = (max_x, max_y),
20              xytext = (max_x+0.01, max_y))
21 plt.plot(score_ori_text_x, score_ori_test_y, 'ks')
22 plt.annotate(score_ori_test_y,
23              xy = (score_ori_text_x, score_ori_test_y),
24              xytext = (score_ori_text_x+0.01, score_ori_test_y))
25 plt.show()

```

1 | Text(1.01, 0.9239, '0.9239')

