

# 模式识别课后作业五

2020E8013282019 李一帆

January 8, 2021

## 1 简述题

### 1.1 问题 1

**问题描述：**请简述 adaboost 算法的设计思想和主要计算步骤。

**解答：**

**设计思想：**Adaboost 从弱学习算法出发，得到一系列弱分类器；然后通过组合这些弱分类器构成一个强分类器。通过改变训练数据的概率分布，针对不同的训练数据的分布，调用弱学习算法来学习一些列分类器。在改变训练数据的权值分布上，Adaboost 通过提高那些被前一轮弱分类器分错的样本的权重，降低已经被正确分类的样本的权重。错分的样本将在下一轮弱分类器中得到更多的关注。于是分类问题被一系列弱分类器“分而治之”。在如何将一系列的弱分类器组合成一个强分类器方面，Adaboost 采用加权（多数）表决的方法。加大分类错误率较小的弱分类器的权重，使其在表决中起更大的作用。

**计算步骤：**输入为给定两类训练数据集  $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ 。

Step1: 初始化训练数据的权值分布  $D_1 = \{w_{11}\}, \{w_{12}\}, \dots, \{w_{1n}\}, w_{1i} = 1/n, i = 1, \dots, n$ 。

Step2: 对  $m = 1, 2, \dots, M$

(a) 使用具有权值分布  $D_m$  的训练数据，学习基本分类器  $G_m(x) : X \rightarrow \{-1, +1\}$ 。

(b) 计算  $G_m(x)$  在训练数据集上的分类错误率（加权）：

$$e_m = P(G_m(x_i) \neq y_i) = \sum_{i=1}^n w_{mi} I(G_m(x_i) \neq y_i)$$

(c) 计算  $G_m(x)$  的贡献系数：

$$\alpha_m = \frac{1}{2} \ln \frac{1 - e_m}{e_m}$$

$\alpha_m$  表示  $G_m(x)$  在最终的分类器中的重要性。当  $e_m \leq 0.5$  时， $\alpha_m \geq 0$ 。同时， $\alpha_m$  将随着  $e_m$  的减小而增大。所以分类误差率越小的基本分类器在最终分类器的作用越大。

(d) 更新训练数据集的权重分布：

$$D_{m+1} = \{w_{m+1,1}, w_{m+1,2}, \dots, w_{m+1,n}\}$$

具体计算如下：

$$\begin{aligned}w_{m+1,i} &= \frac{w_{mi}}{Z_m} \times \begin{cases} \exp(-\alpha_m), & \text{if } G_m(x_i) = y_i \\ \exp(\alpha_m), & \text{if } G_m(x_i) \neq y_i \end{cases} \\&= \frac{w_{mi}}{Z_m} \times \exp(-\alpha_m y_i G_m(x_i))\end{aligned}$$

其中， $Z_m$  是规范化因子，它使  $D_{m+1}$  成为一个概率分布。

$$Z_m = \sum_{i=1}^n w_{mi} \exp(-\alpha_m y_i G_m(x_i))$$

Step3: 构建基本分类器的线性组合。

$$f(x) = \sum_{m=1}^M \alpha_m G_m(x)$$

对于两类分类问题，得到最终的分类器。

$$G(x) = \text{sign}(f(x)) = \text{sign}\left(\sum_{m=1}^M \alpha_m G_m(x)\right)$$

## 1.2 问题 2

**问题描述：**请从混合高斯密度函数估计的角度，简述 k-means 聚类算法的原理（主要用文字描述）；请给出 k-means 聚类算法的计算步骤；请说明哪些因素会影响 k-means 算法的聚类性能。

**解答：**k-means 假设各个聚类出现的概率相等，每个样本以概率 1 属于每个聚类，样本属于哪一类需要计算  $\|x_k - \hat{\mu}_i\|^2$  来判断。因此需要通过迭代来得到  $c$  个高斯成分的均值。

k-means 聚类的计算步骤如下所示。

---

**Algorithm 1:** k-means Algorithm

---

```
1 Function k-means(  $W, k$ ):  
2   begin initialization  $n, c, \mu_1, \mu_2, \dots, \mu_c$ ;  
3   while no change in  $\mu_i$  do  
4     classify  $n$  samples according to nearest  $\mu_i$ ;  
5     re-compute  $\mu_i$ ;  
6   return  $\mu_1, \mu_2, \dots, \mu_c$ 
```

---

影响 k-means 算法的聚类性能的因素主要有聚类数目  $k$ ，初始中心  $\mu_i$  的选取以及聚类数据的分布。

## 1.3 问题 3

**问题描述：**请简述谱聚类算法的原理，给出一种谱聚类算法（经典算法、Shi 算法和 Ng 算法之一）的计算步骤；请指出哪些因素会影响聚类的性能。

**解答：**从图切割的角度，谱聚类就是要找到一种合理的分割图的方法，分割后能形成若干个子图。连接不同子图的边的权重尽可能小，子图内部边权重尽可能大。Ng 聚类算法的原理如 Algorithm 1 所示。

---

**Algorithm 2:** Spectral Clustering-Ng Algorithm

---

**Input:** *similarity matrix*  $W$

**Input:**  $k$  number of clusters

**Output:**  $A_1, \dots, A_k$  clusters

1 **Function** NgSpectralCluster( $W, k$ ):

2     *compute*  $\mathbf{L}_{sym} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$ ;

3     *compute the first  $k$  eigenvectors*  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$  *of*  $\mathbf{L}_{sym}$ ;

4     Let  $\mathbf{U} \in \mathbb{R}^{n \times k}$  be the matrix containing the vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ , namely,

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in \mathbb{R}^{n \times k};$$

5      $max\_time = 0$ ;

6     form the matrix  $\mathbf{T} \in \mathbb{R}^{n \times k}$  from  $\mathbf{U}$  by normalizing the rows to norm 1, namely, set

$$t_{ij} = u_{ij} / \sqrt{\sum_{m=1}^k u_{im}^2};$$

7     for  $i = 1, 2, \dots, n$ , let  $y_i \in \mathbb{R}^k$  be the vector corresponding to the  $i$ -th row of  $\mathbf{T}$ ;

8     cluster the points  $\{y_i\}_{i=1,2,\dots,n}$  in  $\mathbb{R}^k$  with  $k$ -means algorithm into clusters

$$A_1, A_2, \dots, A_k;$$

9     output  $A_1, A_2, \dots, A_k$

---

影响谱聚类性能的因素主要有一下几点。局部连接  $k$  近邻的取值以及  $\epsilon$  半径。如果样本点数目很多，对大型矩阵进行特征值分解的计算仍然不稳定。最后采用  $k$ -means 进行聚类，聚类数目也会影响聚类结果。

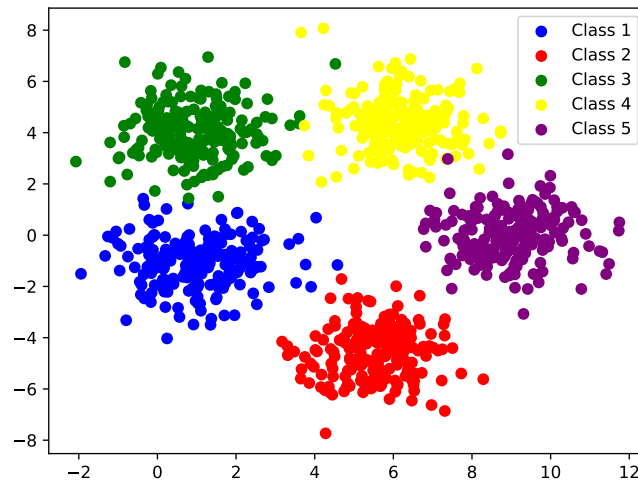


Figure 1: 聚类数据可视化.

## 2 计算机编程

### 2.1 问题 1

将数据点进行可视化后的结果如 Figure 1 所示。

问题描述：

(1) 编写一个程序，实现经典的 k-means 聚类算法。

(2) 令聚类个数等于 5，采用不同的初始值，报告聚类精度、以及最后获得的聚类中心。

并计算所获得的聚类中心与对应的真实分布的均值之间的误差。

解答： k-means 聚类算法的代码如下所示。

```
1 def k_means(data, centers):
2     '''
3     K_means algorithm.
4     :param data: The data that needs to be clustered.
5     :param centers: The centers of cluster.
6     :return: The cluster result and centers.
7     '''
8     N = data.shape[0]
9     class_num = centers.shape[0]
10    standard_ans = np.zeros(N)
11    per_class = int(N / class_num)
12    for i in range(class_num):
13        standard_ans[i * per_class:(i + 1) * per_class] = i
14
15    for step in range(10000):
16        # caculate the distance between cluster centers and data
17        dis = np.zeros([data.shape[0], class_num])
18        for i in range(class_num):
19            data_norm2 = np.linalg.norm(data - centers[i], ord=2,
20                                         axis=1)
21            dis[:, i] = data_norm2
22
23        # find the min index of distance matrix
24        result = np.argmin(dis, axis=1)
25        centers_tmp = np.zeros(centers.shape, dtype='float')
26        for i in range(class_num):
```

```

26         class_inx = np.argwhere(result == i)
27         if class_inx.all():
28             centers_tmp[i] = np.mean(data[class_inx], axis=0)
29     error_matrix = (standard_ans - result)
30     error_matrix[error_matrix < 0] = 1
31     error_matrix[error_matrix > 0] = 1
32     right_num = error_matrix.sum()
33     print("Step %d , acc is : %f" % (step, (N - right_num) / N)
34         )
35     if abs(centers - centers_tmp).sum() == 0:
36         break
37     else:
38         centers = centers_tmp
39     return result, centers

```

采用不同的初始值（随机选择每一聚类中数据点），最终的聚类精度为 98.4%，最终聚类中心与真实分布均值之间的误差如下所示。

1.00000000	1.00000000
0.01610427	0.00131740
0.09496927	0.09183495
0.16539437	0.08684510
0.01165332	0.01481903

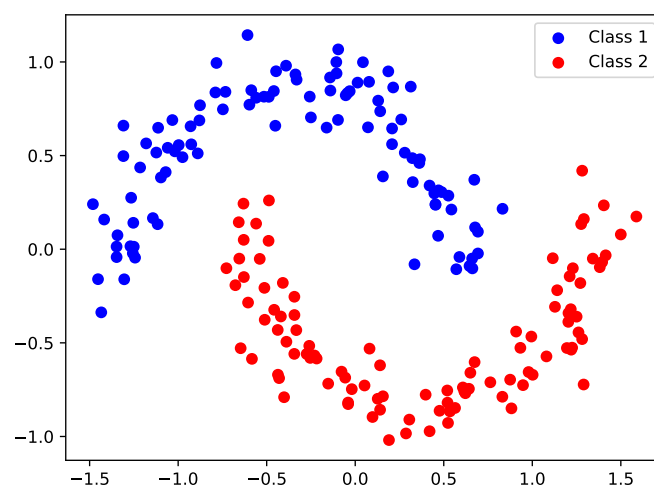


Figure 2: 谱聚类数据可视化.

## 2.2 问题 2

将数据点进行可视化后的结果如 Figure 2 所示。

**问题描述：**请编写一个程序，实现 Ng 谱聚类算法。

**解答：**Ng 算法代码如下所示。

```
1  def spectral_cluster(data, k=10, type='Ng', sim='classical',
2      sigma=5, cluster_num=2):
3      N = data.shape[0]
4      W = np.zeros([N, N])
5      for i in range(N):
6          dis = np.linalg.norm(data - data[i], ord=2, axis=1)
7          dis = np.delete(dis, 0)
8          k_idx = np.argsort(dis)[: k]
9          if sim == 'classical':
10             W[i, k_idx] = 1
11         else:
12             W[i, k_idx] = np.exp(-(dis[k_idx] ** 2) / (2 * (sigma
13                 ** 2)))
14     if sim == 'classical':
15         D = k * np.eye(N)
16         # W = (W.T + W) / 2
17         L = D - W
18     else:
19         D_diag = np.sum(W, axis=0)
20         D = np.diag(D_diag)
21         W = (W.T + W) / 2
22         L = D - W
23     if type == 'Ng':
24         L_sys = ((D * (k ** (-0.5))).dot(L)).dot(D * (k ** (-0.5)))
25         vals, vecs = np.linalg.eig(L_sys)
26         indx = np.argsort(vals)[:cluster_num]
27         U = vecs[:, indx]
28         T = U / np.sqrt(np.sum(U ** 2, axis=0))
29         centers = np.array([T[60], T[150]])
30         result, centers, acc = k_means(T, centers)
31     return result, centers, acc
```

通过 k-means 聚类和谐聚类的结果如下所示。

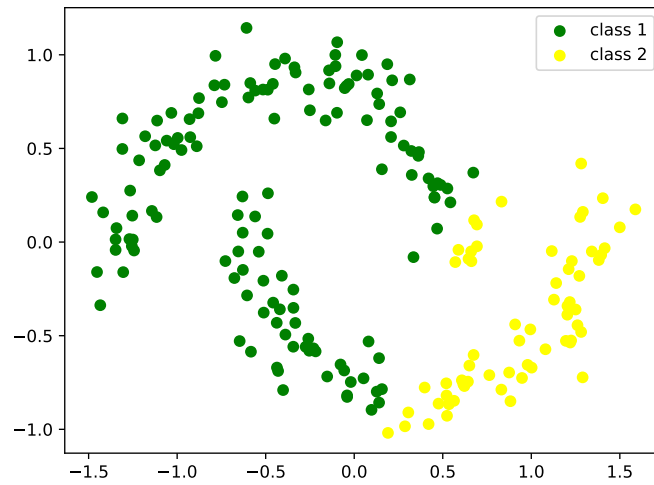


Figure 3: k-means 聚类结果可视化.

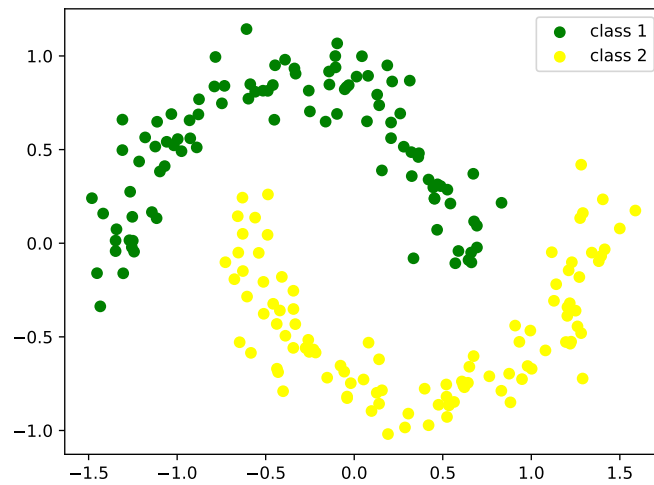


Figure 4: Ng 谱聚类结果可视化.

由上图可以看出，k-means 对于这种分布不服从高斯类型的聚类效果不是很好，而谱聚类则可以很好的解决这个问题。

**问题描述：**设点对亲和性（即边权值）采用如下计算公式。

$$w_{ij} = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right)$$

数据图采用 k-近邻的方法来生成（即对于每个数据点  $x_i$ ，首先在所有样本中找出不包含  $x_i$  的 k 个最近邻的样本点，然后  $x_i$  与每个临近样本点均有一条边相连，从而完成图构造）。注

意，为了保证亲和度矩阵  $W$  是对称矩阵，可以令  $W = (W^T + W)/2$ 。假设已知前 100 个点为一个聚类，后 100 个点为一个聚类，请分析分别取不同的  $\sigma$  值和  $k$  值对聚类结果的影响。

**解答：**在固定  $\sigma = 5$  下，改变  $k$  可以得到的聚类精度变化曲线如 Figure 5 所示。在固定  $k = 10$  下，改变  $\sigma$  可以得到聚类精度变化曲线如 Figure 6 所示。

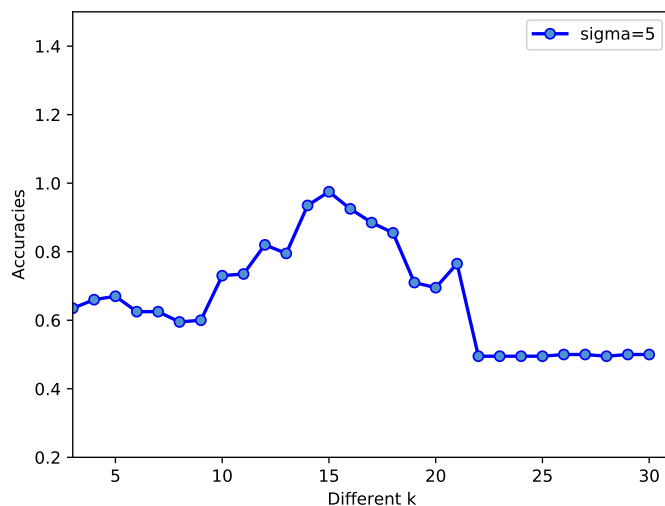


Figure 5: 聚类精度变化曲线 ( $\sigma=5$ ,  $k$  改变) .

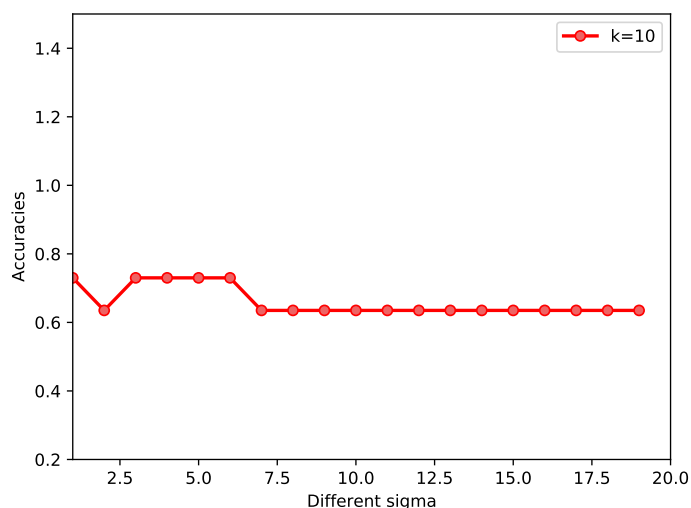


Figure 6: 聚类精度变化曲线 ( $k=10$ ,  $\sigma$  改变) .

由上图可以看出，在固定  $\sigma$  的情况下，不同的  $k$  值会对聚类精度产生不同的影响，并且  $k$  值的选取往往并不是越大越好或是越小越好，最优取值是介于中间的某一个数。在固定  $k$  值的情况下，不同  $\sigma$  的选取也会对聚类精度产生影响，但是影响不如  $k$  值那么明显。



## A 附录——代码 (Python)

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  data_name1 = 'data1.txt'
5  data_name2 = 'data2.txt'
6
7
8  def read_data(data_file):
9      data = []
10     split_string = ' '
11     if data_file == 'data1.txt':
12         split_string = '\t'
13     elif data_file == 'data2.txt':
14         split_string = ' '
15     with open(data_file) as file:
16         for line in file.readlines():
17             data_line = np.array(line.split(split_string),
18                                   dtype='float')
19             data.append(data_line)
20     data = np.array(data)
21     return data
22
23  def k_means(data, centers, print_info=False):
24      '''
25      K_means algorithm.
26      :param data: The data that needs to be clustered.
27      :param centers: The centers of cluster.
28      :return: The cluster result and centers.
29      '''
30      N = data.shape[0]
31      class_num = centers.shape[0]
32      standard_ans = np.zeros(N)
33      per_class = int(N / class_num)
```

```

34     for i in range(class_num):
35         standard_ans[i * per_class:(i + 1) * per_class] = i
36
37     for step in range(10000):
38         # caculate the distance between cluster centers and
           data
39         dis = np.zeros([data.shape[0], class_num])
40         for i in range(class_num):
41             data_norm2 = np.linalg.norm(data - centers[i], ord
           =2, axis=1)
42             dis[:, i] = data_norm2
43
44         # find the min index of distance matrix
45         result = np.argmin(dis, axis=1)
46         centers_tmp = np.zeros(centers.shape, dtype='float')
47         for i in range(class_num):
48             class_inx = np.argwhere(result == i)
49             if class_inx.all():
50                 centers_tmp[i] = np.mean(data[class_inx], axis
           =0)
51         error_matrix = (standard_ans - result)
52         error_matrix[error_matrix < 0] = 1
53         error_matrix[error_matrix > 0] = 1
54         right_num = error_matrix.sum()
55         acc = (N - right_num) / N
56         if print_info:
57             print("Step %d , acc is : %f" % (step, acc))
58         if abs(centers - centers_tmp).sum() == 0:
59             break
60         else:
61             centers = centers_tmp
62     return result, centers, acc
63
64
65 def spectral_cluster(data, k=10, type='Ng', sim='classical',
           sigma=5, cluster_num=2):

```

```

66     N = data.shape[0]
67     W = np.zeros([N, N])
68     for i in range(N):
69         dis = np.linalg.norm(data - data[i], ord=2, axis=1)
70         dis = np.delete(dis, 0)
71         k_idx = np.argsort(dis)[: k]
72         if sim == 'classical':
73             W[i, k_idx] = 1
74         else:
75             W[i, k_idx] = np.exp(-(dis[k_idx] ** 2) / (2 * (
76                 sigma ** 2)))
77     if sim == 'classical':
78         D = k * np.eye(N)
79         # W = (W.T + W) / 2
80         L = D - W
81     else:
82         D_diag = np.sum(W, axis=0)
83         D = np.diag(D_diag)
84         W = (W.T + W) / 2
85         L = D - W
86     if type == 'Ng':
87         L_sys = ((D * (k ** (-0.5))).dot(L)).dot(D * (k **
88             (-0.5)))
89         vals, vecs = np.linalg.eig(L_sys)
90         indx = np.argsort(vals)[:cluster_num]
91         U = vecs[:, indx]
92         T = U / np.sqrt(np.sum(U ** 2, axis=0))
93         centers = np.array([T[60], T[150]])
94         result, centers, acc = k_means(T, centers)
95     return result, centers, acc
96
97 if __name__ == '__main__':
98     # data1 = read_data(data_name1)
99     #####True clusters#####
100     # plt.figure()

```

```

100 # plt.scatter(data1[:200, 0], data1[:200, 1], color='blue',
101               label='Class 1')
102 # plt.scatter(data1[200:400, 0], data1[200:400, 1], color='
103               red', label='Class 2')
104 # plt.scatter(data1[400:600, 0], data1[400:600, 1], color='
105               green', label='Class 3')
106 # plt.scatter(data1[600:800, 0], data1[600:800, 1], color='
107               yellow', label='Class 4')
108 # plt.scatter(data1[800:1000, 0], data1[800:1000, 1], color
109               ='purple', label='Class 5')
110 # plt.legend()
111 #####K-means#####
112 # orig_clusters = np.array([data1[150], data1[350], data1
113                             [550], data1[750], data1[950]], dtype='float')
114 # result, centers, acc = k_means(data1, centers=
115                                 orig_clusters)
116 # plt.figure()
117 # colors = ['blue', 'red', 'green', 'yellow', 'purple']
118 # classes = ['class 1', 'class 2', 'class 3', 'class 4', '
119              class 5']
120 # for i in range(5):
121 #     class_inx = np.argwhere(result == i)
122 #     plt.scatter(data1[class_inx, 0], data1[class_inx, 1],
123                 color=colors[i], label=classes[i])
124 # plt.legend()
125 # centers_mean = np.array([[1, -1], [5.5, -4.5], [1, 4],
126                             [6, 4.5], [9, 0]])
127 # error = abs(orig_clusters - centers_mean).sum()
128 # print("Error between k-means centers and true centers is
129         %f" % error)
130 # print(abs(centers - centers_mean))
131 #
132 #####True Clusters#####
133 plt.figure()
134 data2 = read_data(data_name2)

```

```

124     plt.scatter(data2[:100, 0], data2[:100, 1], color='blue',
125                label='Class 1')
126     plt.scatter(data2[100:200, 0], data2[100:200, 1], color='
127                red', label='Class 2')
128     plt.legend()
129     # np.random.shuffle(data2)
130     # orig_clusters2=np.array([data2[50], data2[150]])
131     # result, centers = k_means(data2, centers=orig_clusters2)
132     # result, centers, acc = spectral_cluster(data2, sim='
133     others')
134     # colors = ['green', 'yellow']
135     # classes = ['class 1', 'class 2']
136     # plt.figure()
137     # for i in range(2):
138     #     class_inx = np.argwhere(result == i)
139     #     plt.scatter(data2[class_inx, 0], data2[class_inx, 1],
140     #                 color=colors[i], label=classes[i])
141     # plt.legend()
142     # plt.show()
143
144     #####分析谱聚类参数的影响#####
145     sigmas = range(1, 20)
146     ks = range(3, 31)
147     accuracies = np.zeros([len(sigmas), len(ks)])
148     for sigma in sigmas:
149         for k in ks:
150             result, centers, acc = spectral_cluster(data2, sim=
151             'others' ,sigma=sigma, k=k)
152             accuracies[sigma - 1, k - 3] = acc
153     plt.figure()
154     sigma = sigmas[5]
155     plt.plot(ks, accuracies[sigma - 1], marker='o', mfc='#4
156             F94CD', label='sigma=5', color='blue', linewidth=2)
157     plt.xlim((3, 31))
158     plt.ylim((0.2, 1.5))
159     plt.xlabel('Different k')

```

```
154 plt.ylabel('Accuracies')
155 plt.legend()
156 plt.figure()
157 k = 10
158 plt.plot(sigmas, accuracies[:, k - 3], marker='o', mfc='#
    EE6363', label='k=10', color='red', linewidth=2)
159 plt.xlim((1, 20))
160 plt.ylim((0.2, 1.5))
161 plt.xlabel('Different sigma')
162 plt.ylabel('Accuracies')
163 plt.legend()
164 plt.show()
```