



南开大学  
Nankai University

南 开 大 学

计 算 机 学 院

机器学习实验报告

---

Lab5 聚类分析

---

马湔怡

学号：2311061

专业：计算机科学与技术

2025 年 12 月 4 日

# 目录

<b>一、 实验原理</b>	<b>1</b>
(一) 层次聚类的基本原理 . . . . .	1
(二) 簇间距离度量 . . . . .	1
1. Single-linkage (单链接 / 最小距离) . . . . .	1
2. Complete-linkage (全链接 / 最大距离) . . . . .	1
3. Average-linkage (平均链接 / 平均距离) . . . . .	1
<b>二、 实验内容</b>	<b>2</b>
(一) 基础任务: Single-linkage 与 Complete-linkage 聚类 . . . . .	2
1. 代码 . . . . .	2
2. 控制台输出结果 . . . . .	3
3. 结果分析 . . . . .	4
(二) 中级任务: Average-linkage 聚类 . . . . .	5
1. 代码 . . . . .	5
2. 控制台输出结果 . . . . .	6
3. 结果分析 . . . . .	6
(三) 提高任务: 算法对比 . . . . .	7
1. 控制台输出结果 . . . . .	7
2. 结果分析 . . . . .	8
(四) 拓展任务: 变换聚类簇数 $k$ . . . . .	10
1. 控制台输出结果 . . . . .	10
2. 结果分析 . . . . .	12

## 一、实验原理

### (一) 层次聚类的基本原理

层次聚类是一种无监督学习方法，用于将数据集中的样本点组织成一个树状结构。它不要求事先指定聚类簇的数量  $k$ ，而是通过迭代地合并或分裂簇来形成层次结构。

本次实验采用**凝聚**方法，其基本步骤如下：

1. **初始化**：将每个样本点视为一个独立的簇。
2. **迭代合并**：重复以下步骤，直到所有样本点合并成一个大簇（或达到预设的簇数  $k$ ）：
  - 计算所有当前簇对之间的距离。
  - 将距离最近的两个簇合并成一个新的簇。
3. **最终划分**：根据用户设定的目标簇数  $k$  或在谱系图上“剪枝”来获得最终的聚类结果。

### (二) 簇间距离度量

层次聚类算法的核心在于定义两个簇  $C_i$  和  $C_j$  之间的距离  $d(C_i, C_j)$ 。本实验实现了三种主要的距离度量方式：

#### 1. Single-linkage (单链接 / 最小距离)

- **定义**：两个簇之间的距离定义为**两个簇中最接近的样本点**之间的距离。
- **公式**：

$$d(C_i, C_j) = \min\{d(a, b) \mid a \in C_i, b \in C_j\}$$

- **特点**：倾向于产生长链状的聚类。它对噪声和离群点高度敏感，容易产生链式效应。

#### 2. Complete-linkage (全链接 / 最大距离)

- **定义**：两个簇之间的距离定义为**两个簇中最远的样本点**之间的距离。
- **公式**：

$$d(C_i, C_j) = \max\{d(a, b) \mid a \in C_i, b \in C_j\}$$

- **特点**：倾向于产生紧凑的球形聚类。它对离群点具有较好的鲁棒性。

#### 3. Average-linkage (平均链接 / 平均距离)

- **定义**：两个簇之间的距离定义为**两个簇中所有样本点对之间距离的平均值**。
- **公式**：

$$d(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{a \in C_i} \sum_{b \in C_j} d(a, b)$$

- **特点**：该方法综合了 Single-linkage 和 Complete-linkage 的特点，通常比前两者更稳健。

## 二、 实验内容

### (一) 基础任务：Single-linkage 与 Complete-linkage 聚类

#### 1. 代码

##### 基础任务

```

1  def _single_linkage_distance(self, cluster1, cluster2, dist_matrix):
2      """
3      Single-linkage: 两个簇之间的【最小】距离
4
5      公式:  $d(C1, C2) = \min\{d(a,b) \mid a \in C1, b \in C2\}$ 
6
7      输入:
8          cluster1: 第一个簇的样本索引列表, 如 [0, 3, 5]
9          cluster2: 第二个簇的样本索引列表, 如 [1, 2]
10         dist_matrix: 距离矩阵, shape=(n_samples, n_samples)
11
12     输出:
13         两个簇之间的最小距离
14
15     TODO: 请完成single-linkage距离计算
16     提示:
17     1. 遍历cluster1中的每个样本i
18     2. 遍历cluster2中的每个样本j
19     3. 使用dist_matrix[i, j]获取i和j的距离
20     4. 返回所有距离中的最小值
21     """
22     min_dist = np.inf
23
24     # TODO: 遍历两个簇中的所有样本对, 找最小距离
25     for i in cluster1:
26         for j in cluster2:
27             # TODO: 比较并更新最小距离
28             current_dist = dist_matrix[i, j]
29             if current_dist < min_dist:
30                 min_dist = current_dist
31
32     return min_dist
33
34 def _complete_linkage_distance(self, cluster1, cluster2, dist_matrix):
35     """
36     Complete-linkage: 两个簇之间的【最大】距离
37
38     公式:  $d(C1, C2) = \max\{d(a,b) \mid a \in C1, b \in C2\}$ 
39
40     输入:
41         cluster1: 第一个簇的样本索引列表

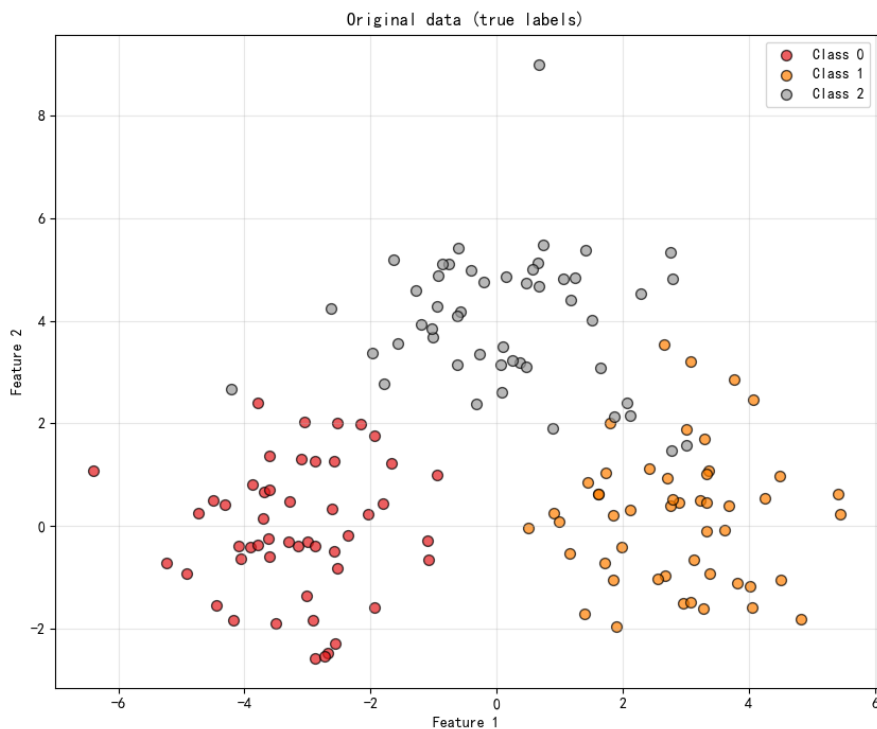
```

```
42         cluster2: 第二个簇的样本索引列表
43         dist_matrix: 距离矩阵
44
45     输出:
46         两个簇之间的最大距离
47     TODO: 请完成complete-linkage距离计算
48     提示: 与single-linkage类似, 但找最大值
49     """
50     max_dist = 0
51
52     # TODO: 遍历两个簇中的所有样本对, 找最大距离
53     for i in cluster1:
54         for j in cluster2:
55             # TODO: 比较并更新最大距离
56             current_dist = dist_matrix[i, j]
57             if current_dist > max_dist:
58                 max_dist = current_dist
59
60     return max_dist
```

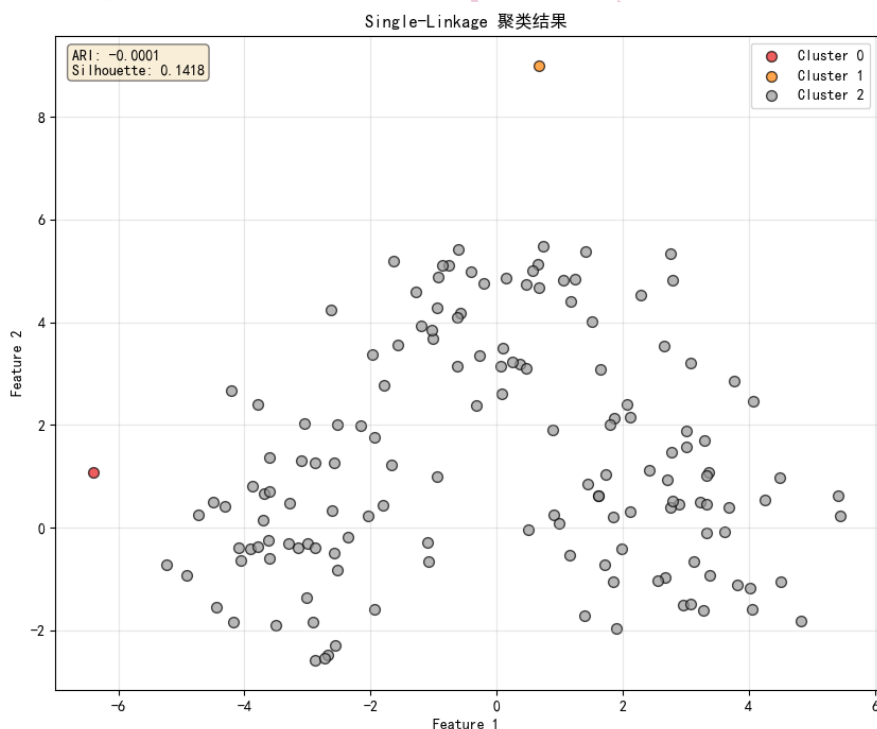
## 2. 控制台输出结果

```
1 [步骤2] 基本要求: Single-linkage 和 Complete-linkage
2
3 Single-linkage聚类...
4     ARI: -0.0001
5 图片已保存: out_exp5\single_linkage.png
6 Complete-linkage聚类...
7     ARI: 0.7481
8 图片已保存: out_exp5\complete_linkage.png
```

### 3. 结果分析

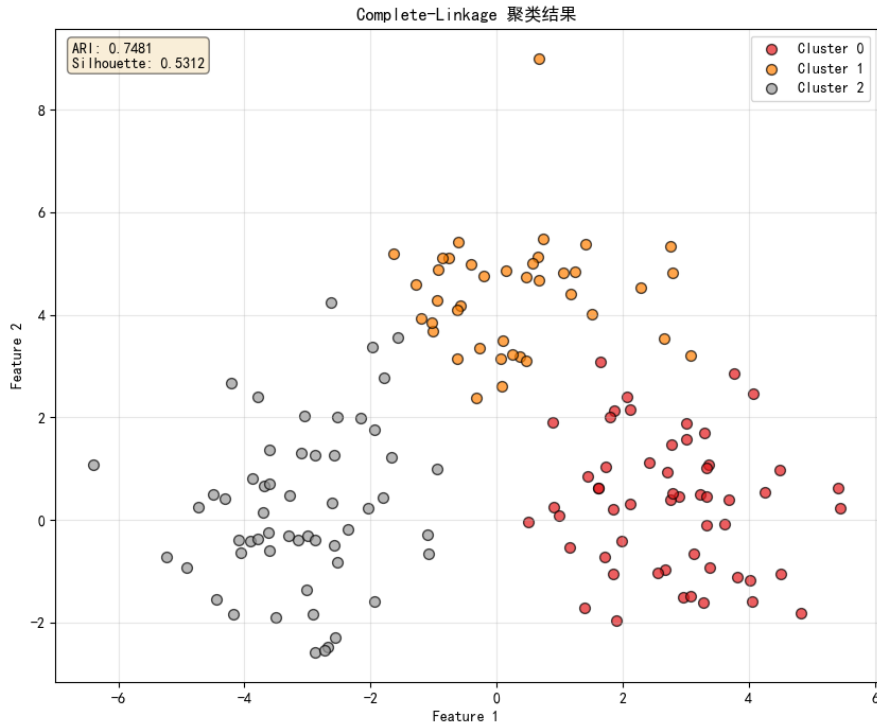


这是原始数据分布。可以看到数据集中有三个大致呈球形的簇（红色 Class 0、橙色 Class 1、灰色 Class 2），且簇之间有一定间隔。



ARI: -0.0001, Silhouette: 0.1418

可以证明这是个失败的聚类。图中显示大部分样本（灰色 Cluster 2）被合并成了一个巨大的、链状的簇，而另外两个簇（红色 Cluster 0 和橙色 Cluster 1）各只有一两个样本。这说明了 Single-linkage 的“链式效应”：由于它只关注簇间最近距离，只要两个点足够近，整个簇就会被连接，导致它无法正确分割出球形、密度均匀的簇。性能指标极差，证明它不适合这类数据。



ARI: 0.7481, Silhouette: 0.5312

可以证明这是个较好的聚类。Complete-linkage 倾向于产生紧凑的球形簇。图中它成功地将三个原始簇大致分离，划分的簇边界清晰。但与真实标签相比，它在簇的边界（例如红色和灰色簇的重叠区域）上有一些误分类，因此 ARI 略低于最优值。

## (二) 中级任务：Average-linkage 聚类

### 1. 代码

#### 中级任务

```

1 def _average_linkage_distance(self, cluster1, cluster2, dist_matrix):
2     """
3     Average-linkage: 两个簇之间的【平均】距离
4
5     公式:  $d(C1, C2) = (1/(|C1| \times |C2|)) \times \sum \sum d(a, b)$ 
6
7     输入:
8         cluster1: 第一个簇的样本索引列表
9         cluster2: 第二个簇的样本索引列表
10        dist_matrix: 距离矩阵
11
12    输出:
13        两个簇之间的平均距离
14    TODO: 请完成average-linkage距离计算
15    提示:
16        1. 累加所有样本对的距离
17        2. 除以样本对的数量 ( $|C1| \times |C2|$ )
18    """

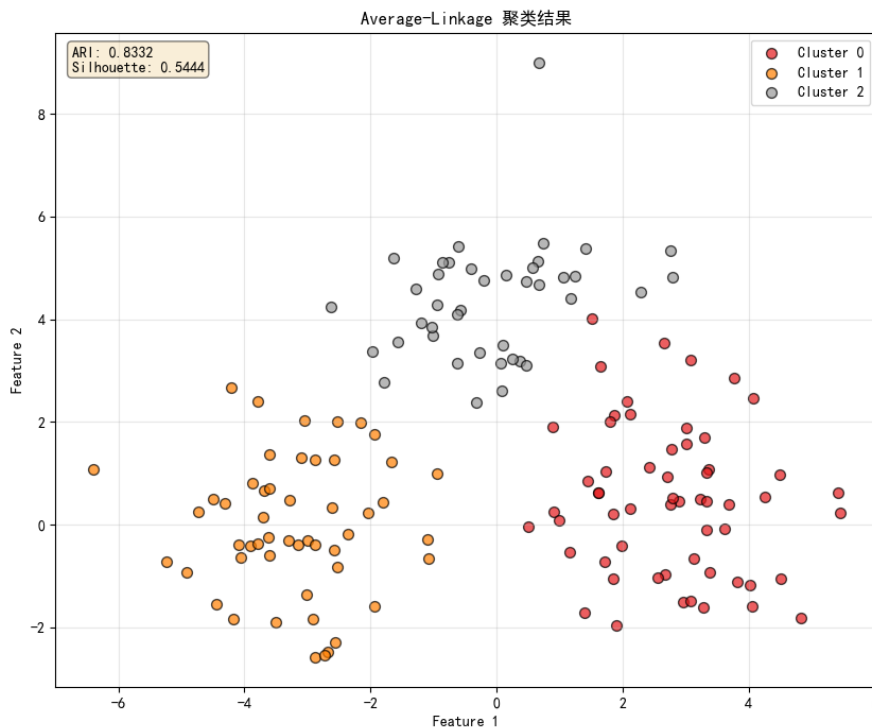
```

```
19     total_dist = 0
20     count = 0
21
22     # TODO: 遍历两个簇中的所有样本对，累加距离
23     for i in cluster1:
24         for j in cluster2:
25             # TODO: 累加距离并计数
26             total_dist += dist_matrix[i, j]
27             count += 1
28
29     # TODO: 返回平均距离
30     if count == 0:
31         return 0 # 避免除以零
32     return total_dist / count
```

## 2. 控制台输出结果

```
1 [步骤3] 中级要求: Average-linkage
2     ARI: 0.8332
3     图片已保存: out_exp5\average_linkage.png
```

## 3. 结果分析



ARI: 0.8332, Silhouette: 0.5444

相比基础任务的两个聚类，这是最佳聚类。Average-linkage 表现最好，其 ARI 和 Silhouette Score 均是三种方法中最高的。这表明在所有样本对的平均距离基础上进行合并，能够提供最稳健、最准确的划分。图中簇的划分比 Complete-linkage 更接近真实边界。



### (三) 提高任务：算法对比

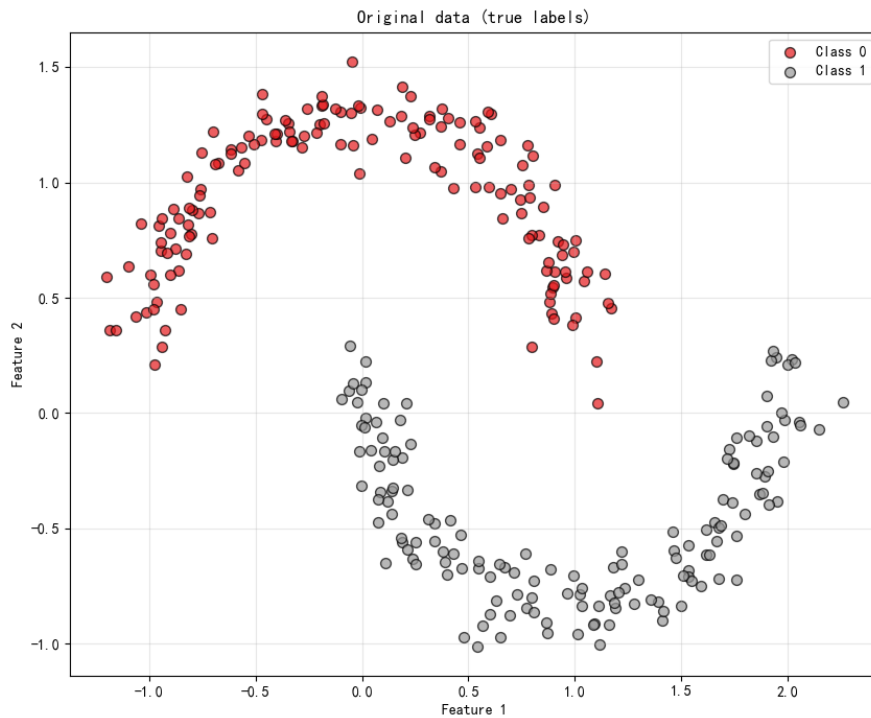
#### 1. 控制台输出结果

```

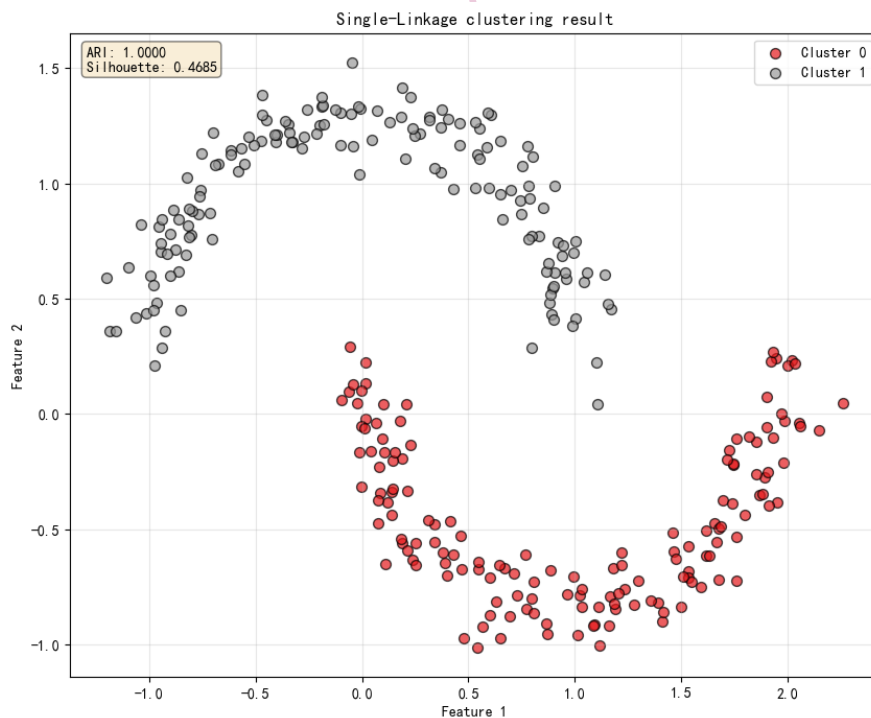
1 [步骤4] 提高要求：算法对比 (make_blobs 数据集)
2
3 [算法对比]
4
5 linkage          ARI          Silhouette
6
7 single           -0.0001       0.1418
8 complete         0.7481        0.5312
9 average          0.8332        0.5444
10
11
12 结论 (make_blobs 数据集):
13   - 在当前数据集上, average-linkage 表现最好
14   - Single-linkage: 容易产生链式效应
15   - Complete-linkage: 倾向于产生紧凑的簇
16   - Average-linkage: 两者的折中
17
18
19 [moons-2] 三种链接方式算法对比 (k=2, make_moons 数据集)
20
21 [算法对比]
22
23
24 linkage          ARI          Silhouette
25
26
27
28
29 single           1.0000        0.4685
30
31 图片已保存: out_exp5_moons\clustering_single.png
32
33 complete         0.6823        0.4887
34
35 图片已保存: out_exp5_moons\clustering_complete.png
36
37 average          0.5073        0.4760
38
39 图片已保存: out_exp5_moons\clustering_average.png

```

## 2. 结果分析

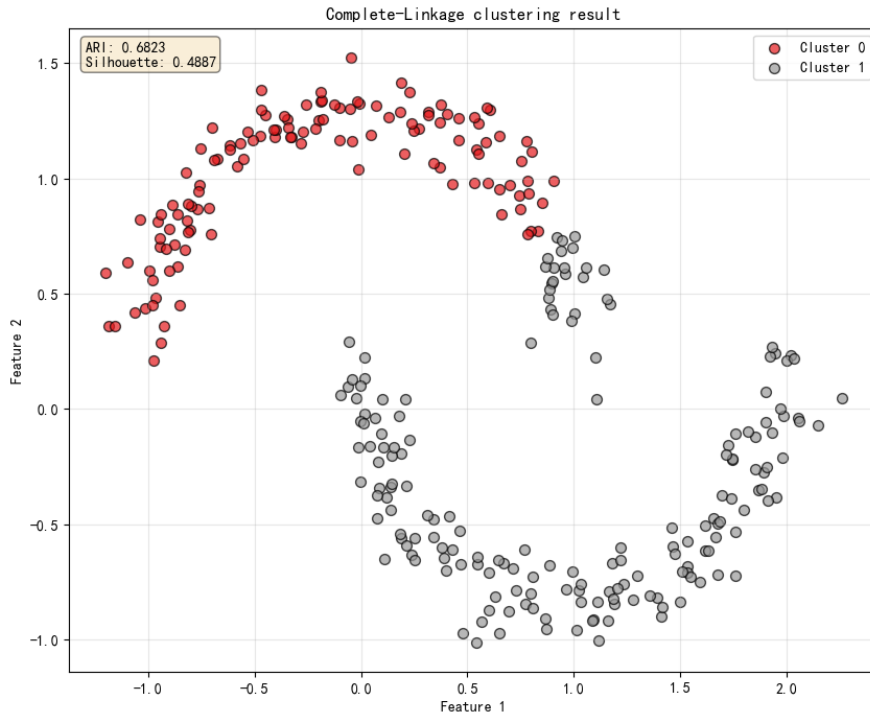


这是原始数据分布。可以看到数据集中包含两个稍微嵌入在一起的弯月形簇（红色 Class 0 和灰色 Class 1），这是典型的非凸形状数据。



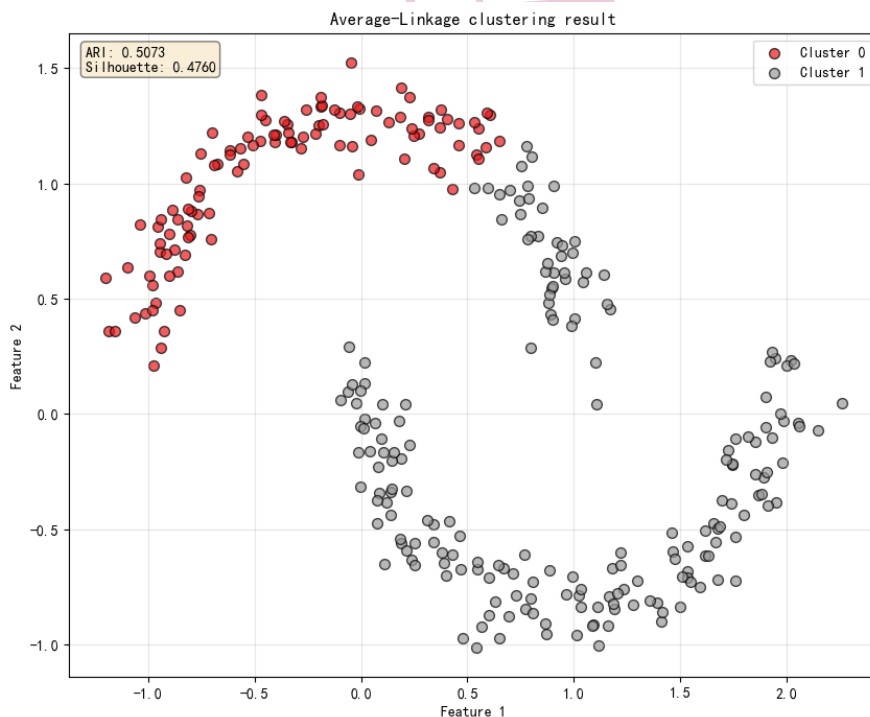
ARI: 1.0000, Silhouette: 0.4685

这是个满分的划分。Single-linkage 达到了  $ARI = 1.0000$ ，这说明聚类结果与真实标签完全匹配。对于非凸形状数据，Single-linkage 效果最佳，因为它只关注簇间最近距离，能够沿着弯月形的曲线进行合并，正确地将两个交织的簇分离。



ARI: 0.6823, Silhouette: 0.4887

这个划分表现一般。Complete-linkage 倾向于发现紧凑的球形簇，因此它在弯月形簇的中间部分进行了“垂直切断”，而不是沿着形状边界划分。结果是，两个簇都包含了来自原始两个弯月形簇的样本，准确度显著下降。



ARI: 0.5073, Silhouette: 0.4760 Average-linkage 的性能比 Complete-linkage 更差一些，结果接近随机划分（即 ARI 接近 0.5）。它处理非凸形状表现较差，错误地将靠近中心的样本混合，表明它更适合球形或凸形状簇。

## (四) 拓展任务：变换聚类簇数 k

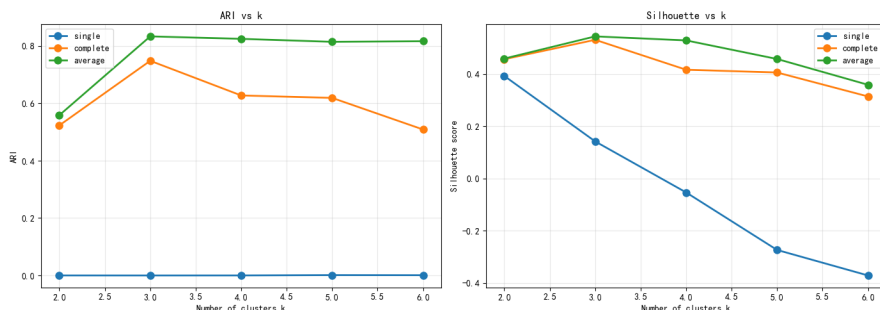
### 1. 控制台输出结果

```
1 [步骤5] 拓展要求：测试不同聚类数k (make_blobs 数据集)
2
3 [不同聚类数k的性能测试]
4
5 k=2:
6
7     single      : ARI=0.0000, Silhouette=0.3928
8
9     complete    : ARI=0.5229, Silhouette=0.4568
10
11    average     : ARI=0.5584, Silhouette=0.4593
12
13 k=3:
14
15    single      : ARI=-0.0001, Silhouette=0.1418
16
17    complete    : ARI=0.7481, Silhouette=0.5312
18
19    average     : ARI=0.8332, Silhouette=0.5444
20
21 k=4:
22
23    single      : ARI=0.0000, Silhouette=-0.0537
24
25    complete    : ARI=0.6273, Silhouette=0.4166
26
27    average     : ARI=0.8246, Silhouette=0.5290
28
29 k=5:
30
31    single      : ARI=0.0012, Silhouette=-0.2742
32
33    complete    : ARI=0.6188, Silhouette=0.4060
34
35    average     : ARI=0.8143, Silhouette=0.4577
36
37 k=6:
38
39    single      : ARI=0.0008, Silhouette=-0.3719
40
41    complete    : ARI=0.5086, Silhouette=0.3142
42
43    average     : ARI=0.8164, Silhouette=0.3588
44
45
```

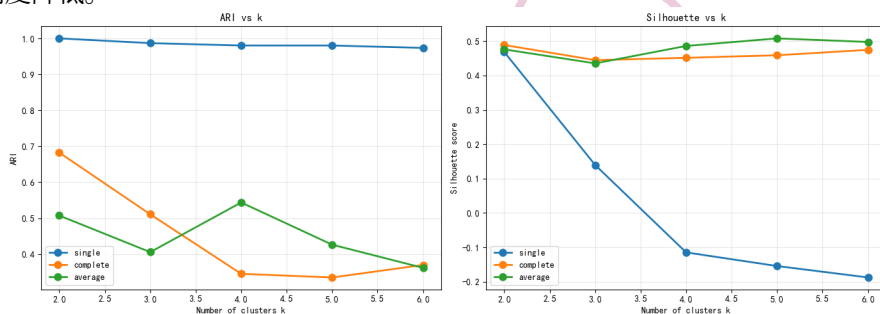
```
46
47 [moons-3] 拓展：测试不同聚类数k (make_moons 数据集)
48 [不同聚类数k的性能测试]
49
50 k=2:
51
52     single      : ARI=1.0000, Silhouette=0.4685
53
54     complete    : ARI=0.6823, Silhouette=0.4887
55
56     average     : ARI=0.5073, Silhouette=0.4760
57
58 k=3:
59
60     single      : ARI=0.9868, Silhouette=0.1394
61
62     complete    : ARI=0.5108, Silhouette=0.4447
63
64     average     : ARI=0.4057, Silhouette=0.4352
65
66 k=4:
67
68     single      : ARI=0.9802, Silhouette=-0.1148
69
70     complete    : ARI=0.3453, Silhouette=0.4515
71
72     average     : ARI=0.5437, Silhouette=0.4860
73
74 k=5:
75
76     single      : ARI=0.9802, Silhouette=-0.1545
77
78     complete    : ARI=0.3352, Silhouette=0.4589
79
80     average     : ARI=0.4264, Silhouette=0.5080
81
82 k=6:
83
84     single      : ARI=0.9736, Silhouette=-0.1877
85
86     complete    : ARI=0.3700, Silhouette=0.4747
87
88     average     : ARI=0.3619, Silhouette=0.4976
89
90 分析 (make_moons 数据集):
91
92     - make_moons 数据集为非凸形状簇，更考验聚类算法的鲁棒性
93
```

- 94 - 在 make\_moons 数据集上, single-linkage 的 ARI 最高
- 95
- 96 - 可以观察不同链接方式在非球形簇上的优劣

## 2. 结果分析



上图是三种链接方式在  $k = 2$  到  $k = 6$  范围内的 ARI 和 Silhouette Score 曲线对比。所有算法的 ARI 都在  $k = 3$  处达到峰值, 这证明了数据集中真实的簇数量是 3。当  $k \neq 3$  时, 性能显著下降 (尤其是 Complete-linkage), 表明强行合并 ( $k = 2$ ) 或强行拆分 ( $k > 3$ ) 会导致聚类准确度降低。



这是三种链接方式在  $k = 2$  到  $k = 6$  范围内的 ARI 和 Silhouette Score 曲线对比 (针对 make\_moons)。Single-linkage 的 ARI 始终保持在 1.0 附近, 再次确认其对该数据集的优秀适应性。Complete 和 Average-linkage 这两种方法的 ARI 随  $k$  的增加而剧烈波动和下降, 表明它们对非凸形状和不正确的  $k$  值非常敏感。Single-linkage 的 Silhouette Score 随着  $k$  的增大急剧下降, 并在  $k \geq 4$  之后变为负值。这表明虽然其 ARI 仍高, 因为链状结构容易在  $k$  增大时保持整体连通性, 但强行拆分非凸形状会导致簇内样本距离过大, 划分质量差。