

层次聚类算法报告

算法设计

在算法设计部分，我采用了一种改进的层次聚类方法，该方法利用了优先队列（最小堆）以优化簇合并过程中距离计算的效率。下面详细描述了代码实现的过程：

- 初始距离计算:** 利用 `scipy.spatial.distance.pdist` 函数计算所有数据点之间的欧氏距离，这一步骤是全对的距离计算的基础。随后，`scipy.spatial.distance.squareform` 函数被用来将距离向量转换成一个对称的距离矩阵，使得矩阵的每个元素 `distances[i, j]` 代表点 `i` 和点 `j` 之间的距离。
- 优先队列初始化:** 使用Python标准库中的 `heapq` 模块初始化一个最小堆。此步骤中，将所有的簇对距离作为元素添加到堆中，其中每个元素是一个三元组 `(distance, i, j)`，`distance` 是簇 `i` 与簇 `j` 之间的距离，`i` 和 `j` 是簇的索引。这样设计的目的是在合并过程中快速找到距离最小的两个簇。
- 迭代合并簇:** 在迭代过程中，算法不断从优先队列中弹出最小距离的簇对进行合并。合并操作涉及将两个簇的索引和数据点合并成一个新簇，并更新相关簇的距离信息。为了保持优先队列的更新，每当两个簇合并成一个新簇时，需要重新计算这个新簇与其他所有簇之间的平均距离，并将这些新计算的距离作为新的元素添加到优先队列中。
- 终止条件:** 这个迭代过程一直持续到簇的数量减少到用户指定的目标簇数量 `n_clusters`。此时，算法终止，得到的簇划分即为最终的聚类结果。
- 聚类结果处理与可视化:** 为了直观展示聚类结果，我们使用了主成分分析（PCA）来降维数据至二维空间，这一步骤是为了可视化目的。随后，根据PCA降维后的结果，利用 `matplotlib` 库绘制了每个簇的散点图，不同的簇用不同的颜色表示，以直观显示聚类的效果。

通过一系列的步骤，层次聚类算法不仅在聚类分析中有效地识别了数据中的自然群组，而且通过优先队列优化了合并过程，提高了算法的效率。尽管时间复杂度和空间复杂度仍较高，但这种方法为处理中等规模数据集提供了一种相对可行的解决方案。

时间和空间复杂度处理

在提供的层次聚类算法中，通过一系列策略对时间复杂度和空间复杂度进行了优化。这些优化措施显著提高了算法处理大数据集的能力。下面详细描述这些优化及其对复杂度的影响。

时间复杂度的优化

- 优先队列的使用:**
 - 原始方法:** 在每次迭代中寻找最小距离的两个簇通常需要遍历整个距离矩阵，其时间复杂度为 $O(n^2)$ ，其中 n 是簇的数量。在聚类过程中，这个操作需要重复执行多次，导致效率低下。
 - 优化方法:** 引入优先队列（最小堆）存储簇对的距离，每次迭代只需 $O(\log n)$ 的时间即可找到最小距离的簇对。这种方法大大减少了查找时间，尤其是在初始迭代时，当簇的数量较多的情况下。
 - 效果:** 通过这种优化，整个聚类过程的时间复杂度从原来的 $O(n^3)$ 或更高降低到了 $O(n^2 \log n)$ 。
- 距离的增量更新:**
 - 原始方法:** 合并两个簇后，重新计算新簇与其他所有簇之间的距离通常需要从头开始，涉及到重复的距离计算。

- **优化方法:** 在合并簇后，只更新涉及新簇的距离，利用已知的距离信息来减少计算量。例如，新簇与其他簇之间的距离可以通过现有簇之间的距离计算得到，而不是重新计算。
- **效果:** 这种方法减少了不必要的计算，进一步提高了算法的运行效率。

空间复杂度的优化

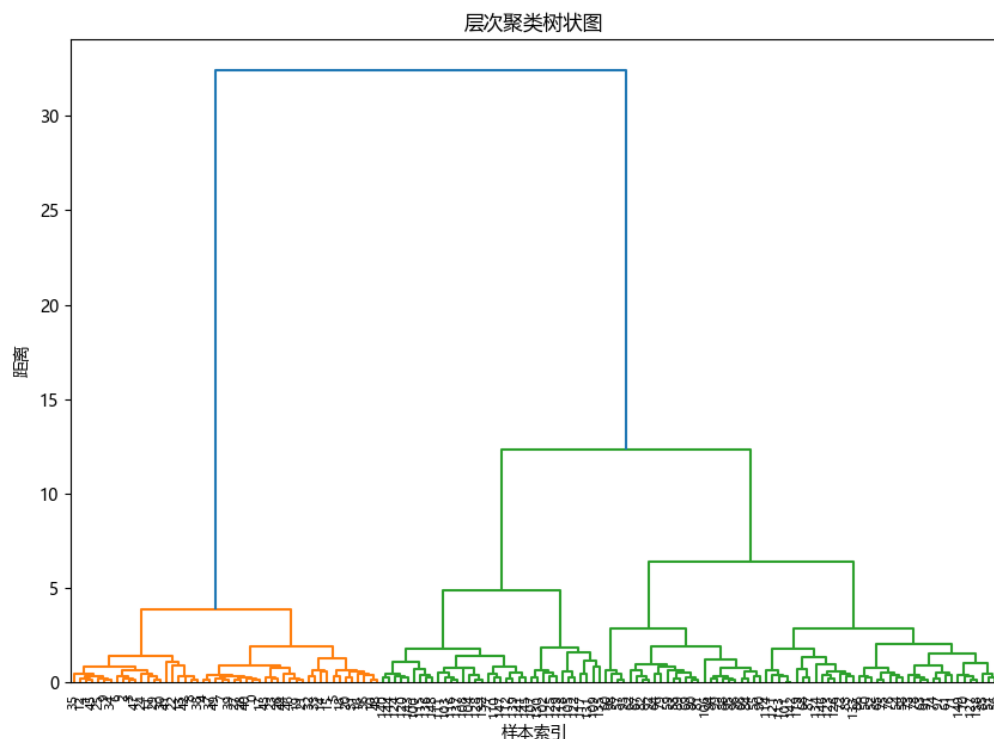
1. 动态距离矩阵更新:

- **原始方法:** 存储完整的距离矩阵需要 $O(n^2)$ 的空间，这在数据点较多时会非常耗费空间。
- **优化方法:** 通过优先队列，我们不再需要存储完整的距离矩阵。在每次合并簇后，只需更新涉及新簇的距离信息，并在优先队列中维护这些信息。这样，空间复杂度主要由优先队列的大小决定，而不是完整的距离矩阵。
- **效果:** 这种策略显著减少了算法的空间需求，使其能够处理更大的数据集。

通过上述优化，算法在处理大规模数据集时的效率和可行性得到了显著提升。特别是使用优先队列的策略，它不仅减少了查找最小距离簇对的时间，也降低了存储距离信息所需的空间，从而优化了整个层次聚类算法的时间复杂度和空间复杂度。

结果展示

树状图



结果分析

PCA的主要目的是将具有多个变量的数据转换成少数几个主要变量（主成分），同时尽可能保留原始数据的变异性。这种转换是通过识别数据中的方差最大的方向，并将这些方向作为新的空间的基，从而达到降维的目的。每一个新的方向（或主成分）都是原始数据在某个特定方向上的投影，且这些方向彼此正交（或互相独立）。

在散点图中，每个点代表一个数据样本，而点的位置则根据样本在主成分上的值来确定。

1. **PCA降维**：使用PCA将数据从原始的高维空间降至二维空间。这意味着，无论原始数据有多少维（特征），通过PCA处理后，每个样本都被映射到了一个由两个主成分构成的新空间中。第一主成分（Principal Component 1）和第二主成分（Principal Component 2）分别对应于数据变异性最大和次大的方向。

2. **散点图解释**：

- **X轴**：第一主成分。它捕获了数据中最大的方差，通常可以解释为数据集中最重要的模式或趋势。
- **Y轴**：第二主成分。它垂直于第一主成分，并捕获了数据中第二大的方差，提供了除了第一主成分以外的新信息。
- **颜色**：用来区分不同的簇。在聚类分析后，每个数据点被分配到一个簇中，散点图上用不同的颜色表示不同的簇，以便观察哪些点被归为同一个簇，从而了解数据中的群集结构。

3. 根据PCA分析的结果，具体化写出各个主成分对应的特征组合：

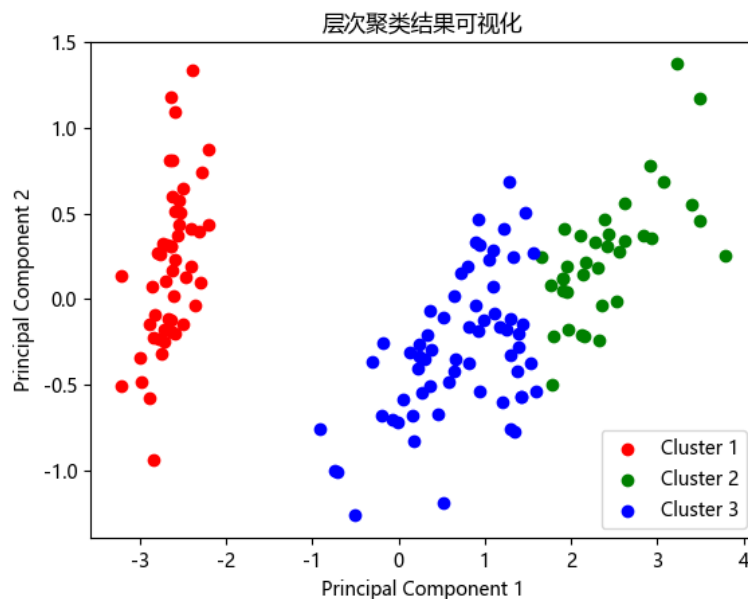
第一主成分

- 萼片长度 (Sepal Length): 0.3614
- 萼片宽度 (Sepal Width): -0.0845
- 花瓣长度 (Petal Length): 0.8567
- 花瓣宽度 (Petal Width): 0.3583

第二主成分

- 萼片长度 (Sepal Length): 0.6566
- 萼片宽度 (Sepal Width): 0.7302
- 花瓣长度 (Petal Length): -0.1734
- 花瓣宽度 (Petal Width): -0.0755

第一主成分主要由花瓣长度的正向贡献和萼片宽度的轻微负向贡献组成，表明花瓣长度是影响数据变异性的主要因素。相反，第二主成分则显示出萼片长度和萼片宽度的强正向贡献，但花瓣的尺寸对它的贡献较小，这可能指向了萼片尺寸在区分数据点上的重要性。



复杂度分析

- **时间复杂度:** 主要开销在于距离矩阵的计算和更新，总体时间复杂度为 $O(n^2 \log n)$ 。
- **空间复杂度:** 主要由初始距离矩阵和优先队列的存储决定，总体空间复杂度为 $O(n^2)$ 。