



模式识别

作业（一） 实现C-Means算法

姓 名	熊恪峥
学 号	22920202204622
日 期	2022年2月28日
学 院	信息学院
课程名称	模式识别

作业（一） 实现C-Means算法

目录

1 实现和运行结果	3
1.1 选择特征和实现	3
1.2 运行结果	3
2 误差统计和分析	3
A 附录： Iris数据集多变量图	6
B 附录： 核心部分源代码	7

1 实现和运行结果

1.1 选择特征和实现

首先做出Iris数据集两两变量的散点图和单个变量的分布曲线图，如图3（见附录：Iris数据集多变量图），可以发现萼片长款参数单独出现对聚类结果没有较好的贡献，重合部分较多。为了充分利用数据，再计算萼片面积和花瓣面积。发现组合“花瓣面积-萼片面积”、“花瓣面积-萼片长度”等都有较好的区分，但是**考虑到K-Means算法团块状聚类的特点**，我选择做出散点图最符合团状特点的“花瓣面积-萼片面积”作为用于聚类的两个特征，但是该特征选择依然难以避免iris-virginica和iris-versicolor两类有所重叠。实现见附录：核心部分源代码中的代码1。

该实现使用了以下第三方库：

- **Pandas** 用于读取CSV数据集文件、计算平均数和方差等统计量
- **NumPy** 用于加速数值计算

1.2 运行结果

聚类结果如图1，图上除了聚类结果还显示了初始时选择的点、聚类中心。可见C-Means给出了较好的聚类结果。

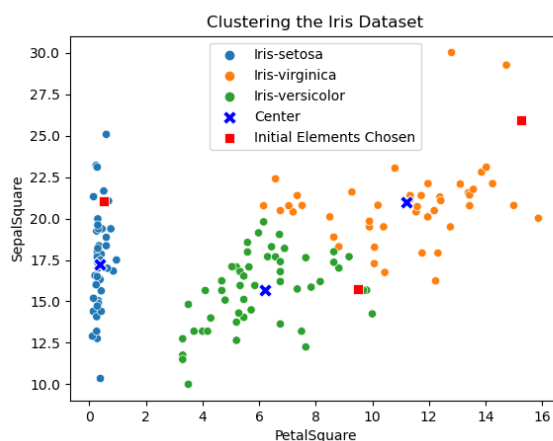


图 1: 聚类结果可视化

2 误差统计和分析

表格 1: 误差统计

物种	数量	聚类数量	错误量	错误率
Iris-virginica	50	46	9	0.2
Iris-versicolor	50	54	13	0.24
Iris-setosa	50	50	0	0

为了进一步判断聚类的效果，将聚类结果和Iris数据集中的标签对照，统计类的大小、错误量、错误率，得到表1和图2。可以发现C-Means较好地区分出了iris-setosa，但是iris-virginica和iris-versicolor中存在一定的错误，这和选择特征时的观察相符合，即iris-virginica和iris-versicolor的分布有重合的部分。

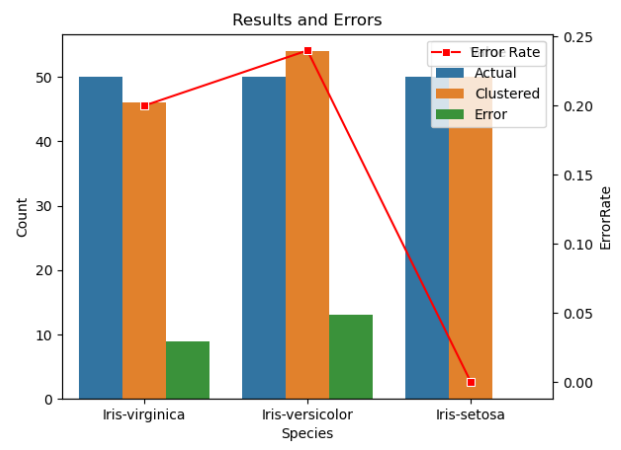


图 2: 结果和错误统计

但是，考虑到错误率相当低，可以认为C-Means算法给出了很好的聚类效果。

A 附录： Iris数据集多变量图

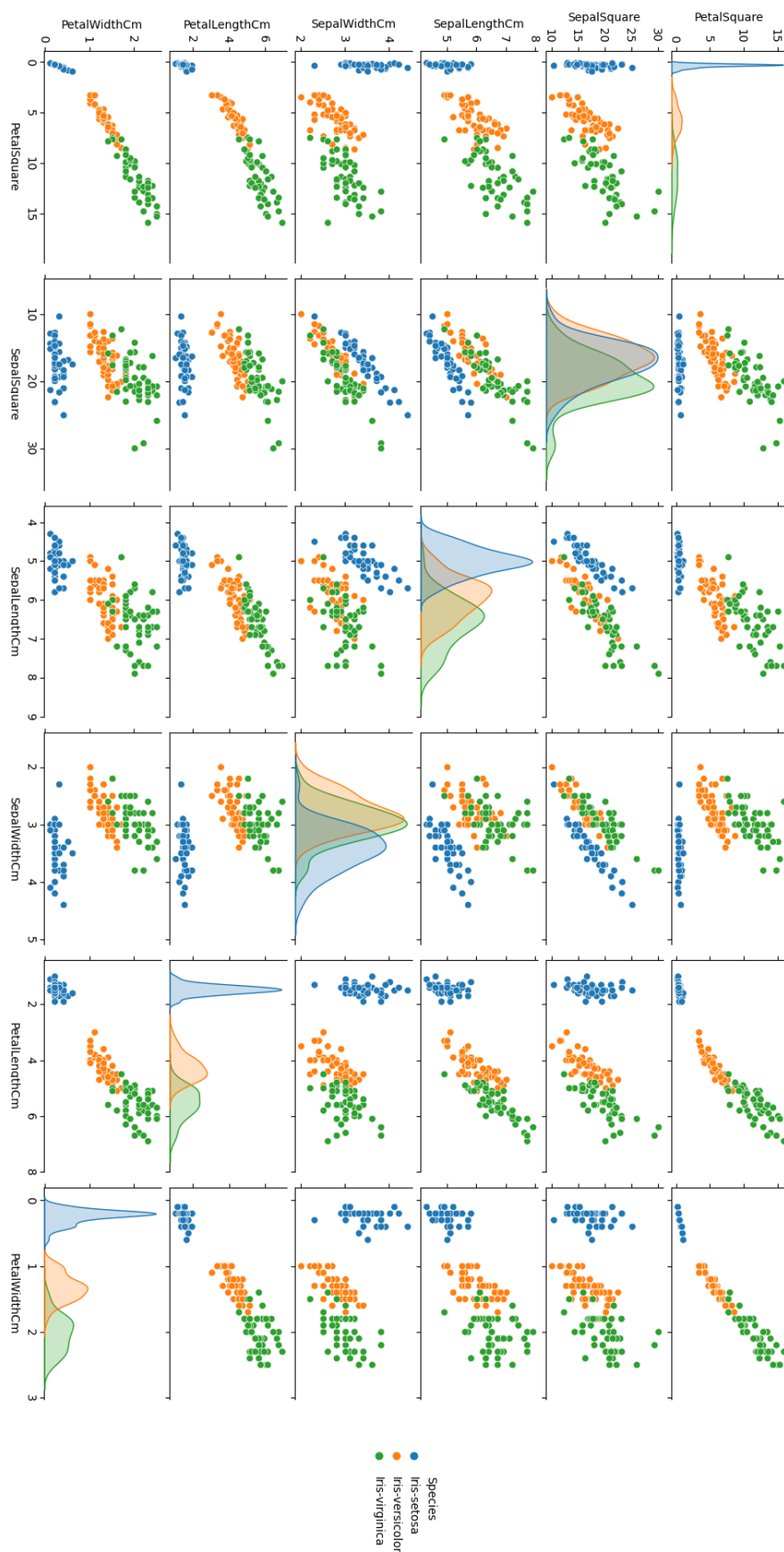


图 3: Iris数据集多变量图(Pair Plot)

B 附录：核心部分源代码

代码 1: 瑞利分布生成

```
1 class KMeans:
2     def __init__(self, feats: pd.DataFrame, k: int):
3         self.tries = 0
4         self.feats = feats
5         self.k = k
6
7     def _wcss(self, centroids, cluster) -> float:
8         ret = 0.0
9         for i, val in enumerate(self.feats.values):
10            ret += np.sqrt(
11                (centroids[int(cluster[i]), 0] - val[0]) ** 2 + (centroids[
12                    int(cluster[i]), 1] - val[1]) ** 2)
13            return ret
14
15    def cluster(self, max_tries: int = 32767):
16        self.tries = 0
17
18        cluster = np.zeros(self.feats.shape[0])
19
20        centroid_indexes, centroids = self.feats.sample(n=self.k).index, self.feats.sample(n
21            =self.k).values
22
23        while self.tries < max_tries:
24            self.tries += 1
25
26            for id, row in enumerate(self.feats.values):
27                min_dist = float('inf')
28                for cid, centroid in enumerate(centroids):
29                    dist = np.sqrt((centroid[0] - row[0]) ** 2 + (centroid[1] -
30                        row[1]) ** 2)
31                    if dist < min_dist:
32                        min_dist, cluster[id] = dist, cid
33                clustered_centroids = self.feats.copy().groupby(by=cluster).
34                    mean().values
35                if np.count_nonzero(centroids - clustered_centroids) == 0:
36                    break
37                else:
38                    centroids = clustered_centroids
39
40            return centroid_indexes, centroids, cluster, self._wcss(centroids, cluster)
41
42    def get_tries(self) -> int:
43        return self.tries
```