

# K-means实验

PB18071477 敖旭扬

## 原理

给定数据集  $D = \{X_1, \dots, X_m\}$  和  $K$ , **K-means** (K均值) 算法针对聚类所得簇划分  $C = \{C_1, C_2, \dots, C_K\}$  最小化平方误差

$$E = \sum_{k=1}^K \sum_{\mathbf{x} \in C_k} \|\mathbf{x} - \boldsymbol{\mu}_k\|_2^2 \quad (1)$$

将数据集D划分成  $K$  个簇, 其中

$$\boldsymbol{\mu}_k = \frac{1}{|C_k|} \sum_{\mathbf{x} \in C_k} \mathbf{x} \quad (2)$$

在本实现中,  $\boldsymbol{\mu}$  的初始值是从数据集  $D$  中随机抽取  $K$  个样本点得到。之后把每个样本点归类为欧氏距离最近的  $\boldsymbol{\mu}_k$  所在的那个簇, 再重新计算  $\boldsymbol{\mu}$ , 重复迭代直到样本分类类别标签不再发生变化时结束训练。

性能度量使用**DBI**, 定义为

$$DBI = \frac{1}{K} \sum_{i=1}^K \max_{j \neq i} \left( \frac{avg(C_i) + avg(C_j)}{d_{cen}(C_i, C_j)} \right) \quad (3)$$

其中

$$avg(C_k) = \frac{1}{|C_k|} \sum_{\mathbf{x} \in C_k} dist(\boldsymbol{\mu}_k, \mathbf{x}) \quad (4)$$

为簇内样本平均距离

$$d_{cen}(C_i, C_j) = dist(\boldsymbol{\mu}_i, \boldsymbol{\mu}_j) \quad (5)$$

为簇中心点距离

一般 **DBI** 越小, 聚类划分效果越好。

## 编程实现

矩阵运算使用 python 的 **numpy** 库实现。

**K-means** 算法如下:

```
1 class KMeans:
2     def fit(self, raw_x, show=False, normalize=True):
3         x = copy.deepcopy(raw_x)
4         m = x.shape[0] # 样本数
5         d = x.shape[1] # 维度
6
7         # 归一化数据
8         if normalize:
9             for i in range(d):
```

```

10         max_xi = np.max(np.abs(X[:, i]))
11         X[:, i] = X[:, i]/max_xi
12
13     # 初始向量设置
14     sample = random.sample(range(m), self.k)
15     self.mu = copy.deepcopy(X[sample])
16     raw_mu = copy.deepcopy(raw_X[sample])
17
18     # 迭代优化
19     self.saved = []
20     self.label = np.zeros(m, dtype=int)
21     while True:
22         old_label = copy.deepcopy(self.label)
23         for i in range(m):
24             self.label[i] = np.argmin(self.dist(X[i], self.mu, 2))
25         if show:
26             label = copy.deepcopy(self.label)
27             self.saved.append((label, raw_mu))
28             raw_mu = self.calculate_mu(raw_X, self.label)
29             self.mu = self.calculate_mu(X, self.label)
30             if (old_label == self.label).all():
31                 # 样本分类类别标签不再发生变化时终止迭代
32                 break
33         self.mu = self.calculate_mu(raw_X, self.label)
34         if show:
35             self.show_dynamic(raw_X, self.saved)
36     return self.label, self.mu

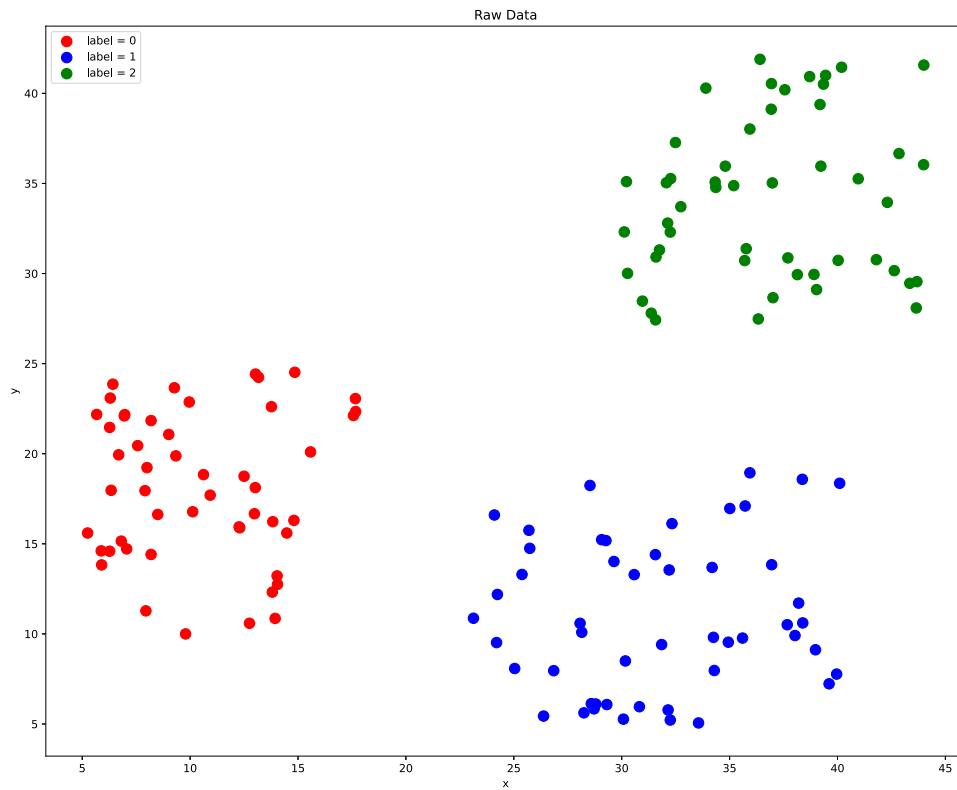
```

完整实验源码见压缩包中的[KMeans.py](#)。

## 运算结果

### 实例

本次实验使用的数据分布如下：



原始数据和理想结果

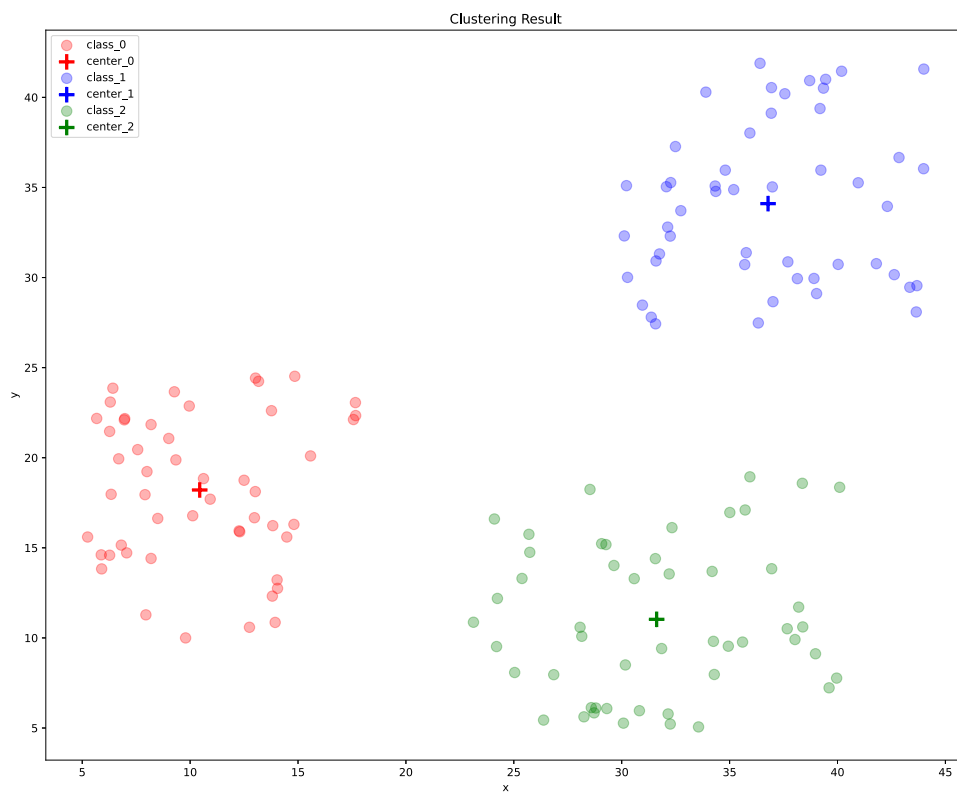
在主函数中调用下面的实例

```
1 # 训练模型
2 kmeans = KMeans(k=3)
3 label, mu = kmeans.fit(X, show=True)
```

命令行输出结果为

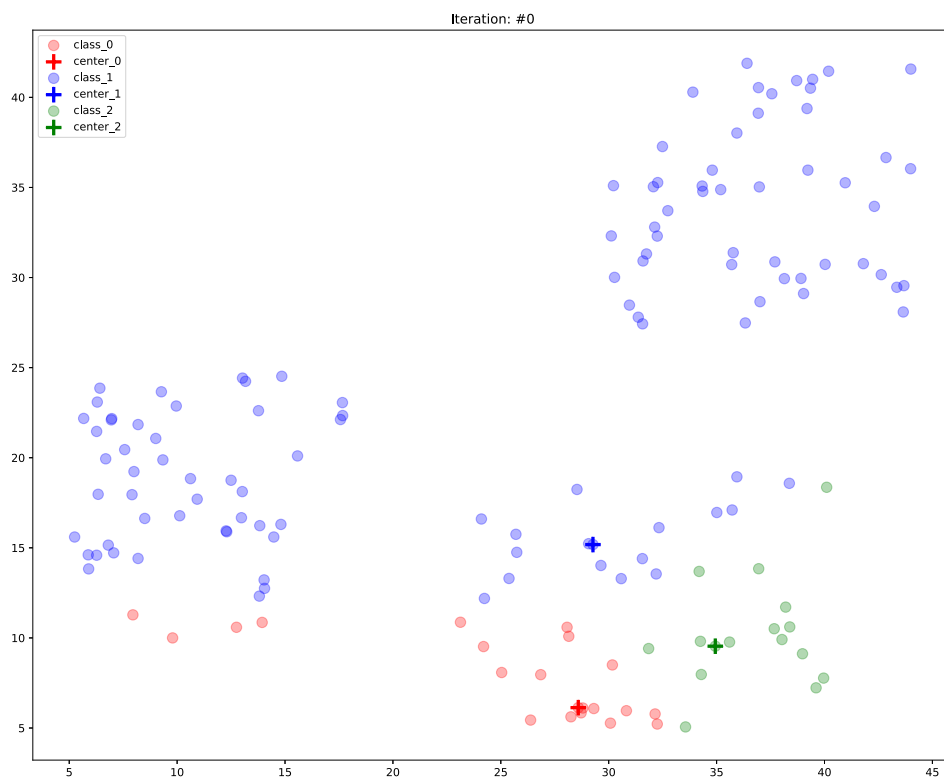
```
1 簇中心坐标:
2 [[10.444  18.2102]
3  [36.782  34.1022]
4  [31.6182 11.0314]]
5 DBI = 0.4936
```

## 训练结果

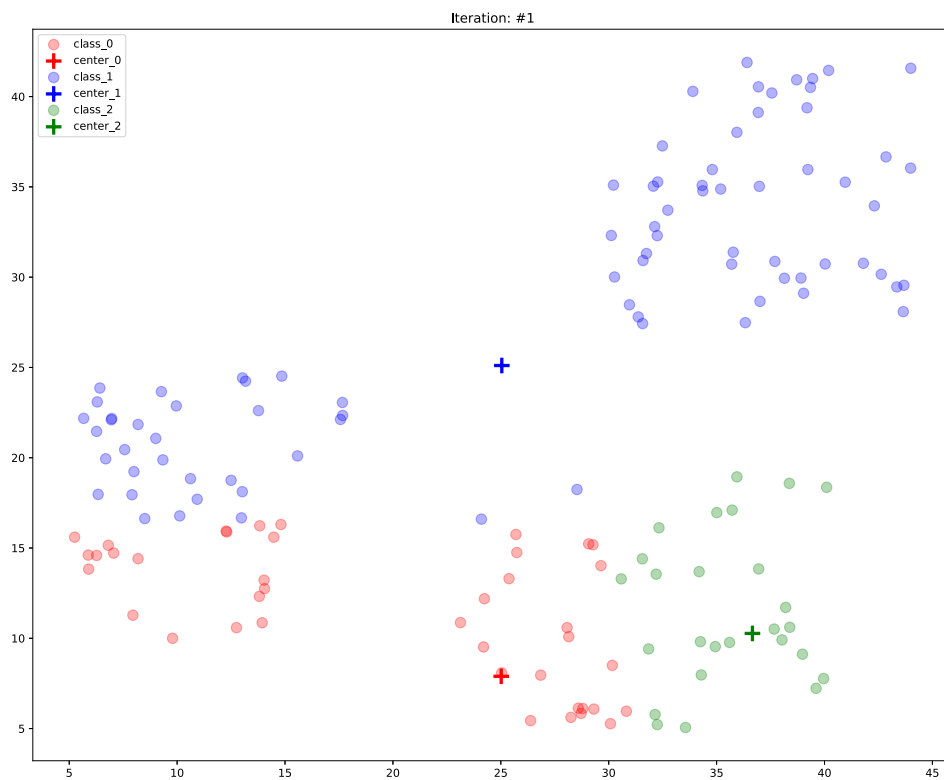


聚类结果

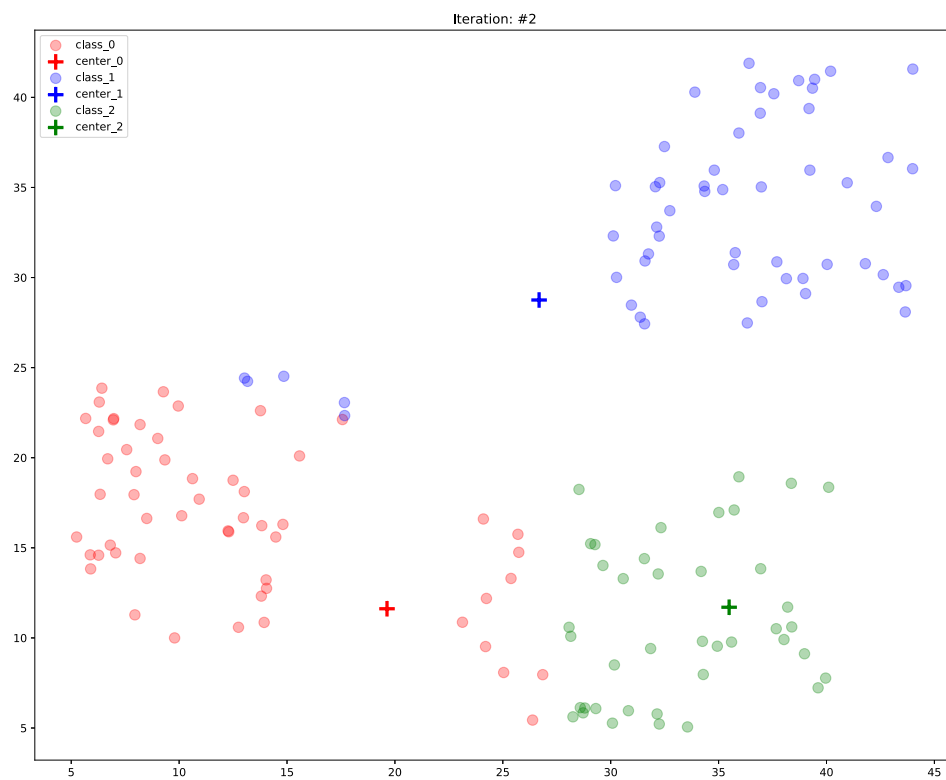
训练迭代过程（某一次的结果）



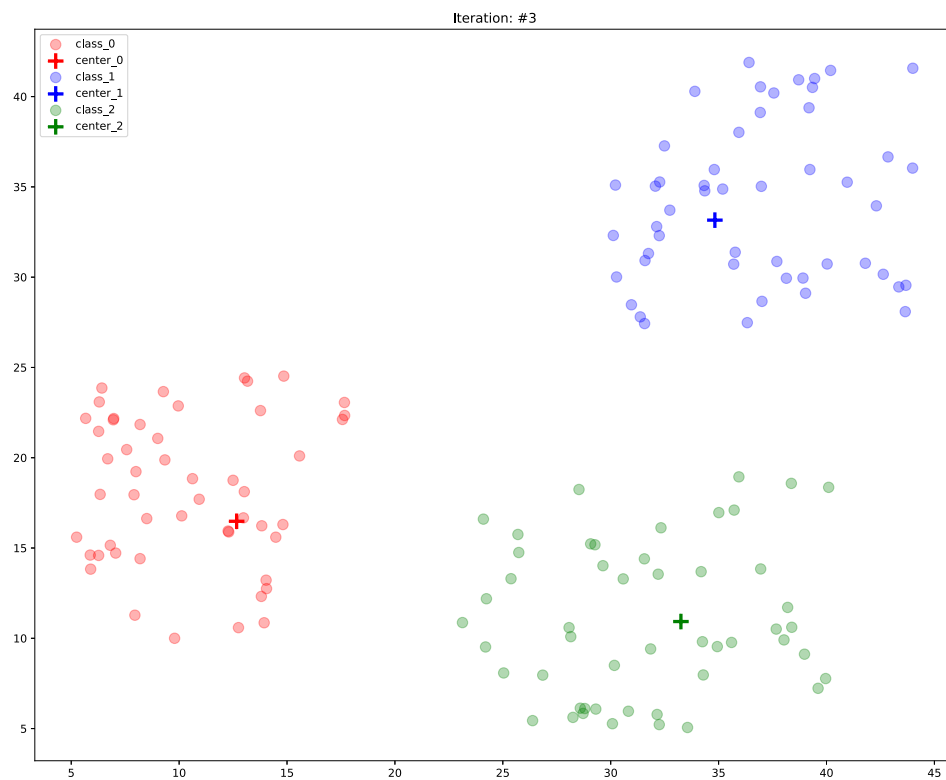
第0次迭代后



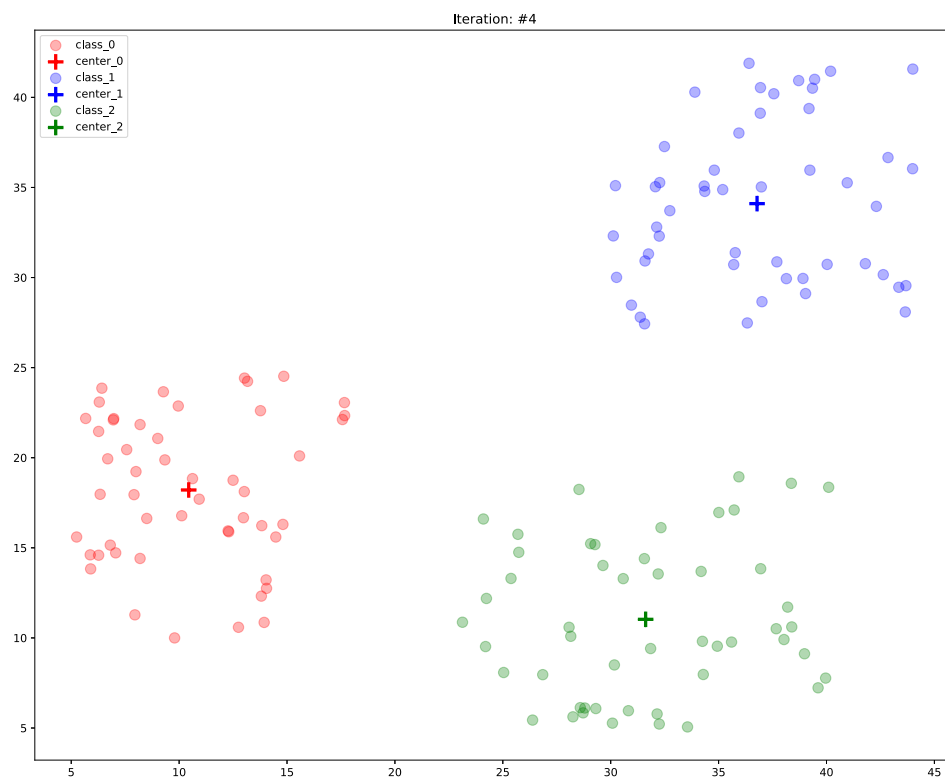
第1次迭代后



第2次迭代后



第3次迭代后



第4次迭代后，样本分类类别标签不变

## 总结

实验要求的 Baseline 为

- 1) 输出聚类后簇中心点坐标
- 2) DBI值小于5

我训练出的结果中，聚类后簇中心点坐标为 (10.444, 18.2102), (36.782, 34.1022), (31.6182, 11.0314), DBI值为 0.4936, 效果良好。