

非监督学习

- PB17111568
- 郭雨轩

数据预处理

将除掉第一维的其余每列数据都做了 `col = (col - mean(col)) / std(col)` 的变换。

PCA

算法思想

对输入的数据，我们希望对其进行线性变换，使得类间散度最大化，我们先构造协方差矩阵，之后逐个选取能使投影后方差最大的方向，直到维度满足我们的要求，这个过程也就是从大到小逐个选取协方差矩阵的特征值对应的特征向量，具体实现如下。

代码实现

训练代码：

```
1  def fit(self, train_x):
2      self.train_x = np.mat(train_x.copy()).T  # m * n m个样本，列向量
3
4      # 计算均值
5      mean_vec = np.mean(self.train_x, axis=1)
6      # 计算误差矩阵
7      var_mat = self.train_x - np.tile(mean_vec, self.train_x.shape[1])
8      # 计算协方差矩阵
9      cov_mat = var_mat * var_mat.T
10     # 计算协方差矩阵的特征向量和特征值
11     e_val, e_vec = np.linalg.eig(cov_mat)
12     # 按照特征值排序排序
13     tmp = sorted(zip(e_val, e_vec.T), reverse=True)
14     e_val, e_vec = zip(*tmp)
15     # 若根据threshold选取特征向量
16     if self.use_threshold:
17         total = sum(e_val)
18         tmp = 0
19         W = []
20         # 逐个选取特征向量，知道大于threshold
```

```

21         for i in range(len(e_val)):
22             tmp += e_val[i]
23             W.append(list(np.squeeze(np.array(e_vec[i]))))
24             if tmp / total >= self.threshold:
25                 break
26         else:
27             # 否则选择前k大特征值对应的特征向量
28             W = [list(np.squeeze(np.array(e_vec[i]))) for i in
range(self.first_k)]
29
30         # w为投影矩阵
31         self.W = np.mat(np.array(W)).T
32         print('DIM:', self.W.shape)
33         self.mean_vec = mean_vec
34         # 计算降维后训练集
35         proj_train = self.train_x.T * self.W
36         return proj_train

```

KMeans

算法思路

初始化时随机选取k个中心点，每轮迭代时不断的将所有的点划分到新的中心点对应的类别处，同时更新中心点，不断迭代直到中心点收敛。

代码实现

```

1  def fit(self, train_x):
2      self.train_x = np.mat(train_x.copy())
3      m, n = np.shape(self.train_x)
4      center = []
5      # 初始化，随机选取中心点
6      for i in random.sample(range(m), self.K):
7          center.append(list(np.squeeze(np.array(self.train_x[i, :]))))
8      center = np.mat(center)
9      error = 10000000
10     label = None
11     # 当中心点更新的幅度大于threshold
12     while error > self.threshold:
13         # 记录旧的中心
14         center_old = center.copy()
15         label = []
16         # 按照现有的中心点，对每个点进行分类，这个点到哪个中心点最近就属于哪一类

```

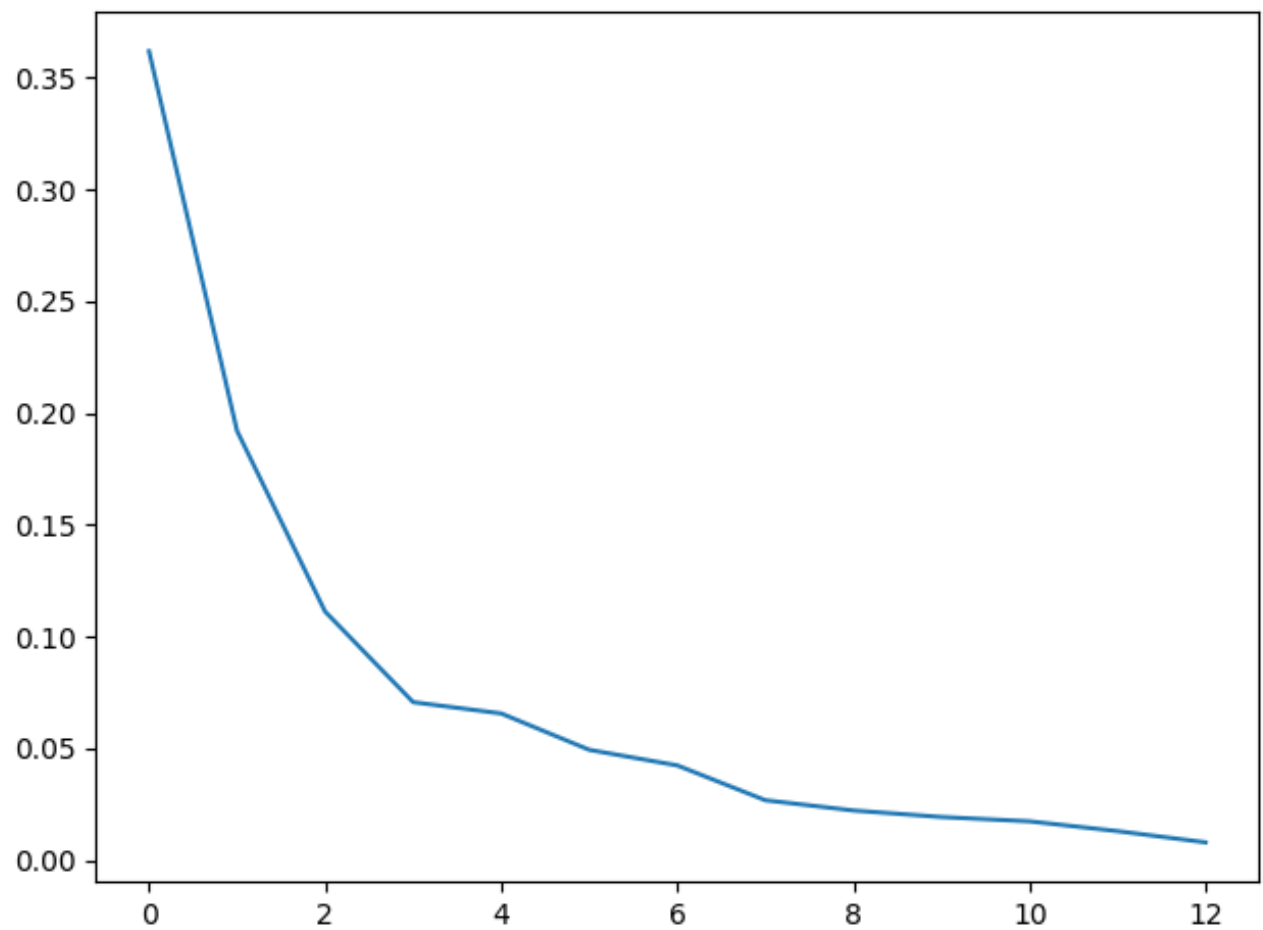
```

17         for i in range(m):
18             min_label = -1
19             min_distance = 10000000
20             for j in range(self.K):
21                 dis_i_j = distance(self.train_x[i], center[j])
22                 if dis_i_j < min_distance:
23                     min_distance = dis_i_j
24                     min_label = j
25             label.append(min_label)
26
27         label = np.array(label)
28         # 重新计算center
29         center = np.array([np.mean(self.train_x[label == i, :], axis=0) for i
in range(self.K)])
30         center = np.mat(np.squeeze(center))
31         # 计算较于旧的center到变化
32         error = np.sum(np.linalg.norm(center-center_old, ord=2, axis=1))
33         return center, np.expand_dims(label, 1)

```

实验部分

PCA中主成分的占比



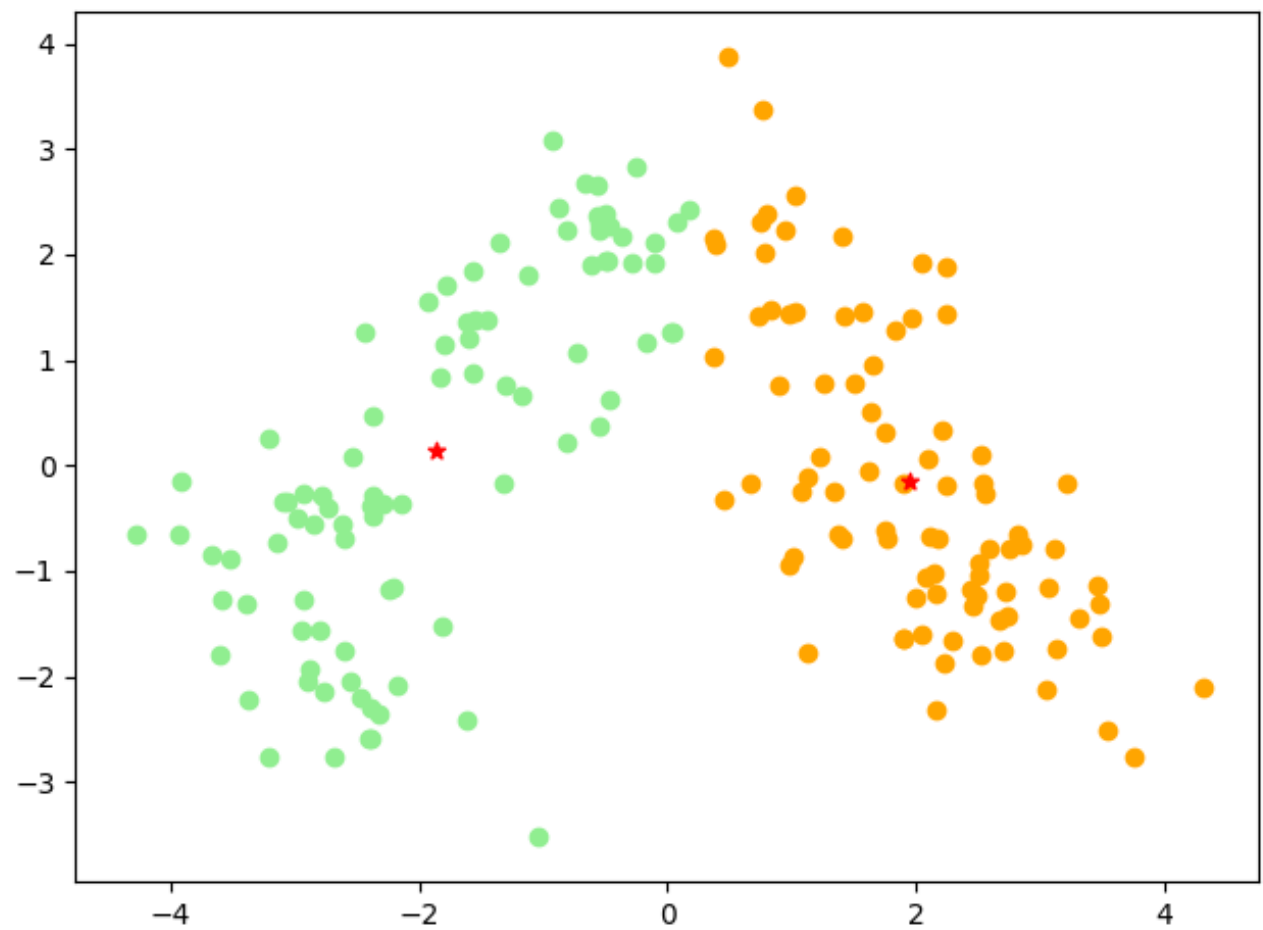
类别数=2

不同程度降维的Kmeans精度

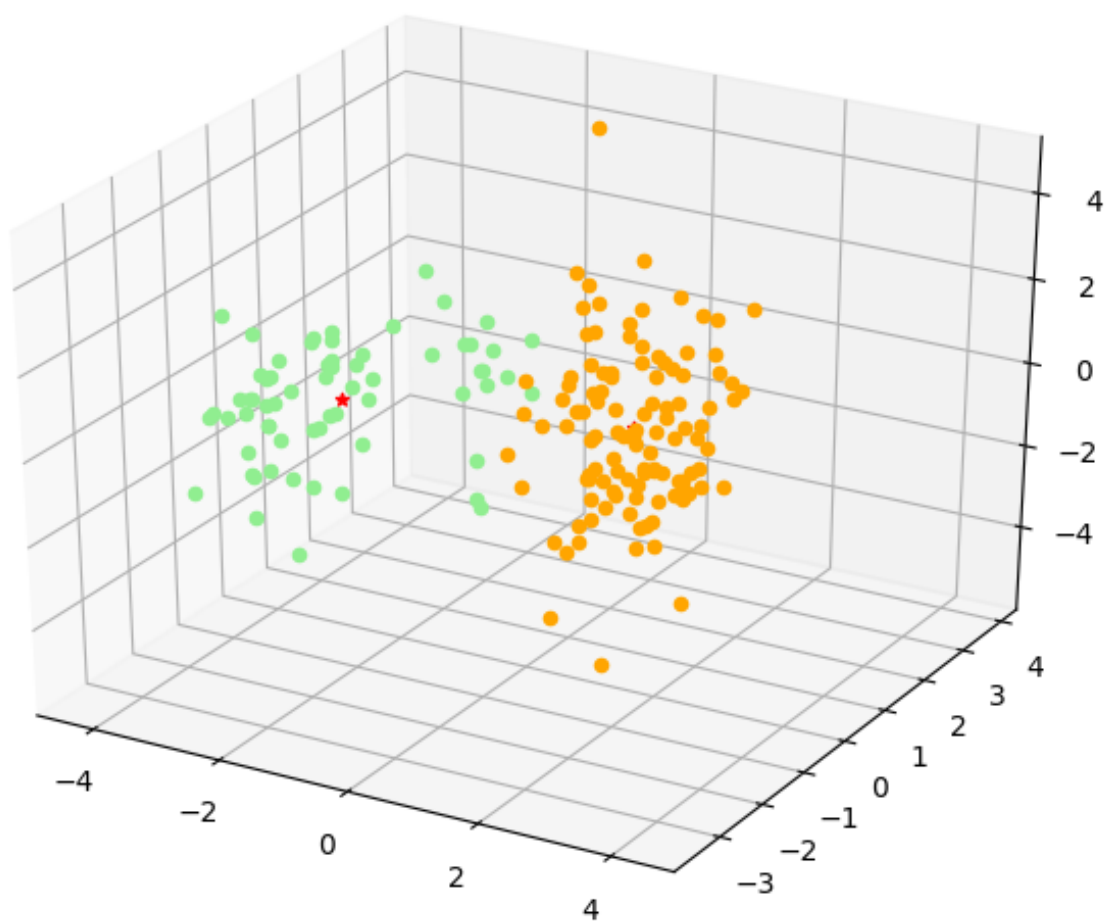
DIM	轮廓系数	兰德系数
2	0.4649	0.6876
3	0.4002	0.6876
4	0.3239	0.6876
5	0.3674	0.6876
6	0.3163	0.6876
7	0.2918	0.6912
8	0.2893	0.6800
9	0.2854	0.6876
10	0.2697	0.6912
11	0.2737	0.6876
12	0.2615	0.6932
13	0.2596	0.6932

降维可视化

降维至二维时：



降维至3维时：



类别数=3

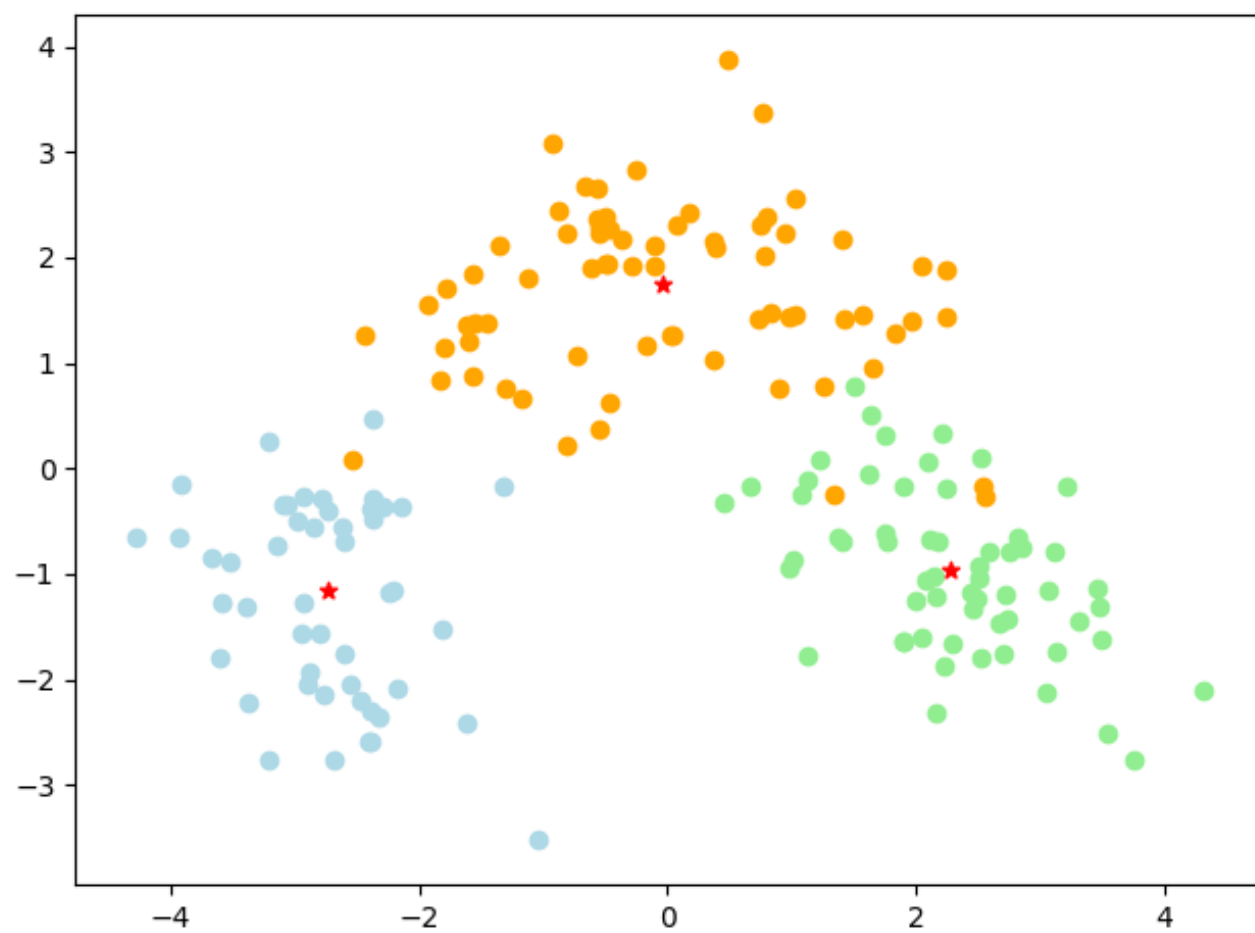
不同程度降维的Kmeans精度

DIM	轮廓系数	兰德系数
2	0.5601	0.9310
3	0.4537	0.9542
4	0.4050	0.9318
5	0.3674	0.9310
6	0.3472	0.9620
7	0.3267	0.9382
8	0.3149	0.9542
9	0.3068	0.9542
10	0.2996	0.9620
11	0.2919	0.9542
12	0.2880	0.9542
13	0.2806	0.9466

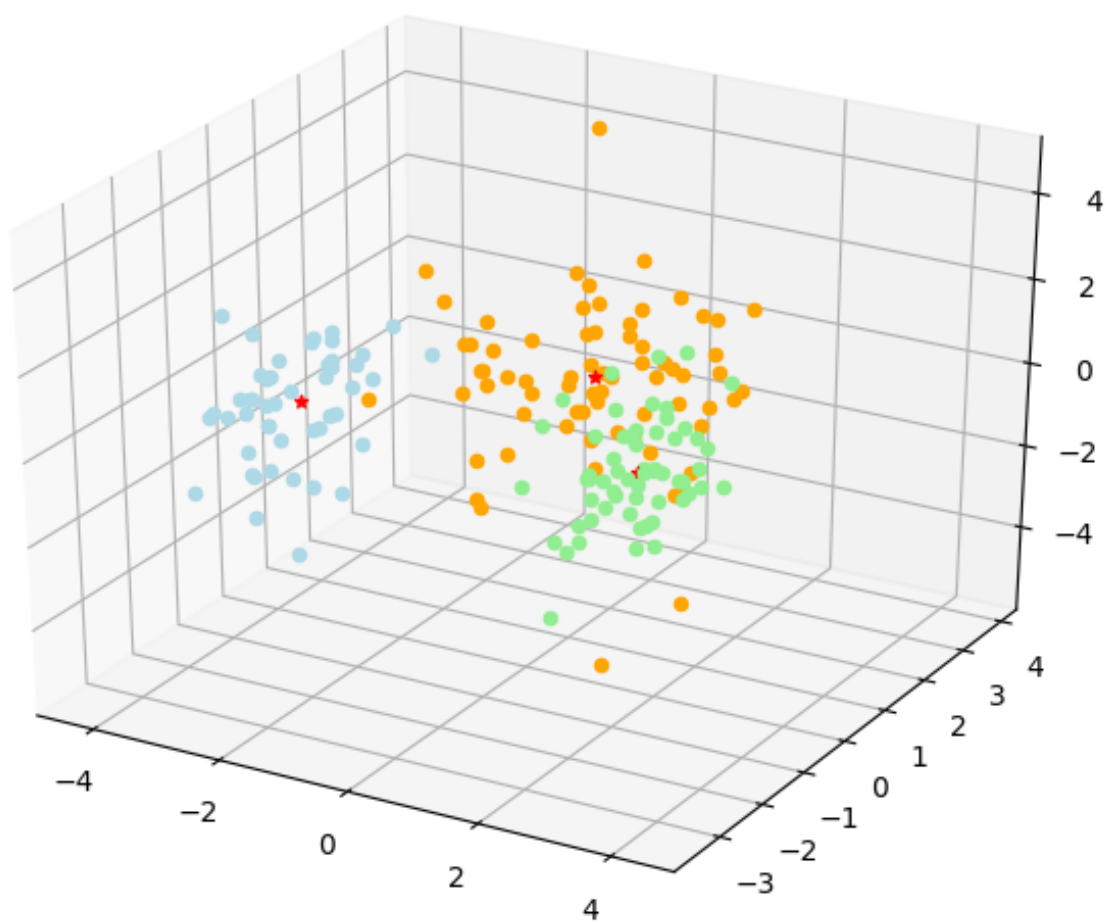
降维可视化

(此图中展示的是label的类别情况而非聚类预测的结果，可以看到聚类的准确性较高)

降维至二维时：



降维至三维时：



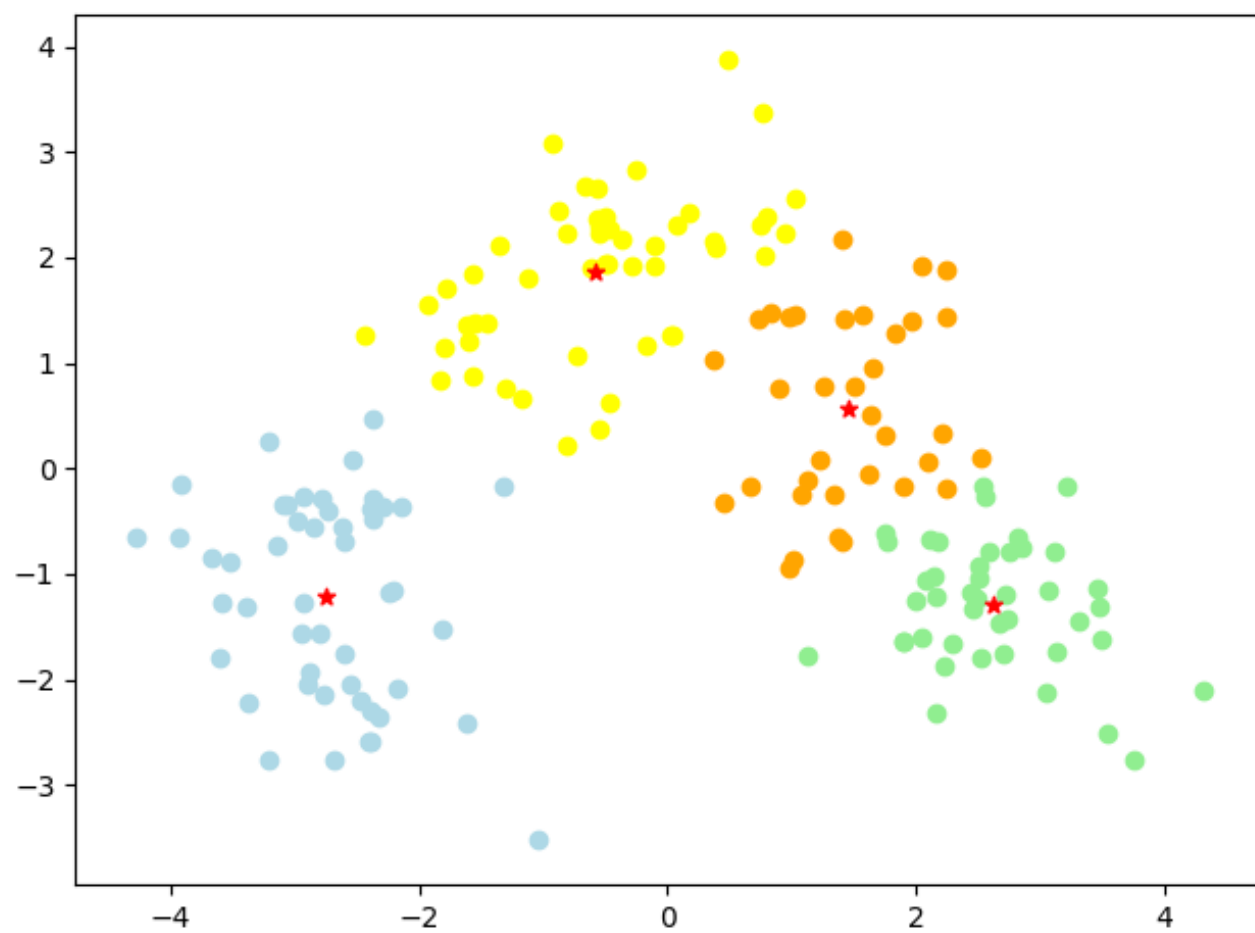
类别数=4

不同程度降维的Kmeans精度

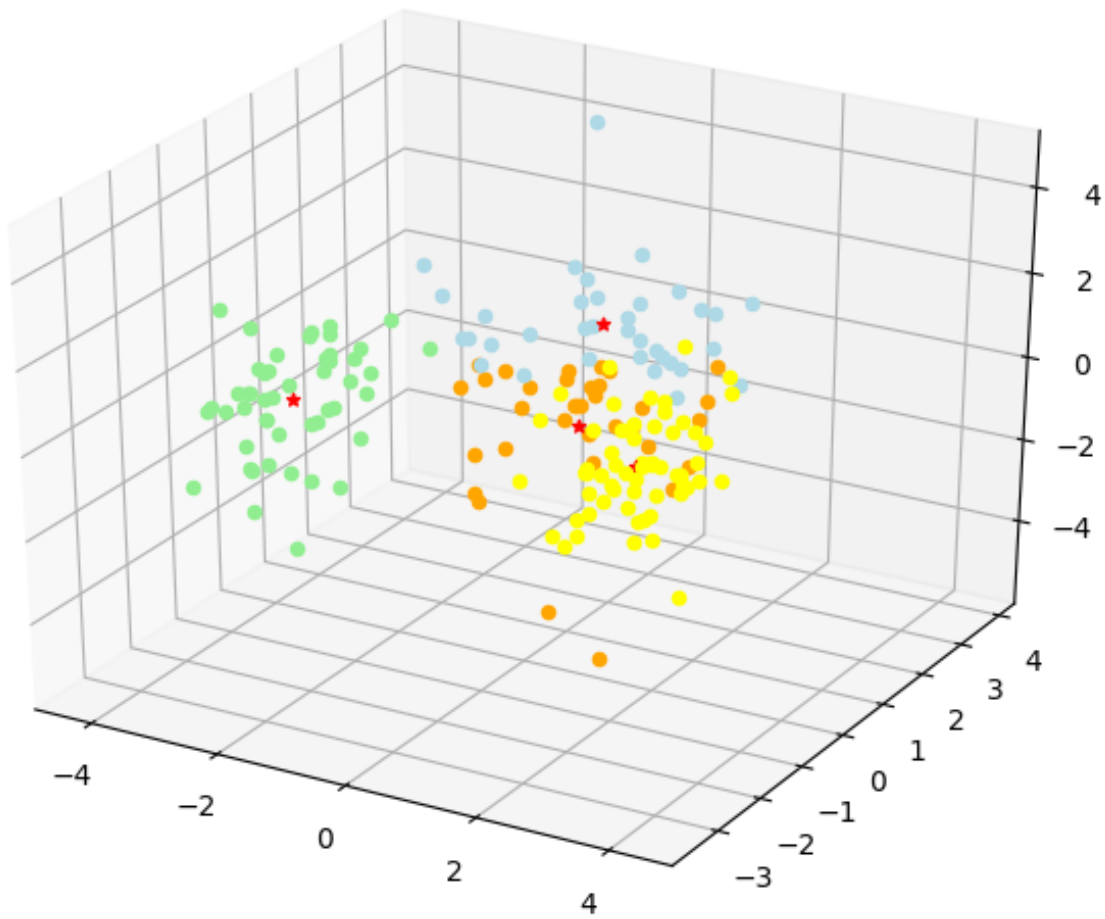
DIM	轮廓系数	兰德系数
2	0.4855	0.8573
3	0.4106	0.8861
4	0.3856	0.9074
5	0.3007	0.9197
6	0.2958	0.8906
7	0.2710	0.8892
8	0.2915	0.9016
9	0.3150	0.9577
10	0.3065	0.9494
11	0.2622	0.9134
12	0.2510	0.8921
13	0.2146	0.8831

降维可视化

降维至二维时：



降维至三维时：



结果分析

- 比较降维维度和聚类个数后，应选择降维维度为2，聚类类数为3，此时轮廓系数最大
- PCA降维的维度对KMeans的精度基本无影响，这是因为数据的维度本身就很小，PCA基本没有降维的必要
- 可视化结果中，当降维到2维时，可以看到绝大部分的点都可以被正确聚类，但是有部分的离群点不能被正确的分类，可以考虑对数据集更细致的统计分析，改进规范化的方法，提升精度