

Q3 (DBSCAN)

姓名：王依睿

学号：1552651

a) DBSCAN

1. 引入 sklearn 和 DBSCAN

```
import sklearn
from sklearn.cluster import DBSCAN
```

2. DBSCAN聚类并计算其silhouette系数

```
# DBSCAN聚类
def dbscan(eps):
    X = StandardScaler().fit_transform(data.T)
    db = DBSCAN(eps = eps).fit(X)
    labels = db.labels_
    if judge(labels):
        silhouette_score = metrics.silhouette_score(X, labels)
        return silhouette_score
    else:
        return -1
```

3. 选择不同的eps进行DBSCAN聚类并计算其silhouette系数

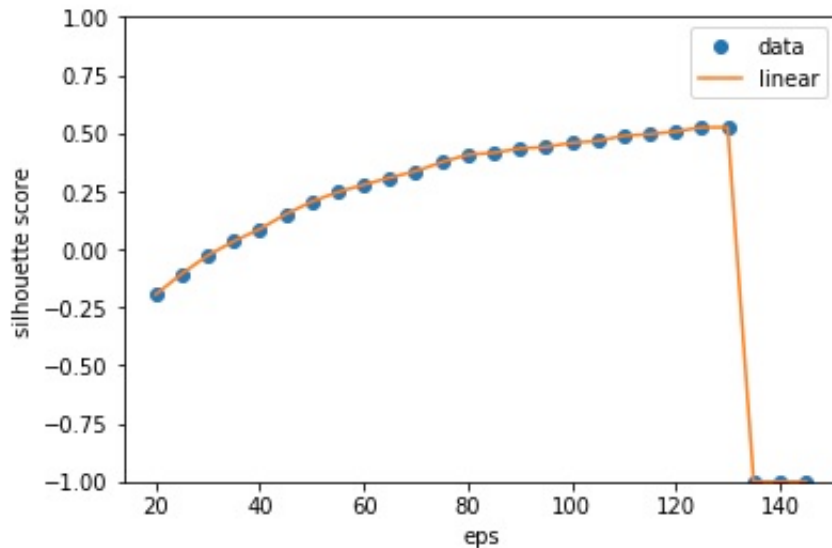
```
# 选择不同的eps聚类并计算其silhouette系数
def vip_dbscan():
    vip_dbscan_df = pd.DataFrame(columns = ['silhouette_score'], index
= np.arange(20, 150, 5))
    for eps in np.arange(20, 150, 5):
        silhouette_score = dbscan(eps)
        vip_dbscan_df['silhouette_score'][eps] = silhouette_score
    return vip_dbscan_df
```

结果

silhouette_score

20	-0.194231
25	-0.10558
30	-0.0265472
35	0.034153
40	0.0870003
45	0.150627
50	0.204767
55	0.245794
60	0.27737
65	0.308034
70	0.336027
75	0.374431
80	0.406719
85	0.414988
90	0.433585
95	0.440235
100	0.457811
105	0.465831
110	0.488395
115	0.495219
120	0.506584
125	0.524478
130	0.524478
135	-1
140	-1
145	-1

Silhouette analysis for DBSCAN clustering on reco data

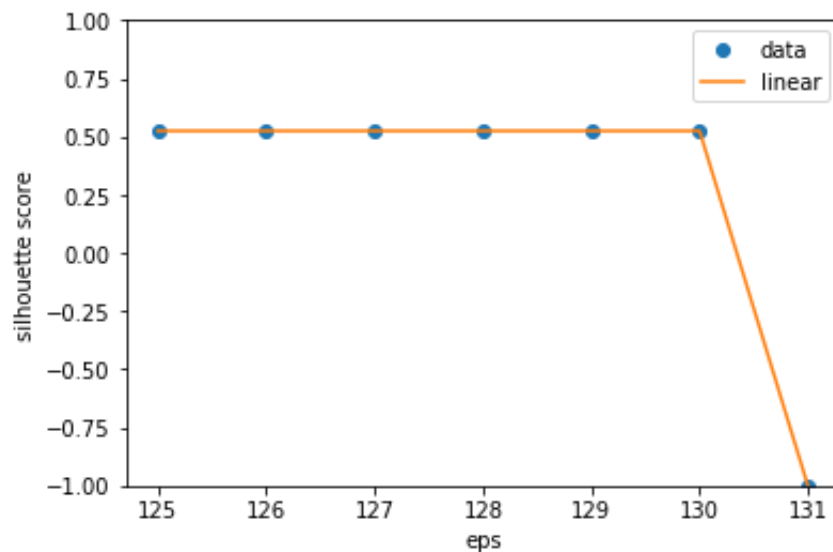


2. 在小范围内查看silhouette score以选出最佳

```
def detail_vip_dbscan():  
    vip_dbscan_df = pd.DataFrame(columns = ['silhouette_score'], index  
= np.arange(125, 132, 1))  
    for eps in np.arange(125, 132, 1):  
        silhouette_score = dbscan(eps)  
        vip_dbscan_df['silhouette_score'][eps] = silhouette_score  
    return vip_dbscan_df  
# 在小范围内查看silhouette score以选出最佳的  
detail_vip_dbscan_df = detail_vip_dbscan()
```

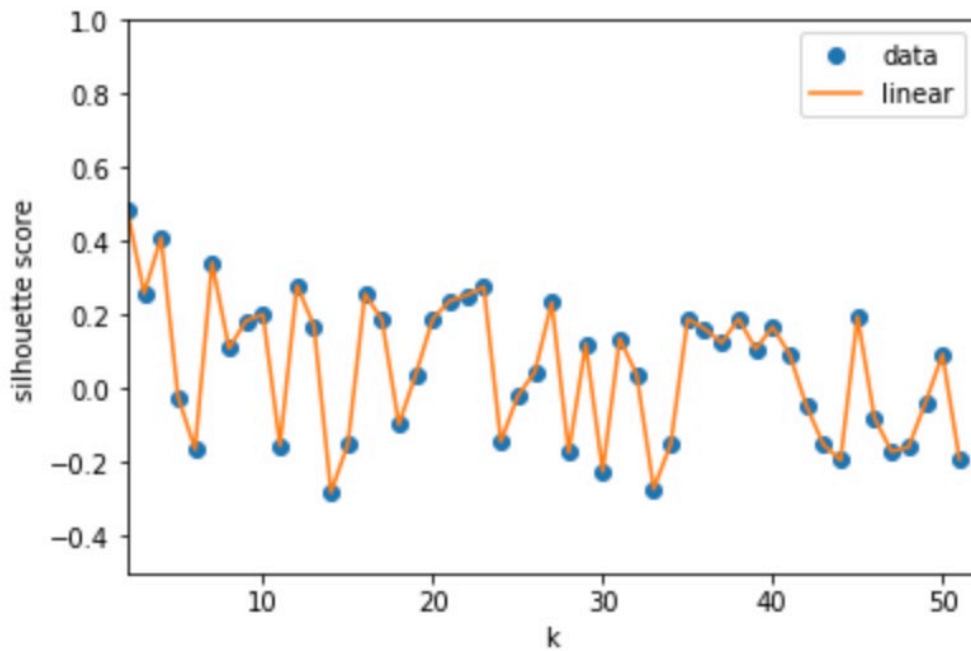
结果

Silhouette analysis for DBSCAN clustering on reco data

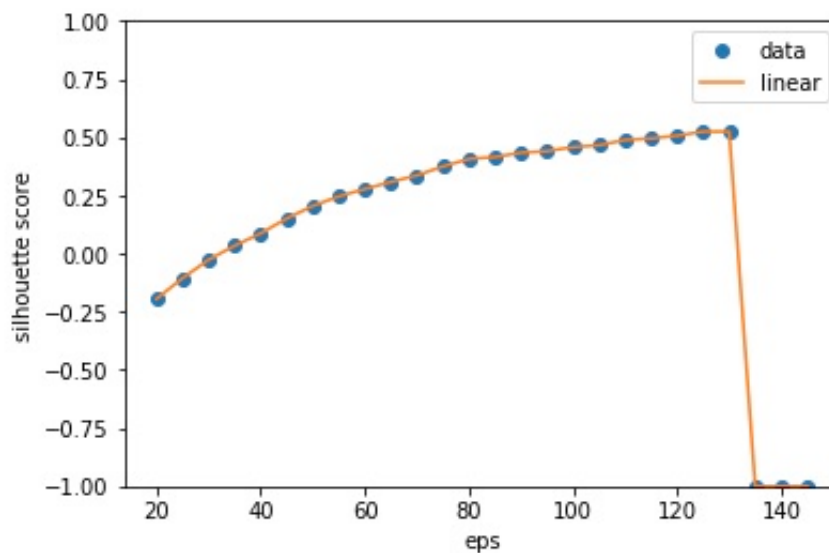


当eps在125到130之间时，Silhouette系数稳定在0.524478，且当eps大于130时，只能得到1类。

Silhouette analysis for KMeans clustering on reco data



Silhouette analysis for DBSCAN clustering on reco data



1. 从Silhouette系数来看，当k-means算法k取2时，取到最大值，为0.48722（标准化后）；当DBSCAN算法eps取130时，取到最大值，为0.524478（标准化后）。两者Silhouette系数相差不大，DBSCAN稍高一些。
2. k-means算法和DBSCAN算法都倾向于将数据点分为两类，并将极大多数点分为一类
3. k-means需要人为输入要聚的类数k，而DBSCAN可以发现任意形状的数据集且不用输入类别数K
4. DBSCAN可区分核心对象、边界点和噪音点，因此对噪声的过滤效果好

DBSCAN原理

基于密度的定义，可将所有样本点分为三类：

- 稠密区域内部的点（核心点）：在半径 Eps 内含有超过 $MinPts$ 数目的点。
- 稠密区域边缘上的点（边界点）：在半径 Eps 内点的数量小于 $MinPts$ （即不是核心点），但是其邻域内包含至少一个核心点。
- 稀疏区域中的点（噪声或背景点）：任何不是核心点或边界点的点。

用图形表示如下图：

