

数据挖掘作业q5

作者：王星洲

学号：1652977

这次实验我使用了DBSCAN的方式，探讨基于密度的聚类能否解决这个问题。

步骤一

```
1 # 读入数据所序列
2 df = pd.read_csv("../trade_new.csv", usecols=["vipno", "pluno", "amt"])
3 # 把pluno列取为第四级商品编号
4 df["pluno"] = (df["pluno"]/1000).astype(int)
5 # 分组，求和
6 df = df.groupby(["vipno", "pluno"]).sum()
```

步骤二

```
1 def jaccard_dist(a, b):
2     fenzi = 0
3     fenmu = 0
4     for i in range(a.size):
5         fenzi += min(a[i], b[i])
6         fenmu += max(a[i], b[i])
7     return 1-(fenzi/fenmu)
```

步骤三

```
1 # visitlist类用于记录访问列表
2 # unvisitedlist记录未访问过的点
3 # visitedlist记录已访问过的点
4 # unvisitednum记录访问过的点数量
5 class visitlist:
6     def __init__(self, count=0):
7         self.unvisitedlist=[i for i in range(count)]
8         self.visitedlist=list()
9         self.unvisitednum=count
10
11     def visit(self, pointId):
12         self.visitedlist.append(pointId)
13         self.unvisitedlist.remove(pointId)
14         self.unvisitednum -= 1
15
16
17 def my_dbscan1(dataSet, eps, minPts):
18     # numpy.ndarray的 shape属性表示矩阵的行数与列数
19     nPoints = dataSet.shape[0]
20     # (1)标记所有对象为unvisited
21     # 在这里用一个类vPoints进行实现
```

```

22 vPoints = visitlist(count=nPoints)
23 # 初始化簇标记列表C,簇标记为 k
24 k = -1
25 C = [-1 for i in range(nPoints)]
26 while(vPoints.unvisitednum > 0):
27     # (3)随机上选择一个unvisited对象p
28     p = random.choice(vPoints.unvisitedlist)
29     # (4)标记p为visited
30     vPoints.visit(p)
31     # (5)if p的 $\epsilon$ -邻域至少有MinPts个对象
32     # N是p的 $\epsilon$ -邻域点列表
33     N = [i for i in range(nPoints) if
jaccard_dist(dataSet.values[i,:], dataSet.values[p,:])<= eps]
34     # print(len(N))
35     if len(N) >= minPts:
36         # (6)创建个新簇C, 并把p添加到C
37         # 这里的C是一个标记列表, 直接对第p个结点进行赋值
38         k += 1
39         C[p]=k
40         # (7)令N为p的 $\epsilon$ -邻域中的对象的集合
41         # N是p的 $\epsilon$ -邻域点集合
42         # (8) for N中的每个点p'
43         for p1 in N:
44             # (9) if p'是unvisited
45             if p1 in vPoints.unvisitedlist:
46                 # (10)标记p'为visited
47                 vPoints.visit(p1)
48                 # (11) if p'的 $\epsilon$ -邻域至少有MinPts个点, 把这些点
添加到N
49                 # 找出p'的 $\epsilon$ -邻域点, 并将这些点去重添加到N
50                 M=[i for i in range(nPoints) if
jaccard_dist(dataSet.values[i,:], \
51                 dataSet.values[p1,:]) <= eps]
52                 if len(M) >= minPts:
53                     for i in M:
54                         if i not in N:
55                             N.append(i)
56                 # (12) if p'还不是任何簇的成员, 把P'添加到C
57                 # C是标记列表, 直接把p'分到对应的簇里即可
58                 if C[p1] == -1:
59                     C[p1]= k
60             # (15)else标记p为噪声
61         else:
62             C[p]=-1
63
64     # (16)until没有标t已为unvisitedl内对象
65     return C
66
67
68 def getCentroid(dataSet):
69     div = len(dataSet)
70     return sum(dataSet)/div
71
72
73 def CP(label, dataSet):
74     centroids = []
75     k = max(label) + 1
76     if k > 1:

```

```

77     print(k)
78     for i in range(1, k):
79         data = []
80         for j in range(len(label)):
81             if label[j] == i:
82                 data.append(dataSet.values[j,:])
83             centroids.append(getCentroid(data))
84     cpnum = 0
85     for i in range(1, k):
86         distance = 0
87         num = 0
88         for j in range(len(label)):
89             if label[j] == i:
90                 distance +=
jaccard_dist(dataSet.values[j,:],centroids[i - 1])
91                 num += 1
92             cpnum += distance/num
93     return cpnum / (k - 1)
94 else:
95     return 0
96
97
98 def getSC(dataSet, label):
99     sum_number = 0
100    k = len(label)
101    for i in range(k):
102        ai = 0
103        bi = 0
104        anum = 0
105        bnum = 0
106        for j in range(k):
107            if label[i]==label[j]:
108                ai +=
jaccard_dist(dataSet.values[i,:],dataSet.values[j,:])
109                anum += 1
110            else:
111                bi +=
jaccard_dist(dataSet.values[i,:],dataSet.values[j,:])
112                bnum += 1
113            if anum == 0:
114                ai = 0
115            else:
116                ai = ai / anum
117            if bnum == 0:
118                bi = 0
119            else:
120                bi = bi / anum
121            sum_number += (bi - ai) / max(ai, bi)
122    return sum_number / k
123
124
125    #数据准备
126    df = df.sort_values(by=["vipno"])
127    vipno_series = df["vipno"].drop_duplicates()
128    vipno_series = vipno_series.reset_index(drop=True)
129    print(vipno_series)
130    df = df.sort_values(by=["pluno"])
131    pluno_series = df["pluno"].drop_duplicates()

```

```

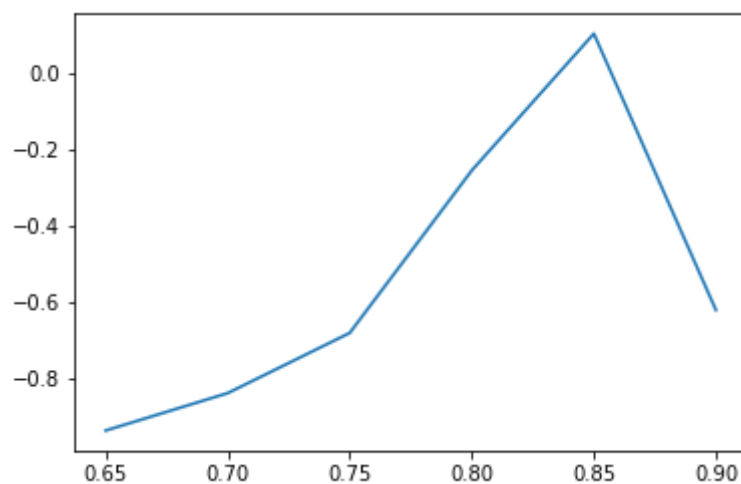
132 pluno_series = pluno_series.reset_index(drop=True)
133 print(pluno_series)
134 data = DataFrame(0, columns=pluno_series, index=vipno_series)
135 # print(data)
136 for i in df.index:
137     vipno = df['vipno'][i]
138     pluno = df['pluno'][i]
139     amt = group_data[vipno][pluno]
140     if math.isnan(data[pluno][vipno]):
141         data[pluno][vipno] = amt
142     else:
143         data[pluno][vipno] += amt
144 print(data)
145
146
147 # 调节eps值
148 sc_cp_label = []
149 for i in np.arange(0.65, 0.9, 0.05):
150     label = my_dbscan1(data, i, 5)
151     print(label)
152     sc = getSC(data, label)
153     cp = CP(label, data)
154     sc_cp_label.append((sc, cp))

```

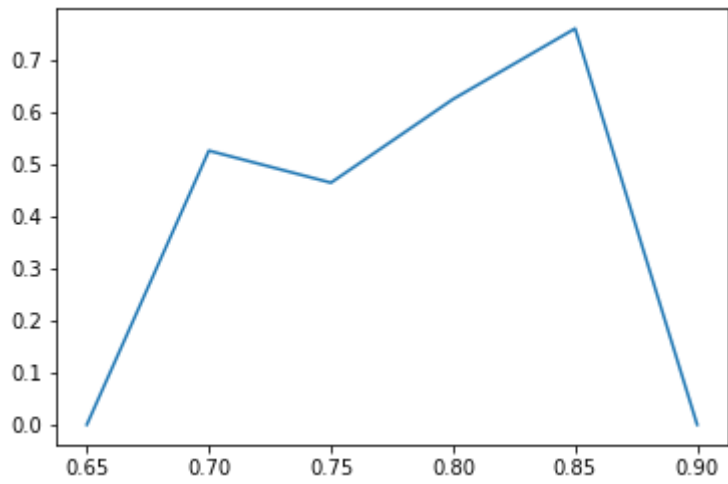
运行结果:

eps	SC	1
0.65	-0.938505	0.000000
0.70	-0.840310	0.524794
0.75	-0.682267	0.463483
0.80	-0.255707	0.623508
0.85	0.103311	0.758701
0.90	-0.622185	0.000000

SC趋势



CP趋势



结论

首先我们可以看出eps在0.85附近取得比较好的效果。

其次可以看出DBSCAN在解决这个问题时的效果也比较一般。分析原因的话：

1. 不确定的eps与minPoint。DBSCAN需要指定eps与minPoint两个参数，设定的是否合理会对结果产生很大的影响
2. 采取的data并没有经过精细的加工，使用了与方法一相同的data，这样的距离并不太可取，导致DBSCAN的半径必须很大，否则就无法进行聚类，这样就会导致效果变差。采用FTC_Tree可能会很好的解决这个问题。
3. DBSCAN的核心是基于密度，顾客购买的物品可能更趋向于散点分布，而没有紧密连在一起的情况，所以效果会不如基于距离的。