

数据挖掘作业q2

作者：王星洲

学号：1652977

步骤一

```
1  # 读入数据所序列
2  row_data = pd.read_csv("F:/CourseData/数据挖掘/datamining20/trade_new.csv")
3  df1 = row_data[["vipno", "pluno", "amt"]]
4  df2 = row_data[["vipno", "pluno", "amt"]]
5  df3 = row_data[["vipno", "pluno", "amt"]]
6  df4 = row_data[["vipno", "pluno", "amt"]]
7  df1["pluno"] = (df1["pluno"]/1000000).astype(int)
8  df2["pluno"] = (df2["pluno"]/100000).astype(int)
9  df3["pluno"] = (df3["pluno"]/10000).astype(int)
10 df4["pluno"] = (df4["pluno"]/1000).astype(int)
11 df1 = df1.sort_values(by=["vipno"])
12 vipno_series = df1["vipno"].drop_duplicates()
13 vipno_series = vipno_series.reset_index(drop=True)
14 print(vipno_series)
15 df1 = df1.sort_values(by=["pluno"])
16 pluno_series1 = df1["pluno"].drop_duplicates()
17 pluno_series1 = pluno_series1.reset_index(drop=True)
18 print(pluno_series1)
19 df2 = df2.sort_values(by=["pluno"])
20 pluno_series2 = df2["pluno"].drop_duplicates()
21 pluno_series2 = pluno_series2.reset_index(drop=True)
22 print(pluno_series2)
23 df3 = df3.sort_values(by=["pluno"])
24 pluno_series3 = df3["pluno"].drop_duplicates()
25 pluno_series3 = pluno_series3.reset_index(drop=True)
26 print(pluno_series3)
27 df4 = df4.sort_values(by=["pluno"])
28 pluno_series4 = df4["pluno"].drop_duplicates()
29 pluno_series4 = pluno_series4.reset_index(drop=True)
30 print(pluno_series4)
31 # 分组，求和
32 group_data1 = df1.groupby(["vipno", "pluno"])["amt"].sum()
33 group_data2 = df2.groupby(["vipno", "pluno"])["amt"].sum()
34 group_data3 = df3.groupby(["vipno", "pluno"])["amt"].sum()
35 group_data4 = df4.groupby(["vipno", "pluno"])["amt"].sum()
36 print(group_data1)
```

运行结果：

vipno	pluno	
781924	10	12.30
	11	72.20
	14	63.47
	15	151.70
	23	18.00
	34	15.00
13325038116	10	1272.00
	14	5.60
	15	98.50
	22	27.54
	23	125.00
	24	12.68
13854627199	27	28.81
	14	190.87
	22	196.03
	23	28.00

步骤二

```

1 #记录：用户486个，plu一级18个，plu二级94个，plu三级329个，plu四级979个
2 def jaccard_dist(a, b):
3     fenzi = 0
4     fenmu = 0
5     for i in range(0,18):
6         fenzi += min(a[i],b[i])
7         fenmu += max(a[i],b[i])
8     sim1 = fenzi/fenmu
9     fenzi = 0
10    fenmu = 0
11    for i in range(18,18+94):
12        fenzi += min(a[i],b[i])
13        fenmu += max(a[i],b[i])
14    sim2 = fenzi/fenmu
15    fenzi = 0
16    fenmu = 0
17    for i in range(18+94,18+94+329):
18        fenzi += min(a[i],b[i])
19        fenmu += max(a[i],b[i])
20    sim3 = fenzi/fenmu
21    fenzi = 0
22    fenmu = 0
23    for i in range(18+94+329,18+94+329+979):
24        fenzi += min(a[i],b[i])
25        fenmu += max(a[i],b[i])
26    sim4 = fenzi/fenmu
27    return 1-(sim1+sim2+sim3+sim4)/4

```

步骤三

```

1 def initCentroids(dataset, k):#dataset-数据点数组 k-设置的质点数

```

```

2     #初始化质心
3     numSamples, dim = dataSet.shape#numSample-数据点个数 dim-数据点维数
4     #shape返回一个关于数组长宽的数组
5     centroids = np.zeros((k, dim))#centroids-存放质心的数组
6     index = random.sample(range(0, numSamples), k)#index-在零到数据点个数间的
    随机数
7     print(index)
8     for i in range(len(index)):
9         centroids[i, :] = dataSet.values[index[i], :]
10    #将随机质心存储入centroids
11    return centroids
12
13    def CP(label, k, centroids, dataSet):
14        cpnum = 0
15        for i in range(k):
16            distance = 0
17            num = 0
18            for j in range((len(label))):
19                if label[j] == i:
20                    distance +=
jaccard_dist(dataSet.values[j, :], centroids[i, :])
21                    num += 1
22            cpnum += distance/num
23        return cpnum/k
24
25
26    def getCentroid(dataSet):
27        div = len(dataSet)
28        return sum(dataSet)/div
29
30
31    def getSC(dataSet, label):
32        sum_number = 0
33        k = len(label)
34        for i in range(k):
35            ai = 0
36            bi = 0
37            anum = 0
38            bnum = 0
39            for j in range(k):
40                if label[i]==label[j]:
41                    ai +=
jaccard_dist(dataSet.values[i, :], dataSet.values[j, :])
42                    anum += 1
43                else:
44                    bi +=
jaccard_dist(dataSet.values[i, :], dataSet.values[j, :])
45                    bnum += 1
46            ai = ai / anum
47            bi = bi / bnum
48            sum_number += (bi - ai) / max(ai, bi)
49        return sum_number / k
50
51
52    def kmeans(dataSet, k):
53        #k-means算法的核心函数
54        numSamples = dataSet.shape[0]#数据点个数为数据点数组的行数
55        label=np.zeros(dataSet.shape[0])

```

```

56     clusterChanged = True#clusterChanged-表示是否需要重新分组的布尔值判定量
57
58     centroids = initCentroids(dataSet, k)#初始化质心
59
60     while clusterChanged:#需要重新分组时
61         clusterChanged = False#重置判定量为假
62         for i in range(numSamples):#遍历所有数据点
63             minDist = 100000.0#minDist-最小的数据点与质心的距离
64             minIndex = 0#minIndex-最小的链接地址
65             for j in range(k):
66                 #计算每个数据点到哪个质心的距离最小，及记录是哪一个质心
67                 distance = jaccard_dist(centroids[j, :], dataSet.values[i,
68 :])#distance-暂时存放数据点到质心的距离，这里是jaccard距离
69                 if distance < minDist:
70                     minDist = distance
71                     minIndex = j
72                 if label[i] != minIndex:#当该数据点所隶属的质心与最小链接地址不同时更
新点中的数据
73                     clusterChanged = True#重置判定量为真
74                     label[i] = minIndex#该数据点的第二列变为一个数组
75                     for j in range(k):#由新的隶属关系中更新质心位置
76                         pointsInCluster = []
77                         for m in range(len(label)):
78                             if label[m]==j:
79                                 pointsInCluster.append(dataSet.values[m, :])
80                                 centroids[j, :] = getCentroid(pointsInCluster)
81                     print(label)
82                 print("分类完成")
83                 #这里计算SC
84                 silhouette_score = getSC(data, label)
85                 #这里计算CP
86                 compactness_score = CP(label,k,centroids,dataSet)
87                 print("sc:" + str(silhouette_score))
88                 print("cp:" + str(compactness_score))
89                 return silhouette_score,compactness_score
90
91 #数据准备
92 data1 = DataFrame(0, columns=pluno_series1, index=vipno_series)
93 # print(data)
94 for i in df1.index:
95     vipno = df1['vipno'][i]
96     pluno = df1['pluno'][i]
97     amt = group_data1[vipno][pluno]
98     if math.isnan(data1[pluno][vipno]):
99         data1[pluno][vipno] = amt
100     else:
101         data1[pluno][vipno] += amt
102 data2 = DataFrame(0, columns=pluno_series2, index=vipno_series)
103 # print(data)
104 for i in df2.index:
105     vipno = df2['vipno'][i]
106     pluno = df2['pluno'][i]
107     amt = group_data2[vipno][pluno]
108     if math.isnan(data2[pluno][vipno]):
109         data2[pluno][vipno] = amt
110     else:
111         data2[pluno][vipno] += amt

```

```

112 data3 = DataFrame(0, columns=pluno_series3, index=vipno_series)
113 # print(data)
114 for i in df1.index:
115     vipno = df3['vipno'][i]
116     pluno = df3['pluno'][i]
117     amt = group_data3[vipno][pluno]
118     if math.isnan(data3[pluno][vipno]):
119         data3[pluno][vipno] = amt
120     else:
121         data3[pluno][vipno] += amt
122 data4 = DataFrame(0, columns=pluno_series4, index=vipno_series)
123 # print(data)
124 for i in df4.index:
125     vipno = df4['vipno'][i]
126     pluno = df4['pluno'][i]
127     amt = group_data4[vipno][pluno]
128     if math.isnan(data4[pluno][vipno]):
129         data4[pluno][vipno] = amt
130     else:
131         data4[pluno][vipno] += amt
132 data = pd.concat([data1, data2, data3, data4],axis=1)
133
134
135 silhouette_score_array = []
136 for i in range(2,40):#从k为2到k为39, 尝试一下
137     silhouette_score_array.append(kmeans(data, i))

```

运行结果:

```

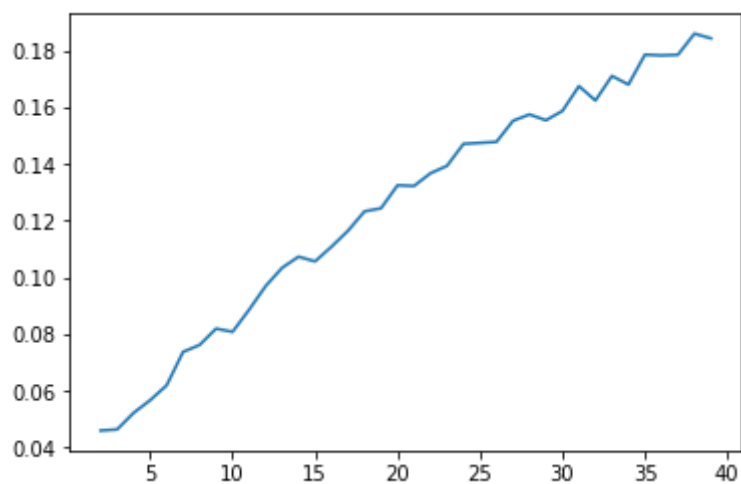
1 [(0.046009060242664455, 0.8703314634660553), (0.04639800140986168,
0.8507745782471007), (0.05223778588281045, 0.8456376368168257),
(0.0566935590430981, 0.8431493097538644), (0.06192274837863743,
0.8381479902059875), (0.07376180474905145, 0.8213263839704223),
(0.07619498835114803, 0.8275969889541755), (0.08196899003053335,
0.8183985455356022), (0.08084097548653246, 0.8227424299647964),
(0.08860153147888537, 0.823448835207196), (0.09692827131098074,
0.8059556649567884), (0.10340349963964743, 0.804504663043267),
(0.10734198666238734, 0.795502089742313), (0.10567994045519696,
0.8000799221307303), (0.11088282728752352, 0.7987122184817311),
(0.11648661299222518, 0.7925162285447176), (0.12333206200997812,
0.7843122079991169), (0.12441300679296502, 0.7877305175020182),
(0.13252647848121082, 0.7833612522715192), (0.13229130338814887,
0.7103942757034106), (0.13677147198394154, 0.7354396374295498),
(0.13937959430553268, 0.7729682405904812), (0.14715264498363456,
0.7618160925415621), (0.147492899268557, 0.7668674654821456),
(0.14784667225537074, 0.7691730476012262), (0.1552782580506372,
0.7601830177832005), (0.15751860731467987, 0.7589987317722463),
(0.15546415700476252, 0.7367407975129647), (0.15877971943254068,
0.7579979619890335), (0.16751464692397214, 0.7347685734583977),
(0.16244209603561277, 0.7228220673242257), (0.1711096423479383,
0.7272031881809556), (0.16806317873181773, 0.7289864361421701),
(0.17860849824890354, 0.7394203609998274), (0.1783595849139852,
0.7344277601077048), (0.1785540606823688, 0.7393803836497161),
(0.18602698564752815, 0.7397692607809426), (0.1843704923390654,
0.7173760381444888)]

```

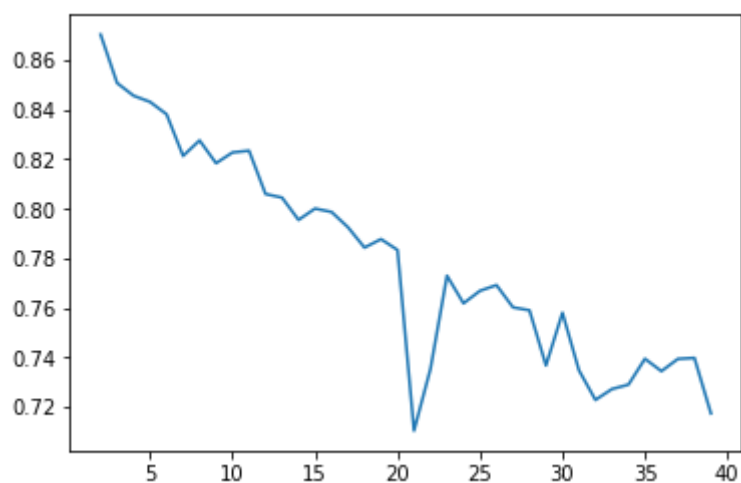
K	SC	CP
2	0.046009	0.870331
3	0.046398	0.850775
4	0.052238	0.845638
5	0.056694	0.843149
6	0.061923	0.838148
7	0.073762	0.821326
8	0.076195	0.827597
9	0.081969	0.818399
10	0.080841	0.822742
11	0.088602	0.823449
12	0.096928	0.805956
13	0.103403	0.804505
14	0.107342	0.795502
15	0.105680	0.800080
16	0.110883	0.798712
17	0.116487	0.792516
18	0.123332	0.784312
19	0.124413	0.787731
20	0.132526	0.783361
21	0.132291	0.710394
22	0.136771	0.735440
23	0.139380	0.772968
24	0.147153	0.761816
25	0.147493	0.766867
26	0.147847	0.769173
27	0.155278	0.760183
28	0.157519	0.758999
29	0.155464	0.736741
30	0.158780	0.757998
31	0.167515	0.734769

K	SC	CP
32	0.162442	0.722822
33	0.171110	0.727203
34	0.168063	0.728986
35	0.178608	0.739420
36	0.178360	0.734428
37	0.178554	0.739380
38	0.186027	0.739769
39	0.184370	0.717376

SC图像



CP图像



结论

可以看到SC和CP的走向和第一种方法差不多，不过由于这次的计算时间要超过第一次，所以我就选取了 $k=2 \sim k=39$ ，结论一样，随着 k 的增大，聚类的效果也在变好。于是我又尝试着寻找拐点：

```

1  #in
2  test_sc_cp = kmeans(data, 200)
3  #out
4  分类完成
5  sc:0.5262420540508803
6  cp:0.3665441372058841
7
8  #in
9  test_sc_cp = kmeans(data, 300)
10 #out
11 分类完成
12 sc:0.7012315838613726
13 cp:0.20827258099201829

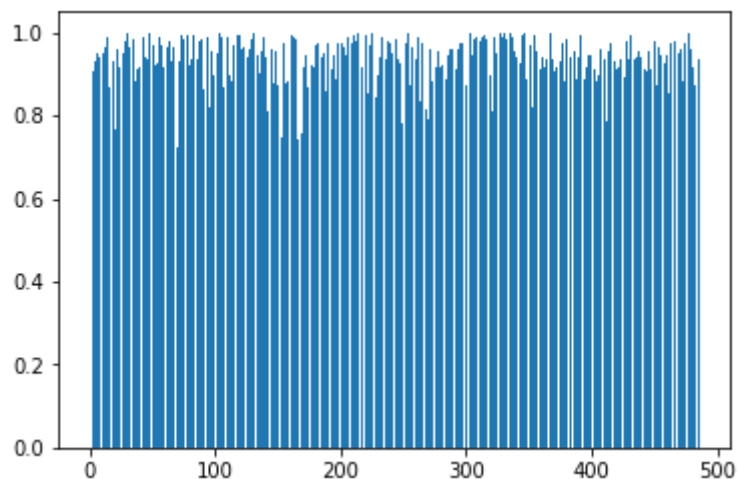
```

可以看到，在k等于300时，仍然没有发现拐点，而聚类将类别划分得这么细已经没有实际意义。所以我尝试着寻找原因，同样使用了距离分布表：

```

1  distance = []
2  for i in range(len(data)):
3      distance.append(jaccard_dist(data.values[0,:], data.values[i,:]))
4  plt.bar(range(len(data)), distance)

```



大体上和第一种方法差不多，大部分距离都分布在1附近，不过已经照比1方法强了很多，可以看到还是有不少点间距离达到0.8以下的。不过作为kmeans算法的距离，太过于敏感了，这么大的距离就意味着聚在一起会严重影响评价，所以效果虽然强于第一种，但仍然不推荐。