# 数据挖掘作业q3

作者：王星洲

学号：1652977

## 步骤一

```python
#省略读取数据，先进行数据准备
df = df.sort_values(by=["sldatime"])
time_series = df["sldatime"].drop_duplicates()
time_series = time_series.reset_index(drop=True)
print(time_series)
df = df.sort_values(by=["vipno"])
vipno_series = df["vipno"].drop_duplicates()
vipno_series = vipno_series.reset_index(drop=True)
print(vipno_series)
df = df.sort_values(by=["pluno"])
pluno_series = df["pluno"].drop_duplicates()
pluno_series = pluno_series.reset_index(drop=True)
print(pluno_series)
df1 = df1.sort_values(by=["pluno"])
pluno_series1 = df1["pluno"].drop_duplicates()
pluno_series1 = pluno_series1.reset_index(drop=True)
df2 = df2.sort_values(by=["pluno"])
pluno_series2 = df2["pluno"].drop_duplicates()
pluno_series2 = pluno_series2.reset_index(drop=True)
df3 = df3.sort_values(by=["pluno"])
pluno_series3 = df3["pluno"].drop_duplicates()
pluno_series3 = pluno_series3.reset_index(drop=True)
df4 = df4.sort_values(by=["pluno"])
pluno_series4 = df4["pluno"].drop_duplicates()
pluno_series4 = pluno_series4.reset_index(drop=True)
# 分组，求和
print(df)
group_data1 = df1.groupby(["vipno","sldatime","pluno"])["qty"].sum()
group_data1
group_data2 = df2.groupby(["vipno","sldatime","pluno"])["qty"].sum()
group_data2
group_data3 = df3.groupby(["vipno","sldatime","pluno"])["qty"].sum()
group_data3
group_data4 = df4.groupby(["vipno","sldatime","pluno"])["qty"].sum()
group_data4
group_data = df.groupby(["vipno","sldatime","pluno"])["qty"].sum()
print(group_data)
#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  购买时间最晚7月31日，最早2月1日

#评价时间和对应的显示层级
def judgeTimeLevel(a):
    delta = datetime.datetime(2016,7,31)-
datetime.datetime(int(a[0:4]),int(a[5:7]),int(a[8:10]))
```

```
43  #       0全部显示
44      if delta < datetime.timedelta(days=30):
45          return 0
46  #       1显示到四级
47      elif delta < datetime.timedelta(days=60):
48          return 1
49  #       2显示到3级
50      elif delta < datetime.timedelta(days=120):
51          return 2
52  #       3显示到2级
53      else:
54          return 3
```

## 步骤二

```
1   #为每一位用户搭建用户树
2   class Point1:
3       def __init__(self, kind, qty):
4           self.kind = kind
5           self.qty = qty
6
7
8   tree_array = []
9   for i in vipno_series:
10      tree = Tree()
11      tree.create_node("Root","root")
12      data = group_data[i]
13      point_array = []
14      for p,q in data.items():
15          time = p[0]
16          kind = p[1]
17          kind1 = int(kind/1000000)
18          kind2 = int(kind/100000)
19          kind3 = int(kind/10000)
20          kind4 = int(kind/1000)
21          qty = q
22          showLevel = judgeTimeLevel(time)
23          print("time: " + str(time) + " level: " + str(showLevel))
24          if showLevel == 0:
25              if kind1 not in point_array:
26                  point_array.append(kind1)
27
    tree.create_node(kind1,kind1,parent="root",data=Point1(kind1,qty))
28              else:
29                  tree.nodes[kind1].data.time = time
30                  tree.nodes[kind1].data.qty += qty
31              if kind2 not in point_array:
32                  point_array.append(kind2)
33
    tree.create_node(kind2,kind2,parent=kind1,data=Point1(kind2,qty))
34              else:
35                  tree.nodes[kind2].data.time = time
36                  tree.nodes[kind2].data.qty += qty
37              if kind3 not in point_array:
38                  point_array.append(kind3)
39
    tree.create_node(kind3,kind3,parent=kind2,data=Point1(kind3,qty))
```

```
40                    else:
41                        tree.nodes[kind3].data.time = time
42                        tree.nodes[kind3].data.qty += qty
43                    if kind4 not in point_array:
44                        point_array.append(kind4)

       tree.create_node(kind4,kind4,parent=kind3,data=Point1(kind4,qty))
46                    else:
47                        tree.nodes[kind4].data.time = time
48                        tree.nodes[kind4].data.qty += qty
49                    if kind not in point_array:
50                        point_array.append(kind)

       tree.create_node(kind,kind,parent=kind4,data=Point1(kind,qty))
52                    else:
53                        tree.nodes[kind].data.time = time
54                        tree.nodes[kind].data.qty += qty
55            elif showLevel == 1:
56                    if kind1 not in point_array:
57                        point_array.append(kind1)

       tree.create_node(kind1,kind1,parent="root",data=Point1(kind1,qty))
59                    else:
60                        tree.nodes[kind1].data.time = time
61                        tree.nodes[kind1].data.qty += qty
62                    if kind2 not in point_array:
63                        point_array.append(kind2)

       tree.create_node(kind2,kind2,parent=kind1,data=Point1(kind2,qty))
65                    else:
66                        tree.nodes[kind2].data.time = time
67                        tree.nodes[kind2].data.qty += qty
68                    if kind3 not in point_array:
69                        point_array.append(kind3)

       tree.create_node(kind3,kind3,parent=kind2,data=Point1(kind3,qty))
71                    else:
72                        tree.nodes[kind3].data.time = time
73                        tree.nodes[kind3].data.qty += qty
74                    if kind4 not in point_array:
75                        point_array.append(kind4)

       tree.create_node(kind4,kind4,parent=kind3,data=Point1(kind4,qty))
77                    else:
78                        tree.nodes[kind4].data.time = time
79                        tree.nodes[kind4].data.qty += qty
80            elif showLevel == 2:
81                    if kind1 not in point_array:
82                        point_array.append(kind1)

       tree.create_node(kind1,kind1,parent="root",data=Point1(kind1,qty))
84                    else:
85                        tree.nodes[kind1].data.time = time
86                        tree.nodes[kind1].data.qty += qty
87                    if kind2 not in point_array:
88                        point_array.append(kind2)

       tree.create_node(kind2,kind2,parent=kind1,data=Point1(kind2,qty))
```
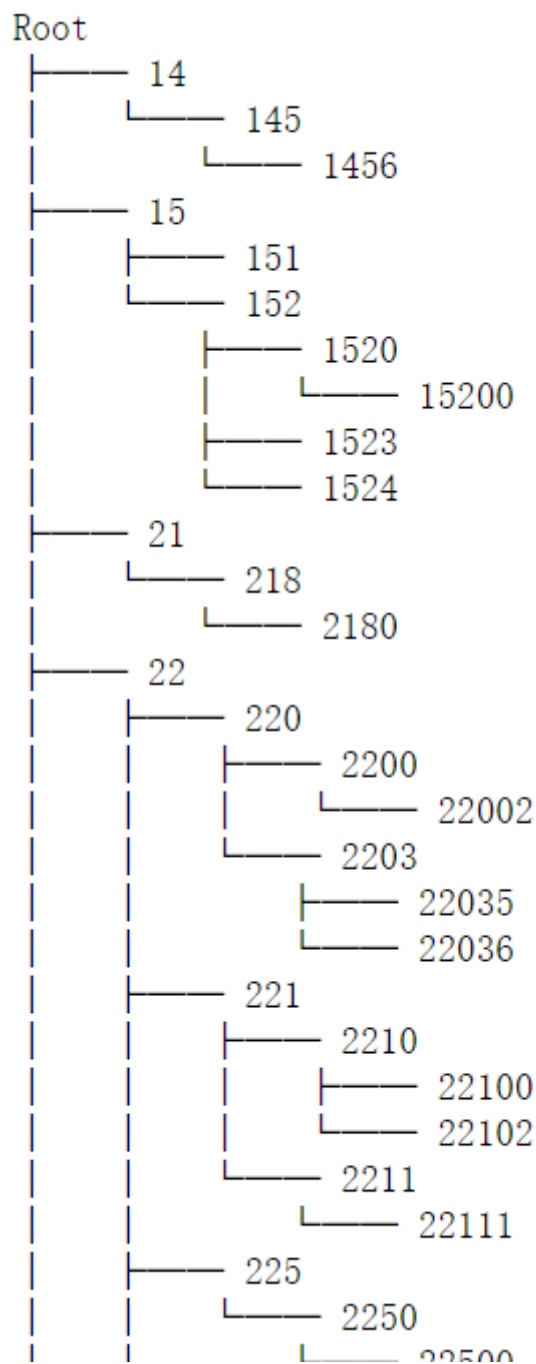
```python
            else:
                tree.nodes[kind2].data.time = time
                tree.nodes[kind2].data.qty += qty
            if kind3 not in point_array:
                point_array.append(kind3)

                tree.create_node(kind3,kind3,parent=kind2,data=Point1(kind3,qty))
            else:
                tree.nodes[kind3].data.time = time
                tree.nodes[kind3].data.qty += qty
        elif showLevel == 3:
            if kind1 not in point_array:
                point_array.append(kind1)

                tree.create_node(kind1,kind1,parent="root",data=Point1(kind1,qty))
            else:
                tree.nodes[kind1].data.time = time
                tree.nodes[kind1].data.qty += qty
            if kind2 not in point_array:
                point_array.append(kind2)

                tree.create_node(kind2,kind2,parent=kind1,data=Point1(kind2,qty))
            else:
                tree.nodes[kind2].data.time = time
                tree.nodes[kind2].data.qty += qty
    tree_array.append(tree)
    tree.show()
```

运行结果:

```
Root
├── 14
│       └── 145
│               └── 1456
├── 15
│       ├── 151
│       └── 152
│               ├── 1520
│               │       └── 15200
│               ├── 1523
│               └── 1524
├── 21
│       └── 218
│               └── 2180
├── 22
│       ├── 220
│       │       ├── 2200
│       │       │       └── 22002
│       │       └── 2203
│       │               ├── 22035
│       │               └── 22036
│       ├── 221
│       │       ├── 2210
│       │       │       ├── 22100
│       │       │       └── 22102
│       │       └── 2211
│       │               └── 22111
│       ├── 225
│       │       └── 2250
│       │               └── 22500
```

## 步骤三

```python
#接下来是实现树之间的距离计算
def IR(tree1, tree2):
    result = Tree()
    result.create_node("Root","root")
    for i in tree1.nodes:
        for j in tree2.nodes:
            if i==j and i!="root":

 result.create_node(i,i,parent=tree1.nodes[i].predecessor(tree1.identifier
),data=Point1(i,tree1.nodes[i].data.qty + tree2.nodes[i].data.qty))
    return result


def UR(tree1, tree2):
    result = Tree()
```

```
14        result.create_node("Root","root")
15        i_cluster = []
16        ii_cluster = []
17        jj_cluster = []
18        for i in tree1.nodes:
19            if i != "root":
20                ii_cluster.append(i)
21        for j in tree2.nodes:
22            if j != "root":
23                if j in ii_cluster:
24                    i_cluster.append(j)
25                    ii_cluster.remove(j)
26                else:
27                    jj_cluster.append(j)
28        for i in i_cluster:
29
     result.create_node(i,i,parent=tree1.nodes[i].predecessor(tree1.identifier
    ),data=Point1(i,tree1.nodes[i].data.qty + tree2.nodes[i].data.qty))
30        for i in ii_cluster:
31
     result.create_node(i,i,parent=tree1.nodes[i].predecessor(tree1.identifier
    ),data=Point1(i,tree1.nodes[i].data.qty))
32        for i in jj_cluster:
33
     result.create_node(i,i,parent=tree2.nodes[i].predecessor(tree2.identifier
    ),data=Point1(i,tree2.nodes[i].data.qty))
34        return result
35
36
37  def FTC_dist(tree1, tree2):
38        i_tree = IR(tree1, tree2)
39  #      i_tree.show()
40        u_tree = UR(tree1, tree2)
41  #      u_tree.show()
42      v1 = 0
43      n1 = 0
44      v2 = 0
45      n2 = 0
46      v3 = 0
47      n3 = 0
48      v4 = 0
49      n4 = 0
50      v5 = 0
51      n5 = 0
52      div_root = 0
53      for i in u_tree.nodes["root"].successors(u_tree.identifier):
54          div_root += u_tree.nodes[i].data.qty
55      for i in i_tree.nodes:
56          if i != "root":
57              dep = len(str(i)) - 1
58  #              print("i: " + str(i) + "dep: " + str(dep))
59              if dep == 1:
60                  div = div_root
61  #                  print("div: " + str(div))
62                  if div != 0:
63                      v1 += i_tree.nodes[i].data.qty / div
64  #                      print("qty: " + str(i_tree.nodes[i].data.qty) +
    "v1_now: " + str(v1))
```

```python
                else:
                    v1 = 0
                n1 += 1
#                   print("n1_now: " + str(n1))
            elif dep == 2:
                div = 0
                for j in
u_tree.get_node(u_tree.nodes[i].predecessor(u_tree.identifier)).successors
(u_tree.identifier):
                    div += u_tree.nodes[j].data.qty
#                   print("div: " + str(div))
                if div != 0:
                    v2 += i_tree.nodes[i].data.qty / div
#                       print("qty: " + str(i_tree.nodes[i].data.qty) +
"v2_now: " + str(v2))
                else:
                    v2 = 0
                n2 += 1
#                   print("n2_now: " + str(n2))
            elif dep == 3:
                div = 0
                for j in
u_tree.get_node(u_tree.nodes[i].predecessor(u_tree.identifier)).successors
(u_tree.identifier):
                    div += u_tree.nodes[j].data.qty
#                   print("div: " + str(div))
                if div != 0:
                    v3 += i_tree.nodes[i].data.qty / div
#                       print("qty: " + str(i_tree.nodes[i].data.qty) +
"v3_now: " + str(v3))
                else:
                    v3 = 0
                n3 += 1
#                   print("n3_now: " + str(n3))
            elif dep == 4:
                div = 0
                for j in
u_tree.get_node(u_tree.nodes[i].predecessor(u_tree.identifier)).successors
(u_tree.identifier):
                    div += u_tree.nodes[j].data.qty
#                   print("div: " + str(div))
                if div != 0:
                    v4 += i_tree.nodes[i].data.qty / div
#                       print("qty: " + str(i_tree.nodes[i].data.qty) +
"v4_now: " + str(v4))
                else:
                    v4 = 0
                n4 += 1
#                   print("n4_now: " + str(n4))
            else:
                div = 0
                for j in
u_tree.get_node(u_tree.nodes[i].predecessor(u_tree.identifier)).successors
(u_tree.identifier):
                    div += u_tree.nodes[j].data.qty
#                   print("div: " + str(div))
                if div != 0:
                    v5 += i_tree.nodes[i].data.qty / div
```

```python
112 #                          print("qty: " + str(i_tree.nodes[i].data.qty) +
    "v5_now: " + str(v5))
113                     else:
114                         v5 = 0
115                     n5 += 1
116 #                      print("n5_now: " + str(n5))
117 #     print("v1: " + str(v1))
118 #     print("v2: " + str(v2))
119 #     print("v3: " + str(v3))
120 #     print("v4: " + str(v4))
121 #     print("v5: " + str(v5))
122 #     print("n1: " + str(n1))
123 #     print("n2: " + str(n2))
124 #     print("n3: " + str(n3))
125 #     print("n4: " + str(n4))
126 #     print("n5: " + str(n5))
127     if n1 != 0:
128         sim1 = v1/n1
129 #         print(sim1)
130     else:
131         sim1 = 0
132     if n2 != 0:
133         sim2 = v2/n2
134 #         print(sim2)
135     else:
136         sim2 = 0
137     if n3 != 0:
138         sim3 = v3/n3
139 #         print(sim3)
140     else:
141         sim3 = 0
142     if n4 != 0:
143         sim4 = v4/n4
144 #         print(sim4)
145     else:
146         sim4 = 0
147     if n5 != 0:
148         sim5 = v5/n5
149 #         print(sim5)
150     else:
151         sim5 = 0
152     if sim2 == 0:
153         result = 1 - sim1
154     elif sim3 == 0:
155         result = 1 - (sim1/3 + sim2*2/3)
156     elif sim4 == 0:
157         result = 1 - (sim1/6 + sim2*2/6 + sim3*3/6)
158     elif sim5 == 0:
159         result = 1 - (sim1/10 + sim2*2/10 + sim3*3/10 + sim4*4/10)
160     else:
161         result = 1 - (sim1/15 + sim2*2/15 + sim3*3/15 + sim4*4/15 +
    sim5*5/15)
162     return result
163
164
165 #测试
166 tree1 = Tree()
167 tree2 = Tree()
```

```
168  tree1.create_node("Root","root")
169  tree2.create_node("Root","root")
170  tree1.create_node(10,10,parent="root",data=Point1(10, 2))
171  tree1.create_node(101,101,parent=10,data=Point1(101, 2))
172  tree1.create_node(1012,1012,parent=101,data=Point1(1012, 2))
173  tree1.create_node(10128,10128,parent=1012,data=Point1(10128, 1))
174  tree1.create_node(10129,10129,parent=1012,data=Point1(10129, 1))
175  tree1.create_node(12,12,parent="root",data=Point1(12, 1))
176  tree1.create_node(122,122,parent=12,data=Point1(122, 1))
177  tree2.create_node(10,10,parent="root",data=Point1(10, 2))
178  tree2.create_node(101,101,parent=10,data=Point1(101, 2))
179  tree2.create_node(1012,1012,parent=101,data=Point1(1012, 1))
180  tree2.create_node(10128,10128,parent=1012,data=Point1(10128, 1))
181  tree2.create_node(1013,1013,parent=101,data=Point1(1013, 1))
182  tree2.create_node(13,13,parent="root",data=Point1(13, 1))
183  test_tree_array = [tree1, tree2]
184  FTC_dist(tree1,tree2)
```

运行结果:

```
1  0.2416666666666667
```

## 步骤四

```
1   #实现质心算法
2   def UTREE(tree_array):
3       print("正在合并大并集树")
4       utree = tree_array[0]
5       for i in range(len(tree_array) - 1):
6           utree = UR(utree, tree_array[i+1])
7       print("合并完成")
8   #     utree.show()
9       return utree
10
11
12  def updateTree(tree, freq):
13      flag = []
14      for i in tree.nodes:
15          if i != "root":
16              if tree.nodes[i].data.qty < freq:
17                  flag.append(i)
18      for i in flag:
19          if tree.get_node(i):
20              tree.remove_node(i)
21      return tree
22
23
24  def GetCT(tree_array):
25      print("正在计算树的质心")
26      ct = Tree()
27      utree = UTREE(tree_array)
28      max_freq = 0
29      sum_freq = 0
30      freq = 1
31      mindist = 999999
32      num_avg_nodes = 0
```

```
33        for q in tree_array:
34            num_avg_nodes += len(q.all_nodes())
35        num_avg_nodes /= len(tree_array)
36        for i in utree.nodes:
37            if i != "root":
38                dep = len(str(i)) - 1
39                if dep == 1:
40                    if utree.nodes[i].data.qty > max_freq:
41                        max_freq = utree.nodes[i].data.qty
42                sum_freq += utree.nodes[i].data.qty
43        avg_freq = sum_freq / (len(utree.nodes) - 1)
44        while freq <= max_freq:
45            utree = updateTree(utree, freq)
46            if len(utree.all_nodes()) <= num_avg_nodes:
47                break
48            dist = 0
49            for i in tree_array:
50                dist += FTC_dist(i,utree)
51            if dist < mindist:
52                mindist = dist
53                ct = utree
54            freq = freq + avg_freq
55        print("计算完成")
56        return ct
```

## 步骤五

```
1   #设计针对FTC树的Kmeans算法
2   def initCentroids(dataSet, k):#dataSet-数据点数组 k-设置的质心数
3       #初始化质心
4       print("正在选择随机质心")
5       centroids = []
6       index = random.sample(range(0, len(dataSet)), 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.append(dataSet[index[i]])
10          #将随机质心存储入centroids
11      print("选择完成")
12      return centroids
13
14
15  def kmeans(dataSet, k):
16      #k-means算法的核心函数
17      numSamples = len(dataSet)#数据点个数为数据点数组的行数
18      label = np.zeros(numSamples)
19      clusterChanged = True#clusterChanged-表示是否需要重新分组的布尔值判定量
20
21      centroids = initCentroids(dataSet, k)#初始化质心
22      step = 0
23  #    print(centroids)
24      while clusterChanged:#需要重新分组时
25          clusterChanged = False#重置判定量为假
26          for i in range(numSamples):#遍历所有数据点
27              minDist = 100000.0#minDist-最小的数据点与质心的距离
28              minIndex = 0#minIndex-最小的链接地址
29              for j in range(k):
```

```python
30                        #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
31                    distance = FTC_dist(centroids[j], dataSet[i])#distance-暂时
     存放数据点到质心的距离，这里是FTC距离
32                        if distance < minDist:
33                            minDist = distance
34                            minIndex = j
35                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
36                        clusterChanged = True#重置判定量为真
37                        label[i] = minIndex
38            print(label)
39  #           print(label)
40            for j in range(k):#由新的隶属关系中更新质心位置
41                pointsInCluster = []
42                for m in range(numSamples):
43                    if label[m] == j:
44                        pointsInCluster.append(dataSet[m])
45  #                   print("m: " + str(m) + " label[m]: " + str(label[m]))

46  #               print(pointsInCluster)
47                centroids[j] = GetCT(pointsInCluster)
48  #           print(centroids)
49            step += 1
50            if step >= 40:
51                break
52        print("分类完成")
53        tree_array_array = []
54        for i in range(k):
55            tree_array = []
56            for j in range(numSamples):
57                if label[j] == i:
58                    tree_array.append(dataSet[j])
59            tree_array_array.append(tree_array)
60        return tree_array_array
```

## 步骤六

```python
1   # 设计BIC算法
2   # Ni该簇用户数，k簇数量，ct质心树，C该簇树列表
3   def variance(Ni,ct,k,C):
4       print("正在计算方差")
5       result = 1/(Ni-k)
6       sum_dist = 0
7       for i in C:
8           sum_dist += math.pow(FTC_dist(i,ct),2)
9       result *= sum_dist
10      print("variance: " + str(result))
11      return result
12
13
14  # N用户数，k簇数，D该簇商品数，ct质心，C该簇树列表
15  def L(N,D,C,ct,k):
16      print("正在计算L")
17      sum_num = 0
18      for i in range(k):
```

```
19        sum_num += N[i] * math.log(N[i])- N[i]*math.log(sum(N)) -
      N[i]/2*math.log(2*math.pi) - N[i]*D/2* math.log(variance(N[i],ct[i],k,C)) -
      (N[i] - k)/2
20      result = sum_num
21      print("计算完成")
22      return result
23

24
25  def BIC(C, array):
26      print("BIC启动")
27      #k=1
28      N = [len(C)]
29      D = len(UTREE(C).leaves())
30      ct = [GetCT(C)]
31      ct[0].show()
32      l = L(N,D,C,ct,1)
33      sub = 1/2*(D+1)*math.log(1)
34      print("N: " + str(N) + " D: " + str(D) + " l: " + str(l) + " sub: " +
      str(sub))
35      result1 = l - sub
36      print("result1: " + str(result1))
37      # k=2
38      tree_array_array = kmeans(C, 2)
39      tree_array1 = tree_array_array[0]
40      tree_array2 = tree_array_array[1]
41      N1 = [len(tree_array1),len(tree_array2)]
42      ct = [GetCT(tree_array1),GetCT(tree_array2)]
43      l1 = L(N1,D,C,ct,2)
44      sub1 = (D+1)*math.log(2)
45      result2 = l1 - sub1
46      print("result1: " + str(result1))
47      print("result2: " + str(result2))
48      if result2 <= result1:
49          print("无需分组")
50          array.append(C)
51      else:
52          BIC(tree_array1, array)
53          BIC(tree_array2, array)
54

55
56  result_array = []
57  a = kmeans(tree_array, 2)
58  for i in a:
59      BIC(i, result_array)
60  k = len(result_array)
61  k
```

```
1  #out
2  2
```

## 步骤七

```
1  #设计聚类效果的评价函数
2  def getSC(tree_array, result_array):
3      sum_number = 0
4      for i in range(len(tree_array)):
```

```
 5                ai = 0
 6                bi = 0
 7                anum = 0
 8                bnum = 0
 9                for j in range(len(tree_array)):
10                    flag = False
11                    for m in range(len(result_array)):
12                        if tree_array[i] in result_array[m] and tree_array[j] in
   result_array[m]:
13                            ai += FTC_dist(tree_array[i], tree_array[j])
14                            anum += 1
15                            flag = True
16                            break
17                    if flag == False:
18                        bi += FTC_dist(tree_array[i], tree_array[j])
19                        bnum += 1
20                ai = ai / anum
21                bi = bi / bnum
22                sum_number += (bi - ai) / max(ai, bi)
23        return sum_number / len(tree_array)
24
25
26  def getCP(tree_array, result_array):
27      k = len(result_array)
28      cpnum = 0
29      for i in result_array:
30          distance = 0
31          ct = GetCT(i)
32          num = 0
33          for j in range(len(tree_array)):
34              if tree_array[j] in i:
35                  distance += FTC_dist(tree_array[j], ct)
36                  num += 1
37          cpnum += distance/num
38      return cpnum/k
39
40
41  sc = getSC(tree_array, result_array)
42  cp = getCP(tree_array, result_array)
43  print(sc)
44  print(cp)
```

运行结果:

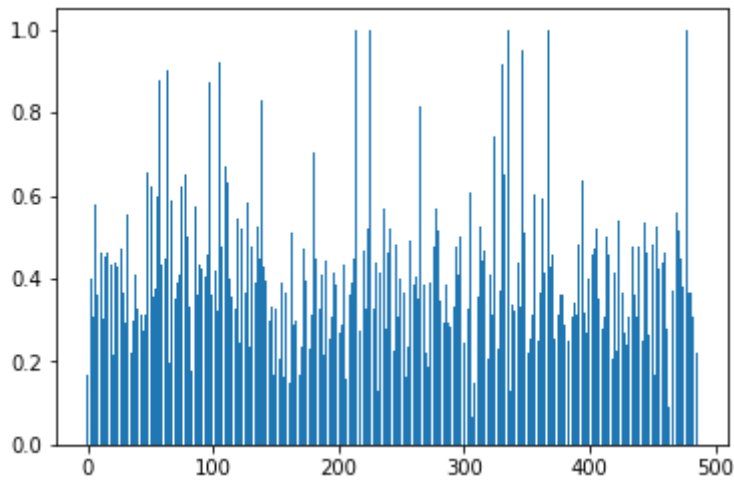```
1  sc:0.03001777457340002
2  cp:0.28558352013812394
```

## 结论

使用FTC树进行分析，在计算用户之间距离的时候有一个非常好的参照，这样计算出来的距离很有道理并且数据也很好，接下来我同样展示一下距离分布:

```
1  distance = []
2  for i in range(len(tree_array)):
3      distance.append(FTC_dist(tree_array[0], tree_array[i]))
4  plt.bar(range(len(tree_array)), distance)
```

可以看到，数据分布很合理，很完美，这样的用户间距离对聚类分析是很不错的。

但是我们的算法仍然出现了问题，我来总结一下算法分析出的结论：

```
1   #以k=2为初始，利用kmeans方法分出的两个簇分别进行BIC运算，再对可以继续分的簇进行递归运算
2   #结果是：
3   k = 2
4   sc = 0.03001777457340002
5   cp = 0.28558352013812394
```

也就是说，除了开始时的主动使用kmeans分为两类以外，BIC公式的结果是不再进行分类，这还是很出乎我的意料的，因为如果利用这个距离分布，使用kmeans算法，最后的SC指标不会这么低，也不会只分为两类，那么问题出在哪里呢？

这里我找到了两个可能的原因：

1. 获取质心算法时，我在结点数小于平均结点数时跳出了循环，避免了最后全部结点消失，但是也带来了问题，就是最后有多个结点，不满足只剩下一个结点

```
1   num_avg_nodes = 0
2   for q in tree_array:
3       num_avg_nodes += len(q.all_nodes())
4   num_avg_nodes /= len(tree_array)
5   if len(utree.all_nodes()) <= num_avg_nodes:
6       break
```

这样就导致了质心选取不准确，同样的数据，质心也可能出现偏差，所以会导致数据出错。

2. 由于上述的原因，kmeans算法同样变得不稳定，导致迟迟不会收敛，因为质心跳跃式改变，每次都会有点的归属发生变化。我不得不进行强制手动收敛：

```
1   step = 0
2   for j in range(k):
3       step += 1
4   if step >= 40:
5       break
```

由于kmeans算法最后没有自动收敛，可能导致分类效果不好，BIC因此给出了不分类的判断。

总体来说，FTC树提供了一个测量用户间距离的很好的方式，这样即便直接使用FTC与kmeans结合，也可以得到很好的聚类效果。而将FTC与BIC结合进行聚类划分提供了一个不需要手动设置k值的，基于距离的优秀的聚类方式，不再依赖于k值的设定，唯一的缺点可能就是仍然受到初始质心选择的影响，但这也是kmeans算法逃不开的。

如果我使用一个其他的数据或者能够在选取质心的时候想到一个不会剪掉所有结点，又能稳定剩下一个结点的方式的话，可能效果会非常好。同样的，在选择初始质心的时候采用非随机，而是尽量较远的范围选点，效果可能会非常不错。