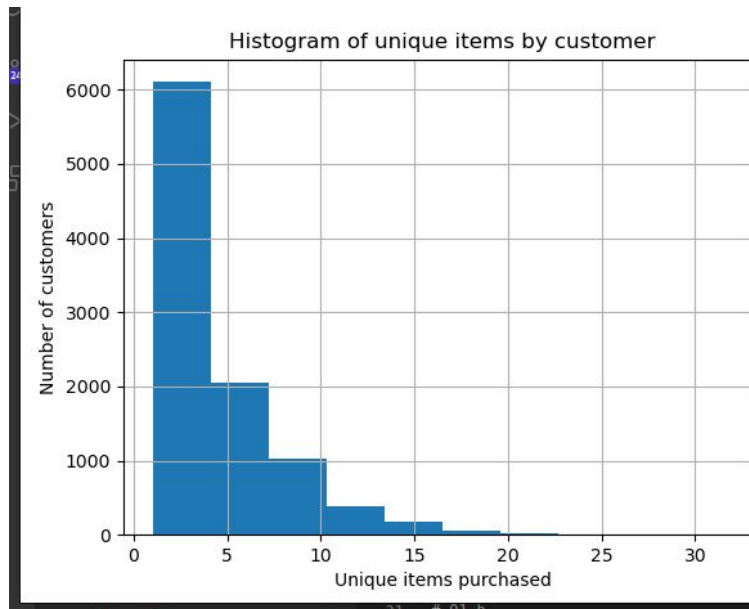


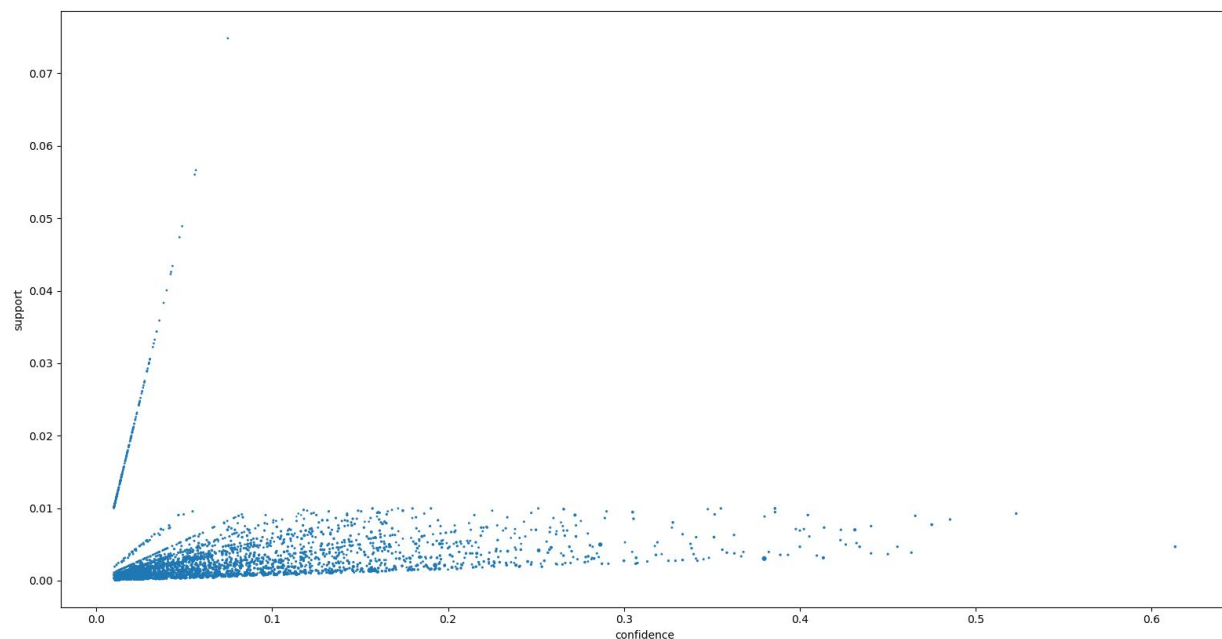
1. See the attached python file q1.py
 - a. See histogram and 5 number summary (q0-4 are the 0-100 percentiles)



```
40000010/20331213, ordered_statistics=[OrderedStatistic(items=frozenset({'Instant food products'}), items_add=frozenset({'Instant food products'}), lift=1.5136008618367895), OrderedStatistic(items=frozenset({'UHT-milk'}), items_add=frozenset({'Instant food products'}), lift=1.5136008618367895)]])
Q1.d
Q1.e
RelationRecord(items=frozenset({'rice', 'whole milk'}), items_add=frozenset({'rice', 'whole milk'}), lift=1.5136008618367895, ordered_statistics=[OrderedStatistic(items=frozenset({'rice', 'whole milk'}), items_add=frozenset({'rice', 'whole milk'}), confidence=0.642711)])
[tate@archbook assignment2]$ ^C
[tate@archbook assignment2]$ ^C
[tate@archbook assignment2]$ python q1.py
Q1.a
Five number summary:
q0 = 1
q1 = 2
q2 = 3
q3 = 6
q4 = 32
```

- b. There are 102 item-sets and the most popular item was purchased by 2513 unique customers
- c.
 - i. Using a min support of 0.0001 I found 13177 association rules with a confidence above 1%.
 - ii. When I limited the discovery to only look for two items (which makes most sense in this case) I found 4877 rules
 - iii. I used (ii) for the rest of the tests as performance was significantly improved

d.



e. Rice -> whole milk

2. As you can see here the three centroids are closest to the tea/coffee, milk, and soda/water

```
2 import pandas as pd
3 import sklearn.cluster as cluster
4
5
6 df = pd.read_csv('cafe.csv')
7 X = pd.DataFrame({'x' : list(df['Frequency'])})
8 myCluster = cluster.KMeans(n_clusters = 3, random_state = 0).fit(X)
9
10 print('Cluster Assignment:', myCluster.labels_)
11 print('Cluster Centroid 0:', myCluster.cluster_centers_[0])
12 print('Cluster Centroid 1:', myCluster.cluster_centers_[1])
13 print('Cluster Centroid 1:', myCluster.cluster_centers_[2])
14
15
```

```
Terminal - tate@archbook:~/Des
Fil Rediger Vis Terminal Faneblade Hjælp
[tate@archbook assignment2]$ python q2.py
Cluster Assignment: [1 2 2 2 0 0 1]
Cluster Centroid 0: [7.5]
Cluster Centroid 1: [80.]
Cluster Centroid 1: [28.33333333]
[tate@archbook assignment2]$
```

3. Struggling with the proof, I'm taking linear algebra in parallel with this course and haven't taken statistics yet, this problem is beyond my abilities
4. I was unable to understand the notation in the powerpoint, the book, or wikipedia. When I went to TA office hours the TA was unable to help and seemed just as lost as me. Regardless, I managed to find the Silhouette and Davies-Bouldin indecies for the

clusters.

```

2
3 import pandas
4 import sklearn.cluster
5 import sklearn.metrics
6
7 print('Loading clusters...')
8 X = pandas.DataFrame({'x' : [-2, -1, 1, 2, 3, 4, 5, 7, 8] })
9 cs = sklearn.cluster.KMeans(n_clusters = 2, random_state = 0).fit(X)
10
11 # Verfiy cluster assignments are correct
12 print('\tCluster Assignment:', cs.labels_)
13
14 ###
15 # Q4.a
16 ###
17 print('Q4.a')
18 # Compute sillouette width
19 sil = sklearn.metrics.silhouette_score(X, cs.labels_)
20 print('\t', sil)
21
22
23 ###
24 # Q4.c
25 ###
26 print('Q4.c')
27
28 # Compute davies
29 dbs = sklearn.metrics.davies_bouldin_score(X, cs.labels_)
30 print('\t', dbs)

```

Terminal - tate@archbook:~/Desktop/cs

File Rediger Vis Terminal Faneblade Hjælp

[tate@archbook assignment2]\$ python q4.py

Loading clusters...

Cluster Assignment: [0 0 0 0 0 1 1 1 1]

Q4.a

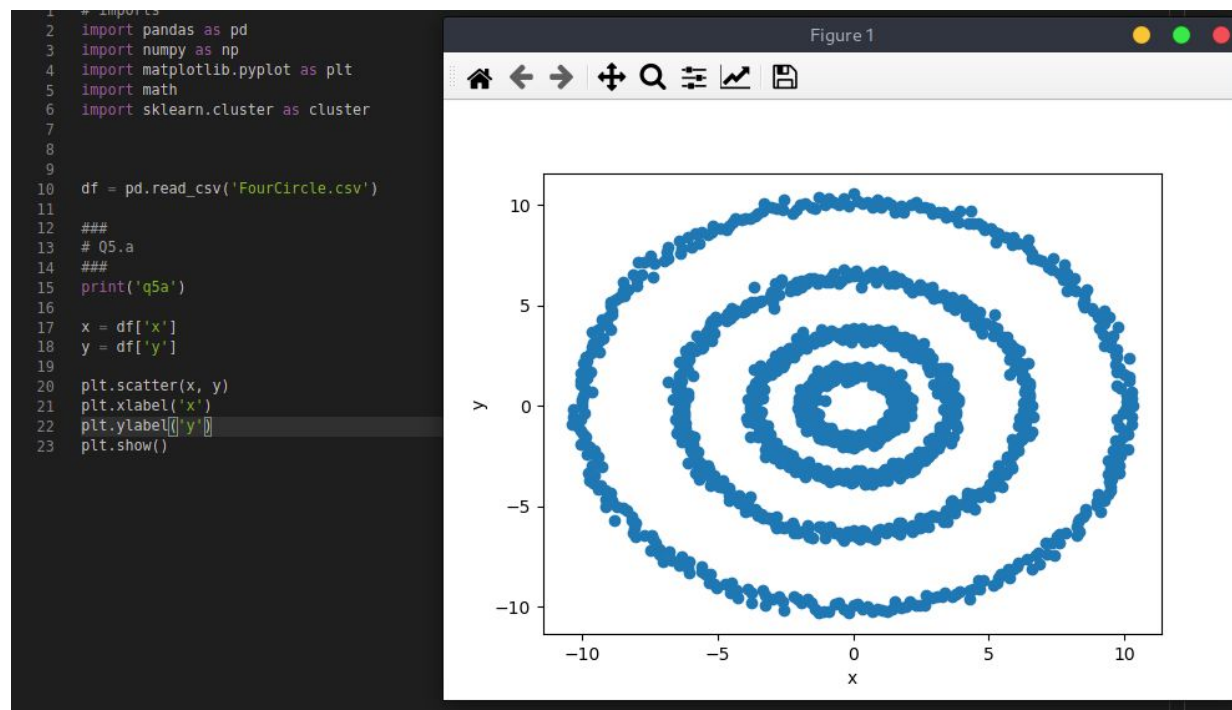
0.4847097336067924

Q4.c

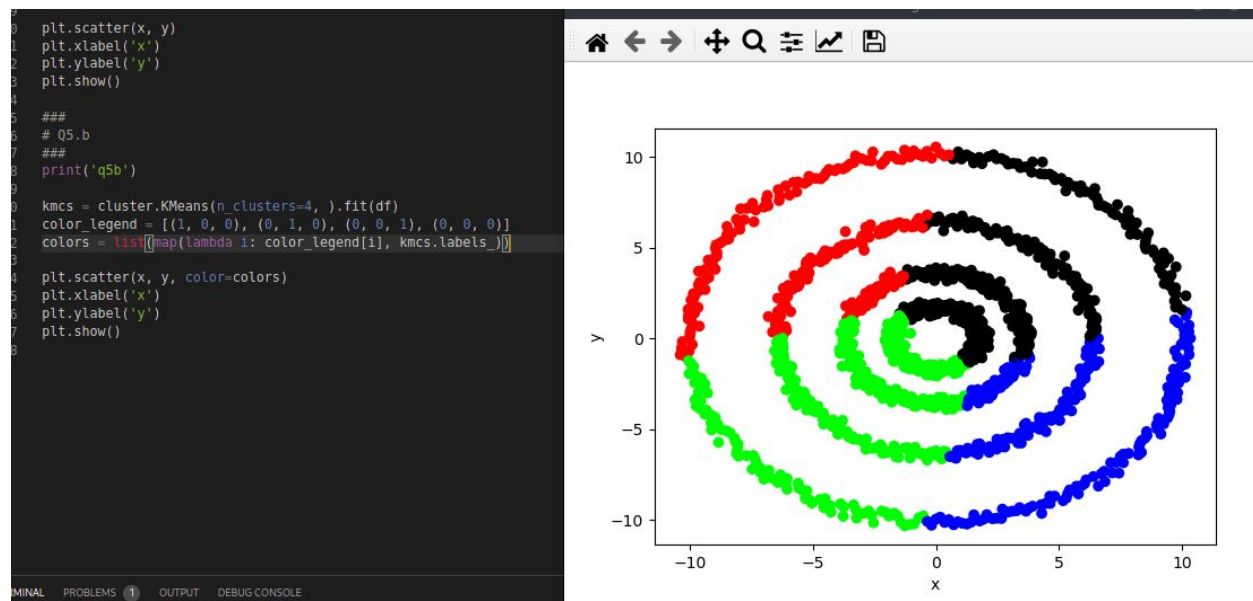
0.5888888888888889

[tate@archbook assignment2]\$

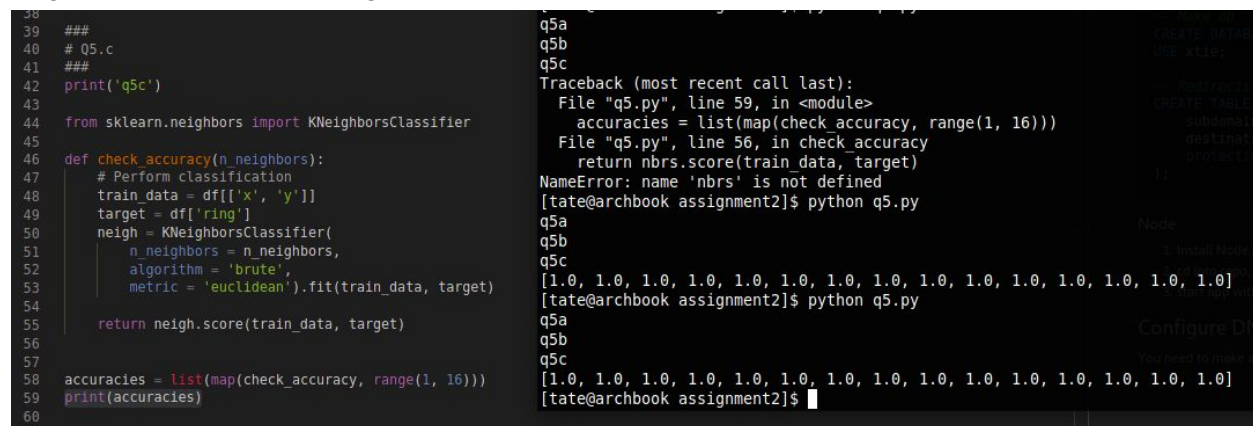
5. See program file q5.py
 - a. It appears there are 4 custers in the dataset



- b. As you can see from the screenshot, the algorithm separates them by quadrant instead of by continuity



- c. It seems that we only need one nearest neighbor to achieve needed accuracy. This makes sense in the case of the problem as any point will have a near neighbor on the relevant ring and it



- d. And e. Because the math here goes well over what I've covered in coursework, I used the spectral clustering class that's built into sklearn to produce the correct

result

```
61  
62 ###  
63 # QS.de  
64 ###  
65 print('q5e-d')  
66  
67 # Math go brrr  
68 sccs = sklearn.cluster.SpectralClustering(  
69     n_neighbors = 1,  
70     n_clusters = 4,  
71     n_components= 3,  
72 ).fit(df)  
73 color_legend = [(1, 0, 0), (0, 1, 0), (0, 0, 1), (0, 0, 0)]  
74 colors = list(map(lambda i: color_legend[i], sccs.labels_))  
75  
76 plt.scatter(x, y, color=colors)  
77 plt.xlabel('x')  
78 plt.ylabel('y')  
79 plt.show()
```

