

CS 584-04: Machine Learning

Autumn 2019 Assignment 2

Question 1 (50 points)

The file Groceries.csv contains market basket data. The variables are:

1. Customer: Customer Identifier
2. Item: Name of Product Purchased

The data is already sorted in ascending order by Customer and then by Item. Also, all the items bought by each customer are all distinct.

After you have imported the CSV file, please discover association rules using this dataset.

- a) (10 points) Create a dataset which contains the number of distinct items in each customer's market basket. Draw a histogram of the number of unique items. What are the median, the 25th percentile and the 75th percentile in this histogram?

Ans :

median (quartile 2): 3.0

25% (quartile 1): 2.0

75% (quartile 3): 6.0

Customer

4918 1

3937 1

3932 1

1819 1

1820 1

..

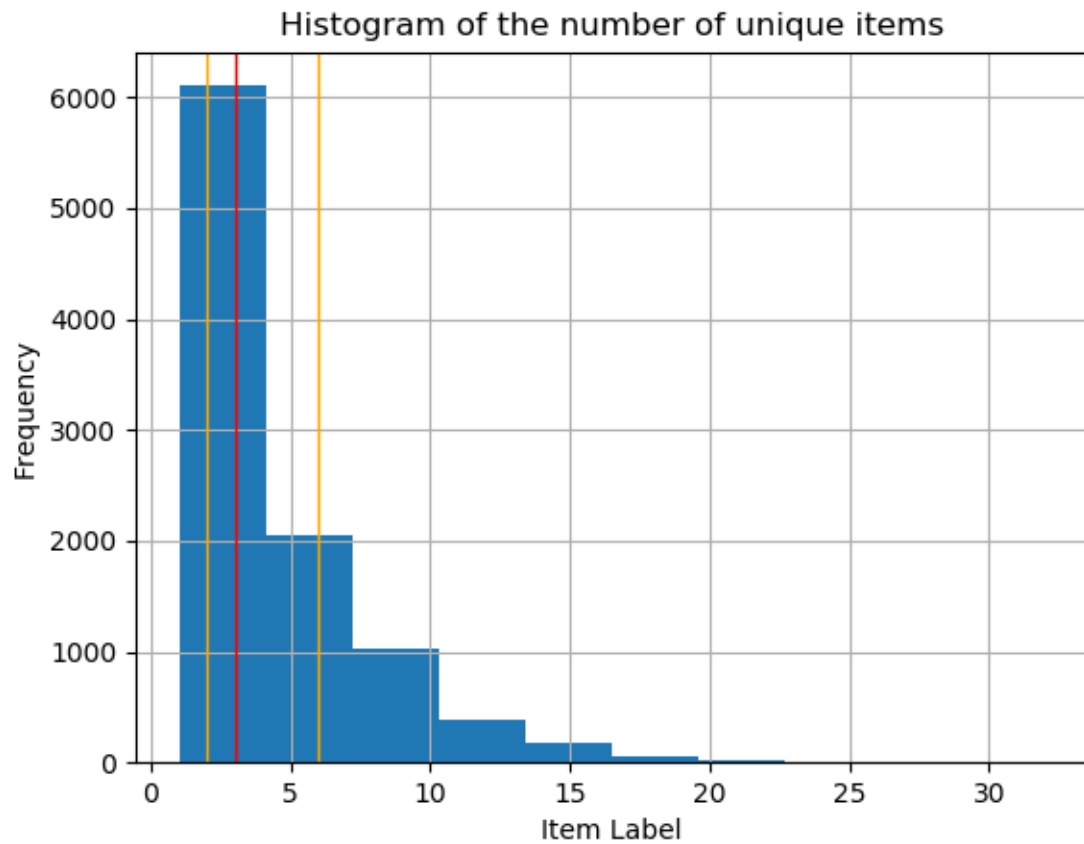
5611 28

9002 29

2974 29

2939 29

1217 32



- b) (10 points) If you are interested in the k -itemsets which can be found in the market baskets of at least seventy five (75) customers. How many itemsets can you find? Also, what is the largest k value among your itemsets?

Ans :

total number of item sets found = 522

the largest value of $k = 3$

frequent item sets :

- | | |
|-----|------------------------------------|
| 0 | (Instant food products) |
| 1 | (UHT-milk) |
| 2 | (baking powder) |
| 3 | (beef) |
| 4 | (berries) |
| ... | |
| 517 | (soda, whole milk, tropical fruit) |

- 518 (soda, yogurt, whole milk)
- 519 (whipped/sour cream, whole milk, tropical fruit)
- 520 (yogurt, whole milk, tropical fruit)
- 521 (whipped/sour cream, whole milk, yogurt)

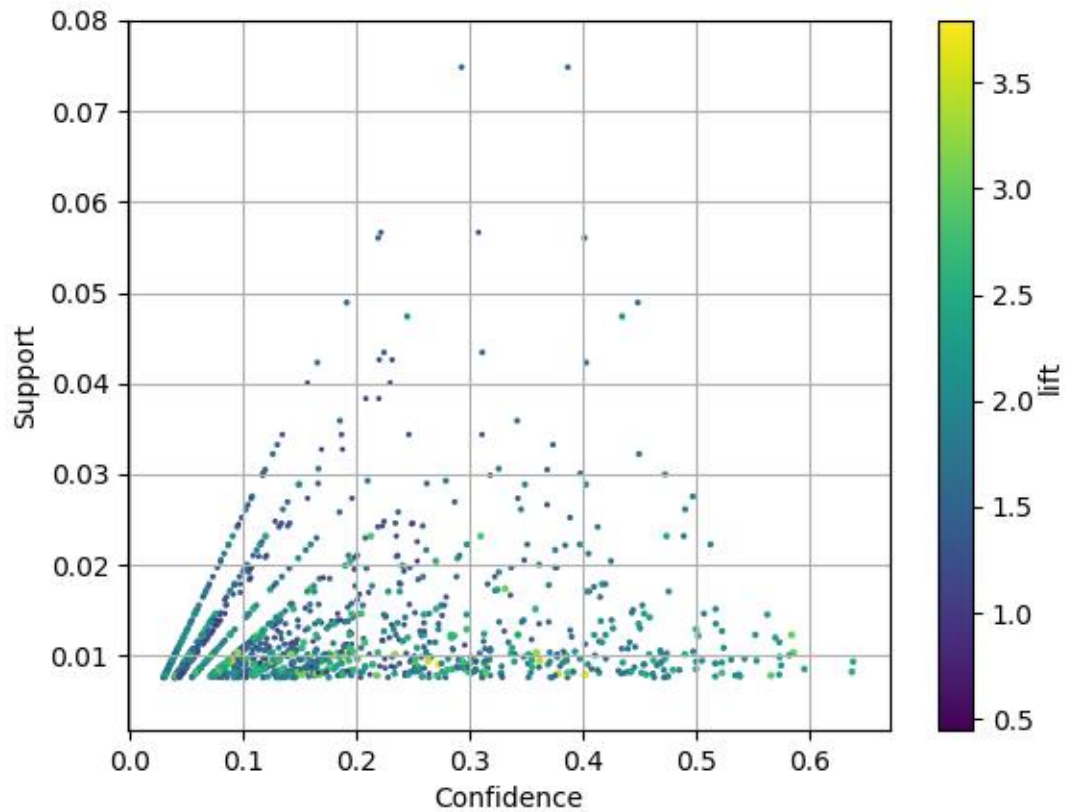
c) **(10 points)** Find out the association rules whose Confidence metrics are at least 1%. How many association rules have you found? Please be reminded that a rule must have a non-empty antecedent and a non-empty consequent. Also, you do not need to show those rules.

Ans :

Total number of association rules found = 1200

d) **(10 points)** Graph the Support metrics on the vertical axis against the Confidence metrics on the horizontal axis for the rules you found in (c). Please use the Lift metrics to indicate the size of the marker.

Ans :



e) (10 points) List the rules whose Confidence metrics are at least 60%. Please include their Support and Lift metrics.

Ans :

association rules :

antecedents consequents antecedent support consequent support support confidence lift
leverage conviction

0 (butter, root vegetables) (whole milk) **0.012913** 0.255516 0.008236 0.637795
2.496107 0.004936 2.055423

1 (yogurt, butter) (whole milk) **0.014642** 0.255516 0.009354 0.638889
2.500387 0.005613 2.061648

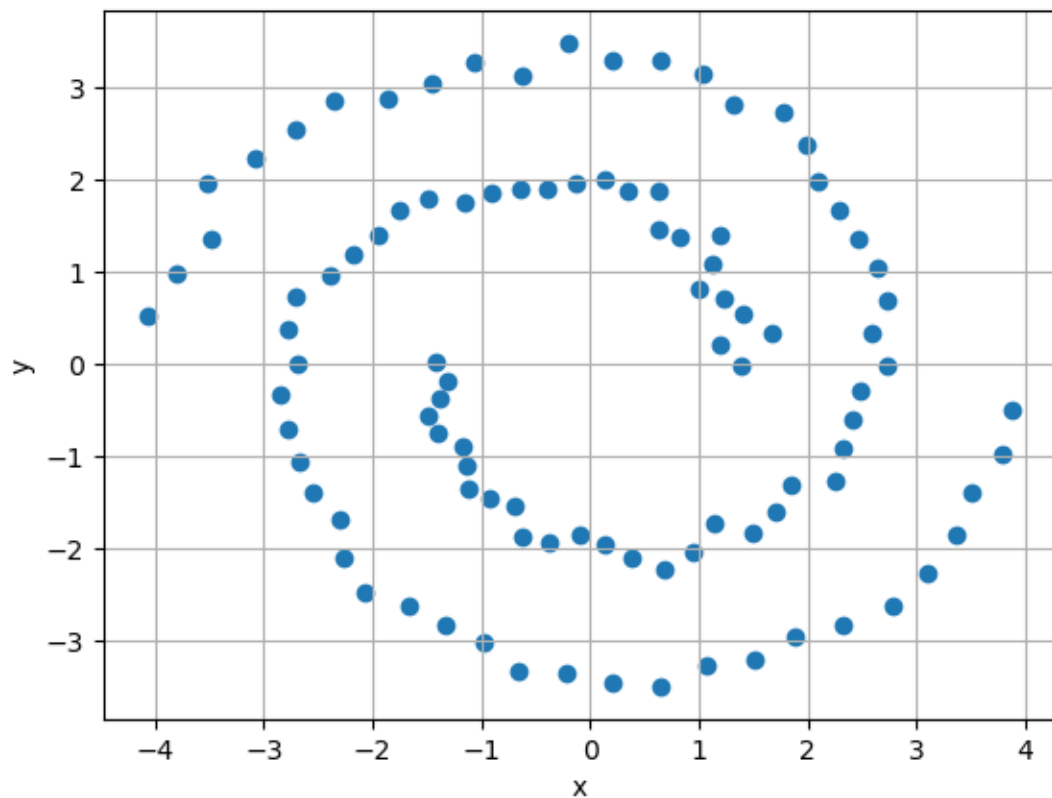
Question 2 (50 points)

Apply the Spectral Clustering method to the Spiral.csv. Your input fields are x and y. Wherever needed, specify `random_state = 60616` in calling the KMeans function.

- a) (10 points) Generate a scatterplot of y (vertical axis) versus x (horizontal axis). How many clusters will you say by visual inspection?

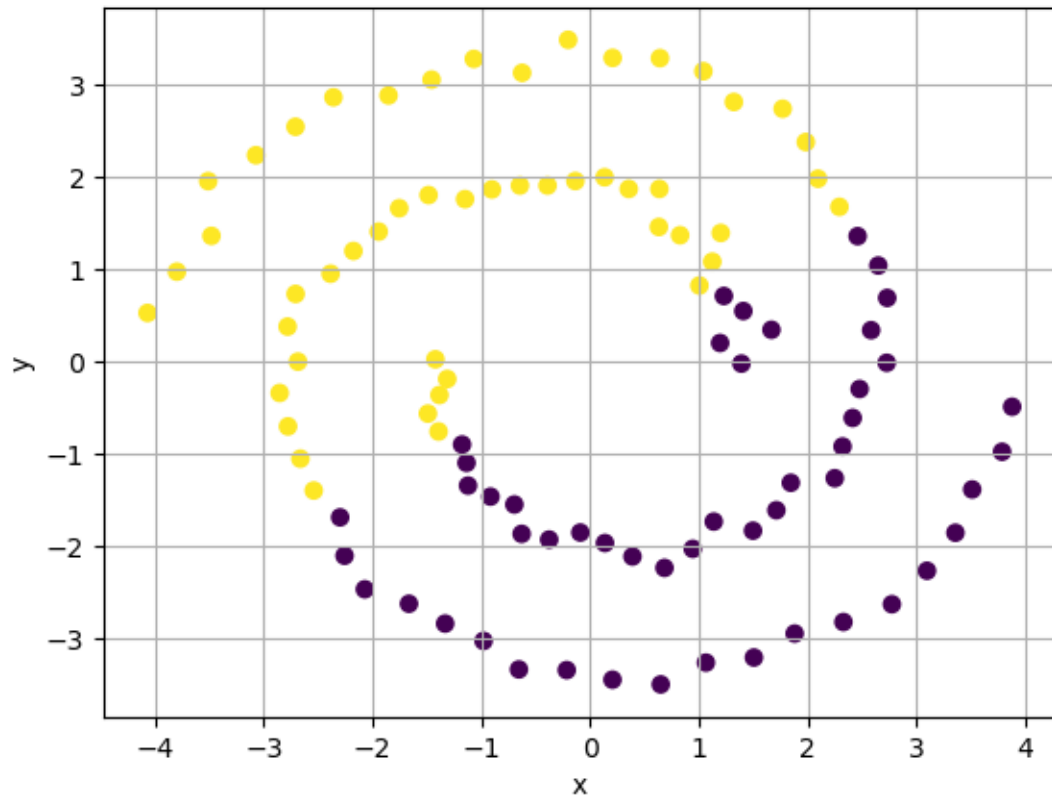
Ans :

I can say 2 clusters by visual inspection.



- b) (10 points) Apply the K-mean algorithm directly using your number of clusters that you think in (a). Regenerate the scatterplot using the K-mean cluster identifier to control the color scheme?

Ans :



- c) (10 points) Apply the nearest neighbor algorithm using the Euclidean distance. How many nearest neighbors will you use? Remember that you may need to try a couple of values first and use the eigenvalue plot to validate your choice.

Ans :

I can see three nearest neighbors.

- d) (10 points) Retrieve the first two eigenvectors that correspond to the first two smallest eigenvalues. Display up to ten decimal places the means and the standard deviation of these two eigenvectors. Also, plot the first eigenvector on the horizontal axis and the second eigenvector on the vertical axis.

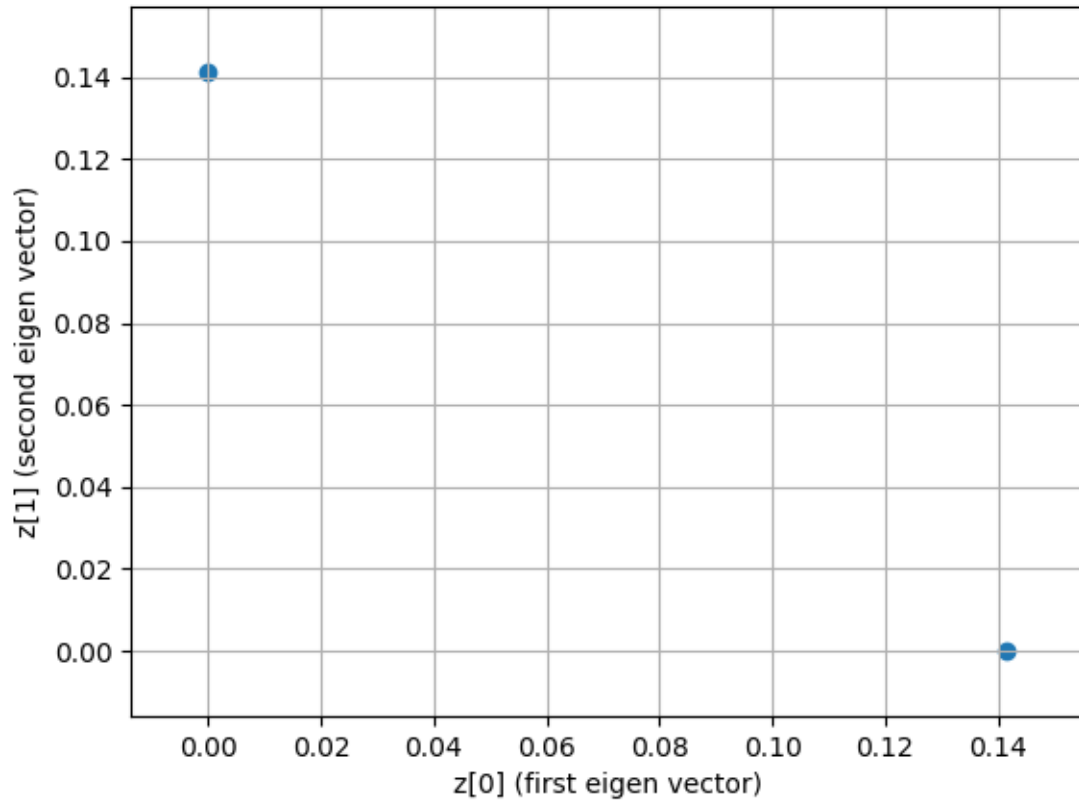
Ans :

Mean of first eigenvectors : 0.0707106781

Standard deviation of first eigenvectors: 0.0707106781

Mean of second eigenvectors : 0.0707106781

Standard deviation of second eigenvectors: 0.0707106781



- e) (10 points) Apply the K-mean algorithm on your first two eigenvectors that correspond to the first two smallest eigenvalues. Regenerate the scatterplot using the K-mean cluster identifier to control the color scheme?

Ans :

