

Ans 1  
Set 1 a) Support = 40% So 4 out of 10 Transaction

C<sub>1</sub>

Items	Support
Bread	7/10 = 70%
Coffee	4/10 = 40%
Sugar	3/10 = 30%
Milk	7/10 = 70%
Butter	4/10 = 40%
Cookies	2/10 = 20%
Eggs	3/10 = 30%

Support > 40%

L<sub>1</sub>

Items	Support
Bread	70%
Coffee	40%
Milk	70%
Butter	40%

Set of 2

C<sub>2</sub>

Items	Support
Bread, Coffee	1 = 10%
Bread, Milk	4 = 40%
Bread, Butter	4 = 40%
Coffee, Milk	3 = 30%
Coffee, Butter	0 = 0%
Milk, Butter	2 = 20%

Support > 40%

L<sub>2</sub>

Items	Support
Bread, Milk	40%
Bread, Butter	40%

C<sub>3</sub>

Items	Support
Bread, Milk, Butter	2 = 20%

No frequent itemset from C<sub>3</sub>

b) Association Rules

Association Rule	Support	Confidence	Confidence %	Valid
Bread → Milk	40%	4/7 = 0.57	57%	True
Milk → Bread	40%	4/7 = 0.57	57%	True
Bread → Butter	40%	4/7 = 0.57	57%	True
Butter → Bread	40%	4/4 = 1	100%	True

ent

E'

variance Matrix

can value it  
all prod

on  
p = 2

aid	locations
B	C
0.25	0.30
2.358	0.711
0.358	0.711

id	location
C	
0.30	
0.85	
0.851	

if  $\Sigma$  is diagonal then a perfect MVE.

$V$  - (Covariance Matrix)

$(x - \mu)^T \Sigma^{-1} (x - \mu) = K$  is Mahalanobis distance of vector  $\vec{x}$  from mean value  $\mu$   
 $\Rightarrow$  when plotted = an ellipse or a  $p$  dimensional ellipsoid  
 and this surface will be centered at  $\mu$ .  
 in MVE

$$[X = [x_1, x_2] \text{ so } p=2]$$

### Question 3

a)

iter	cluster assignments of data points (into A, B & C)								centroid locations		
	0.090	0.172	0.310	0.335	0.429	0.640	0.642	0.851	A	B	C
0	A	A	B	B	B	C	C	C	0.15	0.25	0.90
1	A	A	B	B	B	C	C	C	0.131	0.358	0.711
2	A	A	B	B	B	C	C	C	0.131	0.358	0.711

The algorithm converges after 2 iterations

b)

iter	cluster assignments of data points (into A, B & C)								centroid location		
	0.090	0.172	0.310	0.335	0.429	0.640	0.642	0.851	A	B	C
0	A	A	B	B	B	B	B	C	0.10	0.45	0.90
1	A	A	B	B	B	B	B	C	0.131	0.4712	0.851
2	A	A	B	B	B	B	B	C	0.131	0.4712	0.851

The algorithm converges after 2 iterations.

b)  
After



c) Sum of squared error for (A)

$$\begin{aligned} \text{iter 1 } SSE &= (0.15 - 0.090)^2 + (0.15 - 0.172)^2 + (0.25 - 0.310)^2 + (0.25 - 0.429)^2 \\ &\quad + (0.25 - 0.335)^2 + (0.90 - 0.640)^2 + (0.90 - 0.642)^2 + (0.90 - 0.851)^2 \\ &= 0.0036 + 0.000484 + 0.0036 + 0.032 + 0.067 + 0.066 \\ &\quad + 0.0072 + 0.002 \\ &= 0.1818 \end{aligned}$$

$$\begin{aligned} \text{iter 2 } SSE &= (0.131 - 0.090)^2 + (0.131 - 0.172)^2 + (0.258 - 0.310)^2 + (0.258 - 0.429)^2 + \\ &\quad + (0.258 - 0.335)^2 + (0.711 - 0.640)^2 + (0.711 - 0.642)^2 + (0.711 - 0.851)^2 \\ &= 0.0016 + 0.0016 + 0.0023 + 0.0050 + 0.0050 + 0.0047 + 0.0136 + \\ &\quad + 0.0005 \\ &= 0.0403 \end{aligned}$$

$$\text{iter 3 } SSE = 0.0403 \quad (\text{As algorithm converges})$$

d) Sum of squared error for (b)

$$\begin{aligned} \text{iter 1 } SSE &= (0.10 - 0.090)^2 + (0.10 - 0.172)^2 + (0.45 - 0.310)^2 + (0.45 - 0.429)^2 + (0.45 - 0.640)^2 \\ &\quad + (0.45 - 0.642)^2 + (0.90 - 0.851)^2 + (0.45 - 0.335)^2 \\ &= 0.0001 + 0.0051 + 0.0196 + 0.0004 + 0.0361 + 0.0368 + 0.0024 \\ &\quad + 0.0132 \\ &= 0.1157 \end{aligned}$$

$$\begin{aligned} \text{iter 2 } SSE &= (0.131 - 0.090)^2 + (0.131 - 0.172)^2 + (0.4712 - 0.310)^2 + (0.4712 - 0.429)^2 + \\ &\quad + (0.4712 - 0.335)^2 + (0.4712 - 0.640)^2 + (0.4712 - 0.642)^2 + (0.851 - 0.851)^2 \\ &= 0.0016 + 0.0016 + 0.0253 + 0.0014 + 0.0284 + 0.0294 + 0 + 0.0185 \\ &= 0.1068 \end{aligned}$$

$$\text{iter 3 } SSE = 0.1068 \quad \text{as algorithm converges}$$

$\Rightarrow$  Solution (A) with initial centroids as (0.15, 0.25, 0.90) gives a better solution in terms of SSE

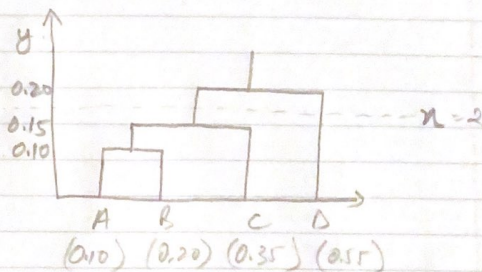
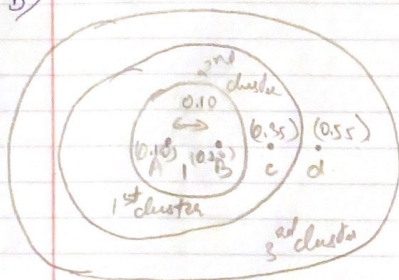


4) Points  $\rightarrow$   $\begin{matrix} \text{Point (A)} & (B) & (C) & (D) \\ [0.10, 0.20, 0.35, 0.55] \end{matrix}$

5) Pair wise Euclidean distance  $\rightarrow$

	A	B	C	D
A	0	0.10	0.25	0.45
B	0.10	0	0.15	0.35
C	0.25	0.15	0	0.20
D	0.45	0.35	0.20	0

6)

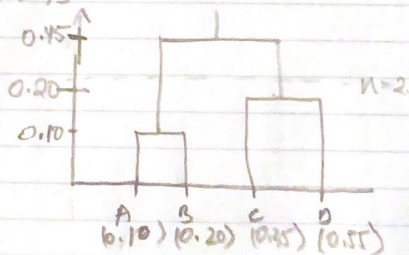
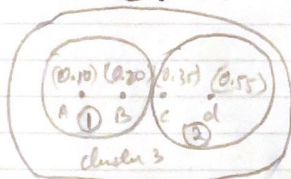


$$SSE_{n=2} = (0.216 - 0.10)^2 + (0.716 - 0.20)^2 + (0.716 - 0.35)^2 + (0.55 - 0.55)^2$$

$$= 0.0134 + 0.0002 + 0.0149$$

$$= \boxed{0.0315}$$

7)



$$SSE_{n=2} = (0.15 - 0.10)^2 + (0.15 - 0.20)^2 + (0.45 - 0.35)^2 + (0.45 - 0.55)^2$$

$$= 0.0025 + 0.0025 + 0.01 + 0.01 \Rightarrow \boxed{0.025}$$

8)

iter	
0	
1	
2	

$$SSE_{n=2} =$$

The Max  
K-means

9)

Algorithm A  
Cluster 1  
Cluster 2

$$\sum_{i=1}^n (w_i) =$$

$$\sum_{i=1}^n (a_i) =$$

$$\sum_{i=1}^n (b_i) =$$

$$ARI_A = 1.55 - 1$$

$$ARI_A = \frac{1}{2} [2.475 +$$

Algorithm



280      2462.5-  
 1237.5      1240      4950.  
 3.03      2265      -16.96      1225      380  
    2762.5-

d)

iter	cluster assignment of data points (A, B, C, D)				centroid location	
	A(0.10)	B(6.20)	C(0.35)	D(0.55)	$\mu_1$	$\mu_2$
0	$C_1$	$C_1$	$C_2$	$C_2$	0.15	0.5
1	$C_1$	$C_1$	$C_2$	$C_2$	0.15	0.45
2	$C_1$	$C_1$	$C_2$	$C_2$	0.15	0.45

$$SSE_{k=2} = (0.15-0.10)^2 + (0.15-0.20)^2 + (0.45-0.35)^2 + (0.45-0.55)^2 = 0.025$$

The Hare Method in (C) produces similar results and Hare method and K-means with Euclidean distance produce a solution with lowest SSE

S)

Algorithm A	ground truth	
	class 1	class 2
cluster 1	10	35
cluster 2	40	15

Algorithm B	ground truth	
	class 1	class 2
cluster 1	35	40
cluster 2	15	10

$$\begin{aligned} \sum_{j=1}^n (u_{ij}) &= \left(\frac{10}{2}\right) + \left(\frac{35}{2}\right) + \left(\frac{40}{2}\right) + \left(\frac{15}{2}\right) \\ &= 45 + 59.5 + 780 + 105 = 1555 \\ \sum_{i=1}^n (q_i) &= \left(\frac{45}{2}\right) + \left(\frac{59.5}{2}\right) \\ &= 900 + 1485 = 2475 \\ \sum_{j=1}^n (b_j) &= \left(\frac{50}{2}\right) + \left(\frac{50}{2}\right) \\ &= 1225 + 1225 = 2450 \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n (u_{ij}) &= \left(\frac{35}{2}\right) + \left(\frac{40}{2}\right) + \left(\frac{15}{2}\right) + \left(\frac{10}{2}\right) \\ &= 152.5 \\ \sum_{i=1}^n (q_i) &= \left(\frac{75}{2}\right) + \left(\frac{125}{2}\right) \\ &= 2775 + 3000 = 3075 \\ \sum_{j=1}^n (b_j) &= \left(\frac{50}{2}\right) + \left(\frac{50}{2}\right) \\ &= 1225 + 1225 = 2450 \end{aligned}$$

$$ARI_A = 1555 - [2475 \times 2450] / \binom{100}{2}$$

$$ARI_B = 1525 - [3075 \times 2450] / \binom{100}{2}$$

$$ARI_A = \frac{\frac{1}{2}[2475 + 2450] - [2475 \times 2450]}{\binom{100}{2}} = 0.2424$$

$$ARI_B = \frac{\frac{1}{2}[3075 + 2450] - [3075 \times 2450]}{\binom{100}{2}} = 0.0024$$

Algorithm A is better in terms of ARI