

Q1

Q1

$$C_1 = \left[ \frac{2+4+2}{3}, \frac{5+3+3}{3} \right] = (2.66, 3.66)$$

$$C_2 = \left[ \frac{6+3+4+3}{4}, \frac{2+5+1+4}{4} \right] = (4, 3)$$

Distance of  $ID_1$  from  $C_1 = 12.66 - 61 + 13.66 - 21 = 5$  | 2.34, 0.66

Distance of  $ID_2$  from  $C_2 = 14 - 61 + 13 - 21 = 3$

Distance of  $ID_3$  from  $C_1 = 0.66 + 0.66 = 1.32$  | 2.68

Distance of  $ID_4$  from  $C_2 = 2 + 0 = 2$

Q2 We will move  $ID_5$  from cluster 2 to 1 because it's closer to  $C_1$ .

$$C_1 = \left( \frac{14}{4}, \frac{15}{4} \right) = (3.5, 3.25)$$

$$C_2 = \left( \frac{10}{3}, \frac{10}{3} \right) = (3.33, 3.33)$$

for  $ID_5$  from  $C_1 = 0.5, 1.75$

$ID_6$  from  $C_2 = 3.33 - 3, 3.33 - 5 = (0.33, 1.67)$

$ID_7$  from  $C_1 = 3.5 - 2, 3.25 - 3 = (1.5, 0.25)$

$ID_8$  from  $C_2 = 3.33 - 2, 3.33 - 5 = (1.33, 0.33)$

Q2

Q.2 1) Minimum distance  $D(3,4) E(4,6) = 13-4 + 14-6 = 23$

2) Average distance

$D(a,e) = 4$   $d(b,e) = 5$

$D(a,g) = 7$   $d(b,g) = 6$

$D(a,f) = 7$   $d(b,f) = 8$

$d(c,e) = 5$   $d(d,e) = 3$

$d(c,f) = 4$   $d(d,f) = 4$

$d(c,g) = 8$   $d(d,g) = 6$

$\frac{67}{12} \approx 5.58$

3)  $C_1 = \left( \frac{5}{2}, \frac{11}{2} \right)$

$C_2 = \left( \frac{24}{5}, \frac{18}{5} \right) = (4.8, 3.6)$

SSE of  $C_1$

$(P_1, C_1) = 0.5^2 + \left( \frac{15}{2} \right)^2 = 2.5$

$(P_2, C_1) = 0.5^2 + \left( \frac{15}{2} \right)^2 = 2.5$

$2.5 + 2.5 = 5$

SSE of  $C_2$

$(C_1, C_2) = 140 = 1$

$(C_2, C_2) = 1 + 1 = 2$

$(C_2, C_2) = 1$

$1 + 1 + 2 = 4$

$C_1, C_2 = (5.8, 5.8)$

$A, C = 14.44 + 3.24 = 17.68$

$B, C = 7.84 + 1.44 = 9.28$

$C, C = 1.44 + 1.44 = 2.88$

$D, C = 10.25 + 1.44 = 11.69$

$E, C = 4.84 + 1.35 = 6.19$

SSE = 47.36

Before merge,  $C_1, C_2 = 9$

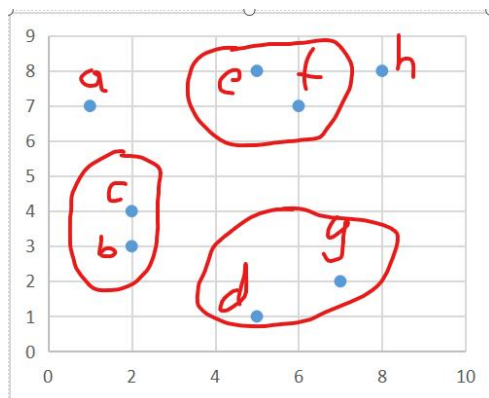
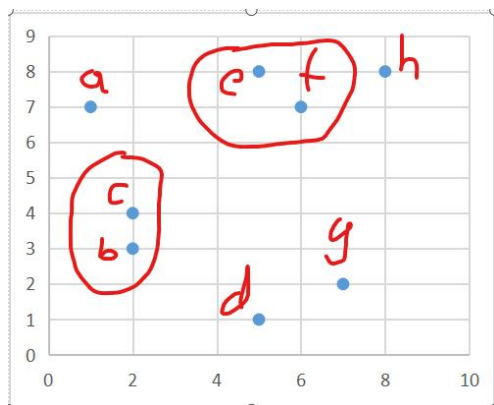
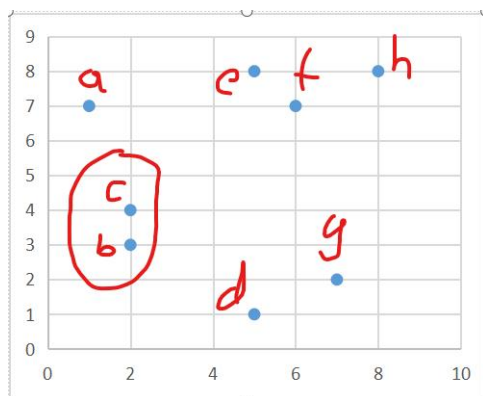
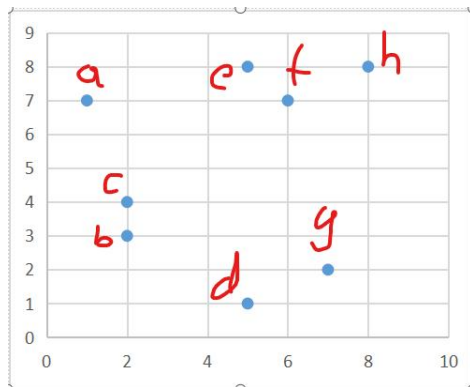
After = 47.36

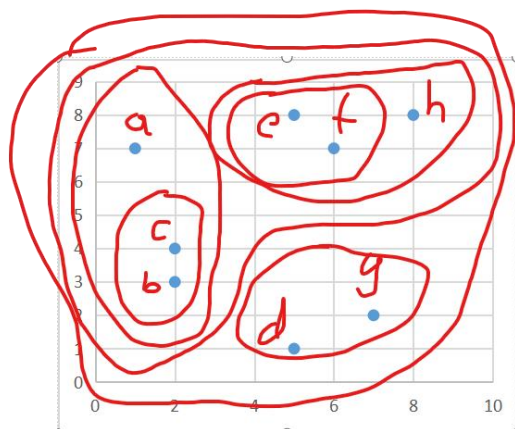
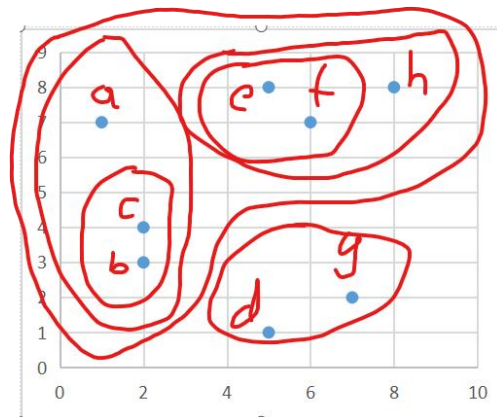
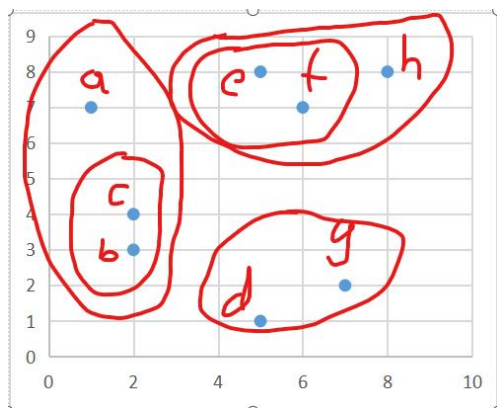
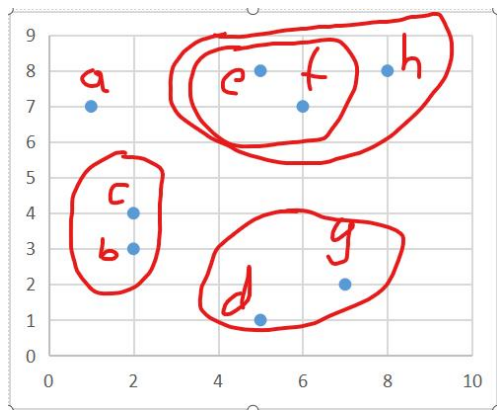
$47.36 - 9 = 38.36$

Increase = 38.36

Cost = 38.36

Q3





(a)(b)(c)(d)(e)(f)(g)(h)

(a)(b,c)(d)(e)(f)(g)(h)

(a)(b,c)(d)(e,f)(g)(h)

(a)(b,c)(d,g)(e,f)(h)

(a)(b,c)(d,g)(e,f,h)

(a,b,c)(d,g)(e,f,h)

Q4

(1)

K=2

Number of iterations: 6  
Within cluster sum of squared errors: 33.30854661169795

Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster#	
		0	1
	(506.0)	(202.0)	(304.0)
A1	68.5749	37.8658	88.9803
A2	356.674	389.2163	335.0506

K=3

kMeans  
=====

Number of iterations: 6  
Within cluster sum of squared errors: 11.529148066305805

Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25  
Cluster 2: 95.6,60.72

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster#		
		0	1	2
	(506.0)	(191.0)	(275.0)	(40.0)
A1	68.5749	36.3461	87.7833	90.41
A2	356.674	388.9075	376.937	63.4515

K=4

Number of iterations: 17  
Within cluster sum of squared errors: 6.895127380017731

Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25  
Cluster 2: 95.6,60.72  
Cluster 3: 95.6,396.9

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Full Data	Cluster#			
		0	1	2	3
	(506.0)	(125.0)	(120.0)	(40.0)	(221.0)
A1	68.5749	27.3688	60.735	90.41	92.1864
A2	356.674	388.1318	390.291	63.4515	373.6995

K=5

Number of iterations: 13  
Within cluster sum of squared errors: 5.355515572616415

Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25  
Cluster 2: 95.6,60.72  
Cluster 3: 95.6,396.9  
Cluster 4: 73.3,385.91

Missing values globally replaced with mean/mode

Final cluster centroids:						
Attribute	Cluster#					
	Full Data	0	1	2	3	4
	(506.0)	(125.0)	(32.0)	(36.0)	(196.0)	(117.0)
A1	68.5749	27.3688	94.0031	89.3444	91.8112	60.3274
A2	356.674	388.1318	278.3478	50.1542	385.6866	390.1994

K=6

KMeans  
=====

Number of iterations: 22  
Within cluster sum of squared errors: 3.923262095570305

Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25  
Cluster 2: 95.6,60.72  
Cluster 3: 95.6,396.9  
Cluster 4: 73.3,385.91  
Cluster 5: 83.2,390.7

Missing values globally replaced with mean/mode

Final cluster centroids:							
Attribute	Cluster#						
	Full Data (506.0)	0 (91.0)	1 (28.0)	2 (36.0)	3 (165.0)	4 (93.0)	5 (93.0)
A1	68.5749	22.8352	94.0357	89.3444	94.303	47.6624	72.8914
A2	356.674	388.692	270.8661	50.1542	383.7804	389.0449	389.3692

K=7

Number of iterations: 37  
Within cluster sum of squared errors: 3.051532898594132

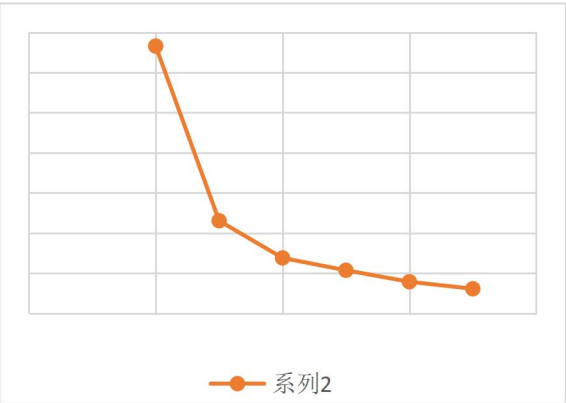
Initial starting points (random):

Cluster 0: 61.5,391.34  
Cluster 1: 92.1,393.25  
Cluster 2: 95.6,60.72  
Cluster 3: 95.6,396.9  
Cluster 4: 73.3,385.91  
Cluster 5: 83.2,390.7  
Cluster 6: 58,396.9

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute	Cluster#						
	Full Data	0	1	2	3	4	5
	(506.0)	(84.0)	(28.0)	(36.0)	(154.0)	(72.0)	(82.0)
A1	68.5749	36.0667	94.0357	89.3444	94.9896	56.7236	76.9122
A2	356.674	389.7713	270.8661	50.1542	383.4029	390.3379	388.7613



optimal number of clusters:k=2

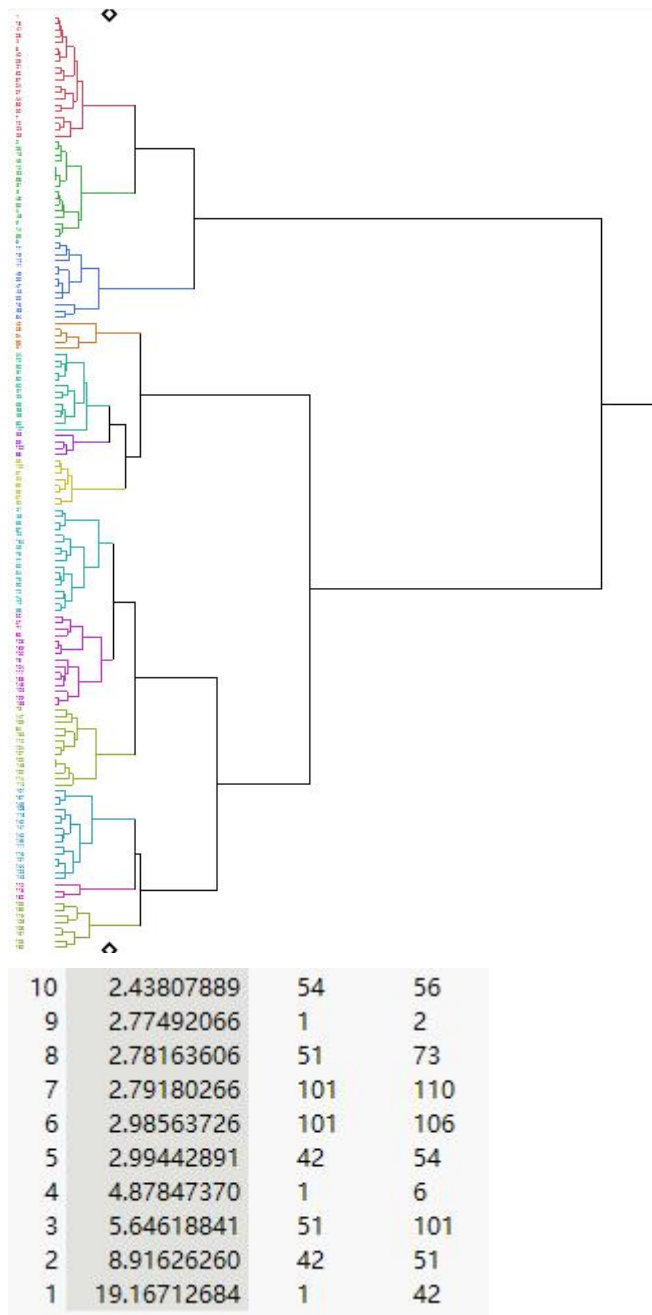
(2)

		1	2	3	4
Calories	Mean	156.5629	214.8264	171.026	167.5224
	STD	+/-11.1385	+/-15.5421	+/-18.191	+/-8.2984
	MAX	211.073064	196.0569717	211.073064	24.10206
	Min	135.9771213	135.9771213	166.0206963	135.9771213
Fiber	Mean	3.9537	7.2813	24.8333	0.225
	STD	+/-4.3567	+/-4.3536	+/-6.3456	+/-1.0665
	MAX	24.5	24.5	8.5	35.5
	Min	-0.5	-0.5	-0.5	-0.5
Sugars	Mean	10.9923	26.7789	26.0188	30.7703
	STD	+/-4.2619	+/-7.5774	+/-10.1438	+/-6.3388
	MAX	38.97474109	38.97474109	28.27701184	42.97815422
	Min	7.150685103	11.95246912	14.32746912	7.150685103
Potassium	Mean	81.2963	146.875	206.6667	46.75
	STD	+/-32.1233	+/-44.7962	+/-36.697	+/-24.4559
	MAX	230	190	100	260
	Min	20	30	35	35

So, cluster 1 has the healthiest cereals

Q5

(1)



I will chose No.4 because since No.4, the distance has a huge change

(2)



