CSCI3220 2018-19 First Term Assignment 5

I declare that the assignment here submitted is original except for source material explicitly acknowledged, and that the same or closely related material has not been previously submitted for another course. I also acknowledge that I am aware of University policy and regulations on honesty in academic work, and of the disciplinary guidelines and procedures applicable to breaches of such policy and regulations, as contained in the following websites.

University Guideline on Academic Honesty: http://www.cuhk.edu.hk/policy/academichonesty/

Student Name: ZHANG Chongzhi

Student ID: 1155077072

<i>]</i> .	X	s_1	<i>S</i> ₂	S 3	<i>S</i> 4	S 5	S 6
a),	g_1	4	6	5	1	2	8
	g_2	1	4	6	8	1	3
	g_3	2	9	3	5	1	6
	g_4	8	5	2	1	3	4
	g_5	7	1	1	3	2	9

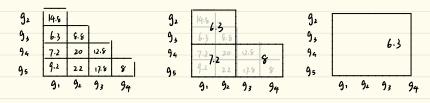
d(g,,g2) =	$\frac{\sum_{k=1}^{l} (x_{ik} - x_{jk})^2}{6}$	9+4+1+49+1+ = 6	- 25 - =
s the same meth			
l(g,,g;) = 6.3	d(91, 94) = 7.2	d(g1,g5)= 9.2	
Lig2, gz) = 8.8	• •	d(92,95)= 22	
, ,-	• •		

d(95,95) = 17.8 d(94,95) = 8

14.8

Produce the distance matrix.

d	9,	د9	9,	94	95
9,	0	14.8	6.3	7.2	9.2
92	14.8	0	ş. ş	20	2.2
9,	6.3	છે.	0	12.8	17.8
94	7.2	20	12.6	0	8
95	9.2	22	17.8	8	0



(C) From the initial quad tree, we knows that g, and g, are marged first.

1). single - Link:
$$d(\{g_1,g_2\},\{g_k\}) = min(d(\{g_3\},\{g_k\}),d(\{g_3\},\{g_k\}))$$

92	94	9s	92	8.8				8.8					
8.8	7.2	9.2		∞	œ			06	§.⊊ <i>⊙</i> ≎				_
	•	<u> </u>	94	7.2	20	œ		7.2	20	∞	C		/. 2 .
			95	9.2	22	8	8	9.2	7.2	00	8		
			•	9.93	92	•	94						

2). Average - link:
$$d(\{g,g\},\{g\}) = \frac{d(\{g\},\{g\}) \cdot 1 \cdot 1 + d(\{g\},\{g\}\},\{g\}) \cdot 1 \cdot 1}{2 \times 1}$$

92	94	Gs	92	11.6	l			11.8	<u> </u>				Г	 	
11.8	10.0	13.5		00	oc		.1	00	/).g			_	1	,	
			94	lo	20	<u>۵</u>	<u> </u>	10	10	-00	ç	ŀ	-	/	υ
			95	13.5	22	∞	8	13.5	2.2	oG.	8		L	 	
				9,9,	92		94								

92	94	95	92	14.8				14.8	14.8]					
14.8	12.8	17.8		00	œ			00	œ			_		_	
			94	12.8	20	03		12.8	12.8	06	ç			ъ	
			95	17.8	22	∞	8	17.8	2.2	o0	8				
				9.93	92		94					_			

3. (b). Every time When re-cleberaining representative, the cluster that contain the outlier will generate a representative draged to the outlier point. Thus next time, the cluster that contain that point may have less points, which means the representative will more closer to the outlier point.

Thursfore, if the outlier point is for enough, it will become a cluster itself and all others points mugad to the other cluster.



(b). The column has a much larger magnitude, which wears that the representation will be very close to the center of that column. Zach times do the iteration. the cluster has a trend to have less point outside that column.

Then fac. the column and the points that close to its center will be marged into one cluster.

(c). Advantage: It uses a ratio rather than a real distance, the difference will be increase for close point and cherease from fer-away points. It will generate clusters with more "valuatively-closed pines Disadvantage: More colculation required.

1d). First let's just make the $k \leftarrow 2$, and nuke the k increase by the variance.

Each time when we get k clusters,

for each cluster, do check the variance of all the points in that cluster, if the

variance is greater than some given value threshold. Let the cluster split: $k \leftarrow k+1$.

if all the cluster has a variance less than threshold, return k.