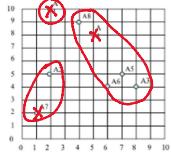


Assignment 5

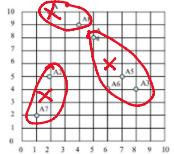
Wednesday, November 24, 2021 2:50 PM

1. [20 points] By considering the following 8 2D data points below:
- [15 points] Group the points into 3 clusters by using k-means algorithm with Euclidean distance. Show the intermediate clusters (by drawing ellipses on this 10 by 10 space) and centroids (by drawing marks like X on this 10 by 10 space) in each iteration until convergence. Consider the initial centroids as: C1 = A1, C2 = A4, and C3 = A7.



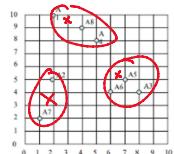
1 st iteration								
Centroid: (C1, C2, C3)	A1	A2	A3	A4	A5	A6	A7	A8
C1 dist.	0	5	4.9	—	7.07	7.31	—	7.23
C2 dist.	—	4.24	5	0	3.61	5.12	—	4.91
C3 dist.	—	3.18	7.18	—	6.7	5.9	0	7.61
Cluster Assigned	U ₁	U ₃	U ₂	U ₂	U ₂	U ₂	U ₃	U ₂

$$C_1 = (1, 1) \quad C_2 = \left[\frac{(4+5+6+7+8)}{5}, \frac{(9+8+5+4+4)}{5} \right] = (6, 6) \quad C_3 = \left[\frac{(1+2)}{2}, \frac{(2+5)}{2} \right] = (1.5, 3.5)$$



2 nd iteration								
Centroid: (C1, C2, C3)	A1	A2	A3	A4	A5	A6	A7	A8
C1 dist.	0	5	4.4	3.6	7.02	7.21	8.66	7.75
C2 dist.	—	4.12	7.42	2.33	1.91	2	6.40	3.6
C3 dist.	—	1.58	6.51	6.3	6.3	4.56	1.79	6.0
Cluster Assigned	U ₁	U ₃	U ₂	U ₂	U ₂	U ₂	U ₃	U ₁

$$C_1 = (2+4)/2, ((0+9)/2) \quad C_2 = \left[\frac{(5+6+7+8)}{4}, \frac{(9+5+4+4)}{4} \right] = (6.5, 5.25) \quad C_3 = (1.5, 3.5)$$



3 rd iteration								
Centroid: (C1, C2, C3)	A1	A2	A3	A4	A5	A6	A7	A8
C1 dist.	1.11	4.6	7.13	3.5	6.02	6.26	7.36	1.11
C2 dist.	—	1.58	—	2.13	—	—	1.58	—
C3 dist.	—	—	—	—	—	—	—	—
Cluster Assigned	U ₁	U ₃	U ₂	U ₁	U ₂	U ₂	U ₃	U ₁

$$C_1 = \frac{7+4+5}{3}, \frac{10+9+8}{3} \quad C_2 = \frac{6+7+8}{3}, \frac{5+4+4}{3} \quad C_3 = (1.5, 3.5) \\ = (3.7, 9) \quad = (7, 4.53)$$

- b. [5 points] Calculate the SSE (Sum of Square Errors) of the final clustering.

$$SSE = \sum_{i=1}^3 \sum_{x \in C_i} dist^2(m_i, x)$$

$$\text{Sum of } C_1 = 1.97^2 + .3^2 + 1.64^2 = 6.66$$

$$\text{Sum of } C_2 = 1.05^2 + .67^2 + 1.05^2 = 2.65$$

$$\text{Sum of } C_3 = 1.58^2 + 1.58^2 = 4.99$$

$$= 14.3$$

2. [20 points] Use the distance matrix below to perform the following operations:

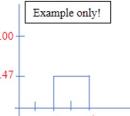
- a. [14 points] Group the points by using single link (MIN) hierarchical clustering. Show your results by informing the updated distance matrix after each merging step and by drawing the corresponding dendrogram that should clearly present the order in which the points are merged.

	p1	p2	p3	p4	p5
p1	0.00	0.10	0.41	0.55	0.35
p2	0.10	0.00	0.64	0.47	0.98
p3	0.41	0.64	0.00	0.44	0.85
p4	0.55	0.47	0.44	0.00	0.76
p5	0.35	0.98	0.85	0.76	0.00

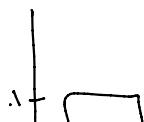
Solution format:

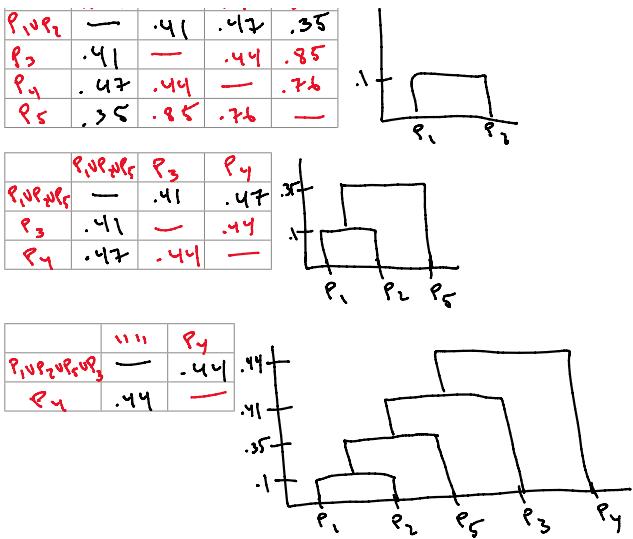
(1st iteration) Suppose the first two points to be merged are p2 and p4, then:

	p1	p2	p2 U p4	p3	p5
p2 U p4	0.00	?	?	0.41	0.35
p3	0.41	?	?	0.00	0.85
p5	0.35	?	?	0.85	0.00

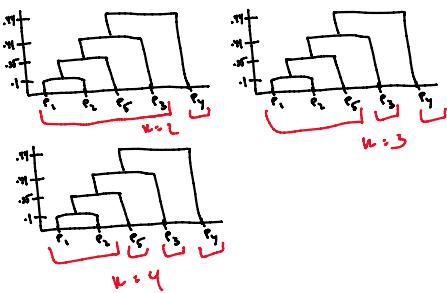


	P1 U P2	P3	P4	P5
P1 U P2	—	.41	.47	.35
P3	.41	—	.44	.85
P4	.47	.44	—	.76





b. [6 points] Show the clusters when $k = 2$, $k = 3$, and $k = 4$.



3. <https://github.com/NootCode/CS4210-Assignment-5/blob/master/Clustering/clustering.py>

4. [10 points] The dataset below presents the user ratings on a 1-3 scale for 6 different rock bands.

	Bon Jovi	Metallica	Scorpions	AC/DC	Kiss	Guns n' Roses
Fred	1	3	1.5	3	1	3
Lillian	3	1.5	2	2	3	1
Cathy	2	2	2	3	1.5	2
John	3	2	2	2	?	?

a. [5 points] Apply user-based collaborative filtering on the dataset to decide about recommending the bands Kiss and Guns n' Roses to John. You should make a recommendation when the predicted rating is greater than or equal to 2.0. Use cosine similarity, a neutral value (1.5) for missing values, and the top 2 similar neighbors to build your model.

$$\text{Cos}(John, Fred) = 0.857$$

$$\text{Cos}(John, Lillian) = .995 \rightarrow \text{Closest 2}$$

$$\text{Cos}(John, Cathy) = .982$$

$$\overline{A_{avg}}_{John} = \frac{3+6}{2} = 4.5 / 4 = 2.25$$

$$2.25 + \frac{[.995 * (3 - 2.25) + .982 * (1.5 - 2.25)]}{(.995 + .982)}$$

$$2.25 + .068 = 2.31 \leftarrow \text{Make the recommendation}$$

Guns n' roses

$$2.25 + \frac{[.995 * (1 - 2.25) + .982 * (2 - 2.25)]}{(.995 + .982)}$$

$$2.25 - .69 = 1.56 \leftarrow \text{Do not recommend}$$

- b. [5 points] Now, apply **item-based** collaborative filtering to make the same decision. Use the same parameters defined before to build your model.

	Bon Jovi	Metallica	Scorpions	AC/DC	Kiss	Guns n' Roses
Fred	1	3	1.5	3	1	3
Lillian	3	1.5	2	2	3	1
Cathy	2	2	2	3	1.5	2
John	3	2	2	2	?	?

Kiss

$$\begin{aligned}\cos(\text{Kiss}, \text{B}) &= .993 \\ \cos(\text{Kiss}, \text{M}) &= .768 \rightarrow 2 \text{ NN} \\ \cos(\text{Kiss}, \text{S}) &= .937 \\ \cos(\text{Kiss}, \text{A}) &= .922\end{aligned}$$

$$\text{Avg Kiss} = 1.833$$

$$\begin{aligned}1.833 + \frac{[.993 * (3 - 2) + .937 * (2 - 1.833)]}{(.993 + .937)} \\ = 1.833 + .597 = \boxed{2.43}\end{aligned}$$

Guns
n' roses

$$\begin{aligned}\cos(\text{G}, \text{B}) &= .71 \\ \cos(\text{G}, \text{M}) &= .99 \rightarrow 2 \text{ NN} \\ \cos(\text{G}, \text{S}) &= .88 \\ \cos(\text{G}, \text{A}) &= .97\end{aligned}$$

$$\begin{aligned}\text{Avg} \\ \text{GNR} &= 2\end{aligned}$$

$$2 + \frac{[.99 * (2 - 2.17) + .97 * (2 - 2.62)]}{.99 + .97}$$

$$\begin{aligned}2 - .417 \\ \boxed{= 1.58}\end{aligned}$$

5. [20 points] Consider the following transaction dataset.

Transaction ID	Items Bought
1	{a, b, c, d, e}
2	{b, c, d}
3	{a, b, c, d}
4	{a, b, d, e}
5	{b, c, d, e}
6	{b, d, e}
7	{c, d}
8	{a, b, c}
9	{a, d, e}
10	{b, d}

Suppose that minimum support is set to 30% (*minsup*) and minimum confidence is set to 60%.

- a. [5 points] Rank all frequent itemsets according to their support (list their support values).

$$\begin{aligned}\{\text{d}\} &= .9 \\ \{\text{b}\} &= .7 \\ \{\text{c}\} &= .6 \\ \{\text{a}\} &= .5 \\ \{\text{e}\} &= .5 \\ \{\text{a, d}\} &= .4 \\ \{\text{b, c}\} &= .4 \\ \{\text{a, c}\} &= .4 \\ \{\text{a, b}\} &= .3 \\ \{\text{b, e}\} &= .4 \\ \{\text{b, c}\} &= .3 \\ \{\text{c, d}\} &= .4 \\ \{\text{a, e}\} &= .5\end{aligned}$$

	a	b	c	d	e
1	██████████			██████████	██████████
2		██████████		██████████	██████████
3	██████████	██████████		██████████	██████████
4	██████████	██████████	██████████	██████████	
5		██████████	██████████	██████████	██████████
6			██████████	██████████	██████████
7			██████████	██████████	██████████
8		██████████	██████████	██████████	
9	██████████	██████████	██████████	██████████	██████████
10		██████████		██████████	██████████

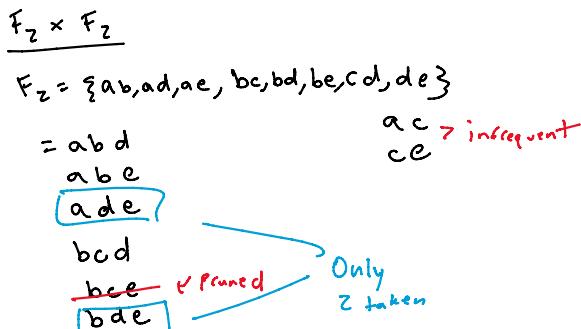
- b. [5 points] For all frequent 3-itemsets, rank all association rules - according to their confidence values - which satisfy the requirements on minimum support and minimum confidence (list their confidence values).

$$\begin{aligned}\{\text{a, d, e}\} \\ \{\text{a, d, e}\} \rightarrow \{\text{e}\} \quad c = \frac{\sigma(\{\text{a, d, e}\})}{\sigma(\{\text{a, d}\})} = \frac{4}{4} = 1\end{aligned}$$

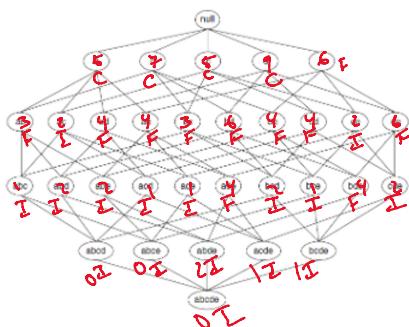
$$\begin{aligned} & \{a, d, e\} \\ & \{a, d\} \rightarrow \{e\} \quad c = \frac{\sigma(\{a, d, e\})}{\sigma(\{a, d\})} = \frac{4}{4} = 1 \\ & \{a, e\} \rightarrow \{d\} \quad c = \frac{\sigma(\{a, d, e\})}{\sigma(a, e)} = \frac{4}{4} = 1 \\ & \{d, e\} \rightarrow \{a\} \quad c = \frac{4}{5} = .8 \\ & \{a\} \rightarrow \{d, e\} \quad c = 4/5 = .8 \\ & \{e\} \rightarrow \{a, d\} \quad c = 4/6 = .667 \end{aligned}$$

$$\begin{aligned} & \{b, d, e\} \\ & \{b, d\} \rightarrow \{e\} \quad c = 3/4 = .75 \\ & \{b, e\} \rightarrow \{d\} \quad c = 3/4 = .75 \\ & \{d, e\} \rightarrow \{b\} \quad c = 3/5 = .6 \end{aligned}$$

- c. [5 points] Show how the 3-itemsets candidates can be generated by the $F_{k-1} \times F_{k-1}$ method and if these candidates will be pruned or not.



- d. [5 points] Consider the lattice structure given below. Label each node with the following letter(s): M if the node is a maximal frequent itemset, C if it is closed frequent itemset, F if it is frequent but neither maximal nor closed, and I if it is infrequent.



	a	b	c	d	e
1	M			M	
2		M			M
3	C	C			C
4	M			M	
5					
6					
7					
8					
9					
10					

Supp by Count(7)
 $\{d\}$
 Supp by Count(8)

6. https://github.com/NootCode/CS4210-Assignment-5/blob/master/ARM/association_rule_mining.py