

EXAM

732A61 and TDDD41
Data Mining –
Clustering and Association Analysis

732A75 Advanced Data Mining

May 10, 2020, kl 8-12

Teachers: Patrick Lambrix, José M Pena

This is an individual exam. No help from others is allowed. No uploading or downloading solutions is allowed. No communication regarding the exam is allowed.

You are allowed to use the course literature, the course slides and your own notes.
Answers to the exam questions may be sent to Urkund.

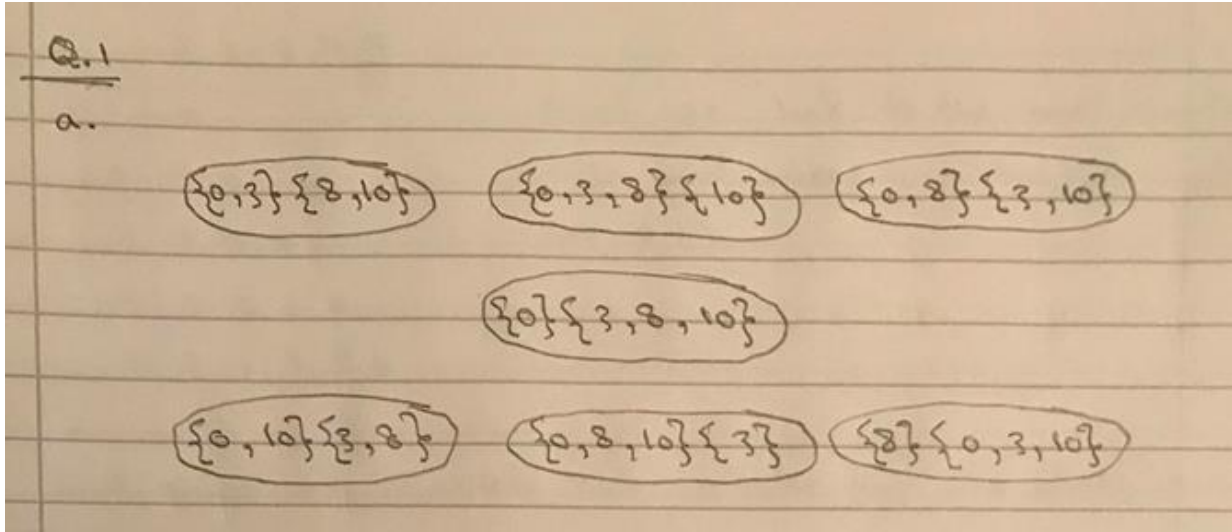
Instructions:

1. Start each question at a new page.
2. Write at one side of a page.
3. Write clearly.
4. If you make assumptions about a question, that are not explicitly stated, you need to write these down. (These assumptions cannot change the exercise or question.)
5. Hand in a pdf file. (You could use Word for text; take pictures of drawings/calculations and insert in a Word file; and then export to pdf.)
6. Hand in via LISAM before 12:00. If you have problems handing in via LISAM, send your answers via e-mail to Patrick.lambrix@liu.se latest 12:05 and keep trying to upload afterwards. Handing in after this time will be considered as not handed in.

GOOD LUCK!

1. Clustering by partitioning (1+2+2=5p)

a. Given the data set $\{0, 3, 8, 10\}$. Assume we use Euclidean distance and $k = 2$. Draw the graph representation of the clustering problem.



b. Start at an arbitrary node and show one iteration of the PAM algorithm on the graph. Give all steps in the computation and show at what node that iteration ends.

Starting at $\{0,3\}\{8,10\}$ the PAM calculates the total distance to the medoids which is 5 and jump to the next node and compare the total distance if it is more it will skip it to the next one.. the best node is $\{0,3\}\{8,10\}$ because it have the least total distance to the medoids (total cost)

c. For each of the questions below, answer yes/no and explain why.

- Does PAM guarantee to find a global optimum of the clustering problem?

Yes, because it takes all the data and it keep iterating until it reaches the global optimum.

- Does PAM guarantee to find a local optimum of the clustering problem?

No, unless the local optimum is the global optimum, it will keep iterating until it converges to the global optimum.

- Does CLARA guarantee to find a local optimum of the original clustering problem?

Yes, because it only takes some samples from the whole data so it will converges to a local optimum.

- Does CLARANS guarantee to find a local optimum of the clustering problem?

Yes, CLARANS can find the best local optimum because it keep randomly searching the whole graph as well as part of the graph until it converges to the best local optimum given the number of local minima & maximum number of neighbors to compare that it has provided with in the beginning.

2. Hierarchical clustering (3+2=5p)

a. (i) Show the different steps of the Agglomerative Hierarchical Clustering algorithm using the dissimilarity matrix below and *complete* link clustering. Give partial results after each step.

	1	2	3	4	5
1	0				
2	5	0			
3	9	10	0		
4	3	2	6	0	
5	7	1	4	8	0

Q.2

a. (i)

	1	2	3	4	5
1	0				
2	5	0			
3	9	10	0		
4	3	2	6	0	
5	7	1	4	8	0

First we look to the most similar pair and merge them together which are (2,5)

and since it is complete link we take the max distance with other points

$$d_{(2,5),1} = \max(d_{2,1}, d_{5,1}) = \max(5, 7) = 7$$

$$d_{(2,5),3} = \max(d_{2,3}, d_{5,3}) = \max(10, 4) = 10$$

$$d_{(2,5),4} = \max(d_{2,4}, d_{5,4}) = \max(2, 8) = 8$$

	1	(2,5)	3	4
1	0			
(2,5)	7	0		
3	9	10	0	
4	3	8	6	0

now we merge (1,4) because they are the most similar

$$d_{(1,4), (2,5)} = \max(d_{1,(2,5)}, d_{4,(2,5)}) = \max(7, 8) = 8$$

$$d_{(1,4), 3} = \max(d_{1,3}, d_{4,3}) = \max(9, 6) = 9$$

	(1,4)	(2,5)	3
(1,4)	0		
(2,5)	8	0	
3	9	10	0

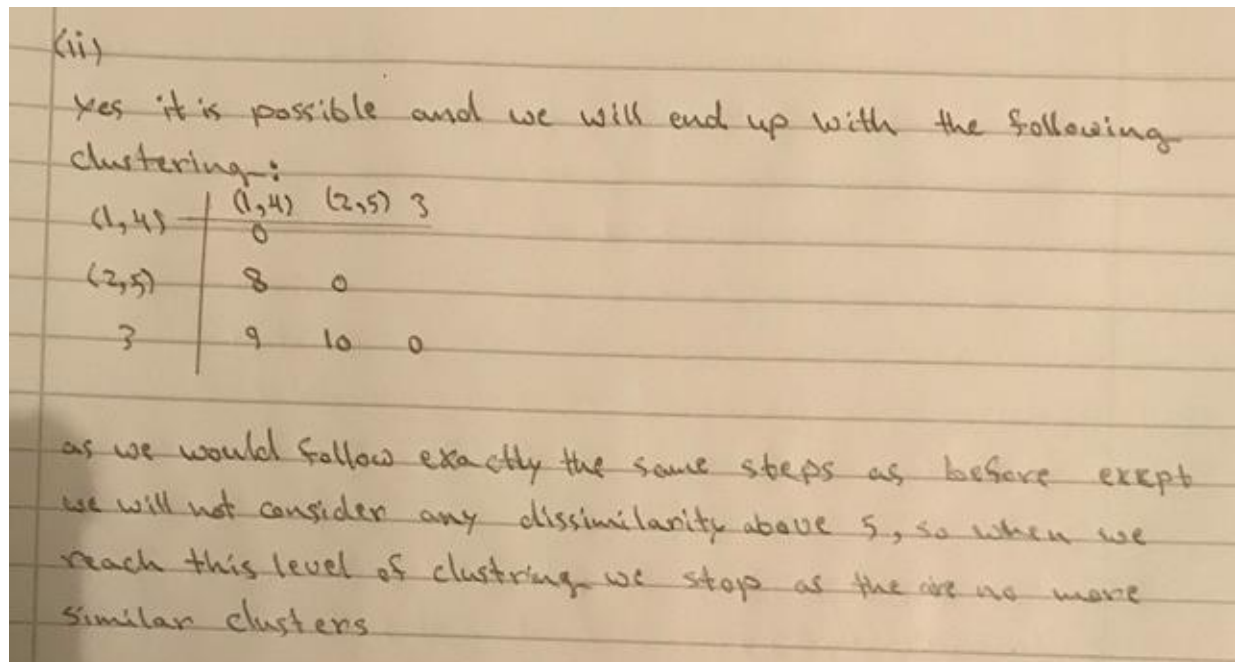
now we merge (1,4) and (2,5) because they are the most similar

$$d_{(1,2,4,5), 3} = \max(d_{(1,4), 3}, d_{(2,5), 3}) = \max(9, 10) = 10$$

	(1,2,4,5)	3
(1,2,4,5)	0	
3	10	0

now all the points (1,2,3,4,5) can be in one cluster

(ii) Would it be possible to optimize the computation above if you know that the threshold is 5? What is the result if the threshold is 5?



b. For the ROCK algorithm:

(i) Given the similarity matrix below. What is $\text{link}(A,B)$ if the threshold is 0.6?

	A	B	C	D	E
A	1				
B	0.9	1			
C	0.8	0.7	1		
D	0.1	0.2	0.5	1	
E	0.2	0	0.3	0.4	1

The link is (AB), (AC) and (BC)

(ii) If the number of elements in a cluster is n and the number of neighbors for each element in the cluster is m , what is the contribution of an object in the cluster to the expected link for the cluster? Explain why.

3. Density-based clustering (2p)

For the following statements say whether they are true or false.

If a statement is true, then prove it. (Observe that an example is not a proof.)

If a statement is false, then give a counterexample.

- If p and q are density connected wrt ϵ and MinPts , then q and p are density connected wrt ϵ and MinPts .
TRUE
- If p is density reachable from q wrt ϵ and MinPts , then q is density reachable from p wrt ϵ and MinPts .
- If p is directly density reachable from q wrt ϵ and MinPts , then p is density reachable from q wrt ϵ and MinPts
- If p is directly density reachable from q wrt ϵ and MinPts , then p and q are density connected wrt ϵ and MinPts .

4. Different types of data and their distance measures (2+2=4p)

a. What is the distance between Item K and Item L? (no normalization needed)

	A	B	C	D	E	F	G
Item K	(50,500)	(2,1,0)	Y	N	Y	N	77
Item L	(50,505)	(1,3,0)	Y	Y	N	N	no-value-available

Attribute A is interval-based and Euclidean distance is used.

Attribute B is interval-based and Manhattan distance is used.

Attributes C and D are binary symmetric variables.

Attributes E and F are binary asymmetric variables.

Attribute G is interval-based.

Handwritten calculation of the distance between Item K and Item L:

$$d(K, L) = (1) \sqrt{(50-50)^2 + (500-505)^2} + (1) (|2-1| + |1-3| + |0-0| + (1)(0) + (1)(1) + (1)(1) + (0)(0))$$

$$= \frac{5 + 3 + 1 + 1}{5} = 2$$

$S_{K,L}^E$ and $S_{K,L}^G$ are 0 because one is asymmetric binary attribute and the other is missing value

b. Can the formula for distance for objects with variables of mixed types (that you used in a) also be used for objects with only asymmetric variables? If no, explain why. If yes, state whether you would get the same results as with the method using contingency tables and explain why or why not.

b. Yes, it can be used and it will result in similar results because in the nominator we only sum the dissimilarities of the asymmetric items and in the denominator we take the sum of all similarity and dissimilarity except for similarity in zero.

if we have contingency table like below:

	object i		sum
	1	0	
	a	b	a+b
object j	0	c	d
	a+c	b+d	P

$$d(i, j) = \frac{b+c}{a+b+c} \text{ which is equal to } \frac{\sum_{k \in S(i, j)} d(i, j)}{\sum_{k \in S(i, j)} 1}$$

don't count asymmetric values that are 0 (or d is in the contingency table)

5. Apriori algorithm ($2p+2p+2p+2p+1p=9p$)

- a. Run the Apriori algorithm on the following transactional database with minimum support equal to two transactions. Explain step by step the execution.

Transaction id	Items
1	A,B
2	B,C
3	B,A
4	C,B

Ans

database	C_1	sup	L_1	C_2	L_2	C_3	sup	L_3
A,B	A	2	A	A,B	2	A,B,C	0	NULL
B,C	B	4	B	A,C	0	B,C	2	
B,A	C	2	C	B,C	2			
C,B								

- First we set the input of the algorithm with 2 parameters
D the data set and minsup the minimum support
- we start the algorithm by setting L_1 large item set with all possible combinations
- start a loop 1 which will call candidate generation function every iteration
- the candidate generation function when called it will run 2 loops one after another the first generate all the possible combinations of the previous large item set that inputted as parameter to this function and the second loop prune this list from item sets that does not meet the minimum sup
- after candidate list is generated it will be passed to second nested loop (nested in loop 1) to look for all candidates that exist in database, and if item exists it will be added to the next large item set

- b. Run the Apriori algorithm on the following transactional database with minimum support equal to two transactions, and the following additional constraint: Find the frequent itemsets that contain the item D. Explain step by step the execution. Make clear when and how the constraint is used. Incorporate the constraint into the algorithm, i.e. do not simply run the algorithm and afterwards consider the constraint.

Transaction id	Items
1	A, B, C, D
2	B, C, D, E

b.

- First we enter the data base as well as the minimum sup. and constraints as parameters

1 A,B,C,D minsup = 2

2 B,C,D,E constraints = D ∈ L;

- we run loop1 and call candidate generation function each iteration.

- in the first iteration we get the following C_1 which has all the combinations of one item that exist in the data base:

A 1

B 2

C 2

D 2

E 1

- this C_1 goes then to another nested loop in the main func. that look for items that satisfy both minsup and the constraints which yields L, large item set since it is under above constrain.

D 2

- this then goes to the candidate generation function but it will stop since there are no more items to join

- c. Run the Apriori algorithm on the last transactional database with minimum support equal to two transactions, and the following two additional constraints: Find the frequent itemsets that (1) contain the item D and (2) do not contain the item B. Explain step by step the execution. Make clear when and how the constraint is used. Incorporate the constraint into the algorithm, i.e. do not simply run the algorithm and afterwards consider the constraint.

c.

- A, B, C, D minsup = 2
 B, C, D, E cons₁ = $D \in L_i$; cons₂ = $B \notin L_i$
- A 1 first C_i is generated
 B 2
 C 2
 D 2
 E 1
- when passed in the nested loop of main function it will check for both minsup and anti monotone constraints ($D \in L_i$), so we get:
 D 2
- before we iterate the main loop again we check for monotone constraints this time $B \notin L_i$
 so D 2 will be passed
- second iteration starts and candidate generation function is called, since there is only D 2 in the input it will stop

- d. Apply the rule generation algorithm to the frequent itemset ABC on the database below in order to produce association rules with confidence greater or equal than 50 %. Explain step by step the execution.

Transaction id	Items
1	C, B, A
2	D, C, A
3	A, B
4	A, B
5	A, D
6	A, D

- e. In the light of the previous result, propose an improvement to the rule generation algorithm.

6. FP grow algorithm (3p)

Run the FP grow algorithm on the following transactional database with minimum support equal to one transaction. Explain step by step the execution.

Transaction id	Items
1	C, B, A
2	D, C, A
3	A, B
4	A, B
5	A, D
6	A, D

Q6

1. C, B, A
2. D, C, A
3. A, B
4. A, B
5. A, D
6. A, D

minsup = 1

- we count the support for each item

A	B	C	D
6	3	2	3

- all items meet min sup

- sort them in descending order

A, B, C
A, D, C
A, B
A, B
A, D
A, D

- create a FP tree

- we run a recursive insert-tree function that increase the count for each item whenever the item exist in database

```
graph TD
    Root["{}"] --> A["A (4)"]
    Root --> B["B (1)"]
    A --> C1["C (1)"]
    A --> D1["D (1)"]
    B --> C2["C (1)"]
```

- a conditional data base is created

A: - - -

B: A:3

D: A:3

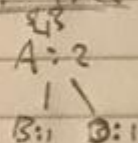
C: A,B:1, A,D:1

- C- conditional data base

1. A,B → sorted A, B

2. A,D A, D

⑥ FP tree



- B- conditional data base

1. A A

2. A → sorted A →

3. A A

{3}

A:3

- D conditional data base

1. A A

2. A → sorted A

3. A A

{3}

A:3

- A conditional data base

1. -

{3}

7. Constraints (1p+1p+1p=3p)

- a. Prove that the constraint "the average price of a set of items is greater than 10" is neither monotone nor antimonotone.
- b. Prove that the previous constraint is convertible monotone. Check the slides for the definition of convertible monotone constraint.
- c. Prove that the previous constraint is convertible antimonotone as well. Check the slides for the definition of convertible antimonotone constraint.