

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

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

- (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

- (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 .
- 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.

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.

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

- 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

- 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.
- 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

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

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