

EXAM

732A61 and TDDD41
Data Mining –
Clustering and Association Analysis

732A75 Advanced Data Mining

June 9, 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 (except with the teachers mentioned above).

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 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.

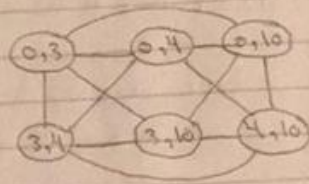
GOOD LUCK!

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

- a. Given the data set $\{0, 3, 4, 10\}$. 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. Assume $\text{numlocal} = 1$ and $\text{maxneighbor} = 2$. Start at the same node as in question b and show one iteration of the CLARANS algorithm on the graph. Give all steps in the computation and show at what node that iteration ends.

Q.1 $\{0, 3, 4, 10\}$ $K=2$

a.



each node represents a potential solution

b.

- Start at node $\{0, 3\}$ arbitrarily and assign data points to closest centroid and calculate the cost

$$= \sqrt{(0-0)^2} + \sqrt{(3-3)^2} + \sqrt{(3-4)^2} + \sqrt{(3-10)^2} = 8$$

- calculate the total swapping cost TC_{ij} for each pair of non-selected object j and selected object i

	cluster 1	cluster 2	
$TC_{0,4}$	$\{0, 3\}$	$\{3, 4, 10\}$	$= 7 - 8 = -1$
$TC_{0,10}$	$\{0, 3, 4\}$	$\{10\}$	$= 7 - 8 = -1$
$TC_{3,4}$	$\{0, 3\}$	$\{4, 10\}$	$= 9 - 8 = 1$
$TC_{3,10}$	$\{0, 3, 4\}$	$\{10\}$	$= 4 - 3 = -1$

PAM will choose the node with the minimum swapping cost in this case $\{3, 4\}$ and new iteration starts

c.

- start at $\{0, 3\}$
- select random neighbor $\{0, 4\}$
- check cost differential
- if $C_{\{0, 3\}} > C_{\{0, 4\}} \rightarrow \{0, 4\}$ current node
- else increase the neighbor counter
- select another random neighbor
- * in this case the condition is true
- re-set the neighbor counter to 1 and select a new node
- when neighbor counter reach 2, it will check if we reached numLocal - if true \rightarrow out put

- in this case next node is $\{3, 10\}$

$c\{3, 10\} > c\{0, 4\} \Rightarrow \{3, 10\}$ is current node

- check if numLocal is $> 1 \rightarrow$ true \rightarrow output

best node $\{3, 10\}$

2. Hierarchical clustering (2+2+1=5p)

a. Show the different steps of the Agglomerative Hierarchical Clustering algorithm using the dissimilarity matrix below and *single* 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

b. For the ROCK algorithm:

Assume the similarity matrix below. Assume that cluster C1 contains the data objects A, B and E. Assume that cluster C2 contains the data objects C and D. What is $\text{Link}(C1, C2)$? Show your computations.

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

c. In BIRCH, given a cluster with the data points (1,2), (1,3) and (2,2), what is its cluster feature?:

Q.2 1 2 3 4 5 1 {2,5} 3 4

a.	1	0				1	0			
	2	5	0			{2,5}	5	0		
	3	9	10	0		3	9	4	0	
	4	3	2	6	0	4	3	2	6	0
	5	7	1	4	8	0				

- take the most similar pair and merge them together

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

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

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

	1	{2,5,4}	3	
1	0			
{2,5,4}	3	0		
3	9	4	0	

$$d_{(2,5,4),1} = \min(d_{(2,5),1}, d_{4,1}) = \min(5, 3) = 3$$

$$d_{(2,5,4),3} = \min(d_{(2,5),3}, d_{4,3}) = \min(4, 6) = 4$$

	{2,5,4,1}	3	
{2,5,4,1}	0		
3	3	3	0

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

b. Link is the number of common neighbors between C_1 and C_2

- for C_1 :

- Neighbors of A:

$$\text{sim}(A, A) = 1 \geq 0.6, \text{sim}(A, B) = 0.9 \geq 0.6$$

$$\text{sim}(A, C) = 0.8 \geq 0.6$$

\therefore neighbors of A are: A, B, C

- Neighbors of B: B, A, C

- " " E: E

- for C_2 :

- neighbors of C: A, B, C, D

" " D: C, D

$$\text{Link}(C_1, C_2) = \frac{C_1 \cap C_2}{C_1 \cup C_2} = \frac{3}{5} = 0.6$$

c. $CF = (N, LS, SS)$

where:

$N = \text{number of data points} = 3$

$$LS = \sum_{i=1}^N \vec{X}_i = (4, 7)$$

$$SS = \sum_{i=1}^N X_i^2 = (6, 17)$$

$$\Rightarrow CF = (3, (4, 7), (6, 17))$$

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 is density reachable from p 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 density reachable from q wrt ϵ and MinPts , then p and q are density connected 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 .

Q-3

1- False

for q to be DR from p there need to be a chain of points p_1, \dots, p_n where $p_1 = q, p_n = p$ such that p_{i+1} is DDR from p_i

but if p and q are only density connected there is a point (a) that both p and q are only DR from

2- False

because p need to be core to be DR from q

3- True

if p is DR from $q \Rightarrow$ there exist a point (a) that both p and q are DR from $\Rightarrow p$ and q are density connected

4- True

if p is DDR from $q \Rightarrow p$ is DR from q
there exist a common point (a) which both p & q are DR from

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	(35,20,4)	(2,5,1)	Y	N	Y	N	77
Item L	(40,20,4)	(1,2,1)	Y	Y	Y	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. Assume we have categorical data. One method to define a distance between two data objects is $(p-m)/p$ where p is the total number of categorical variables and m is the number of categorical variables for which there is a match between the objects. A second method is to introduce a new asymmetric binary variable for each of the possible values for each of the categorical variables. Give a formula for the distance between two objects in the second method in terms of p and m (where p and m have the same meaning as above; i.e. p is the number of categorical variables - not the number of introduced binary variables, and m is the number of matches in the categorical variables). Show how you obtained the formula.

$$Q.4 \quad A = \sqrt{(35-40)^2 + (20-20)^2 + (4-4)^2} = 5$$

a.

$$B = |2-1| + |5-2| + |1-1| = 4$$

$$C = 0$$

$$D = 1$$

$$E = 0$$

$$F = 0$$

$$G = 0$$

$$d(k, l) = \frac{(1)(5) + (1)(4) + (1)(0) + (1)(1) + (1)(0) + (0)(0) + (0)(0)}{5}$$

$$= \frac{10}{5} = 2$$

$$b. \quad \sin = 1 - \frac{p-m}{p} = \frac{p}{p} - \frac{p-m}{p} = \frac{m}{p}$$

$$\begin{array}{cc} & B \\ & \begin{array}{cc} 1 & 0 \\ \frac{m}{p} & \frac{p-m}{p} \\ 0 & \frac{p-m}{p} \end{array} \\ A & \begin{array}{cc} 1 & 0 \\ \frac{m}{p} & \frac{p-m}{p} \\ 0 & \frac{p-m}{p} \end{array} \end{array}$$

$$d(A, B) = \frac{\frac{p-m}{p} + \frac{p-m}{p}}{\frac{m}{p} + 2\left(\frac{p-m}{p}\right)} = \frac{\frac{2(p-m)}{p}}{\frac{2p-m}{p}} = \frac{2p-2m}{2p-m}$$

5. Apriori algorithm (3+2+2=7p)

- 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
2	B,C,D
3	B,D
4	B,C,E
5	C,D
6	C,E
7	F
8	A
9	B,C,D
10	B,D
11	B,C,E
12	C,D
13	C,E

- b. Run the Apriori algorithm on the previous transactional database with minimum support equal to two transactions, and the following additional constraint: Find the frequent itemsets that contain the itemset {B, C}. 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. Run the Apriori algorithm on the previous transactional database with minimum support equal to two transactions, and the following additional constraint: Find the frequent itemsets that do NOT contain the itemset {B, C, 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.

Q.5 minSup = 2

a. First we order the items in alphabetical order, in this case it is already sorted

Generate candidate list size 1, and calculate the support

<u>C₁</u>	<u>Sup</u>	<u>L₁</u>	<u>output</u>
A	2	A	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> A, B, C, D, E </div> F will not go to output because it did not satisfy minSup
B	6	B	
C	8	C	
D	6	D	
E	4	E	
F	1		

<u>C₂</u>	<u>Sup</u>	<u>L₂</u>	<u>output</u>
AB	0	BC	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> BC, BD, BE, CD, CE </div>
AC	0	BD	
AD	0	BE	
AE	0	CD	
BC	4	CE	
BD	4		items AB, AC, AD, AE and CE will not go to output because they did not satisfy minSup
BE	2		
CD	4		
CE	4		
DE	0		

<u>C₃</u>	<u>Sup</u>	<u>L₃</u>	<u>output</u>
BCD	2	BCD	<div style="border: 1px solid black; padding: 5px; display: inline-block;"> BCD, BCE </div>
BCE	2	BCE	

* Items BDE and CDE are not in C₃ because their subset DE are not in L₂, i.e. subset checking

<u>C₁</u>	<u>Sup</u>	<u>L₁</u>	<u>output</u>
BCDE	0	-	-

- stop no more items

b. - First we check if the constraint is monotone or antimonotone
 In this case it is monotone since if it is true for A
 it will be true for B where $A \subseteq B$

- items already ordered

<u>C₁</u>	<u>Sup</u>	<u>Const.</u>	<u>L₁</u>	<u>output</u>
A	2	F	A	
B	6	F	B	
C	3	F	C	
D	6	F	D	
E	4	F	E	
F	1	F		

* because the const. is
 monotone we need to
 keep the items that did
 not satisfy const. but
 will not go to output

<u>C₂</u>	<u>Sup</u>	<u>Const.</u>	<u>L₂</u>	<u>output</u>
AB	0	F	BC	
AC	0	F	BD	<u>BC</u>
AD	0	F	BE	
AE	0	F	CD	
BC	4	T	CE	
BD	4	F		
BE	2	F		
CD	4	F		
CE	4	F		
DE	0	F		

<u>C₃</u>	<u>Sup</u>	<u>Const.</u>	<u>L₃</u>	<u>output</u>
BCD	2	T	BCD	
BCE	2	T	BCE	[BCD, BCE]

* items BDE and CDE did not pass subset checking

<u>C₄</u>	<u>Sup</u>	<u>Const.</u>	<u>L₄</u>	<u>output</u>
BCDE	0	T	-	-

- no items → stop

C

check if the const. is monotone or antimonotone
in this case it is antimonotone since if it is true for A it is true for B where $B \subseteq A$

- items already ordered

<u>C₁</u>	<u>Sup</u>	<u>Const. 1</u>	<u>Const. 2</u>	<u>output</u>
A	2	F	T	A
B	6	F	T	B
C	8	F	T	C
D	6	F	T	D
E	4	F	T	E
F	1	F	T	

<u>C₂</u>	<u>sup</u>	<u>const.1</u>	<u>const.2</u>	<u>L₂</u>	<u>output</u>
AB	0	F	T	BC	
AC	0	F	T	BD	<u>[CE]</u>
AC	0	F	T	BE	
AE	0	F	T	CD	
BC	4	T	T	CE	
BD	4	F	T		
BE	2	F	T		
CD	4	F	T		
CE	4	F	T		
DE	0	F	T		

<u>C₃</u>	<u>sup</u>	<u>const.1</u>	<u>const.2</u>	<u>L₃</u>	<u>output</u>
BCD	2	T	F	BCE	<u>[BCE]</u>
BCE	2	T	T		

* items BDE and CDE are not in C₃ because they did not pass subset checking

- no more items → stop

6. FP grow algorithm (4p)

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	B, A
4	B, A
5	D, A
6	D, A

Q.6 1 C, B, A

2 D, C, A

3 B, A

4 B, A

5 D, A

6 D, A

minsup = 1

we count support for each item

A B C D

6 3 2 3

we take out items that do not satisfy minsup, in this case all items satisfy minsup

Sort the items in support descending order

A B D C

6 3 3 2

1 A, B, C

2 A, D, C

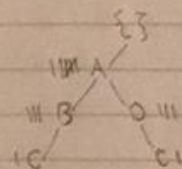
3 A, B

4 A, B

5 A, D

6 A, D

we create FP tree by writing a recursive insert-tree function that increase the count of each item whenever the item exist in database



- Construct conditional database where X-conditional database consists of all the prefix paths leading to X in the FP tree

A -

B A: III

D A: III

C AB: I, AD: I

- A conditional database empty

- B conditional database

1 A

2 A

3 A

- * sort in support descending order and prune items that do not satisfy support, in this case no change

- * output the frequent 1-iter sets by adding B as a suffix

output [AB]

- * construct FP tree and the conditional database

{ } then conditional database
III A / A -

- * re-start the process, no more items within B conditional database

- D conditional database

1 A

2 A

3 A

* After finding frequent 1-itemsets and sorting the transactions in support descending order, we have

- 1 A
- 2 A
- 3 A

* output the frequent 1-itemsets by adding D as a suffix

output [AD]

* construct FP tree and the conditional database

$\{ \}$ item conditional database
 111 A / A -

* re-start the process, no more items within D conditional database

C - conditional database

- 1 A, B
- 2 A, D

* After finding frequent 1-itemsets and sorting accordingly, we have

- 1 A, B
- 2 A, D

* output the frequent 1-itemsets by adding C as a suffix

output [AC, BC, DC]

* construct FP tree and the conditional database

$\{ \}$ item
 111 A / A -
 11 B / B A:1
 11 D / D A:1

* re-start the process
- BC - conditional database
IA

* prune and sort accordingly
IA

* output [ABC]

* FP tree $\{ \}$
IA'

* re-start, no items left

- DC - conditional database
IA

* prune and sort accordingly
IA

* output [ADC]

* FP tree $\{ \}$
IA'

* re-start, no items left - go back, no items - go back, no items
Done

7. Rule generation (4p)

Apply the rule generation algorithm to the frequent itemset {A, B, C} 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

1. C, B, A frequent Itemset: ABC
 2. D, C, A min Conf: 0.5
 3. A, B
 4. A, B
 5. A, D
 6. A, D

genrules(ABC, ABC, 0.5)

$X \rightarrow Y$

* if X doesn't result in a rule with min Conf, so neither does any subset of X, so we start with the largest possible antecedent to take advantage of this rule

$A = \{AB, AC, CB\}$ antecedent subset of size 2

$AB \rightarrow C$

$$\text{Conf}(AB \rightarrow C) = \frac{\text{sup}(ABC)}{\text{sup}(AB)} = \frac{1}{3} < 0.5$$

* does not generate rule we stop and don't do recursive calls to generate for subsets of AB

$AC \rightarrow B$

$$\text{Conf}(AC \rightarrow B) = \frac{\text{sup}(ABC)}{\text{sup}(AC)} = \frac{1}{2} = 0.5 \Rightarrow \boxed{AC \rightarrow B} \text{ out put}$$

we call genrules(ABC, AC, 0.5) to check for subsets of size 1 as antecedents to the rule

$A \rightarrow CB$ $\{A, C\}$

$$\text{Conf}(A \rightarrow CB) = \frac{\text{sup}(ABC)}{\text{sup}(A)} = \frac{1}{6} < 0.5$$

* does not generate rule so we don't call genrules

$C \rightarrow AB$

$$\text{Conf}(C \rightarrow AB) = \frac{\text{sup}(ABC)}{\text{sup}(C)} = \frac{1}{2} = 0.5 \Rightarrow \boxed{C \rightarrow AB} \text{ out put}$$

* C does not have subsets so we stop and go back to antecedent subset size 2

- $BC \rightarrow A$

$$\text{Conf}(BC \rightarrow A) = \frac{\text{Sup}(ABC)}{\text{Sup}(BC)} = \frac{1}{1} > 0.5 \rightarrow \underline{BC \rightarrow A} \text{ output}$$

* we call generate $(ABC, BC, 0.5)$ to look in antecedent subset size 1 of item $BC \rightarrow \{B, C\}$

- $B \rightarrow CA$

$$\text{Conf}(B \rightarrow CA) = \frac{\text{Sup}(ABC)}{\text{Sup}(B)} = \frac{1}{2} < 0.5$$

* does not generate a rule, we stop and go to the next item in the subset

- $C \rightarrow AB$

$$\text{Conf}(C \rightarrow AB) = \frac{\text{Sup}(ABC)}{\text{Sup}(C)} = \frac{1}{2} = 0.5 \rightarrow \underline{C \rightarrow AB} \text{ output}$$

C does not have subsets and all larger subsets have been generated so we stop