

Question 1

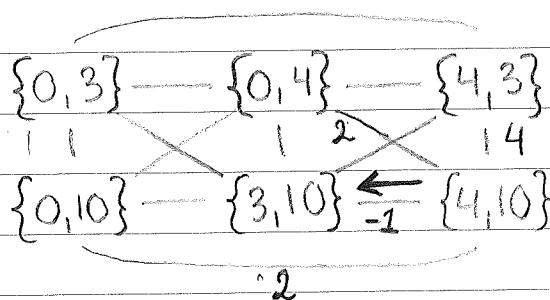
a) i. A node contains k objects and represent a proposal of the cluster medoids.

ii. Two nodes are neighbors when they differ with exactly one object. A node has $k(n-k)$ neighbors

iii. PAM and CLARANS are the algorithms with the graphs containing the most nodes. They search the same graph with the difference that PAM checks all neighbors of a node while CLARANS checks only the set maximum number of neighbors. CLARA takes samples from the data to search and therefore the graph contains fewer nodes.

iv. None of the algorithms promise a global optimum, but PAM guarantees to find a local optimum (which the two others don't)

b)



The node $\{4,10\}$ is chosen as starting point.

Swapping costs to it's neighbors are calculated with

$$TC_{2ih} = \sum C_{ijh} = \sum_i (d(x_i, x_j) - d(x_i, x_h)).$$

They can be seen in the graph.

The algorithm then moves on to the neighbor with lowest swapping cost (as long as it's negative). If there are only positive TC_{2ih} then the algorithm stops.

In this case the algorithm moves on to $\{3,10\}$.

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad number: Sheet number: 2

Question 2

a) i. The BIRCH algorithm consists of two steps:

1 Build a CF-tree → how?

2 Use a cluster algorithm of your choice on the sub clusters in the leaf nodes.

ii. Cluster Feature Vector contains information about the cluster : n = number of observations

LS = sum of the observations

SS = sum of the squared observations

For the case given, the CF is:

$$n = 3$$

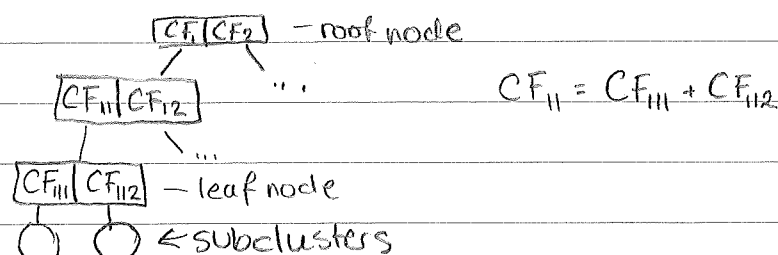
$$LS = (1+1+2=4, 2+3+2=7)$$

$$SS = (1+1+4=6, 4+9+4=17)$$

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

iii. The CF tree is a height balanced tree. It consists of a root-node, middle level nodes and leaf nodes.

The leaf nodes contains information on the data points in the subclusters with the help of CFs. The leaf-nodes are connected to parent nodes which holds information about their children (with combined CFs). In the end the root node contains information about all the data points by combining it's childrens CFs. An example:



iv. Branching factor specifies how many children a parent can have. Threshold states the max. allowed diameter of subclusters.

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN7

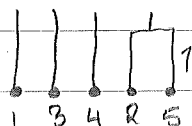
Blad nummer: Sheet number:
3

Question 2

b) All clusters start out as individual clusters. The two clusters that are the least dissimilar (according to complete link) are joined together in each step. In the end, all data points form a single cluster.

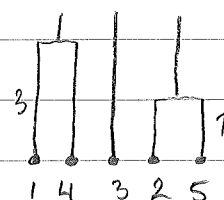
Step 1:

	1	2	3	4	5
1	0				
2	50	0			
3	9	10	0		
4	8	2	6	0	
5	7	4	8	0	0



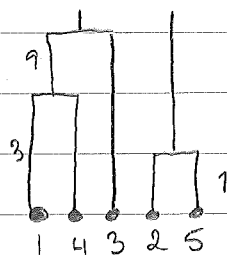
Step 2:

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



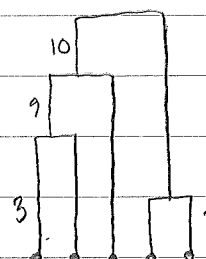
Step 3:

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



Step 4:

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



2

Question 3

A core point is a data point that has at least min.point of data points within its neighborhood, defined by radius ϵ . ϵ and min.point is set by the user. A point p is directly density-reachable from point q if q is a core point and p lies inside its ϵ -neighborhood. The point p is density reachable from q if there exists a chain of points (q as start and p as end) where each point is directly density-reachable from the point one step previous in the chain. The points p and o are density-connected if there exists a point q from which both p and o are density-reachable. In DBSCAN a cluster is formed by density-connected data points. If a point isn't density-connected to any other points it's considered to be an outlier since it's not placed in a dense enough region.

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad nummer: Sheet number: 5

Question 4

a) To compute the distance between K and L the following formula is used:

$$d_{KL} = \frac{\sum_{f=1}^p \delta_f^f \times d_{KL}^f}{\sum_{f=1}^p \delta_f^f}$$

where d_{KL}^f is the distance between the objects for attribute f and $\delta_f^f = 0$ if f has missing value or f is asymmetric binary and is of case 0,0. Otherwise $\delta_f^f = 1$

$$d_{KL}^A = \sqrt{(10-10)^2 + (500-505)^2} = 5 \quad d_{KL}^B = |2-1| + |1-3| + |1-1| = 3$$

$$d_{KL}^C = 0 \quad d_{KL}^D = 1 \quad d_{KL}^E = 1 \quad d_{KL}^F = 0$$

$$d_{KL} = \frac{1 \times 5 + 1 \times 3 + 1 \times 0 + 1 \times 1 + 1 \times 1 + 0 \times 0 + 0 \times d_{KL}^G}{5 \times 1} = \frac{10}{5} = 2$$

b) Create ranking for the observations $r = \{r_1, \dots, r_r\}$ where M_r is the highest rank and r_i is the lowest rank. Then you need to normalize the ranks by

$$z_i = \frac{r_i - 1}{M_r - 1}$$

An interval-based measure can now be used on the normalized z -values.

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad nummer: Sheet number: 6
--

Question 5

a) min. sup = 1

C_1		L_1	C_2		L_2
Item	S		Items	S	
A	6	→	AB	3	→
B	3		AC	2	
C	2		AD	3	
D	3		BC	1	
		↑ output	BD	0	↑ output
			CD	1	

Candidates of size 1
check support

C_3		L_3
Items	S	
ABC	1	→
ABD	—	
ACD	1	
		↑ output

Candidates size 2
created by self join (joined by prefix)
check support

Candidates size 3
created by self join
prune itemsets that have
a subset not included
in L_2 . Check support

Question 5

b)

C_1			L_1	C_2			L_2
Item	C	S		Item	C	S	
A	✓	6	A	AB	✓	3	AB
B	✓	3	B	AC	✓	2	AC
C	✓	-	C	AD	✓	3	AD
D	✓	-	D	BC	✓	1	BC
				BD	✓	0	too low support
				CD	X		constraint not fulfilled

The constraint is antimonotone!

C_3			L_3
Item	C	S	
ABC	✓	1	→ ABC
ABD	-	-	} subsets not in L_2
ACD	-	-	

The constraint is checked after the controll that all subsets of the candidates are frequent (is part of L_{k-1}) and before checking support. If the constraint is false for an itemset this itemset is pruned and there is no need to check support for this itemset.

c) $\min.conf \geq 0.5$

$AB \rightarrow C$

$$c = \frac{\text{Sup}(ABC)}{\text{Sup}(AB)} = \frac{1}{3} \times$$

$AC \rightarrow B$

$$c = \frac{\text{Sup}(ABC)}{\text{Sup}(AC)} = \frac{1}{2} \checkmark$$

$BC \rightarrow A$

$$c = \frac{\text{Sup}(ABC)}{\text{Sup}(BC)} = \frac{1}{1} \checkmark$$

Because this rule has too low confidence we don't need to check the rules $A \rightarrow BC$ and $B \rightarrow AC$

$C \rightarrow AB$

$$c = \frac{\text{Sup}(ABC)}{\text{Sup}(C)} = \frac{1}{2} \checkmark$$

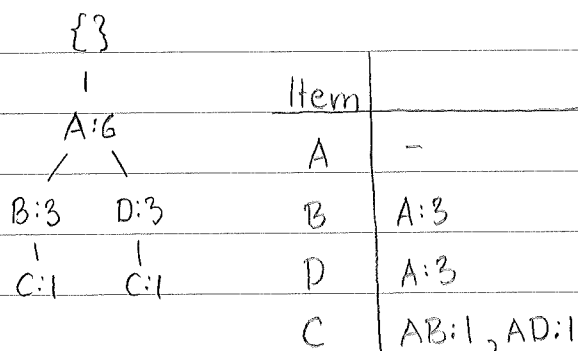
Question 6

a) Sort the items after support and remove those that aren't frequent ($\text{sup} \leq \text{min. sup} = 1$). Reorder transactions.

Item	A	B	D	C	TID	Items
Sup	6	3	3	2	1	ABC
					2	ADC
					3	AB
					4	AB
					5	AD
					6	AD

IN THE OUTPUT

Build the FP-tree based on the reordered transactions



Repeat the steps above for the conditional data bases

A-conditional is
empty, no need
to check

B-conditional

A, A, A

Item A

Sup 3

↓
OUTPUT: AB

{}	Item
A:3	A

AB-cond empty, no need
to check further

AID-number: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad nummer: Sheet number: 9

Continue question 6a)

D-conditional

A, A, A	{}	Item	
Item A	1	A	-
Sup 3	A:3		

↓

AD-cond empty, no need to check

OUTPUT: AD

C-conditional

AB, AD	{}	Item	
Item A B D	A:2	A	-
Sup 2 1 1	B:1 D:1	B	A:1
		D	A:1

↓

OUTPUT: AC, BC, DC

BC-conditional

A	{}
Item A	A:1
Sup 1	

↓

Output: ABC

DC-conditional

A	{}
Item A	1
Sup 1	A:1
	2

↓

Output: ADC

b) Antimonotone constraints are incorporated the same way support checking is. If an item/itemset doesn't fulfill the constraint it is pruned before the reordering of the items in the data base and is thus not part of any resulting tree or conditional data base.

Monotone are checked when support is checked. If not fulfilled the item/itemset is removed from output but not from tree or data base. If the constraint is true for an item/itemset you needn't check it for the conditional data bases that follow that item/itemset.

0.5
2 → missing places

0

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad nummer: Sheet number: 10

Question 6

- c) The main advantage of FP grow over Apriori algorithm is that it doesn't produce candidates and thus saves time (especially for low minimum support) and storage. In addition it only need to scan the original data base once, then it can be discarded (unlike in Apriori algorithm)

AID-nummer: AID-number: 2691	Datum: Date: 2018-03-17
Kurskod: Course code: 732A75	Provkod: Exam code: TEN1

Blad nummer: Sheet number: 11

Question 7

a) A convertible monotone constraint is $\text{avg}(S) \geq p$ where S is the prices of the itemset and items are ordered after ascending price (item with lowest price first). Adding a new item with higher price can only increase average price and that mean that the constraint is convertible monotone.

A convertible antimonotone constraint is $\text{avg}(S) \geq p$ when the items are ordered after descending price. Adding an item with lower price than previous items can only decrease average price and the constraint is convertible monotone.

In both cases it can be shown that the constraint $\text{avg}(S) \geq p$ is neither monotone nor antimonotone since adding an item (when there is no specific order) can both increase and decrease the price - we can't know.

3