

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

Blad nummer:
Sheet number:
1

Q-1 @ PAM

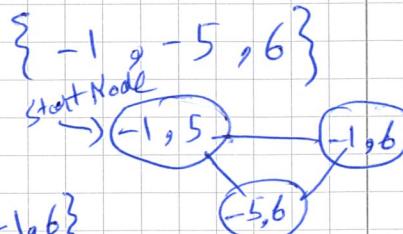
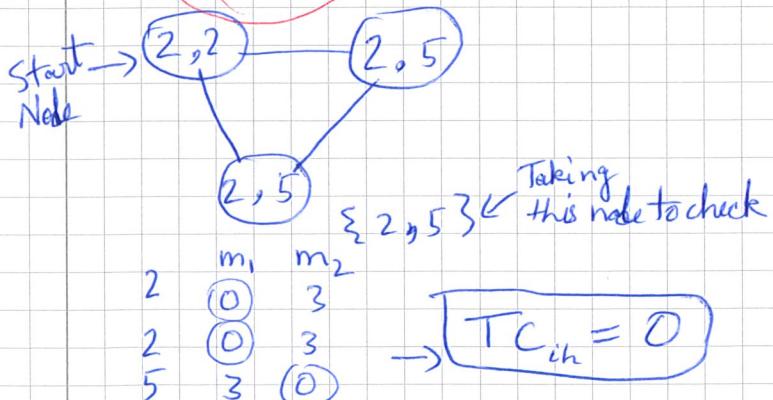
(i) PAM is better than K means as it uses medoids -
It is good for small datasets -
Algo

- ① Select K arbitrary nodes as cluster medoids -
- ② Select ~~n~~ n arbitrary points i - For each selected object i and non selected object h calculate the total swapping cost TC_{ih} , do it for all items -
Repeat -
- ③ Select the ~~other~~ ~~not~~ objects with the minimum swapping cost TC_{ih} -
- ④ If $TC_{ih} < 0$ replace i by h and repeat the process for all the objects in the dataset -
- ⑤ Set the current cost to the minimum cost and current node the best node so far -
Repeat for all items -

(ii) Swapping cost - Is the cost got by replacing the selected item i and the non selected item h in the data set - TC_{ih} -

Is the cost we get if we replace two items we replace items with each other if the cost is minimum or $TC_{ih} < 0$ -

(iii) $\{2, 2, 5\}$?



	m_1	m_2
-1	0	1
-5	-1	1
0	1	0

$$TC_{ih} = -4$$

Strictly Negative

AID-nummer: <i>AID-number:</i>	2203	Datum: <i>Date:</i>	18-06-05
Kurskod: <i>Course code:</i>	732A75	Provkod: <i>Exam code:</i>	TENT

Blad nummer:
Sheet number:

2

Q-1-b

~~PAM, CLARA~~ 0,5

Q-1-c

None

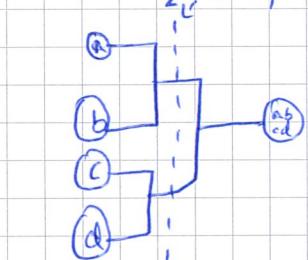
1

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732 AT 5	Provkod: Exam code:	TENT	3

Q - 2 Agglomerative Hierarchical Clustering -

Agglomerative Hierarchical clustering , AGNES is a bottom up hierarchical clustering algorithm which uses dissimilarity matrix to make clusters of the dataset . It combines the items based on the similarity between them to make a cluster until the user specified threshold -

² user specified threshold



It uses different types of links to combine clusters complete link, average link, medoids, centroid etc -

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

- First we check the minimum item in the matrix which is 1 in (1, 4) so we will combine 1 and 4 -

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

$$\text{dist}((1, 4), 2) = \text{Max}(\text{dist}(1, 4), \text{dist}(4, 2)) \\ = \text{Max}(\text{dist}(1, 2), \text{dist}(4, 2))$$

$$= \text{Max}(8, 7)$$

$$= 8$$

$$\text{dist}((1, 4), 3) = \text{Max}(\text{dist}(1, 4), \text{dist}(3)) \\ = \text{Max}(\text{dist}(1, 3), \text{dist}(4, 3)) \\ = \text{Max}(3, 9) = 9$$

: We are using Max because it is complete link clustering.

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	4

$$\begin{aligned}
 \text{dist}((1,4), 5) &= \text{Max}(\text{dist}(1,4), \text{dist}5) \\
 &= \text{Max}(\text{dist}(1,5), \text{dist}(4,5)) \\
 &= \text{Max}(10, 5) = 10
 \end{aligned}$$

② Now again we check for minimum value which is 2 in 2 and 5, so we combine (2, 5)

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

$$\begin{aligned}
 \text{dist}((1,4), (2,5)) &= \text{Max}(\text{dist}(1,4), \text{dist}(2,5)) \\
 &= \text{Max}(\text{dist}(1,2), \text{dist}(4,5)) \\
 &= \text{Max}(\text{dist}((1,4), 2), \text{dist}((1,4), 5)) \\
 &= \text{Max}(8, 10) \\
 &= 10
 \end{aligned}$$

$$\begin{aligned}
 \text{dist}(3, (2,5)) &= \text{Max}(\text{dist}(3, 2), \text{dist}(3, 5)) \\
 &= \text{Max}(4, 6) \\
 &= 6
 \end{aligned}$$

③ Now again checking for minimum value which is 6 so combining $(3, (2,5))$ -

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

$$\begin{aligned}
 \text{Max}(\text{dist}(1,4), \text{dist}(3,2,5)) &= \text{Max}(\text{dist}((1,4), 3), \text{dist}(1,4), (2,5)) \\
 \text{Max}[\text{dist}(1,4), \text{dist}(3,2,5)] &= \text{Max}(9, 10) \\
 &= 9
 \end{aligned}$$

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	5

Q.3

ROCK

Rock algorithm is a clustering algorithm used for categorical data because traditional hierarchical algorithms do not do well with categorical data because categorical data sample space is discrete and we cannot apply distance measure on it-

ROCK algorithm uses cluster neighbors, common neighbors to combine the clusters -

In ROCK if two points ~~have~~ are similar to each other they are neighbors. If points have other elements ^{similar} ~~common~~ to them they are common neighbors -

Let

$$\langle a, b, c, d, e \rangle = \{ \{a, b, c\}, \{a, b, d\}, \{a, b, e\}, \\ \{a, c, d\}, \{a, c, e\}, \{a, d, e\}, \\ \{b, c, d\}, \{b, c, e\}, \{b, d, e\}, \{c, d, e\} \}$$

$$\cancel{\langle a, b, c, d \rangle} = \{ \{f, g, h\}, \{f, g, i\}, \{f, h, i\} \}$$

$$\langle a, b, e, f \rangle = \{ \{a, b, e\}, \{a, b, f\}, \{a, e, f\}, \\ \{b, e, f\} \}$$

Let $T_1 = \langle a, b, c \rangle$ and $T_2 = \langle a, b, e \rangle$

$$\text{Here } T_1 \cup T_2 = a, b, c, e$$

$$T_1 \cap T_2 = a, b.$$

$$\text{sim}(T_1, T_2) = \frac{a, b}{a, b, c, e} = \frac{2}{4} = \frac{1}{2}$$

AID-nummer: AID-number:	2203	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

 Blad nummer:
 Sheet number:
 6

In the example $\{a, b, c\}, \{a, b, d\}, \{a, b, e\}$
 $\{a, b, e\}, \{a, b, f\}$ are

neighbors as they have common a, b in them.

The points/objects are neighbors iff

$$\text{sim}(T_1, T_2) > t$$

$$\text{sim}(T_1, T_2) = \frac{T_1 \cap T_2}{T_1 \cup T_2}$$

Links between the objects is the number
 common neighbors-

$$\text{Goodness Measure (G)} = \frac{\text{Number of links b/w clusters}}{\text{Number of expected links b/w the clusters}}$$

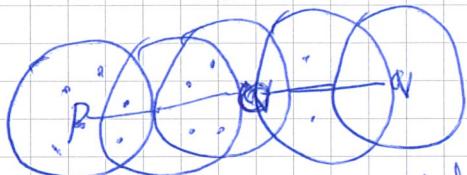
ROCK uses goodness Measure to combine clusters
 the ~~most~~ items with more link are in the same cluster
 If the two elements are neighbors there common
 neighbors are the elements ~~not~~ which are similar to
 the nodes -

3, ✓

AID-nummer: AID-number:	2203	Datum: Date:	18-16-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	F

Q.4

If p is density connected to q w.r.t Eps & Minpts
then



If p is density ~~reachable~~^{connected} from q w.r.t. Eps & Minpts then it is that is p is directly density connected to o and q is directly density connected to point o . Both p and q are directly density connected to point o so they are density connected to each other. As p and q are density connected to each other so they are not density-reachable but they are directly density connected?

For density reachable the points p and q should be density connected with the point o .

There should be chain of points connecting p & q w.r.t. Eps and Minpts.

Q-4-b OPTICS

Optics is used for finding the value of the epsilon - It is used to check how the epsilon will effect the cluster - It uses ordering of the objects - Optics selected epsilon

AID-nummer: <i>AID-number:</i>	2203	Datum: <i>Date:</i>	18 - 06 - 05	Blad nummer: <i>Sheet number:</i>
Kurskod: <i>Course code:</i>	732A75	Provkod: <i>Exam code:</i>	TENT	8

Value is used in the DBSCAN algorithm-

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	9

Q-5		A	B	C	D	E	F	G
Item-K	gold	(0,0)	Y	N	Y	N	silver	
	bronze	(1,1)	N	N	N	N		
	distance	1	1+1=2	1	0	1	0	0
	delta	1	1	1	1	1	0	$M_f = \{1, \dots, M_f\}$

A and G are ordinal variables

Taking Y = 1
and N = 0

$\{gold, silver, bronze\}$ $m=3$

↓ ↓ ↓
1 2 3

$$Gold \Rightarrow \frac{M_f - 1}{M_f - 1} = \frac{1 - 1}{3 - 1} = \frac{0}{2} = 0$$

$$Silver \Rightarrow \frac{2 - 1}{3 - 1} = \frac{1}{2}$$

$$Bronze \Rightarrow \frac{3 - 1}{3 - 1} = \frac{2}{2} = 1$$

E and F asymmetric
 $\delta_{ij} = 0$ if $x=0$ and
 $f=\text{asymmetric}$

or $x = \text{Missing}$

$$\text{So } d(i,j) = \frac{\sum_{i=1}^n s_{ij}}{\sum_{i=1}^n s_i}$$

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

2

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	10

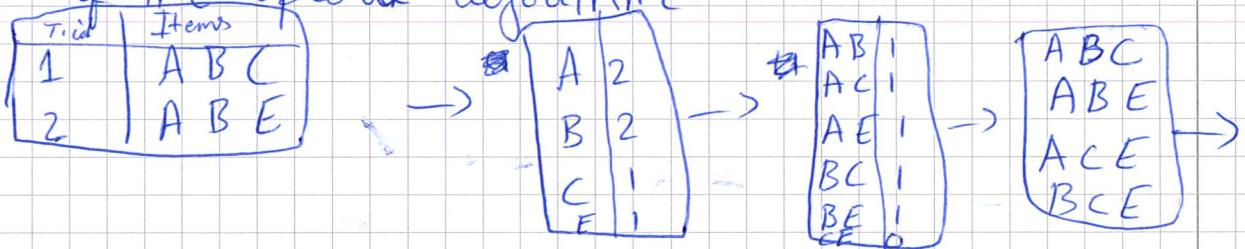
Q-6

(a) Apriori Property:

~~If a subset is infrequent all of its supersets will also be infrequent -~~

It states If a superset is frequent all of its not empty subsets will be frequent -

(b) We produce candidates in the aprori algorithm- by combining the subset with each other ^{now?} after applying descending ordering to the itemsets - We generate candidates at each step or iteration of the aprori algorithm -



(c) At the first step we check if the items satisfy the constraint we prune the items and do not take the item further - If the item do not satisfy the antimonotone constraint we take the item further - Because if an item do not satisfy the constraint none of its superset will satisfy the constraint in antimonotone constraint so we prune the item - and check for the constraint again at the next step -

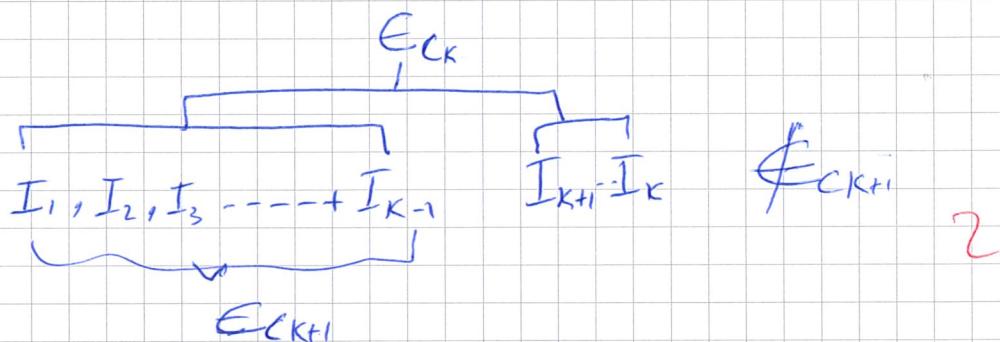
AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

(d) Correctness of Apriori Algorithm =

- ① Let $L_k \subseteq C_k$ ($k=0$) trivial case -
- ② Induction Hypothesis $L_{k+1} \subseteq C_{k+1}$, such that $1 \dots k$
- ③ To prove $L_{k+1} \subseteq C_{k+1}$ -

Assume:

$$L_{k+1} \notin C_{k+1}$$



As self joining is producing $I_1 \dots I_2$ we are not able to prune the itemset so our assumption $L_{k+1} \notin C_{k+1}$ is denied -

Hence

$$L_{k+1} \in C_{k+1} \text{ is True}$$

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

Blad nummer:
Sheet number:
12

Q-7

FP grow Algorithm -

checking the frequency support of the items

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

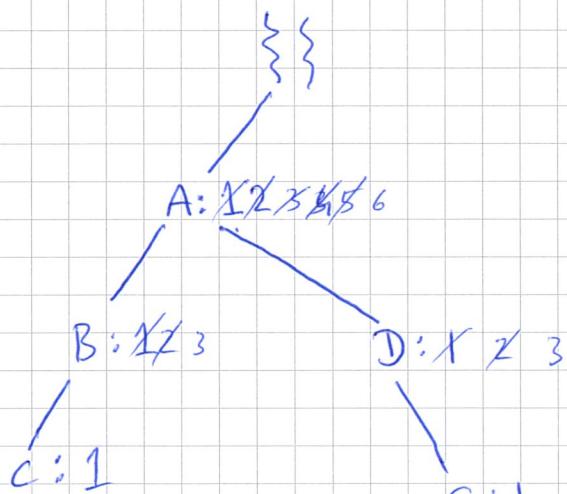
Items	Support
A	6
B	3
C	2
D	3

All of the items fulfill the minsupport requirement > 1 .

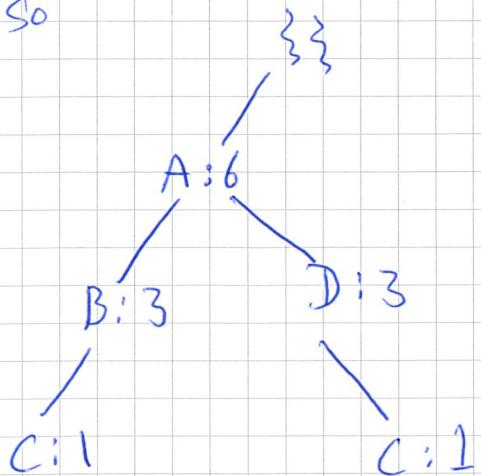
Arranging the items in the transaction database according to frequency -

$$A \rightarrow B \rightarrow D \rightarrow C$$

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



So



AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

Blad nummer:
Sheet number:
13

⇒ Now we will check the
 α =conditional databases -

B-conditional database

Items : A

Support : 1

Cond DB : AB: 3



AB-conditional database

Items : Null

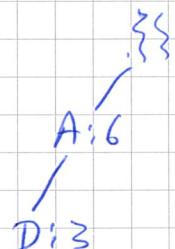
∅

D-conditional database

Items : A

Support : 1

Cond DB : AD: 3



AD-conditional database

∅

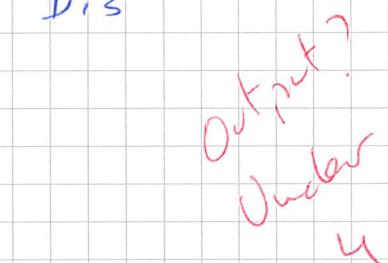
C-conditional database

Items : A, B, D

Support : 2, 1, 1

Cond DB : AB: 1, AD: 1

Set : ABC, ADC



ABC-conditional database

Items : ABC

Support : 1, 1, 1

Cond DB : ∅

ADC-conditional database

Items : ADC

Support : 1, 1, 1

Cond DB : ∅

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

Blad nummer:
Sheet number:
14

A-conditional database =

Items : { }

support ; { }

φ

{ }

/
A:6

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT

Blad nummer:
Sheet number:

15

Q-8 a)

Monotone

$$\text{sum}(s) \geq v \quad \text{no - except first}$$

$$\text{range}(s) \geq v$$

$$\min(s) \leq v$$

Antimonotone

$$\text{sum}(s) \leq v$$

$$\text{range}(s) \leq v$$

$$\min(s) \geq v$$

Convertible Monotone But Not Monotone

$\text{Avg}(s) \geq v$ when items are in ~~increasing~~ order -

Convertible Antimonotone But Not Antimonotone

$\text{Avg}(s) \geq v$ when items are in ~~decreasing~~ order -

2

AID-nummer: AID-number:	2203	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	16

(b) Frequent Itemset ABC

ABC

$$\cancel{ABC} \rightarrow AB \rightarrow C \quad \frac{\text{sup}(ABC)}{\text{sup}(AB)} = \frac{1}{3} = 0.33 = 33\%$$

$$AC \rightarrow B \quad \frac{\text{sup}(ABC)}{\text{sup}(AB)} = \frac{1}{2} = 50\%$$

$$BC \rightarrow A \quad \frac{\text{sup}(ABC)}{\text{sup}(BC)} = \frac{1}{1} = 100\%$$

As $AC \rightarrow B$ and $BC \rightarrow A$ support are 50% and 100% respectively so we will not use them to check for others -

$$C \rightarrow AB \quad \frac{\text{sup}(ABC)}{\text{sup}(C)} = \frac{1}{2} = 0.5 = 50\%$$

So the frequent Item ABC generate 3 association rules that satisfy condition of confidence greater or equal than 50%.

$AC \rightarrow B$

$BC \rightarrow A$

$C \rightarrow AB$

(c) Association rule is a causal rule when the rule using items are independent on each other

Diper \rightarrow Milk

Milk \rightarrow Diper