

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer: Sheet number:
1

## Q. 1 (a) - i

PAM Algorithm :

PAM algorithm can be described as:

1. Select  $k$  arbitrary objects as medoid
2. For each pair of medoid  $i$  and non-medoid object  $h$ , calculate the swapping cost.
3. Select the object with the lowest swapping cost
  - If  $TC_{ih} < 0$   
then swap (update) the current medoid  $i$  with  $h$ .
  - Assign each non-selected objects w.r.t the new medoid.
4. Repeat steps until there is no change.

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN 1

Blad nummer:  
Sheet number:  
2

### Q. 1 (a) - ii

The total cost of swapping an object is called swapping cost. It can also be defined as sum of cost-

i.e

$$TC_{ih} = \sum_{j=1}^n C_{jih}$$

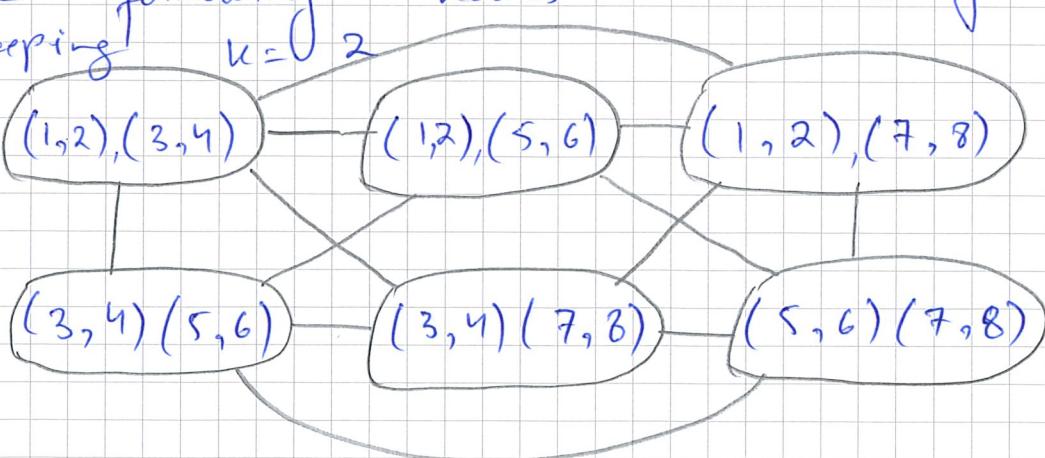
### Q. 1 (a) - iii

Consider a dataset in 2 dimension

i.e

$$\{(1,2), (3,4), (5,6), (7,8)\}$$

The following nodes will be formed keeping  $k=2$



AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
3

Now take initial k medoids  $(3,4), (5,6)$   
 $(3,4), (5,6)$

cluster 1 :  $(3,4) \rightarrow (3,4), (1,2)$

cluster 2 :  $(5,6) \rightarrow (5,6), (7,8)$

Considering Manhattan dist (for reducing complexity of calculation)

Cost of it will be

$$TC_{(3,4),(5,6)} = 0 + 4 + 0 + 4 = 8$$

↳ old cost

Now take another <sup>neighbouring</sup> point to compare its swapping cost with our initially selected medoids

$(1,2)(5,6)$

cluster 1 :  $(1,2) \rightarrow (1,2), (3,4)$

cluster 2 :  $(5,6) \rightarrow (5,6), (7,8)$

$$\textcircled{O} \quad TC_{(1,2)(5,6)} = 0 + 4 + 0 + 4 = 8$$

↳ New cost

$$\begin{aligned} \text{Swapping Cost} &= \text{New cost} - \text{old cost} \\ &= 8 - 8 \\ &= 0 \end{aligned}$$

Thus the swapping is 0 as there is no change in cost

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:

4

## 2<sup>nd</sup> Condition

Now consider  $(1, 2), (3, 4)$  as initial medoids

cluster 1 :  $(1, 2) \rightarrow (1, 2)$

cluster 2 :  $(3, 4) \rightarrow (3, 4), (5, 6), (7, 8)$

$$TC_{(1,2),(3,4)} = 0 + 0 + 4 + 8 = 12$$

$\hookdownarrow$  old-cost

for its neighbour

$(3, 4), (5, 6)$

cluster 1 :  $(3, 4) \rightarrow (1, 2), (3, 4)$

cluster 2 :  $(5, 6) \rightarrow (5, 6), (7, 8)$

$$TC_{(3,4),(5,6)} = 4 + 0 + 0 + 4 = 8$$

$\hookdownarrow$  new-cost

Then,

$$\Rightarrow \text{new-cost} - \text{old-cost}$$

$$= 8 - 12$$

$$= -4$$

Since new medoid is better than the previously chosen, thus the result is strictly negative.

AID-nummer: <i>AID-number:</i>	2198	Datum: <i>Date:</i>	18 - 06 - 05
Kurskod: <i>Course code:</i>	732A75	Provkod: <i>Exam code:</i>	TEN 1

Blad nummer:  
*Sheet number:*

5

Q. 1 (b)

PAM and CLARANS

o

Q. 1 (c)

none

)

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:

6

## Q. 2

### Agglomerative Hierarchical Clustering:

Agglomerative hierarchical clustering is a hierarchical clustering approach in which a data set is divided into smaller clusters.

Agglomerative hierarchical clustering could be complete link or single link.

AGNES is an example of agglomerative hierarchical clustering.

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

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer: Sheet number:
7

Step 1 : we start with the smallest value then we combine the nodes

Step 2 : determine the distance of combined nodes / points to every other point (for complete link we take max dist)

$$\text{dist}((1,4), 2) = \max(d(1,2), d(4,2)) = \max(8, 7) \\ = 8$$

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

$$\text{dist}((1,4), 5) = \max(d(1,5), d(4,5)) = \max(10, 5) \\ = 10$$

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

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN 1

Blad nummer:  
Sheet number:  
8

Now we repeat step 1 and 2 until all nodes are ~~compe~~ combined (or threshold 'n' reached if defined)

$$d((2,5), (1,4)) = \max(d(2, (1,4)), d(5, (1,4))) = \max(8, 10) \\ = 10$$

$$d((2,5), (3)) = \max(d(2, 3), d(5, 3)) = \max(4, 6) \\ = 6$$

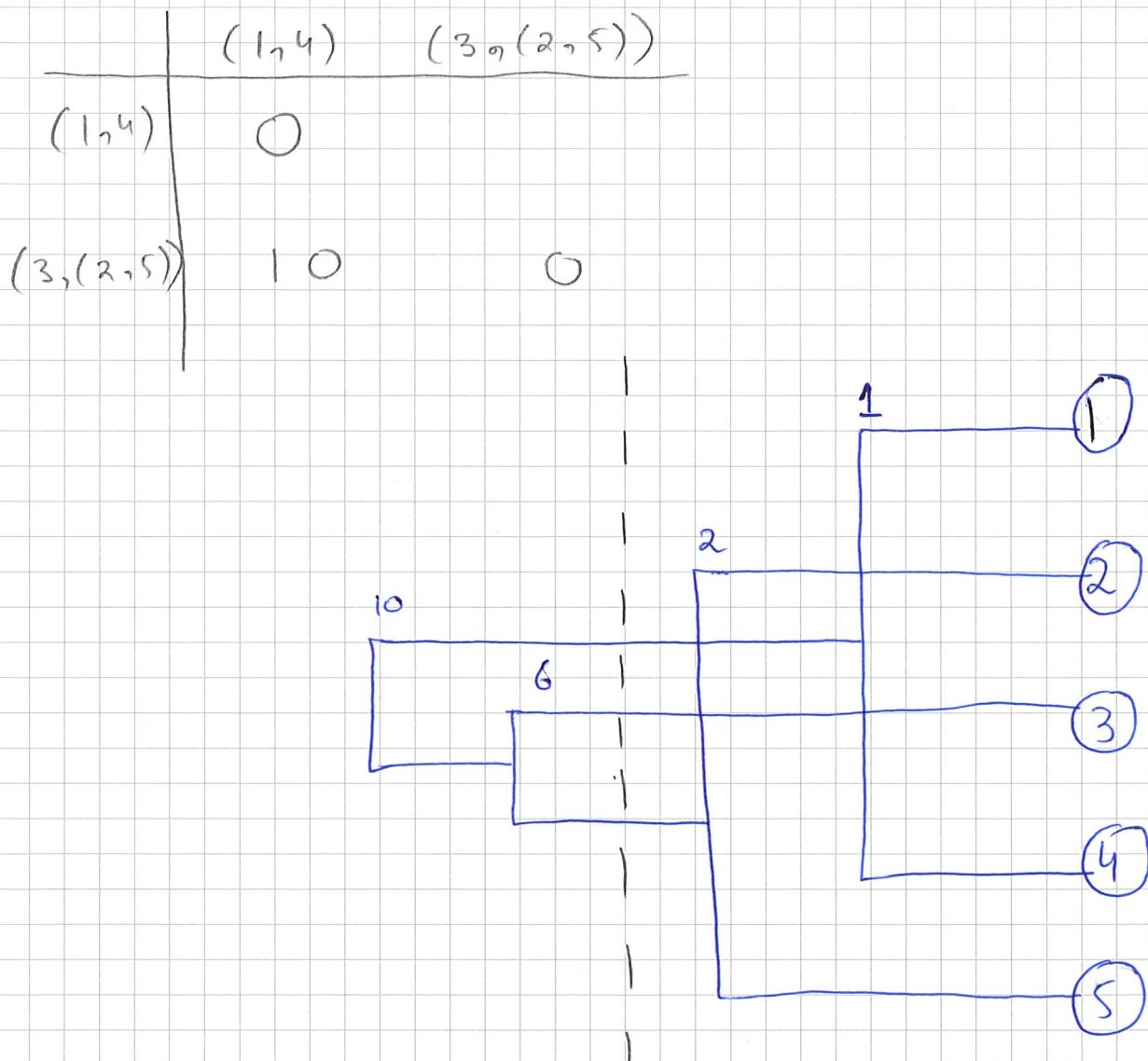
	(1,4)	(2,5)	3
(1,4)	0		
(2,5)	10	0	
3	9	6	0

Combine 3 and (2,5)

$$d((3, (2,5)), (1,4)) = \max(d(3, (1,4)), d((2,5), (1,4))) \\ = \max(9, 10) = 10$$

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
 Sheet number:  
 9



Threshold = 4  $\Rightarrow$  if defined

(Then process would have stopped here)

3

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer: Sheet number:
10

### Q. 3

#### ROCK Algorithm:

ROCK is a hierarchical clustering approach in which clustering is done by checking similarity/proximity between the points.

- It maximizes the similarity - between points i.e. number of links within a cluster.
- Minimizes the link between points in different cluster.

#### Algorithm:

- Take an arbitrary point  $p$ , after scanning finds neighbour of the selected point and determine the links.
- label data on disk is marks data having the highest goodness of measure (or above threshold)

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1	11

Neighbour :

An object is said to be neighbour if it has at most difference of 1 item i.e

$$\text{sim} \geq t \quad (\text{where } t = \text{threshold})$$

it can be explained by the following example :

abc      abd      abe      acd      ace

ade      bcd      bce      cde

abf      abg      afg      bfg

Consider an object abc

its similarity with itsdp can be defined by

$$d(\text{abc}, \text{abc}) = \frac{a}{a+b+c} = \frac{3}{3} = 1$$

a	b	c
1	1	1
1	1	1

Now consider another object abd  
the similarity b/w abc and abd is

$$d(\text{abc}, \text{abd}) = \frac{a}{a+b+c} = \frac{2}{2+1+1} = \frac{2}{4} = 0.5$$

a	b	c	d
1	1	1	0
1	1	1	0

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
12

Thus in the given scenario, an object will be neighbour if

$$\text{sim} \geq 0.5$$

Common Neighbour :

Two objects is said to have common neighbours if neighbour of "A" is also a neighbour of "B" i.e the neighbour must have at most 1 item different from node "A" and node "B".

For example :

If we consider the same data set as in previous example,

abd is a neighbour of abc

If we take another node bcd and determine its similarity with abd then

$$d(bcd, abd) = \frac{2}{2+1+1} = 0.5$$

a	b	c	d
0	1	1	1
1	1	1	0

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
13

which implies that abd is a common neighbour of abc and bcd

Link for Objects:

Link for objects can be defined as number of common neighbours between point  $P_i$  and  $P_j$ .

~~The above example~~

For example:

<u>abc</u>	<u>abd</u>	<u>abe</u>	<u>acd</u>	<u>ace</u>
<u>ade</u>	<u>bcd</u>	<u>bce</u>	<u>cde</u>	bde
<u>abf</u>	<u>abg</u>	aff	bfg	

in above example

abc has the following neighbours  
 abc, abd, abe, acd, ace  
 bcd, bce, abf, abg

and

cde has following neighbours  
 acd, ace, ade, bcd, bce, cde

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer: Sheet number:
14

Here it can be seen that  
 acd, ace, bcd and bce  
 are ~~not~~ common neighbours of abc &  
 cde so link between abc and cde  
 will be given as  
 $\text{link}(\text{abc}, \text{cde}) = 4$

Link for clusters :

Link for clusters is defined as  
 number of cross links between  
 points in  $C_i$  and points in  $C_j$   
 where  $C_i$  and  $C_j$  are two clusters.

$$\text{Link}(C_i, C_j) = \sum_{\substack{P_i \in C_i \\ P_j \in C_j}} \text{link}(P_i, P_j)$$

Goodness Measure :

Goodness measure can be defined  
 as ratio of cross links between  
 points in cluster  $C_i$  with points in  
 cluster  $C_j$  to the expected number  
 of ~~cross~~ cross links b/w ~~the~~ points  
 in cluster  $C_i$  to points in  $C_j$ .

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
15

$$G_1 = \frac{\text{Link}(C_i, C_j)}{\text{Expected # of cross links between } C_i \text{ & } C_j}$$

Expected number of Cross Links Between  
 $C_i$  &  $C_j$  is given by :

- ↳ expected number of links in  $C_i \cup C_j$
- ↳ expected number of links in  $C_i$
- ↳ expected number of links in  $C_j$

$$G_1 = \frac{\text{Link}(C_i, C_j)}{(n_i + n_j)^{1+2f(t)} - n_i^{1+2f(t)} - n_j^{1+2f(t)}}$$

where

expected # of links in  $C_i \cup C_j$  is given by  
 $(n_i + n_j)^{1+2f(t)}$

expected # of links in  $C_i = n_i^{1+2f(t)}$

expected # of links in  $C_j = n_j^{1+2f(t)}$

✓  
4

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN 1	16

Q. 4 (a)

The answer is both YES and NO.

PROVE :

If point 'p' is density connected to  $q$  wrt Eps and Minpts, then it implies that there exist a point 'o' from which 'p' and 'q' are density reachable.

If  $q = o$ , then p will become density reachable from 'q'

Counter Argument:

It might be possible that ' $q$ ' and ' $p$ ' are density reachable from ' $o$ ' but ' $q$ ' is not a core object i.e ' $q$ ' does not have minpts in its ~~Eps~~ neighborhood. In this case  $p$  is not density reachable from ' $q$ ' i.e  $q \neq o$  in any possible way (since it's not a core object).

General core

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
17

## Q. 4 (b)

OPTICS provides ordering <sup>to</sup> clustering structure.  
 OPTICS uses the principles of DBSCAN  
 but ~~as~~ with some extension.

OPTICS uses two extra parameters (as compared to DBSCAN)

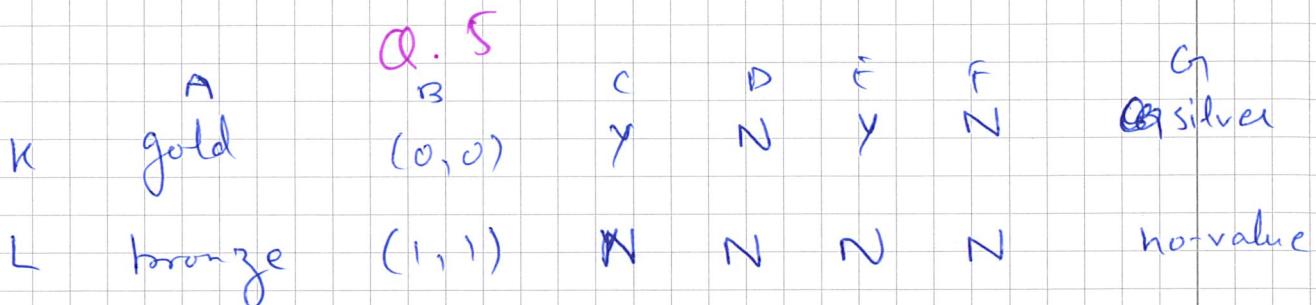
1) Core distance : The minimum ~~as~~  $\epsilon$  value to make an object a core point  
 If not found then undefined.

2) Reachability distance :  
 $\max(\text{core distance}(o), d(o,p))$  i.e  
 maximum of core distance for a point 'o' and the distance of object 'o' to 'p'.

These parameters make optics perform better and consider area of high density while clustering.

OPTICS is often used for determining the  $\epsilon$  parameter for DBSCAN (or other density based algorithm).

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1	18



for  $\text{dist}(K, L)$  for variable A  
ranks

$$\text{gold} = 3$$

$$\text{silver} = 2$$

$$\text{bronze} = 1$$

$$z_f = \frac{r_f - 1}{m_f - 1} \quad \text{where } m_f = \max \text{ rank}$$

$$z_f(\text{gold}) = \frac{3 - 1}{3 - 1} = \frac{2}{2}$$

$$z_f(\text{gold}) = 1$$

$$z_f(\text{silver}) = \frac{2 - 1}{3 - 1} = \frac{1}{2} = 0.5$$

$$z_f(\text{bronze}) = \frac{1 - 1}{3 - 1} = 0$$

$$\text{dist}(\text{gold}, \text{bronze}) = d(1, 0)$$

considering Euclidean distance

$$= \sqrt{(1 - 0)^2} = \sqrt{1^2}$$

$$d(\text{gold}, \text{bronze}) = 1$$

AID-nummer: AID-number:	2198	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN 1

Blad nummer: Sheet number:	19
-------------------------------	----

dist for attribute B :

$$= (10-11) + (10-11)$$

$$= |-1| + |-1| = 1+1$$

$$d(K, L) = 2$$

For binary symmetric variables if value is same (matching),  $d(K, L) = 0$   
else 1

For binary assymetric 0,0 case,  $\delta_{0,0} = 0$   
else 1

so

A      B      C      D      E      F      G

$$d(K, L) = \frac{\sum_{i=1}^N S_i(s) d_i(f)}{\sum_{i=1}^N S_i(s)}$$

$$= \frac{1 \times 1 + 2 \times 1 + 1 \times 1 + 0 \times 1 + 1 \times 1 + 0 \times 0 + 0}{1 + 1 + 1 + 1 + 1}$$

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

$$= \frac{5}{5}$$

$$d(K, L) = 1$$

?

AID-nummer: AID-number:	2198	Datum: Date:	18 - 06 - 05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
20  
89

### Q. 6 (a)

According to apriori property,  
the infrequent itemsets are pruned  
in pruning step.

o

### Q. 6 (b)

Candidates are generated by  
performing self-join operation on  
itemset produced in k step.

Data ~~base~~ is first scanned, then  
candidates are generated using self-join  
That is if  $I_k$  and  $J$  represents  
itemset then

join is applied such that

$$I_1, \dots, I_{k-1} \supseteq J_{k-1}$$

join  
(  
subset  
checking?)

In each k-step, the candidate size  
is increased by one if constraints are  
satisfied.