

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:

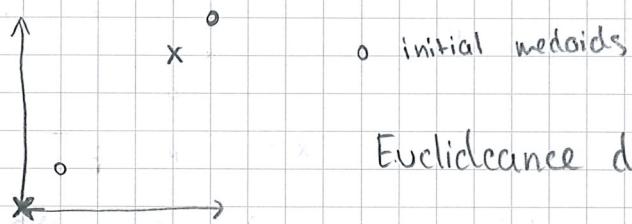
1

1) a) i Pam-algorithm

1. Define  $k$  arbitrary points as your initial medoids.
2. Assign each non-medoid point to the closest medoid to form  $k$  clusters.
3. Calculate the total cost for the current medoid.
4. Calculate the total cost for other possible medoids.
5. Calculate the swapping cost  $\text{New cost} - \text{Old cost}$ .
6. If at least one swapping cost is lower than 0 we switch to the medoids that corresponds to the lowest swapping cost.
7. repeat steps 3-6 until nothing changes

ii Swapping cost =  $T\mathcal{C}_{\text{new}} - T\mathcal{C}_{\text{old}}$

iii Example of when swapping cost is 0:



Euclidean distance is used

$$\begin{array}{ll}
 (5,5) & (1,1) \\
 (0,0) & \sqrt{25} \quad \sqrt{2} \\
 (4,4) & \sqrt{2} \quad \sqrt{9} \\
 \hline
 T\mathcal{C} = \sqrt{2} + \sqrt{2}
 \end{array}$$

our initial medoids are (5,5) and (1,1)

calculate the total cost

other possible medoids: (4,4) and (1,1)

$$\begin{array}{ll}
 (4,4) & (1,1) \\
 (0,0) & \sqrt{16} \quad \sqrt{2} \\
 (5,5) & \sqrt{2} \quad \sqrt{16} \\
 \hline
 T\mathcal{C} = \sqrt{2} + \sqrt{2}
 \end{array}$$

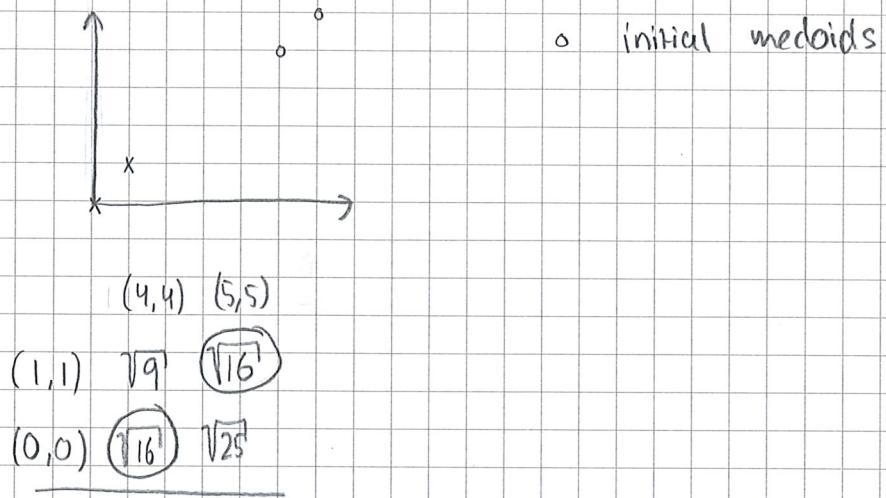
calculate the total cost for (4,4) and (1,1)

The swapping cost is the new cost - old cost

$$(\sqrt{2} + \sqrt{2}) - (\sqrt{2} - \sqrt{2}) = 0$$

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Example of when swapping cost is strictly negative:



$$TC = 8$$

Here we check all neighbours to see if they all result in negative TC

(4,4) (1,1)	(4,4) (0,0)	(5,5) (1,1)	(5,5) (0,0)
(0,0) $\sqrt{16}$ $\sqrt{2}$	$\sqrt{16}$ 0	$\sqrt{25}$ $\sqrt{2}$	$\sqrt{25}$ 0
(1,1) $\sqrt{9}$ 0	$\sqrt{9}$ $\sqrt{2}$	$\sqrt{16}$ 0	$\sqrt{16}$ $\sqrt{2}$
(4,4) 0 $\sqrt{9}$	0 $\sqrt{16}$	$\sqrt{2}$ $\sqrt{9}$	$\sqrt{2}$ $\sqrt{16}$
(5,5) $\sqrt{2}$ $\sqrt{16}$	$\sqrt{2}$ $\sqrt{25}$	0 $\sqrt{16}$	0 $\sqrt{25}$

$TC = \sqrt{2} + \sqrt{2}$			
----------------------------	----------------------------	----------------------------	----------------------------

all of the total cost for the neighbours are smaller than 8.

This will lead to all swapping costs being negative (new cost - old cost)

Therefore in this example the swapping cost is strictly negative.

AID-nummer: <i>AID-number:</i>	2201	Datum: <i>Date:</i>	18-06-05
Kurskod: <i>Course code:</i>	232A75	Provkod: <i>Exam code:</i>	TEN1

Blad nummer:  
*Sheet number:*

3

b)

PAM, CLARANS

0

c)

None

1

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05	Blad nummer: Sheet number:
Kurskod: Course code:	732A75	Provkod: Exam code:	TENT	4

2) Dissimilarity matrix

	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

Agglomerative Clustering

Hierarchical clustering when we start with all the objects in separate clusters and merge clusters that are most similar until all objects are in one cluster.

We use complete-link clustering, which means that the distance between two clusters is the maximum distance between two points in those clusters. Therefore, when we calculate the distance between clusters we use the maximum distance.

First we merge 1 and 4 since they have the shortest distance.

$$\{1,4\} \quad 2 \quad 3 \quad 5$$

$\{1,4\}$	0
2	8
3	9
5	10

$$d_{(1,4),2} = \max(d_{1,2}, d_{4,2}) = \max(8, 7) = 8$$

$$d_{(1,4),3} = \max(d_{1,3}, d_{4,3}) = \max(3, 9) = 9$$

$$d_{(1,4),5} = \max(d_{1,5}, d_{4,5}) = \max(10, 5) = 10$$

threshold

Then we merge 2 and 5 since they have the shortest distance now.

$$\{1,4\} \quad 3 \quad \{2,5\}$$

a. /

$\{1,4\}$	0
3	9

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

$\{2,5\}$	10	6	0
-----------	----	---	---

$$d_{(2,5),1} = \max(d_{2,1}, d_{5,1}) = \max(8, 10) = 10$$

Then we merge 3 and  $\{2,5\}$

$$\{1,4\} \quad \{2,3,5\}$$

$$d_{(1,4),\{2,3,5\}} = \max(d_{(1,4),2}, d_{(1,4),3}, d_{(1,4),5}) = \max(9, 10) = 10$$

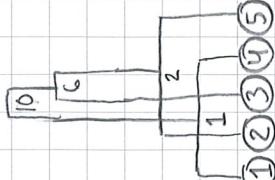
$\{1,4\}$	0
-----------	---

$\{2,3,5\}$	10	0
-------------	----	---

Dendrogram

Then we merge  $\{1,4\}$  and  $\{2,3,5\}$

↓



AID-nummer: AID-number:	Datum: Date:
Kurskod: Course code:	Provkod: Exam code:

Blad nummer:  
Sheet number:

5

3)

### Rock-algorithm

category  
like

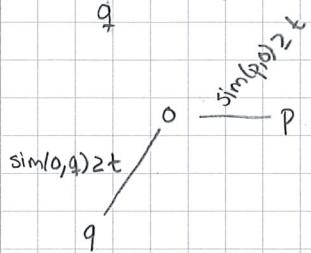
Rock uses a link measure to create clusters.  
We take in a threshold describing the similarity we want for two points to be neighbours

$$P_i \xrightarrow{\text{sim}(p_i, p_j) \geq t} P_j$$

points  $p_i$  and  $p_j$  are neighbours if the sim between the two points are higher or equal to some threshold  $t$ .

points  $p_i$  and  $p_j$  are common neighbours if they have a neighbour  $o$  that is neighbour to both  $p_i$  and  $p_j$

$o$  is the common neighbour



#### link for objects

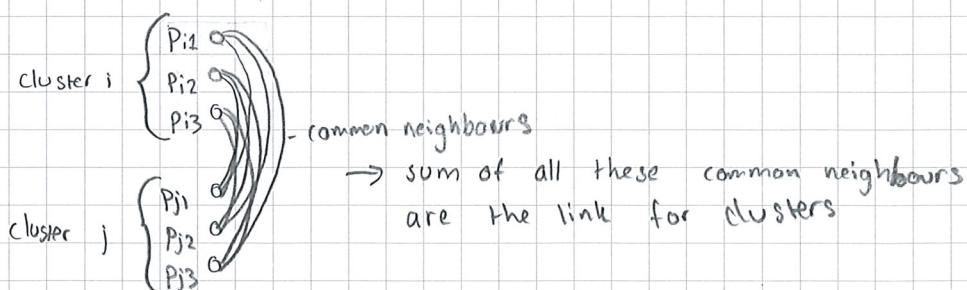
the link measure is the number of common neighbours between two objects. the link between  $p_i$  and  $p_j$  are the number of points that are neighbours to both  $p_i$  and  $p_j$ .

P.

#### link for clusters

$$\sum_{\substack{p_i \in C_i \\ p_j \in C_j}} \text{link}(p_i, p_j)$$

this is the sum of all links where point  $p_i$  is in cluster  $i$  and  $p_j$  is in cluster  $j$ .



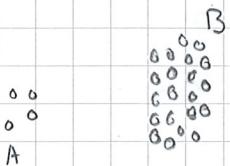
AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEW1

Blad nummer: Sheet number:
6

3)  $G$  = Goodness measure

$$G = \frac{\text{links for clusters}}{\text{Expected number of links for clusters}} = \frac{\sum_{i \in C_i, j \in C_j} (p_i, p_j)}{E \left( \sum_{i \in C_i, j \in C_j} (p_i, p_j) \right)}$$

This is good for example when two clusters consist of different number of points.



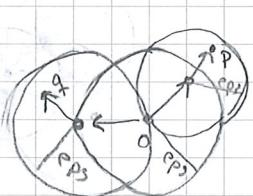
It is quite likely that points in cluster A will find common neighbours in cluster B. therefore the links between those clusters will be more than links between two small clusters of equal sizes. Therefore we use the goodness measure.

3)

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

4) a) DBSCAN

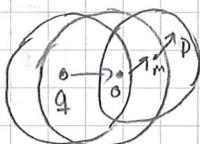
$p$  and  $q$  are density-connected. if there exist a point  $o$  from where both  $p$  and  $q$  are density reachable



Each circle has radius  $\epsilon_{ps}$  and points  $\geq \text{minpts}$ . Each circle has a core point.

here  $p$  is not density-reachable from  $q$  since  $q$  is not a core-point.

a point is density-reachable from a core-point if that core point has a chain of directly-density reachable points to that point.



here  $p$  is density reachable from  $q$ .

In summary,  $p$  does not have to be density-reachable from  $q$ , when they are density-connected because  $q$  doesn't have to be a core point for them to be density connected. However,  $q$  has to be a core point if  $p$  is density-reachable from  $q$ .

b) OPTICS is ran before DBSCAN in order to find a good  $\epsilon$ . This is good since DBScan is sensitive to the  $\epsilon$  parameter.

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:

8

5)

	A	B	C	D	E	F	G
item K	gold	(0,0)	Y	N	Y	N	silver
item L	bronze	(1,1)	N	N	N	N	missing

keep all dummies

Distance A: ordinal variable. Transform to assymmetric binary (like dummies)

item i	item K						
item L	{	gold	0	silver	0	bronze	0
	gold	0	0	silver	0	0	1
	0	1	0	0	1	0	0

$$\text{distance} = \frac{b+c}{a+b+c} = \frac{1}{1} = 1$$

distance =  $\frac{\sum \partial d}{\sum \partial}$  where  $\partial = 0$  if one value is missing or if variable is assymmetric binary and 0, and 1 otherwise

$$\text{distance} = \frac{1+1}{2} = 1$$

Distance B: Manhattan distance  $|x_1 - y_1| + |x_2 - y_2| + \dots + |x_n - y_n|$

$$\text{distance} = 1 + 1 = 2$$

Distance C	binary symmetric	→	item i	item j	distance = $\frac{b+c}{a+b+c+d}$
	item k		a	b	
Yes	No		0	0	
item L	Yes 0	0	1	0	
	No 1	0	0	1	

$$\text{distance} = \frac{0+1}{0+0+1+0} = 1$$

Distance D binary symmetric

item K	Yes	No	distance = $\frac{0+0}{0+0+0+1} = 0$
item L	Yes 0	0	
	No 0	1	

Distance E binary assymetric → item i

item k	Yes	No	item i	j	0	distance = $\frac{b+c}{a+b+c}$
item L	Yes 0	0	a	b		
	No 1	0	0	c	d	

$$\text{distance} = \frac{1+0}{0+0+1} = 1$$

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:

9

Distance F: binary assymetric

item k	
	Yes      No
item l	Yes      0      0
	No      0      1

$$\text{distance} = \frac{0+0}{0+0+0}$$

Distance G: missing value

$$d = \begin{cases} 0 & \text{if missing value or assymetric binary and 0} \\ 1 & \text{otherwise} \end{cases}$$

$$\text{distance between item } k \text{ and item } l = \frac{\sum d}{\text{number of } d}$$

$$= \frac{1+1+1+2+1+1+1+0+1+1+0+0}{1+1+1+1+1+0+0} = \frac{5}{5} = 1$$

The distance between the items is 1.

2

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
Sheet number:  
10

6 a) Apriori property

Subsets of frequent itemsets are also frequent.

b) Candidates  $C_k$  are produced by self-joining  $L_{k-1}$

and prune itemsets that have subsets that are not in  $L_{k-1}$ .

$C_1$  consist of the 1-itemsets in the database

c) An antimonotone constraint is checked for when  $\text{minsup}$  is checked. If an itemset in  $C_k$  does not fulfill the constrained it is not in  $L_k$  or the output. The reason we can prune it altogether is because if an itemset does not fulfill an antimonotone constraint neither will a superset of that itemset. So we don't have to check again.

AID-nummer: AID-number:	Datum: Date:
Kurskod: Course code:	Provkod: Exam code:

2201                    18-06-05

732A75                TEN1

Blad nummer:  
Sheet number:

11

6) d) Apriori algorithm is proved by induction

1.  $L_k \subseteq C_k$  is true for  $k=1$  (seen by 1st line in pseudocode)

2. We assume  $L_{k-1} \subseteq C_{k-1}$  (induction hypothesis)

3. We prove for  $L_k \subseteq C_k$

i. assume contradiction: there exists an itemset  $I$  such that  $I \in L_k, I \notin C_k$  ( $I$  is in  $L_k$ , but not in  $C_k$ )

ii. induction hypothesis  $L_{k-1} \subseteq C_{k-1}$

iii.  $C_k$  is made from self-joining  $L_{k-1}$  (apriori-gen)

since  $I$  is a large itemset  $\in L_k$  it was not pruned from  $C_k$ , therefore  $I \in C_k$

This is a contradiction to our assumption in i, and therefore apriori is correct for

$L_k \subseteq C_k$ .

2

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer:  
 Sheet number:

12

7)

Tid	Items	minsup = 1
1	C B A	
2	D C A	
3	A B	
4	A B	
5	A D	
6	A D	

1) Scan the database to find frequent 1-itemsets

sup

C 2  
 B 3  
 A 6  
 D 3

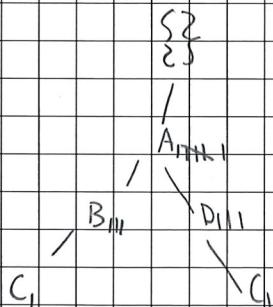
order: A, B, D, C

all have sup ≥ minsup all go in the output

Reorder the transactions in sup descending order

Tid	Items
1	A B C
2	A D C
3	A B
4	A B
5	A D
6	A D

construct Fp-tree



Conditional data bases

A: -  
 B: A:3  
 C: A B:1 | A D:1  
 D: A:3

AID-nummer: AID-number:	2201	Datum: Date:	18-06-05
Kurskod: Course code:	732A75	Provkod: Exam code:	TEN1

Blad nummer: Sheet number:
13

B - conditional database

items A  
support 3  
 $\text{Support} \geq \text{minsup}$  ok!  
output: AB

S<sub>1</sub>  
S<sub>2</sub>  
|  
A<sub>3</sub>

A B - conditional database  
empty, we don't need  
to check further

C - conditional database

items A, B, D  
support 2, 1, 1 dont need to reorder  
 $\text{Support} \geq \text{minsup}$  for all items  
output: AC, BC, DC

S<sub>2</sub>  
S<sub>3</sub>  
|  
A<sub>2</sub>  
B<sub>1</sub> / \ D<sub>1</sub>

AC - conditional database

empty, we don't need to check further

BC - conditional database

items A  
support 1  
 $\text{Support} \geq \text{minsup}$  ok!  
output: ABC

S<sub>1</sub>  
|  
A<sub>1</sub>

ABC - conditional database

empty, we don't need to check further

DC - conditional database

items A  
support 1  
 $\text{Support} \geq \text{minsup}$  ok!  
output: ADC

S<sub>1</sub>  
|  
A<sub>1</sub>

5

ADC - conditional database

empty, we don't need to check further

D - conditional database

items A  
support 3  
 $\text{Support} \geq \text{minsup}$  ok!  
output: AD

S<sub>1</sub>  
|  
A<sub>3</sub>

Output: {A, B<sub>1</sub>, C<sub>1</sub>, D<sub>1</sub>, A B<sub>1</sub>, A C<sub>1</sub>, B<sub>1</sub>, D<sub>1</sub>, A B<sub>1</sub> C<sub>1</sub>, A B<sub>1</sub> D<sub>1</sub>, A C<sub>1</sub> D<sub>1</sub>, A D<sub>1</sub>}

AD - conditional database

empty, we don't need to check further

AID-nummer: AID-number:	Datum: Date:	Blad nummer: Sheet number:
AID-nummer: AID-number:	Datum: Date:	Blad nummer: Sheet number:

14

8) a) Monotone constraint:  
 $\text{sum(prices)} \geq 5$  *valuewise nondecreasing*

Antimonotone constraint:

$$\text{sum(prices)} \leq 5$$

Convertible monotone, but not monotone, constraint:

$$\text{avg(prices)} \geq 5 \quad \text{with respect to decreasing price order}$$

Convertible antimonotone, but not antimonotone, constraint:

$$\text{avg(prices)} \geq 5 \quad \text{with respect to increasing price order}$$

b) ABC rules, confidence  $\geq 50\%$

$$c(A \rightarrow B) = \frac{\text{sup}(AB)}{\text{sup}(A)} = \frac{1}{3} \approx 33\% \quad \text{not a rule, we don't have to check further}$$

$$c(A \rightarrow C) = \frac{\text{sup}(AC)}{\text{sup}(A)} = \frac{1}{1} = 100\% \quad \text{RULE}$$

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

$$\rightarrow c(A \rightarrow CB) = \frac{\text{sup}(ABC)}{\text{sup}(A)} = \frac{1}{6} \quad \text{not a rule}$$

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

$$\rightarrow c(B \rightarrow CA) = \frac{\text{sup}(ABC)}{\text{sup}(B)} = \frac{1}{3} \approx 33\% \quad \text{not a rule}$$

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

RULES:  $AC \rightarrow B$ ,  $BC \rightarrow A$ ,  $C \rightarrow AB$

AID-nummer: <i>AID-number:</i>	2201	Datum: <i>Date:</i>	18-06-05
Kurskod: <i>Course code:</i>	732A75	Provkod: <i>Exam code:</i>	TEW1

Blad nummer:  
*Sheet number:*

15

c) given that  $X$  and  $Y$  are dependent,

$X \rightarrow Y$     $Y \rightarrow X$    is a causal rule

0