

Q1

(a) Yes.

Apriori algorithm for items. the table is like this:

TID	A	B	C	D	E
t ₁	1	0	0	1	0
t ₂	1	0	0	1	1
t ₃	0	0	1	0	0
t ₄	1	0	1	1	1
t ₅	1	0	1	0	1

Apriori algorithm for itemsets

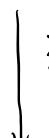
(S₁, S₂, S₃, S₄) means itemsets

SequenceID	S ₁	S ₂	S ₃	S ₄	S ₅
SA	1	0	0	1	0
SB	1	0	0	1	1
Sc	0	0	1	0	0
SD	1	0	1	1	1
SE	1	0	1	0	1

transform item
into itemsets

L ₁ :	itemsets	Support	C ₂ :	Start from here!
	{S ₁ }	4		↓
	{S ₂ }	0	→	
	{S ₃ }	3		
	{S ₄ }	3		
	{S ₅ }	3		

L ₂ :	Sequence	Support	C ₃ :
	{S ₁ , S ₃ }	2	
	{S ₁ , S ₄ }	3	→ {S ₁ , S ₃ , S ₄ }
	{S ₁ , S ₅ }	3	
	{S ₃ , S ₄ }	1	
	{S ₃ , S ₅ }	2	
	{S ₄ , S ₅ }	2	



the result is {{S₁, S₃, S₅}, {S₁, S₄, S₅}}

S₁, S₂, S₃ ... are the smallest units in
this kind of Apriori algorithm.

The Apriori Algorithm is like this:

- ① $k=2$
- ② look for all 2-sequences (which support ≥ 2) and store them in L_2 ;
- ③ Loop
- ④ $k=k+1$;
- ⑤ generate candidate k -sequences from L_{k-1} and store them in C_k
- ⑥ scan the database once and count the support of each candidate in C_k
- ⑦ look for the k -sequences in C_k with support ≥ 2 and store them in L_k
- ⑧ if ($L_k == \text{null}$) break;
- ⑨ Loop
- ⑩ for ($i=2, i \leq k; k++$)
- ⑪ return L_i

And here is the join step:

- ① if there are sequence S' and S'' , $S' = \{s_1, s_2, s_3\}$, $S'' = \{s_2, s_3, s_4\}$, $S'' \text{join } S' = \{s_1, s_2, s_3, s_4\}$
- ② if there are sequence S' and S'' , $S' = \{s_1, s_2, s_3\}$, $S'' = \{s_2, s_3, s_4\}$, $S'' \text{join } S' = \{s_1, s_2, s_3, s_4\}$
for S' and S'' , if $|S'| - \text{first item}| = |S''| - \text{last item}|$ (for example $s_2 = s_2$)

then S'' can join S' $\Rightarrow S''' = S' + \text{the last item of } S''$

The prune step:

Check [the first item of S'] + [the last item of S''] is large

if the answer is large, then store S''' into C_1

Here are the situations

- ① if s_1, s_2, s_3, \dots are all $|s_i| > k$ ($k=2, 3, \dots$) ($i=1, 2, \dots, m$)
then the number of k -sequence with support ≥ 2 is 0
- ② we can count every items/sequence from L_i to answer
if $\sum_{k=1}^m |s_i| = k$ (for example: $\{s_4, s_5\}$ $|s_4| + |s_5| = k$, $\{s_4, s_5\} = k$ -sequence)
then $k_{\max} = |s_1| + |s_2| + |s_3| + |s_4| + |s_5|$
 $k_{\min} = \min(|s_1|, |s_2|, |s_3|, |s_4|, |s_5|)$

for any k , $k_{\min} \leq k \leq k_{\max}$, we can use Apriori algorithm for itemsets to count the number of k -sequences easily.

(3) if $k > k_{\max}$, the number of k -sequence is 0

(b) NO, Apriori Algorithm can't adapting to this case

According to the lecture: Apriori Algorithm has property 1:

If an itemset S is large, then any proper subset of S must be large.

In this case, the property 1 can describe as this: If a k -sequence is large, then any proper sub-sequence of S must be large. However, according to the definition of support of one sequence, the sub-sequences are possibly not large no matter the subsequence is large or not.

For example:

$S = \{A, B, C\}$ Support for $\{A, B\}, \{B, C\}, \{A, C\} = 2$

$S = \{A, B, C\}$ however support for $\{A\} = 1$

Q2

(a) (i) $\{C, D\}$

(ii) $\{C, D\}, \{A, C\}, \{A, D\}$

(iii) $\{A, C, D\}, \{B, E, F\}$

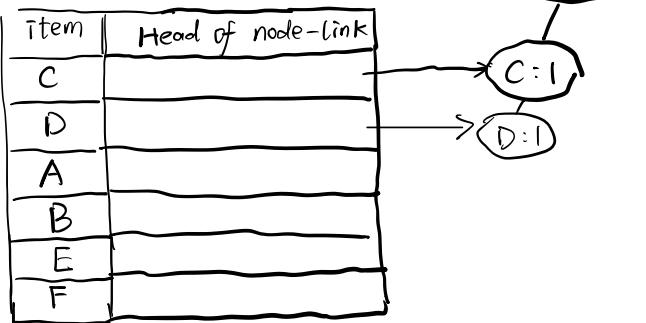
(b) Yes, it can be adapted in this case.

① counting: firstly, the threshold = 1.

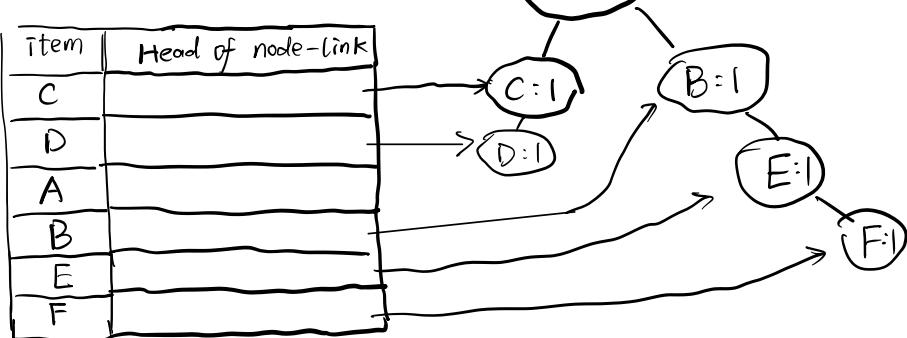
Item	Frequency	Item	Frequency	(ordered) frequent items
A	2	C	3	C, D
B	1	D	3	B, E, F
C	3	A	2	C, D, A
D	3	B	1	
E	1	E	1	
F	1	F	1	C, D, A

→ sort → →

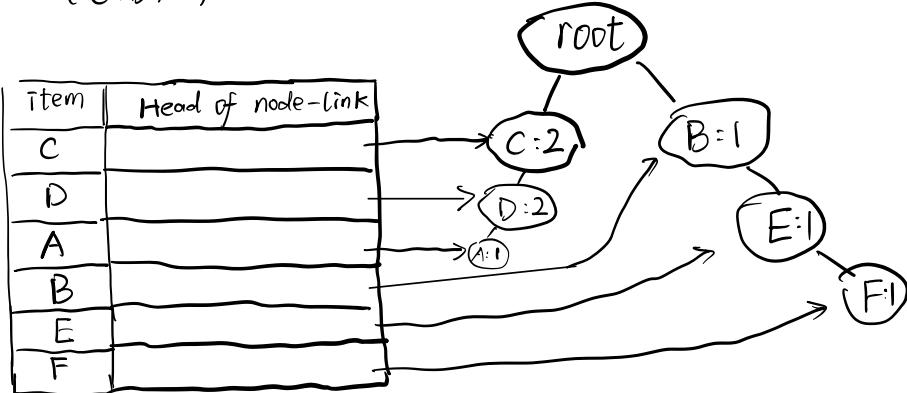
② FP - Tree (C, D)



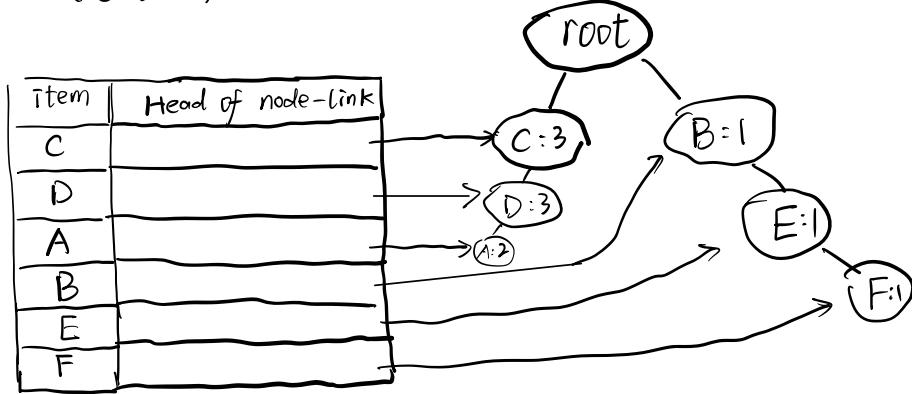
(B, E, F)



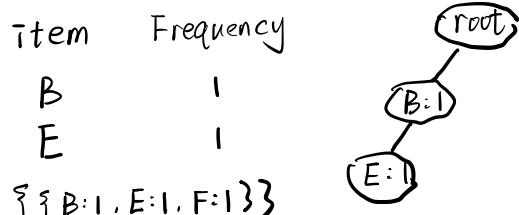
(C, D, A)



(C, D, A)



Conditional FP-tree on "F" ①



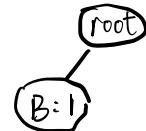
$$S_{1,2} = \{ F:1 \}$$

$$S_{2,2} = \{ BF:1, EF:1 \}$$

$$S_{3,2} = \{ BEF:1 \}$$

Condition FP-tree on "E" ②

$$\{ \{ B:1, E:1 \} \}$$



$$S_{1,2} = \{ E:1, F:1 \}$$

$$S_{2,2} = \{ BF:1, EF:1, BE:1 \}$$

$$S_{3,2} = \{ BEF:1 \}$$

Condition FP-tree on "B" ③

$$\{ \{ B:1 \} \}$$



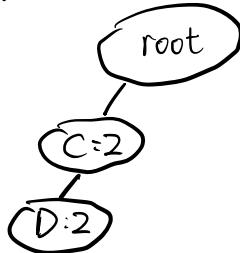
$$S_{1,2} = \{ B:1, E:1, F:1 \}$$

$$S_{2,2} = \{ BF:1, EF:1, BE:1 \}$$

$$S_{3,2} = \{ BEF:1 \}$$

Condition FP-tree on "A" ④

$$\{ \{ A=2, C=2, D=2 \} \}$$



$\because ACD=2$, $AC=AD=2 > 1$, delete BF, EF, BE in $S_{2,2}$

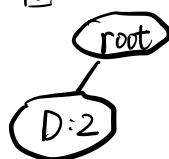
$$S_{1,2} = \{ A:2, B:1, E:1, F:1 \}$$

$$S_{2,2} = \{ AC:2, AD:2 \}$$

$$S_{3,2} = \{ ACD:2, BEF:1 \}$$

Condition FP-tree on "C" ⑤

$$\{ \{ C=3, D=3 \} \}$$



$\because CD=3$, delete B, E, F in L_i

$$S_{1,2} = \{ C=3, A:2 \}$$

$$S_{2,2} = \{ CD:3, AD:2, AC:2 \}$$

$$S_{3,2} = \{ ACD:2, BEF:1 \}$$

Condition FP-tree on "D" ⑫ root
 $\{ \{D:3\} \}$

$S_{1,2} : \{C=3, D=3\}$
 $S_{2,2} : \{CD=3, AD=2, AC=2\}$
 $S_{3,2} : \{ACD=2, BEF=1\}$

$$\therefore S_{1,2} = \{C:3, D:3\}, S_{2,2} = \{CD:3, AD:2, AC:2\}; S_{3,2} = \{ACD:2, BEF:1\}$$

This kind of FP-tree algorithm has "changible" threshold. The threshold set to 1 at first.

$S_{1,2}, S_{2,2}, S_{3,2}$ are also set to 1 at first.

Also we need to three variables to record the the 2nd (if there is just one, then record the 1st) greatest value in the multi-set of the supports of all itemsets in $S_{1,2}, S_{2,2}, S_{3,2}$ respectively. I name them a_1, a_2, a_3 .

In this phrase, we can get the FP-tree with the support threshold of 1.

Then we need to get the FP-conditional tree with the support threshold at the stage.

① Specifically, set the variable a_i mentioned above to the i -th greatest value of the multi-set of the supports of all k -itemsets in $S_{k,1}$.

② If we can find other itemset with support greater than a_i , we add the itemset to $S_{k,1}$.

③ update a_i to be the newest i -th largest value of $S_{k,1}$.

④ Remove all multi-itemset with support smaller than a from $S_{k,1}$.

⑤ do the same things to a_2, a_3

⑥ update the threshold of support = $\min(a_1, a_2, a_3)$

(C) ① We can get top-1 itemsets with greater support, while traditional FP-tree algorithm returns all frequent items with support greater than the fixed threshold

② We can have chagable support threshold, the fixed-size support-threshold is hard to set, we don't know the most appropriate value

D₂

(A) (i) Distance ² to A(50,4)		Distance ² to B(50,70)	belongs	Cluster A	Cluster B
$x_1(55,50)$	106	425	A		
$x_2(43,50)$	130	449	A		
$x_3(55,52)$	146	349	A		
$x_4(43,54)$	218	305	A		
$x_5(58,53)$	208	353	A		
$x_6(41,47)$	117	610	A		
$x_7(50,41)$	0		A		
$x_8(50,70)$		0	B		

(ii) Cluster A Cluster B
 x_2, x_4, x_6, x_7
 $(44.25, 48)$
mean of cluster(44.5, 56.25)

(iii) Cluster A Cluster B Cluster C
 x_1, x_7
 x_2, x_4, x_6
 $(42.33, 50.33)$
mean (52.5, 45.5) (54, 3, 58.3)
(iv) Cluster A Cluster B Cluster C Cluster D
 x_7
 x_2, x_4, x_6
 $(42.33, 50.33)$
mean (50.41) (56, 51, 67) (50.70)

(b) I think the answer is negative.

Here is an example:

$$x_1: (0,0), x_2: (1,1), x_3: (0.5,0), x_4: (0.25,1) \quad \text{meanA:}(0,0.5), \text{meanB:}(1,0.5)$$

the result of k-means: cluster A cluster B

$$\begin{array}{ll} x_1, x_4 & x_2, x_3 \\ \text{mean} & (0.5, 1) \end{array}$$

the result of k-means:

after x_1 is assigned, meanA: (0,0) meanB: (1,0.5)

after x_2 is assigned, meanA: (0,0) meanB: (1,1)

after x_3 is assigned, meanA: (0.375,0) meanB: (1,1)

after x_4 is assigned, meanA: (0.375,0) meanB: (1,0.5)

$$\begin{array}{ll} x_1, x_3 & x_2, x_4 \end{array}$$

different from k-means

This is because the result obtain by sequential k-means on the presentation order of the data points, but the result of k-means are not related to the presentation order.

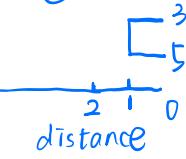
Q4

(a)	1 2 3 4 5 6 7 8
1	0
2	11 0
3	5 13 0
4	12 2 14 0
5	7 17 18 0
6	13 4 15 5 20 0
7	9 15 12 16 15 19 0
8	11 20 12 21 17 22 30 0

group(3,5)

1 2 (3.5) 4 6 7 8
1 0
2 11 0
3 5 13 0
4 12 2 14 0
5 7 17 18 0
6 13 4 15 5 20 0
7 9 15 12 16 15 19 0
8 11 20 12 21 17 22 30 0

①

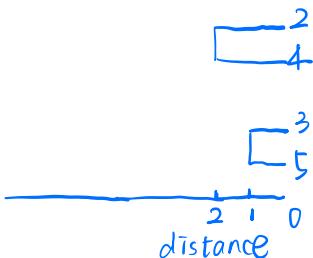


1 (2.4,6) (3.5) 7 8
1 0
(2.4,6) 12 0
(3.5) 6 16,17 0
7 9 16,17 13,5 0
8 11 21 14,5 30 0

group(2,4,6)

1 (2.4) (3.5) 6 7 8
1 0
(2.4) 12 0
(3.5) 6 16,17 0
7 9 16,17 13,5 0
8 11 21 14,5 30 0

②



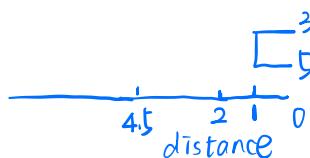
(1,3,5) 0
(2,4,6) 14,7,8 0
7 12 16,17 0
8 13,3 21 30 0

group(1,3,5,7)

1 (2.4,6) (3.5) 6 7 8
1 0
(2.4,6) 14,7,8 0
7 12 16,17 0
8 13,3 21 30 0

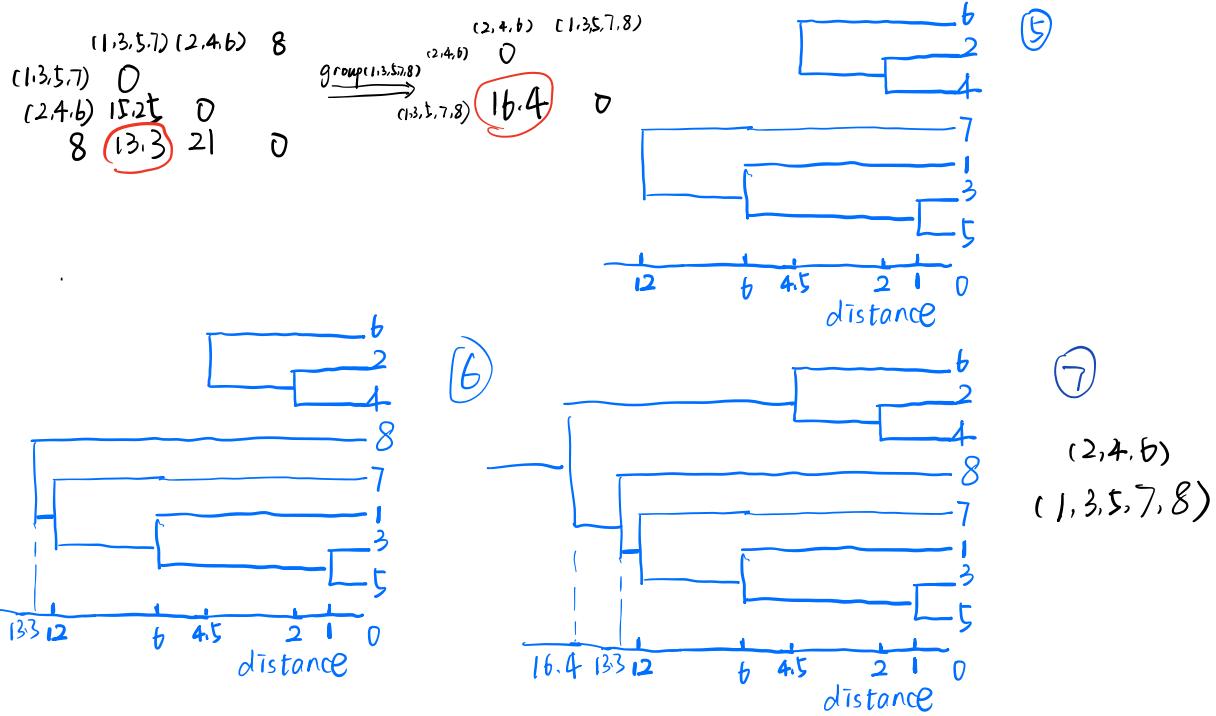


③



distance

distance



(b) from the root, we can get 2 branches $(2,4,6); (1,3,5,7,8)$

$$\because \text{distance}(8, (1,3,5,7)) = 13.3 > \text{distance}(6, (2,4)) = 4.5$$

i. we get 3 branches $(2,4,6) (8) (1,3,5,7)$

$$\text{distance}(7, (1,3,5)) = 12 > \text{distance}(6, (2,4)) = 4.5$$

ii. finally, we get 4 branches: $(2,4,6) (7) (8) (1,3,5)$

(C) (i) the greatest possible number is 8, all of them

(ii) the smallest possible number is 3; $(1,3,5)$

(d) Yes, triangle inequality refer to $||a,b| - |b,c|| \leq |a,c| \leq |a,b| + |b,c|$

let a, b, c now to be cluster centroids and we wish to determine the closest two centroid to join.

if we already know $|a,b|$ and $|b,c|$ and $||a,b| - |b,c|| \geq k$, then we don't need to compute the $|a,c|$, as lower bound on

$|a,c|$ already exceed k , hence a, c can't be joined.

(2) the algorithm is like this:

given that k constraint for agglomerative methods: Two clusters whose geometric centroids are separated by a distance greater than k cannot be joined

① for $i=2$ to $n-1$

② $d_{i,j} = D(C_i, C_j)$

③ for $i=2$ to $n-1$

④ for $j=i+1$ to $n-1$ $\hat{d}_{i,j} = |d_{i,i} - d_{i,j}|$

⑤ if $\hat{d}_{i,j} > k$ then $d_{i,j} = k+1$; i, j cannot be joined
else $d_{i,j} = D(x_i, x_j)$; i, j can be joined

⑥ return $d_{i,j}$.

Q5

(a)

DBSCAN(DB, distFunc, ϵ , MinPts)

```

C = 0      # C means the cluster counter
for each point p in database DB
    if label(p) ≠ undefined then continue      # check if processed in the inner loop
    Neighbors N = RangeQuery(DB, distFunc, P,  $\epsilon$ )      # Find neighbors
    if |N| < MinPts then      # check density
        label(p) = Noise
        continue;
    C ++      # next cluster
    label(p) = C      # Label initial point
    SeedSet S = N \ {p}      # Neighbors to expand
    for each point Q in S      # process every seed point Q
        if label(Q) = Noise then label(Q) = C      # update Noise to border point
        if label(Q) ≠ undefined then continue      # check whether processed before
        label(Q) = C      # label neighbor
        Neighbors N = RangeQuery(DB, distFunc, Q,  $\epsilon$ )      # find neighbors
        if |N| ≥ MinPts then      # check if Q is a core point
            S = S ∪ N      # Add new neighbors to seed set.
    
```

RangeQuery(DB, distFunc, Q, ϵ)

Neighbors N := empty list

```

for each point P in database DB      # Scan in database
    if distFunc(Q, P) ≤  $\epsilon$  then      # Compute distance and check  $\epsilon$ 
        N = N ∪ {P}      # Add to result .
    
```

return N

(b) $\epsilon = 10$, MinPt = 3

$N(1) = \{3, 5, 7\} = 3 = \text{MinPt}$
 $N(2) = \{4, 6\} = 2 < \text{MinPt}$
 $N(3) = \{1, 5\} = 2 < \text{MinPt}$
 $N(4) = \{2, 6\} = 2 < \text{MinPt}$
 $N(5) = \{1, 3\} = 2 < \text{MinPt}$
 $N(6) = \{2, 4\} = 2 < \text{MinPt}$
 $N(7) = \{3\} = 0 < \text{MinPt}$

$$N(8) = \emptyset = 0 < \text{MinPt}$$

i. core point : 1

border point : 3, 5, 7

noise point : 2, 4, 6, 8

(c) DBSCAN visits each point of the database, possibly multiple times.

As for time complexity, we look for neighbors for each points.

$O_{\text{overall}} = n \cdot O(\text{rangeQuery})$. As for query, we can use binary tree to store the data

so $O(\text{rangeQuery})$ can be minimum to be $O(n)$. $O(\text{DBSCAN overall averages}) = O(n \log n)$

and the worst situation is use a list or array to store, $O(\text{DBSCAN worst situation}) = O(n^2)$

As for database, we can easily use a list or array to store, so spatial complexity can be $O(n)$

if we use Q4 situation to store database, the spatial complexity is $O(n^2)$

(d) ① DBSCAN based on density, while hierarchical clustering is based on distance.

② DBSCAN groups together closely-packed points, not every point is part of cluster (noise).
hierarchical clustering don't.

③ DBSCAN can deal with noise and outliers. hierarchical clustering can't.

④ the border points of two clusters may be arbitrarily classified. hierarchical clustering have
specific cluster.

⑤ DBSCAN has difficulty to find clusters of varying densities, but hierarchical clustering can easily
find clusters of different density.