2022 3.

(i)

$\{(0,0),(0,1)\}$ —— $\{(0,0),(4,4)\}$ —— $\{(0,0),(5,3)\}$

$\{(0,1),(4,4)\}$ —— $\{(0,1),(5,3)\}$ —— $\{(4,4),(5,3)\}$

(ii) assume we start from $\{(0,0),(0,1)\}$

cost of $\{(0,0),(0,1)\}$ = $d_{(0,1)-(4,4)} + d_{(0,1)-(5,3)}$

= $7 + 7 = 14$

cost of $\{(0,0),(4,4)\}$ = $d_{(0,0)-(0,1)} + d_{(4,4)-(5,3)}$

= $1 + 2 = 3$

cost of $\{(0,0),(5,3)\}$ = $d_{(0,1)-(0,1)} + d_{(5,3)-(4,4)}$

= $1 + 2 = 3$

cost of $\{(0,1),(4,4)\}$ = $d_{(0,1)-(0,0)} + d_{(4,4)-(5,3)}$

= $1 + 2 = 3$

cost of $\{(0,1),(5,3)\}$ = $d_{(0,1)-(0,0)} + d_{(5,3)-(4,4)}$

= $1 + 2 = 3$

In the PAM, we will calculate all the neighbors and find one with minimum cost, here we have 4 neighboors equally cost 3 , all of them represent a minimum cost, so we can randomly pick one from $\{(0,0),(4,4)\}$, $\{(0,0),(5,3)\}$ , $\{(0,1),(4,4)\}$, $\{(0,1),(5,3)\}$

b.1) No. PAM always follow the local minimum, so it can only find local minimum.

2) Yes. As said above, PAM will follow optimal path in each iteration, so it can find local minimum.

3) No. CLARA is done on subdataset, so it can only find local optimum of subdata.

4) No. CLAPANS doesn't always follow the local minimum in each iteration, so it can not guarantee to find local optimum.

# 2. Hierarchical clustering

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

iteration 1

$\Rightarrow$

|   | (1,5) | 2 | 3 | 4 |
|---|---|---|---|---|
| (1,5) | 0 | | | |
| 2 | 7 | 0 | | |
| 3 | 10 | 9 | 0 | |
| 4 | 8 | ② | 5 | 0 |

iteration 2

$\Rightarrow$

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

iteration 3

$\Rightarrow$

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

iteration 4

In iteration 1, we find $d_{(1,5)}=1$ is minimum.
So we combine (1,5) together

$d_{(1,5)-2} = max\{d_{(2,1)}, d_{(2,5)}\} = 7$

$d_{(1,5)-3} = max[d_{(3,1)}, d_{(3,5)}] = 10$   distance in iter 2

$d_{(1,5)-4} = max[d_{(1,4)}, d_{(5,4)}] = 8$

In iteration 2, we find $d_{(2,4)}=2$ is minimum
so, we combine (2,4) together.

$d_{(1,5)-(2,4)} = max\{d_{(1,5)-2}, d_{(1,5)-4}\} = 8$   distance in iter 3

$d_{(2,4)-3} = max\{d_{(2,3)}, d_{(4,3)}\} = 9$

In iteration3, we find $d_{(1,5)-(2,4)}=8$ is minimum,
so, we combine [(1,5)(2,4)] together.

$d_{[(1,5),(2,4)]-3} = max\{d_{(1,5)-3}, d_{(2,4)-3}\} = 10 \rightarrow$ distance in iter 4



(i) From the graph in the left side, we can see

(i) threhold = 5 :
   cluster 1 :   (1,5)
   cluster 2 :   (2,4)
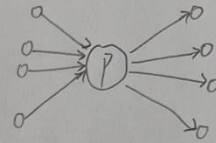   cluster 3 :   (3)

(ii) threhold = 9 :
   cluster 1 :   (1,2,4,5)
   cluster 2 :   (3)

## 3. ROCK

(i) if the threshold = 0.6, then the neighbors of A are (A, B, C). The neighbors of
are (A, B, C), so A, B have 3 common neighbors => link (A, B) = 3

(ii) The expected link is

$$n^{1+2m}$$



We consider the graph in the right side,
for every point P in cluster C, it has $P_j$-th
in-degree and $P_j$-th out-degree. Every point
form in-degree and every point from out degree can
be a pair whose common neighbor must contain P.
So in this graph, P can contribute $P_i * P_j$ -th common
neighbors.

Back to our problems, each neighbor can be
in-degree or out-degree so for each point the
expected link is $n^m \cdot n^m = n^{2m}$. And we have
n-th points here, so the expected links are

$$n^{1+2m}$$

(iii) The expected link here works as a kind of normalization. If there is no such normalization
the greatest goodness will be reached when we put all points into a single cluster. So, it
is necessary to use this normalization.

# 4. Density-based clustering

1) False. If p and q are density connected, then it means there is a core "A" that p and q are density reachable from "A". So, "q is density from p" is not guaranteed.

2) False. "Directly density reachable" requires core point. There is no information showing p is a core point.

3) True. From the statement, we can know q is a core point. Meanwhile, there is a path from q to p : $P_1, P_2, \cdots P_n$ ; $P_1 = q$, $P_n = P$ and $P_{i+1}$ is directly density reachable from $P_i$. So for $P_{n-1}$, both q and P are density reachable from $P_{n-1}$, so p and q are density connected.

4) True. If p and q are density connected, then there is an point "o" where p and q are density reachable from. So, we can see density connected is a symmetric condition.

## 5. Distance Measure

a.

$$d_{KL}(A) = \sqrt{[3-(-2)]^2} = 5 \qquad \delta_{KL}(A) = 1$$

$$d_{KL}(B) = |2-(-1)| + |3-2| = 4 \qquad \delta_{KL}(B) = 1$$

$$d_{KL}(C) = 0 \qquad\qquad \delta_{KL}(C) = 1$$

$$d_{KL}(D) = 1 \qquad\qquad \delta_{KL}(D) = 1$$

$$d_{KL}(E) = 0 \qquad\qquad \delta_{KL}(E) = 1$$

$$d_{KL}(F) = 0 \qquad\qquad \delta_{KL}(F) = 0$$

$$d_{KL}(G) = 0$$

$$d_{KL} = \frac{\sum \delta_{ij}(f) \cdot d_{ij}(f)}{\sum \delta_{ij}(f)} = \frac{5+4+1}{5} = 2$$

b. We will have the same answer.

If we us asymmetric variable only, we may have a table looks like below

| | Y | N |
|---|---|---|
| Y | a | b |
| N | c | d |

then the distance is $\frac{b+c}{a+b+c}$

If using the formula in question (a), the (Y, N) and (N, Y) will be counted into distance.

For denominator, (Y, Y) (Y, N) (N, Y) will contribute. From the table, we have a-th pairs of (Y, Y), b-th pairs of (Y, N) and c-th pairs of (N, Y). So, the result would be

$$\frac{b+c}{a+b+c}.$$

Now, we can see they will have the same answer.

# 6. Aprior

## a.

$C_1$

| items | sup |
|-------|-----|
| A | 4 |
| B | 4 |
| C | 4 |
| D | 4 |
| E | 2 |

$\longrightarrow L_1$

| item | sup |
|------|-----|
| A | 4 |
| B | 4 |
| C | 4 |
| D | 4 |
| E | 2 |

$\longrightarrow C_2$

| item | sup |
|------|-----|
| AB | 3 |
| AC | 3 |
| AD | 3 |
| AE | 2 |
| BC | 3 |
| BD | 3 |
| BE | 2 |
| CD | 3 |
| CE | 1 |
| DE | 1 |

$\longrightarrow L_2$

| item | sup |
|------|-----|
| AB | 3 |
| AC | 3 |
| AD | 3 |
| AE | 2 |
| BC | 3 |
| BD | 3 |
| BE | 2 |
| CD | 3 |

$\longrightarrow C_3$

| item | sup |
|------|-----|
| ABC | 2 |
| ABD | 1 |
| ABE | 2 |
| ACD | 2 |
| ACE | 1 |
| ADE | 1 |
| BCD | 2 |
| BCE | 1 |
| BDE | 1 |

$\longrightarrow L_3$

| item | sup |
|------|-----|
| ABC | 2 |
| ABE | 2 |
| ACD | 2 |
| BCD | 2 |

$\longrightarrow C_4$

| item | sup |
|------|-----|
| ABCE | 1 |
| ABCD | 1 |

$\longrightarrow L_4$

| item | sup |
|------|-----|
| null | null |

output
$$L_1: \{A\}, \{B\}, \{C\}, \{D\}, \{E\}$$
$$L_2: \{A,B\}, \{A,C\}, \{A,D\}, \{A,E\}, \{B,C\}, \{B,D\}, \{B,E\}, \{C,D\}$$
$$L_3: \{A,B,C\}, \{A,B,E\}, \{A,C,D\}, \{B,C,D\}$$

6. b

This is antimonotone constraint.

C.

| item | sup | value |
|------|-----|-------|
| A | 4 | 0 |
| B | 4 | 0 |
| C | 4 | 0 |
| D | 4 | 0 |
| E | 2 | 0 |

→ L₁

| item | sup | range |
|------|-----|-------|
| A | 4 | 0 |
| B | 4 | 0 |
| C | 4 | 0 |
| D | 4 | 0 |
| E | 2 | 0 |

→ C₂

→ L₂

| item | sup | range |
|------|-----|-------|
| AB | 3 | 1 |
| AC | 3 | 2 |
| AD | 3 | 3 |
| AE | 2 | 4 |
| BC | 3 | 1 |
| BD | 3 | 2 |
| BE | 2 | 3 |
| CD | 3 | 1 |
| CE | 1 | 2 |
| DE | 1 | 1 |

→ L₂

| item | sup | range |
|------|-----|-------|
| AB | 3 | 1 |
| AC | 3 | 2 |
| BC | 3 | 1 |
| BD | 3 | 2 |
| CD | 3 | 1 |

By using constraint, we move AD, AE, BE
out of L₂

C₃ →

| item | sup | range |
|------|-----|-------|
| ABC | 2 | 2 |
| ABD | 1 | 3 |
| ACD | 2 | 3 |
| BCD | 2 | 2 |

→ L₃

| item | sup | range |
|------|-----|-------|
| ABC | 2 | 2 |
| BCD | 2 | 2 |

→ C₄

| item | sup | range |
|------|-----|-------|
| ABCD | 1 | 3 |

→ L₄ empty

By applying constraint,
we move ACD out.

So: { {A} {B} {C} {D} {E}
{AB}, {AC}, {BC}, {BD} {CD}
{ABC}, {BCD} }

C. This is monotone constraint

C₁

| item | sup | min(S) |
|------|-----|--------|
| A | 4 | 5 |
| B | 4 | 4 |
| C | 4 | 3 |
| D | 4 | 2 |
| E | 2 | 1 |

→ L₁

| item | sup | min(S) |
|------|-----|--------|
| A | 4 | 5 |
| B | 4 | 4 |
| C | 4 | 3 |
| D | 4 | 2 |
| E | 2 | 1 |

→ C₂

| item | sup | min(S) |
|------|-----|--------|
| AB | 3 | 4 |
| AC | 3 | 3 |
| AD | 3 | 2 |
| AE | 2 | 1 |
| BC | 3 | 3 |
| BD | 3 | 2 |
| BE | 2 | 1 |
| CD | 3 | 2 |
| CE | 1 | 1 |
| DE | 1 | 1 |

→ L₂

| item | sup | min(S) |
|------|-----|--------|
| AB | 3 | 4 |
| AC | 3 | 3 |
| AD | 3 | 2 |
| AE | 2 | 1 |
| BC | 3 | 3 |
| BD | 3 | 2 |
| BE | 2 | 1 |
| CD | 3 | 2 |

[ ]In L₁, by applying constraint, we
can only output {D}, {E}

In L₂, by applying the constraint, we
can only output {A,D}, {A,E}, {B,D}, {B,E}
{C,D}

→ C₃

| item | sup | min(S) |
|------|-----|--------|
| ABC | 2 | 3 |
| ABD | 1 | 2 |
| ABE | 2 | 1 |
| ACD | 2 | 2 |
| ACE | 1 | 1 |
| ADE | 1 | 1 |
| BCD | 2 | 2 |
| BCE | 1 | 1 |
| BDE | 1 | 1 |

→ L₃

| item | sup | min(S) |
|------|-----|--------|
| ABC | 2 | 3 |
| ABE | 2 | 1 |
| ACD | 2 | 2 |
| BCD | 2 | 2 |

→ C₄

| item | sup | min(S) |
|------|-----|--------|
| ABCE | 1 | 1 |
| ABCD | 1 | 2 |

→ L₄ empty

In L₃, by applying constraint,
we output {A,B,E}, {A,C,D}, {B,C,D}

So, the output is  { {D}, {E}
{A,D}, {A,E}, {B,D}, {B,E}, {C,D}
{A,B,E}, {A,C,D}, {B,C,D} }

# 7. FP grow algorithm

a.

| item | sup |
|------|-----|
| T | 8 |
| Q | 6 |
| S | 4 |
| U | 3 |
| R | 5 |
| P | 1 |

(1)

| Tid | items |
|-----|-------|
| 1 | T |
| 2 | T Q S |
| 3 | T Q U |
| 4 | T R |
| 5 | T Q R U |
| 6 | T Q R S |
| 7 | T Q R S U |
| 8 | T Q R S |

```
        {}
         |
        T:8
       /    \
     Q:6    R:1
     / |  \
   S:1 R:4 U:1
        |  \
       U:1  S:3
              \
              U:1
```

b.

Using the table (1) from question a. This is an antimonotone constraint

| Tid | items | Range |
|-----|-------|-------|
| 1 | T | 0 |
| 2 | T Q S | 3 |
| 3 | T Q U | 4 |
| 4 | T R | 2 |
| 5 | T Q R U | 4 |
| 6 | T Q R S | 3 |
| 7 | T Q R S U | 4 |
| 8 | T Q R S | 3 |

After applying constraint, we only have

| Tid | items | Range |
|-----|-------|-------|
| 1 | T | 0 |
| 4 | T R | 2 |

```
  {}
   |
  T:2
   |
  R:1
```

The conditional database is

| item | |
|------|-----|
| T | S1 |
| R | T:1 |

# 8. Rule Generation

$\{\bar{F}, H, J, L\}$

① $FHJ \to L$  $P(F,H,J,L) = \frac{5}{10}$
$\quad P(F,H,J) = \frac{5}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(FHJ)} = 1 > 60\% \quad,\quad \text{output}$

② $\bar{F}H \to JL$  $P(\bar{F}HJL) = \frac{5}{10}$
$\quad P(\bar{F}H) = \frac{6}{10}$  $\text{confidence} = \frac{P(\bar{F}HJL)}{P(\bar{F}H)} = \frac{5}{6} > 60\% \quad,\quad \text{output}$

③ $\bar{F} \to HJL$  $P(\bar{F}HJL) = \frac{5}{10}$
$\quad P(\bar{F}) = \frac{6}{10}$  $\text{confidence} = \frac{P(\bar{F}HJL)}{P(\bar{F})} = \frac{5}{6} > 60\% \quad,\quad \text{output}$

④ $H \to FJL$  $P(FHJL) = \frac{5}{10}$
$\quad P(H) = \frac{9}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(H)} = \frac{5}{9} < 60\% \quad\quad \times$

⑤ $FJ \to HL$  $P(FHJL) = \frac{5}{10}$
$\quad P(FJ) = \frac{5}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(FJ)} = 1 > 60\%,\ \text{output}$

⑥ $F \to HJL$  already discussed ③

⑦ $J \to FHL$  $P(FHJL) = \frac{5}{10}$
$\quad P(J) = \frac{8}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(J)} = \frac{5}{8} > 60\%,\ \text{output}$

⑧ $HJ \to FL$  $P(FHJL) = \frac{5}{10}$
$\quad P(HJ) = \frac{7}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(HJ)} = \frac{5}{7} > 60\%,\ \text{output}$

⑨ $H \to JFL$  already discussed ④

⑩ $J \to HFL$  already discussed ⑦

⑪ $FHL \to J$  $P(FHJL) = \frac{5}{10}$
$\quad P(FHL) = \frac{5}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(FHL)} = 1 > 60\%,\ \text{output}$

⑫ $FH \to JL$  already discussed ②

⑬ $F \to HJL$  already discussed ⑥

⑭ $H \to FJL$  already discussed ④

⑮ $\bar{F}L \to HJ$  $P(FHJL) = \frac{5}{10}$
$\quad P(\bar{F}L) = \frac{5}{10}$  $\text{confidence} = \frac{P(FHJL)}{P(\bar{F}L)} = 1 > 60\%,\ \text{output}$

# 8. Rule Generation

⑯ F→HLJ    already discussed ⑥

⑰ L→FHJ    $P(FHJL) = \frac{5}{10}$   confidence $= \frac{P(FHJL)}{P(L)} = \frac{5}{6} > 60\%$, output
           $P(L) = \frac{6}{10}$

⑱ HL→FJ    $P(FHJL) = \frac{5}{10}$   confidence $= \frac{P(FHJL)}{P(HL)} = \frac{5}{6} > 60\%$, output
          $P(HL) = \frac{6}{10}$

⑲ H→FJL    already discussed ④

⑳ L→FHJ    already discussed ⑰

㉑ FJL→H    $P(FJLH) = \frac{5}{10}$   confidence $= \frac{P(FHJL)}{P(FJL)} = 1 > 60\%$, output
          $P(FJL) = \frac{5}{10}$

㉒ FJ→LH    already discussed ⑤

㉓ F→LHJ    already discussed ③

㉔ J→FLH    already discussed ⑦

㉕ FL→JH    already discussed ⑮

㉖ F→LJH    already discussed ⑥

㉗ L→FJH    already discussed ⑰

㉘ JL→FH    $P(FHJL) = \frac{5}{10}$   confidence $= \frac{P(FHJL)}{P(JL)} = \frac{5}{6} > 60\%$, output
         $P(JL) = \frac{6}{10}$

㉙ J→FHL    already discussed ⑦

㉚ L→FJH    already discussed ⑰

㉛ HJL→F    $P(HJFL) = \frac{5}{10}$   confidence $= \frac{5}{6} > 60\%$, output
         $P(HJL) = \frac{6}{10}$

㉜ HJ→FL    already discussed ⑧ ; ㉝ H→JFL already discussed ④; ㉞ J→FHL already discussed ⑦

㉟ HL→FJ    already discussed ⑱ ; ㊱ H→FJL already discussed ④ ;㊲L→FJH already discussed ⑰

㊳ JL→FH    $P(FHJL) = \frac{5}{10}$   confidence $= \frac{P(FHJL)}{P(JL)} = \frac{5}{6} > 60\%$, output ;㊳J→FHL already discussed ⑦
        $P(JL) = \frac{6}{10}$                          ;㊴ L→FJH already discussed ⑰

Finally, we will output

| | | | |
|---|---|---|---|
| FHJ→L | J→FHL | L→FHJ | HJL→F |
| FH→JL | HJ→FL | HL→FJ | JL→FH |
| F→HJL | FHL→J | FJL→H | |
| FJ→HL | FL→HJ | JL→FH | |