

7/4/23  
Monday

## Module-4 Associations Rule Analysis (AP)

Possibilities of buying two items together

It is also known as market basket analysis.

It is used to find relationship / association b/w large set of data items.

$$x \rightarrow y$$

Support and confidence are used to find association b/w items in a dataset.

Support

The support for  $x \rightarrow y$  is the probability of both  $x$  and  $y$  appearing together, i.e.,  $P(x \cup y)$

$$\text{support}(A \cup B) = P(A \cup B)$$

Confidence

The confidence of  $x \rightarrow y$  is the conditional probability of  $y$  appearing given that  $x$  exists. It is written as  $P(y/x)$  and read as P of  $y$  given  $x$ .

$$\text{confidence}(A \Rightarrow B) = P(B/A) = \frac{\text{support\_count}(A \cup B)}{\text{support\_count}(A)}$$

Find the support value of each itemset

Transaction ID	Items purchased
1	Bread, cheese, Egg, juice
2	Bread, cheese, juice
3	Bread, milk, yogurt.
4	Bread, juice, milk.
5	Cheese, juice, Milk.

Support (item) = Frequency of item / Total no. of brands

Total no. of transaction = 5

$$\text{support}(\text{Bread}) = \frac{4}{5} = 80\%$$

$$\text{support}(\text{cheese}) = \frac{3}{5} = 60\%$$

$$\text{support}(\text{Egg}) = \frac{1}{5} = 20\%$$

$$\text{support}(\text{Juice}) = \frac{4}{5} = 80\%$$

$$\text{support}(\text{Milk}) = \frac{3}{5} = 60\%$$

$$\text{support}(\text{yogurt}) = \frac{1}{5} = 20\%$$

## Methods to discover Association rules

Basic idea behind the method to discover association rule.

Step 1: Find all frequent itemset

Step 2: Generate strong association rules from the frequent itemset

## Frequent itemset

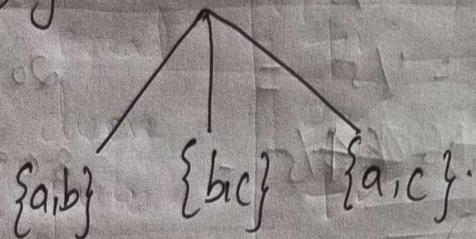
An itemset  $X \subseteq A$  is said to be a frequent itemset with respect to  $\sigma$ , if

$$\text{support}(X) \geq \sigma$$

Downward closure property: Any subset of a frequent set is a frequent set.

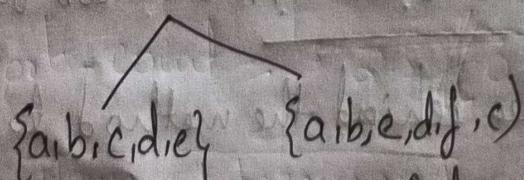
Upward closure property: Any superset of an infrequent set is an infrequent set.

e.g:  $\{a, b, c\} \rightarrow$  frequent



$\{a\}$   $\{b\}$   $\{c\}$

→ If ~~a & b & c~~  $\{a, b, c\}$  is frequent then all



the supersets of  $\{a, b, c\}$  are infrequent.

## APRIORI algorithm

It is also known as level-wise algorithm

Step 1: discover all frequent (single) items that have support  $\geq \sigma$  above the minimum support number

Step 2: use the set of frequent items to generate the association rules that have high enough confidence level.

The candidate generation process and the pruning process are the most important parts of the algorithm.

Candidate set generation method

Given represented as  $C_k$ .

Given  $L_{k-1}$  the set of all frequent  $(k-1)$ -itemset, we

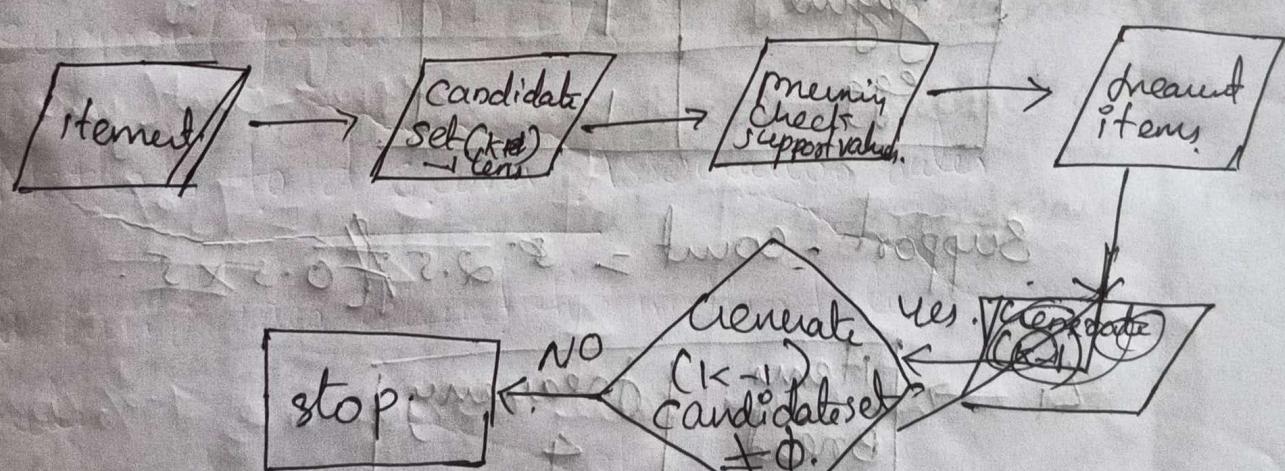
want to generate frequent  $k$ -itemset.

Eg:  $\{1, 2, 3\}$      $\{1, 2, 3, 4\}$

$$C = \{1, 2, 3, 4\}$$

Pruning

It is used to eliminate infrequent items.



Find association rule from the following data with minimum support of 50% & minimum confidence of 75%

Output to be displayed

Transaction ID	Items
100	Bread, cheese, egg, juice
200	Bread, cheese, juice
300	Bread, Milk, yogurt
400	Bread, juice, Milk
500	cheese, juice, Milk

Step 1: Find L1  
 The items that are frequent ( $> 50\%$  of support)

Item	freenency
Bread	4 ✓
cheese	3 ✓
juice	4 ✓
Milk	3 ✓
yogurt	1
egg	1

$$\text{Support - count} = \frac{2.5}{10} \times 5$$

L1 → Items	freenency
Bread	4
cheese	3
juice	4
Milk.	3

Step 2: Find C2 → candidate set of two items.

Items	Frequency
{Bread, cheese}	2.
{Bread, juice}	3
{Bread, Milk}	2.
{cheese, juice}	3
{cheese, milk}	1
{juice, milk}	2.

3. Find  $L_2$ . ( $\geq 2.5$ ).

$L_2 \rightarrow$	Items	Frequency
	{Bread, juice}	3
	{cheese, juice}	3.

4. Find  $C_3$ .

Items. Frequency.

If it is not possible to create subset a candidate set  
 $C_3$  {Bread, cheese, juice} because <sup>all</sup> subset of frequent

subset must be frequent. There is no frequent item  
of the form  
{Bread, cheese}.

Algorithm will stop

generate 4 association rules from the input data.

Bread  $\rightarrow$  juice

juice  $\rightarrow$  Bread

cheese  $\rightarrow$  juice

juice  $\rightarrow$  Bread, cheese.

$$\text{Bread} \rightarrow \text{Juice} \rightarrow \text{confidence} = \frac{\text{support-count}(\text{Bread, Juic})}{\text{support-count}(\text{Bread})}$$

$$= \frac{3}{4} = 75\%$$

$\text{Juice} \rightarrow \text{Bread}$

$$\text{confidence} = \frac{3}{4} = 75\%$$

$\text{Cheese} \rightarrow \text{Juice}$

$$\text{confidence} = \frac{\text{support-count}(\text{cheese, juice})}{\text{support-count}(\text{cheese})}$$

$$= \frac{3}{3} = 100\%$$

$\text{Juice} \rightarrow \text{Cheese}$

$$\text{confidence} = \frac{\text{support-count}(\text{juice, cheese})}{\text{support-count}(\text{juice})}$$

$$= \frac{3}{4} = 75\%$$

Since all items have minimum of 75% confidence.

For

2. Find all AR from the following data using Apriori algorithm with minimum support count measured as 2 & confidence 70%.

TID.	list of item-IDs.
T <sub>100</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub> .
T <sub>200</sub>	I <sub>2</sub> , I <sub>4</sub> .
T <sub>300</sub>	I <sub>2</sub> , I <sub>3</sub> .
T <sub>400</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub> .
T <sub>500</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>600</sub>	I <sub>2</sub> , I <sub>3</sub> .
T <sub>700</sub>	I <sub>1</sub> , I <sub>3</sub> .
T <sub>800</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub> .
T <sub>900</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> .

1. Find L<sub>1</sub>.

Items	frequency	Support
I <sub>1</sub>	6	
I <sub>2</sub>	7	
I <sub>3</sub>	6	
I <sub>4</sub>	2	
I <sub>5</sub>	2	Support-count > 2.

L <sub>1</sub> → Items	frequency
I <sub>1</sub>	6
I <sub>2</sub>	7
I <sub>3</sub>	6
I <sub>4</sub>	2
I <sub>5</sub>	2

2. Find  $C_2$

Items	frequency.
$\{I_1, I_2\}$	4.
$\{I_1, I_3\}$	4
$\{I_1, I_4\}$	1
$\{I_1, I_5\}$	2.
$\{I_2, I_3\}$	4
$\{I_2, I_4\}$	2.
$\{I_2, I_5\}$	2.
$\{I_3, I_4\}$	0
$\{I_3, I_5\}$	1
$\{I_4, I_5\}$	0.

3. Find  $L_2$ .  $\geq 2$

Items	frequency.	
$\{I_1, I_2\}$	4	$\{I_1, I_2, I_3\} \checkmark$
$\{I_1, I_3\}$ .	4	$\{I_1, I_2, I_5\} \times$
$\{I_1, I_5\}$ .	2.	$\{I_1, I_2, I_4\} \checkmark$
$\{I_2, I_3\}$ .	4	$\{I_1, I_3, I_4\} \times$
$\{I_2, I_4\}$ .	2.	$\{I_1, I_3, I_5\} \times$
$\{I_2, I_5\}$ .	2.	$\{I_2, I_3, I_4\} \times$

4. Find  $C_3$

Items      frequency.

$\{I_1, I_2, I_3\}$

2.

$\{I_1, I_2, I_3\}$ .

$\{I_1, I_2\} \quad \{I_1, I_3\}$

$\{I_1, I_2, I_5\}$

2.

lives - happens = evidence

5. Find  $C_3$

Items      frequency.

$\{I_1, I_2, I_3\}$

2

F

$\{I_1, I_2, I_5\}$

2.

$\{I_1, I_5\} \leftarrow C^1$

6. Algorithm stop.

Association rule - happens

$\{I_1, I_2, I_3\}:$

$\{I_1, I_2, I_3\}.$

$\{I_1, I_2\} \rightarrow \{I_3\}$

$\{I_1, I_3\} \rightarrow I_2$

$\{I_2, I_3\} \rightarrow I_1$

$I_1 \rightarrow \{I_2, I_3\}.$

$I_2 \rightarrow \{I_1, I_3\}.$

$I_3 \rightarrow \{I_1, I_2\}.$

$I_1 \rightarrow \{I_2, I_3\}.$

confidence = support-cont ( $I_1 \rightarrow I_2, I_3$ )

support-cont ( $I_1$ ).

$$= \frac{2}{8} = 33.3\%$$

confidence.

$$I_2 \rightarrow \{I_1, I_3\}$$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_3)}{\text{support-count}(I_2)}$$

$$= \frac{2}{7} = 28.2\%$$

$$I_3 \rightarrow \{I_1, I_2\}$$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_3)}{\text{support-count}(I_3)}$$

$$= \frac{2}{6} = 33.3\%$$

$$\{I_1, I_2\} \rightarrow \{I_3\}$$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_3)}{\text{support-count}(I_1, I_2)}$$

$$= \frac{2}{4} = 50\%$$

$$\{I_1, I_3\} \rightarrow I_2$$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_3)}{\text{support-count}(I_1, I_3)}$$

$$= \frac{2}{4} 50\%$$

$\{I_2, I_3\} \rightarrow I_1$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_3)}{\text{support-count}(I_2, I_3)}$$

$$= \frac{2}{4} 50\%$$

$\{I_1, I_2, I_5\} \rightarrow I_3$

$I_1 \rightarrow \{I_2, I_5\}$

$I_2 \rightarrow \{I_1, I_5\}$

$I_5 \rightarrow \{I_1, I_2\}$

$\{I_2, I_5\} \rightarrow I_1$

$\{I_1, I_5\} \rightarrow I_2$

$\{I_1, I_2\} \rightarrow I_5$

~~co~~  $\{I_1 \rightarrow \{I_2, I_5\}\}$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_1)}$$

$$= \frac{2}{6} 33.3\%$$

$I_2 \rightarrow \{I_1, I_5\}$

$$\text{confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_2)}$$

$$I_5 \rightarrow \{I_1, I_2\}$$

$$\text{Confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_5)}$$
$$= \frac{2}{2} = 100$$

$$\{I_1, I_5\} \rightarrow I_1$$

$$\text{Confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_1, I_5)}$$
$$= \frac{2}{2} = 100$$

$$\{I_1, I_5\} \rightarrow I_2$$

$$\text{Confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_1, I_5)}$$
$$= \frac{2}{2} = 100$$

$$\{I_1, I_2\} \rightarrow I_5$$

$$\text{Confidence} = \frac{\text{support-count}(I_1, I_2, I_5)}{\text{support-count}(I_1, I_2)}$$

$$\text{Final o/p} = \frac{2}{7} 50\%$$

$$I_5 \rightarrow \{I_1, I_2\}$$

$$\{I_2, I_5\} \rightarrow I_1$$

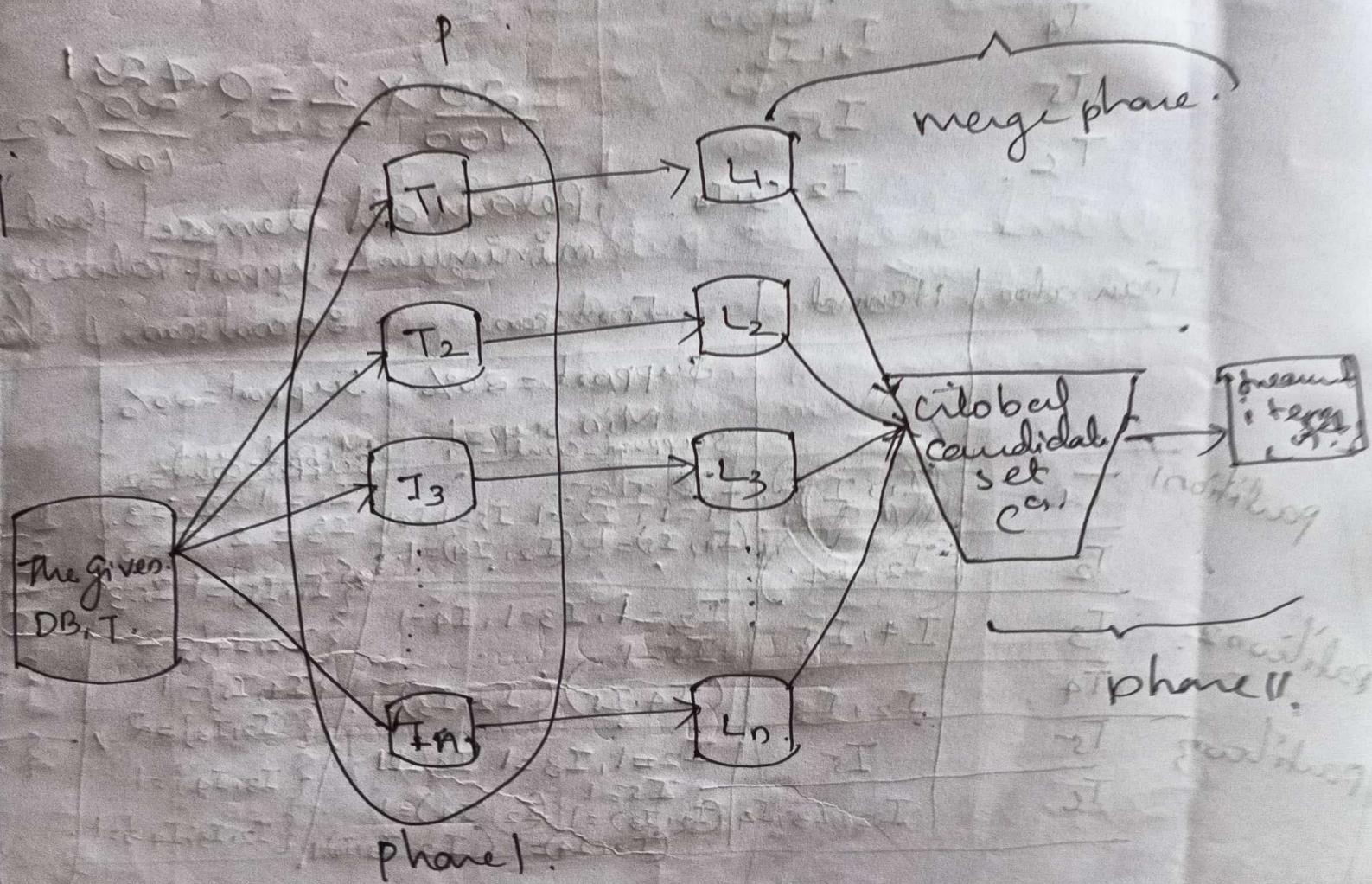
$$\{I_1, I_5\} \rightarrow I_2$$

## Partition Algorithm

Partition algorithm scans the database only twice.

First scan: generate a set of all potentially large itemset.

Second scan: set up counters for each potentially large itemset and compute their actual supports.



Transaction : Item  
 Find frequent itemset from the following data  
 with support value of 20% using partition of size 3

Transaction itemset

T<sub>1</sub> I<sub>1</sub>, I<sub>5</sub>

T<sub>2</sub> I<sub>2</sub>, I<sub>4</sub>

T<sub>3</sub> I<sub>4</sub>, I<sub>5</sub>

T<sub>4</sub> I<sub>2</sub>, I<sub>3</sub>

T<sub>5</sub> I<sub>5</sub>

T<sub>6</sub> I<sub>2</sub>, I<sub>3</sub>, I<sub>4</sub>

$$\text{Support\_count} = \frac{\text{Support}}{\text{Total}} \times 100$$

$$= \frac{100}{505} \times 100$$

$$= \frac{20}{101} \times 100$$

$$= 0.492 \approx 1$$

$$\frac{20 \times 2}{100} = 0.4 \approx 1$$

Retain all itemset that have minimum support value in first scan.

second scan. shortlisted.

Transaction	Itemset	First scan. Support = 20%	Second scan. Support = 20%
T <sub>1</sub>	I <sub>1</sub> , I <sub>5</sub>	Min. sup = 1. Support count = 1.	Support count = 2.
T <sub>2</sub>	I <sub>2</sub> , I <sub>4</sub>	I <sub>1</sub> =1, I <sub>2</sub> =1, I <sub>3</sub> =1, I <sub>4</sub> =1 (I <sub>1</sub> , I <sub>2</sub> )=1, (I <sub>2</sub> , I <sub>4</sub> )=1.	I <sub>1</sub> =1, I <sub>2</sub> =3, I <sub>3</sub> =2, I <sub>4</sub> =3 {I <sub>1</sub> , I <sub>2</sub> } = 1, {I <sub>2</sub> , I <sub>3</sub> } = 2.
T <sub>3</sub>	I <sub>4</sub> , I <sub>5</sub>	I <sub>2</sub> =1, I <sub>3</sub> =1, I <sub>4</sub> =1 I <sub>5</sub> =1	{I <sub>2</sub> , I <sub>4</sub> } = 2
T <sub>4</sub>	I <sub>2</sub> , I <sub>3</sub>	(I <sub>4</sub> , I <sub>5</sub> ), (I <sub>2</sub> , I <sub>3</sub> ) = 1.	{I <sub>4</sub> , I <sub>5</sub> } = 1
T <sub>5</sub>	I <sub>5</sub>	I <sub>2</sub> =1, I <sub>3</sub> =1, I <sub>4</sub> =1 I <sub>5</sub> =1	{I <sub>2</sub> , I <sub>3</sub> } = 2
T <sub>6</sub>	I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub>	(I <sub>2</sub> , I <sub>3</sub> ) = 1, (I <sub>2</sub> , I <sub>4</sub> ) = 1 (I <sub>3</sub> , I <sub>4</sub> ) = 1	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> } = 1

local frequent itemset

Here, all itemsets are satisfying the condition. All items are retained during first scan. These itemsets are local frequent itemset.

max op frequent itemset

$$\{I_2, I_4\} - 2$$

$$\{I_2, I_3\} - 2$$

$$I_2 \rightarrow I_4$$

$$\text{confidence} = \frac{\text{support-count}(I_2, I_4)}{\text{support-count}(I_2)}$$
$$= \frac{2}{1} = 2$$

$$\text{confidence} = \frac{\text{support-count}}{\text{support-count}}$$

20 X  
100

~~$$\frac{20}{100} \times 2$$~~  
$$2 \times 0.4 \approx 1$$

# FP Growth Algorithm → FP Tree

Find AR from the following data using FPgrowth algorithm with min cont support count 4 & 2 & confidence 80%.

Transactions

Itemset

T<sub>100</sub>

I<sub>1</sub>, I<sub>2</sub>, I<sub>5</sub>

T<sub>200</sub>

I<sub>2</sub>, I<sub>4</sub>

T<sub>300</sub>

I<sub>2</sub>, I<sub>3</sub>

T<sub>400</sub>

I<sub>1</sub>, I<sub>2</sub>, I<sub>4</sub> (bel if not)

T<sub>500</sub>

I<sub>1</sub>, I<sub>3</sub>

T<sub>600</sub>

I<sub>2</sub>, I<sub>3</sub>

T<sub>700</sub>

I<sub>1</sub>, I<sub>3</sub>

T<sub>800</sub>

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>, I<sub>5</sub> (bel if not)

T<sub>900</sub>

I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>

freq: I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>, I<sub>5</sub>  
I<sub>4</sub>

1) Find frequent pattern set ( $L$ ): It contains all items greater than minimum support value and these items are arranged in decreasing order.

Item

frequenc

I<sub>2</sub>

7

I<sub>1</sub>

6

I<sub>3</sub>

6

I<sub>4</sub>

2

I<sub>5</sub>

2

Item

frequenc

I<sub>1</sub>

6

I<sub>2</sub>

7

I<sub>3</sub>

6

I<sub>4</sub>

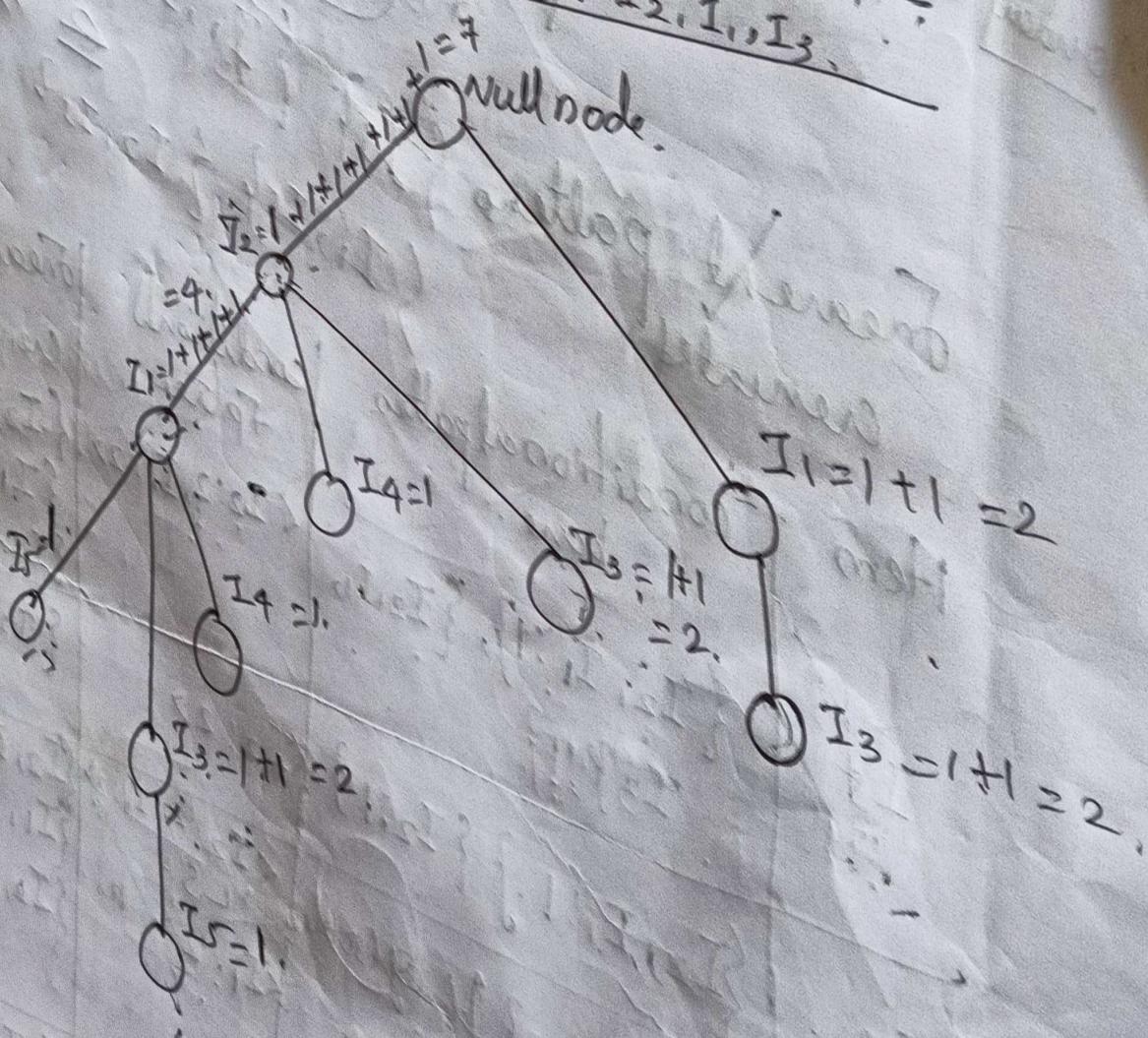
2

I<sub>5</sub>

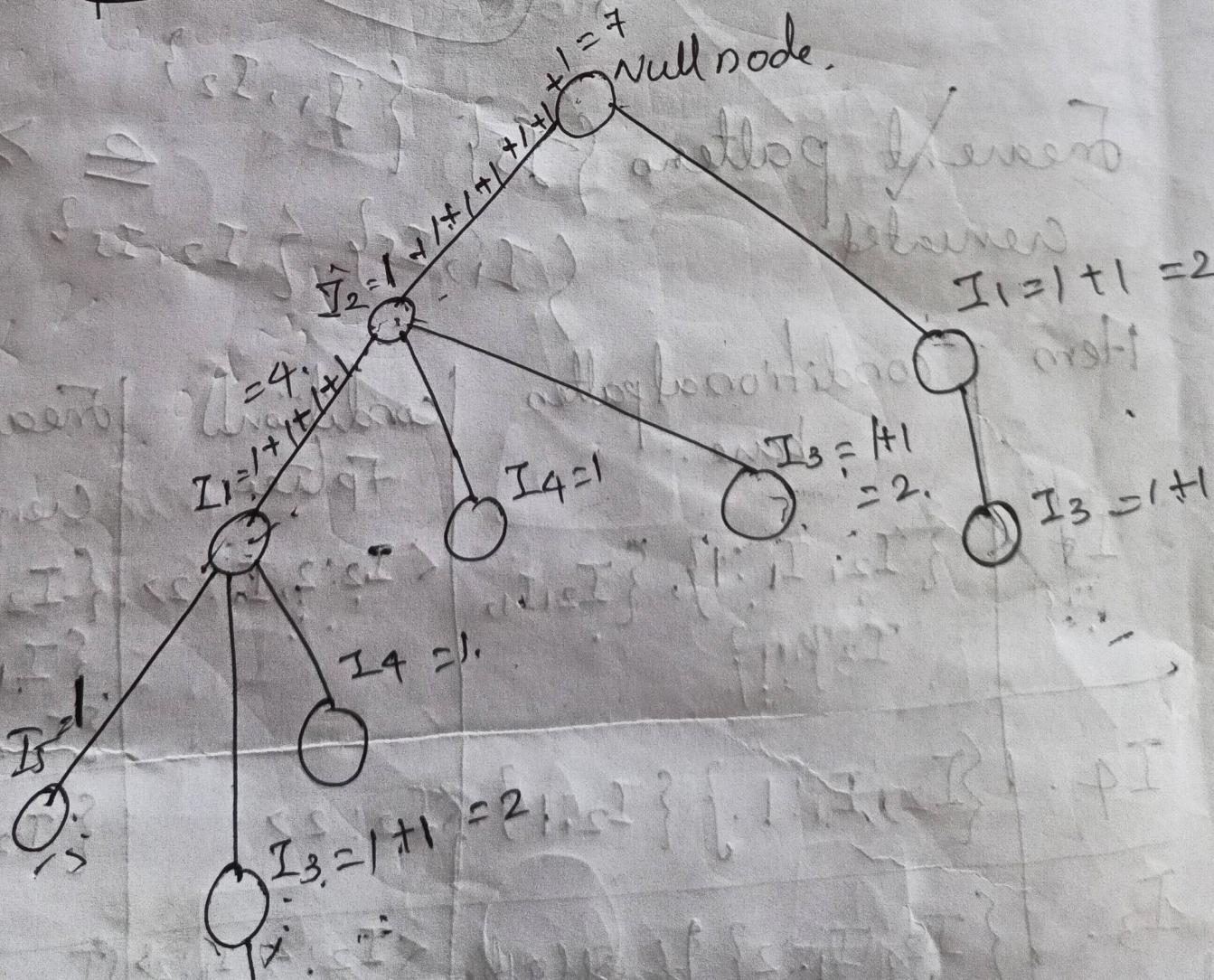
2

2) Find ordered itemset

Transaction	Item-set	Ordered itemset
T <sub>100</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>5</sub>
T <sub>200</sub>	I <sub>2</sub> , I <sub>4</sub>	I <sub>2</sub> , I <sub>4</sub>
T <sub>300</sub>	I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>400</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>4</sub>
T <sub>500</sub>	I <sub>1</sub> , I <sub>3</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>600</sub>	I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>700</sub>	I <sub>1</sub> , I <sub>3</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>800</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub>
T <sub>900</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub>



Transaction	Item-set	ordered itemset
T <sub>100</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>5</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>5</sub>
T <sub>200</sub>	I <sub>2</sub> , I <sub>4</sub>	I <sub>2</sub> , I <sub>4</sub>
T <sub>300</sub>	I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>400</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>4</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>4</sub>
T <sub>500</sub>	I <sub>1</sub> , I <sub>3</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>600</sub>	I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>3</sub>
T <sub>700</sub>	I <sub>1</sub> , I <sub>3</sub>	I <sub>1</sub> , I <sub>3</sub>
T <sub>800</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> , I <sub>5</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> , I <sub>5</sub>
T <sub>900</sub>	I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub>	I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub>



Item	conditional patterns base	conditional FP tree	
I <sub>5</sub>	{I <sub>2</sub> , I <sub>1</sub> : 1}, {I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> : 1}.	< I <sub>1</sub> : 2, I <sub>1</sub> : 2 >	
I <sub>4</sub>	{I <sub>2</sub> , I <sub>1</sub> : 1}, {I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> : 1}	< I <sub>2</sub> : 2, I <sub>4</sub> : 2 >	
I <sub>3</sub>	{I <sub>1</sub> , I <sub>2</sub> : 2}, {I <sub>2</sub> , I <sub>3</sub> : 2}, {I <sub>1</sub> , I <sub>2</sub> , I <sub>3</sub> : 2}.		
I <sub>2</sub>	I <sub>2</sub> : 2 > LS. I <sub>2</sub> : I <sub>1</sub> : 2. I <sub>1</sub> , I <sub>3</sub> : 2 > RS.		
I <sub>1</sub>		{I <sub>1</sub> , I <sub>2</sub> : 2}	
		{I <sub>1</sub> , I <sub>2</sub> } {I <sub>1</sub> , I <sub>3</sub> : 2} > 2	
		{I <sub>1</sub> , I <sub>3</sub> } {I <sub>2</sub> , I <sub>3</sub> : 2} > 2	
Item	conditional patterns base	conditional FP tree	frequent patterns generated.
I <sub>5</sub>	{I <sub>2</sub> , I <sub>1</sub> : 1}, {I <sub>2</sub> , I <sub>1</sub> , I <sub>3</sub> : 1}.	< I <sub>2</sub> : 2, I <sub>1</sub> : 2 >	{I <sub>2</sub> , I <sub>5</sub> : 2}
I <sub>4</sub>	{I <sub>2</sub> , I <sub>1</sub> : 1}, {I <sub>2</sub> : 1}.	KT <sub>2</sub> : 2 >	{I <sub>2</sub> , I <sub>4</sub> : 2}
I <sub>3</sub>	{I <sub>2</sub> , I <sub>1</sub> : 2}, {I <sub>2</sub> , I <sub>3</sub> : 2}, {I <sub>1</sub> , I <sub>3</sub> : 2}.	< I <sub>2</sub> : 4, I <sub>1</sub> : 2 > LS. I <sub>1</sub> : 2 > RS.	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> : 4} {I <sub>1</sub> , I <sub>3</sub> : 4}. {I <sub>4</sub> , I <sub>1</sub> , I <sub>3</sub> : 2}
I <sub>1</sub>	{I <sub>2</sub> : 4}.	< I <sub>2</sub> : 4 >	{I <sub>2</sub> , I <sub>1</sub> : 4}

I<sub>5</sub> & I<sub>3</sub> are highest frequent patterns

2/5/23  
Tuesday

# Dynamic itemset counting algorithm (DIC)

Items are wanted in 4 different ways.

solid box: confirmed frequent itemset.

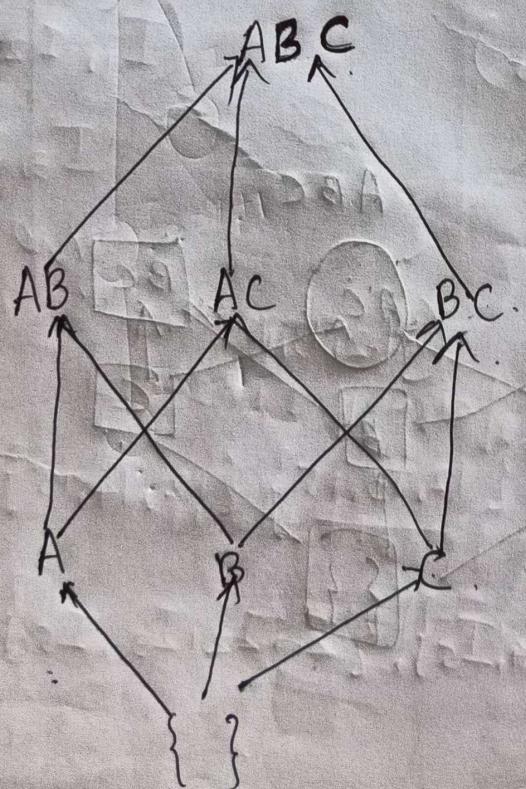
solid circle: confirmed infrequent itemset

dashed box: suspected frequent itemset

dashed circle: suspected infrequent itemset

$$\min \text{Supp} = 25\% \left( S : C = 0.25 \times 4 = \right) : M = 2 : 3$$

TID	A	B	C	
T <sub>1</sub>	1	1	0	AB
T <sub>2</sub>	1	0	0	A
T <sub>3</sub>	0	1	1	BC
T <sub>4</sub>	0	0	0	-

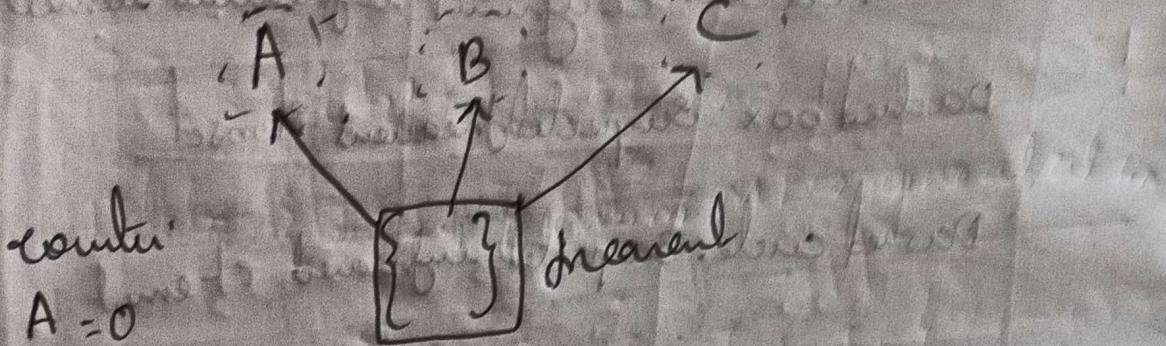


ABC.

AB

AC

BC.



lambdaz

$\lambda_A = 2$

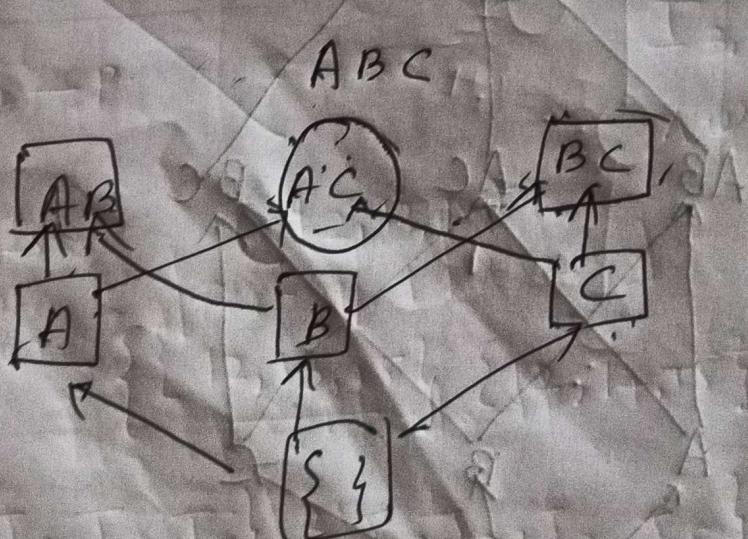
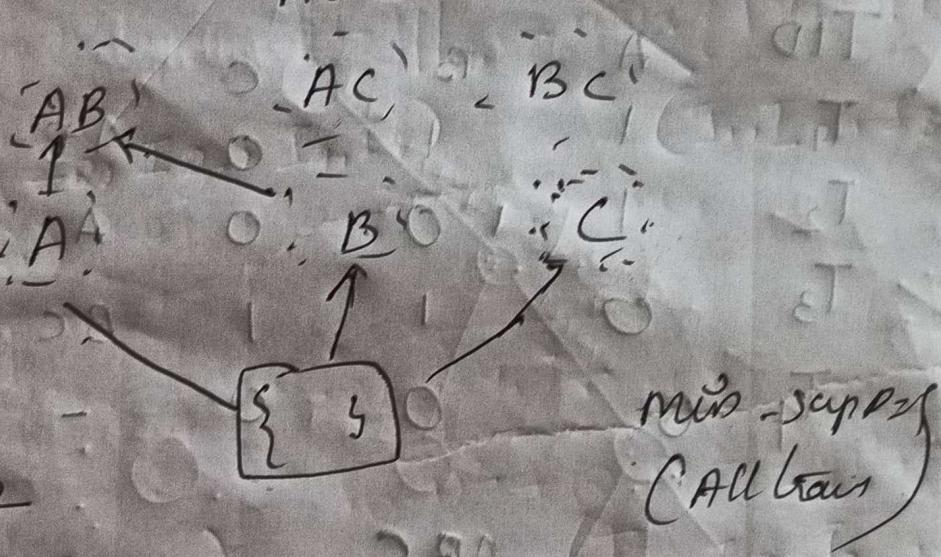
$\lambda_B = 1$

$\lambda_C = 1$

$\lambda_A = 2$

$\lambda_B = 2$

$\lambda_C = 1$ .



3/1/23  
Wednesday Pincer search Algorithm

Bottom up procedure.  
BFS method.  
computations

Item	Count
I <sub>2</sub>	7
I <sub>1</sub>	6
I <sub>3</sub>	6
I <sub>4</sub>	2
I <sub>5</sub>	2

$$C_1 = \{\{I_1\}, \{I_2\}, \{I_3\}, \{I_4\}, \{I_5\}\}$$

$$\text{MFCS} = \{I_1, I_2, I_3, I_4, I_5\}.$$

$$\{I_1, I_2, I_3, I_4, I_5\} - 0$$

LQ:

$$L_1 = \{\{I_1\}, \{I_2\}, \{I_3\}, \{I_4\}, \{I_5\}\}$$

S<sub>1</sub> = ∅. refinement item

Item	Sup-Set	
{I <sub>1</sub> , I <sub>2</sub> }	4	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> }
{I <sub>1</sub> , {I <sub>3</sub> }}	4	{I <sub>2</sub> , {I <sub>4</sub> }} 2
{I <sub>1</sub> , I <sub>4</sub> }	1	{I <sub>2</sub> , I <sub>5</sub> } 2
{I <sub>1</sub> , I <sub>5</sub> }	2	{I <sub>3</sub> , I <sub>4</sub> } 0 —
		{I <sub>3</sub> , I <sub>5</sub> } 1 —
		{I <sub>4</sub> , I <sub>5</sub> } 6 —

Item	Sup-Set	
{I <sub>1</sub> , I <sub>2</sub> }	4	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> }
{I <sub>1</sub> , {I <sub>3</sub> }}	4	{I <sub>2</sub> , {I <sub>4</sub> }} 2
{I <sub>1</sub> , I <sub>4</sub> }	1	{I <sub>2</sub> , I <sub>5</sub> } 2
{I <sub>1</sub> , I <sub>5</sub> }	2	{I <sub>3</sub> , I <sub>4</sub> } 0 —
		{I <sub>3</sub> , I <sub>5</sub> } 1 —
		{I <sub>4</sub> , I <sub>5</sub> } 6 —

Item	Sup-Set	
{I <sub>1</sub> , I <sub>2</sub> }	4	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> }
{I <sub>1</sub> , {I <sub>3</sub> }}	4	{I <sub>2</sub> , {I <sub>4</sub> }} 2
{I <sub>1</sub> , I <sub>4</sub> }	1	{I <sub>2</sub> , I <sub>5</sub> } 2
{I <sub>1</sub> , I <sub>5</sub> }	2	{I <sub>3</sub> , I <sub>4</sub> } 0 —
		{I <sub>3</sub> , I <sub>5</sub> } 1 —
		{I <sub>4</sub> , I <sub>5</sub> } 6 —

Item	Sup-Set	
{I <sub>1</sub> , I <sub>2</sub> }	4	{I <sub>2</sub> , I <sub>3</sub> , I <sub>4</sub> , I <sub>5</sub> }
{I <sub>1</sub> , {I <sub>3</sub> }}	4	{I <sub>2</sub> , {I <sub>4</sub> }} 2
{I <sub>1</sub> , I <sub>4</sub> }	1	{I <sub>2</sub> , I <sub>5</sub> } 2
{I <sub>1</sub> , I <sub>5</sub> }	2	{I <sub>3</sub> , I <sub>4</sub> } 0 —
		{I <sub>3</sub> , I <sub>5</sub> } 1 —
		{I <sub>4</sub> , I <sub>5</sub> } 6 —

Count (MFCs) = 9  
 $L_2 = \{I_1, I_2\} \{I_1, I_3\}; \{I_1, I_5\} \{I_2, I_3\}$   
 $\{I_2, I_4\} \{I_2, I_5\}$

$$S_2 = \{\{I_1, I_4\}\}$$

Delete  $\{I_1, I_4\}$  from MFCs one at a time.

$$\{I_1, I_2, I_3, I_4, I_5\}$$

$$(I_2, I_3, I_4, I_5)$$

$$(I_1, I_2, I_3, I_5).$$

$$(I_2, I_3, I_4, I_5)$$

$$(I_1, I_2, I_3, I_5)$$

$$(I_3, I_4) - P.$$

$I_3, I_4$  - absent  
no change.

$$(I_2, I_4, I_5)$$

$$(I_2, I_3, I_5)$$

$$\text{New MFCs} = \{\{I_2, I_3, I_5\}, \{I_2, I_4, I_5\}, \{I_1, I_2, I_3, I_5\}\}$$

$$= \{\{I_3, I_4, I_5\}, \{I_1, I_2, I_3, I_5\}\}$$

Delete  $\{I_3, I_{s-3}\}$  from each item in mfc's one at a time.

$I_1, I_2, I_3, I_s$

$I_2, I_4, I_5$

$(I_3, I_5) - p.$

$(I_3, I_5) \rightarrow$  absent

$I_3$

$(I_1, I_2, I_3)$

$\rightarrow$  No change

$I_5$

$(I_1, I_2, I_3)$

New MFCs.

$\{ \{I_1, I_2, I_5\}, \{I_1, I_2, I_3\}, \{I_2, I_4, I_5\} \}$

Delete  $\{I_1, I_2, I_5\}$ .

$\{I_1, I_4, I_5\}$

$I_2$

$\{I_4, I_5\}$

$I_5$

$\{I_2, I_5\}$

$I_2, I_5 \subseteq \{I_2, I_4, I_5\}$ .

MFCs.  $\{I_1, I_2, I_5\}, \{I_1, I_2, I_3\}, \{I_2, I_4, I_5\}$

Final O/P.

Stream Path  $\rightarrow \{I_1, I_2, I_5\}$

$\{I_1, I_2, I_3\}$ .

4/5/23

Thursday

clustering

Cluster analysis is the task of grouping set of objects in such a way that objects in the same group are more similar to each other than to those in other groups.

(cluster)

Measures of dissimilarity (2.1.1.1)

Numerical measures of distance between points  
data Euclidean distance =  $\|\bar{x} - \bar{y}\|_2^2 = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$

Squared euclidean distance =  $\|\bar{x} - \bar{y}\|_2^2 = (x_1 - y_1)^2 + \dots + (x_n - y_n)^2$

Manhattan distance =  $\|\bar{x} - \bar{y}\|_1 = |x_1 - y_1| + \dots + |x_n - y_n|$

Maximum distance =  $\|\bar{x} - \bar{y}\|_\infty = \max \{|x_1 - y_1|, \dots, |x_n - y_n|\}$

$$x = (2, 2)$$

$$y = (1, 1)$$

1) Euclidean distance:  $\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$

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

2) squared

3) CB/Manh =  $|x_1 - y_1| + |x_2 - y_2|$

$$= |2-1| + |2-1|$$

$$= \underline{\underline{1+1}} = 2$$

$$\text{Manhattan} = \max \{ |x_1 - y_1|, |x_2 - y_2| \}$$

$$= \max \{ |x_1 - 1|, |x_2 - 1| \} = \max \{ l_1, l_2 \}$$

$$x = (2, 1, 3)$$

$$y = (2, 0, 1)$$

$$d \in D = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

Minkowski distance.

Non-numerical data

Levenshtein distance.

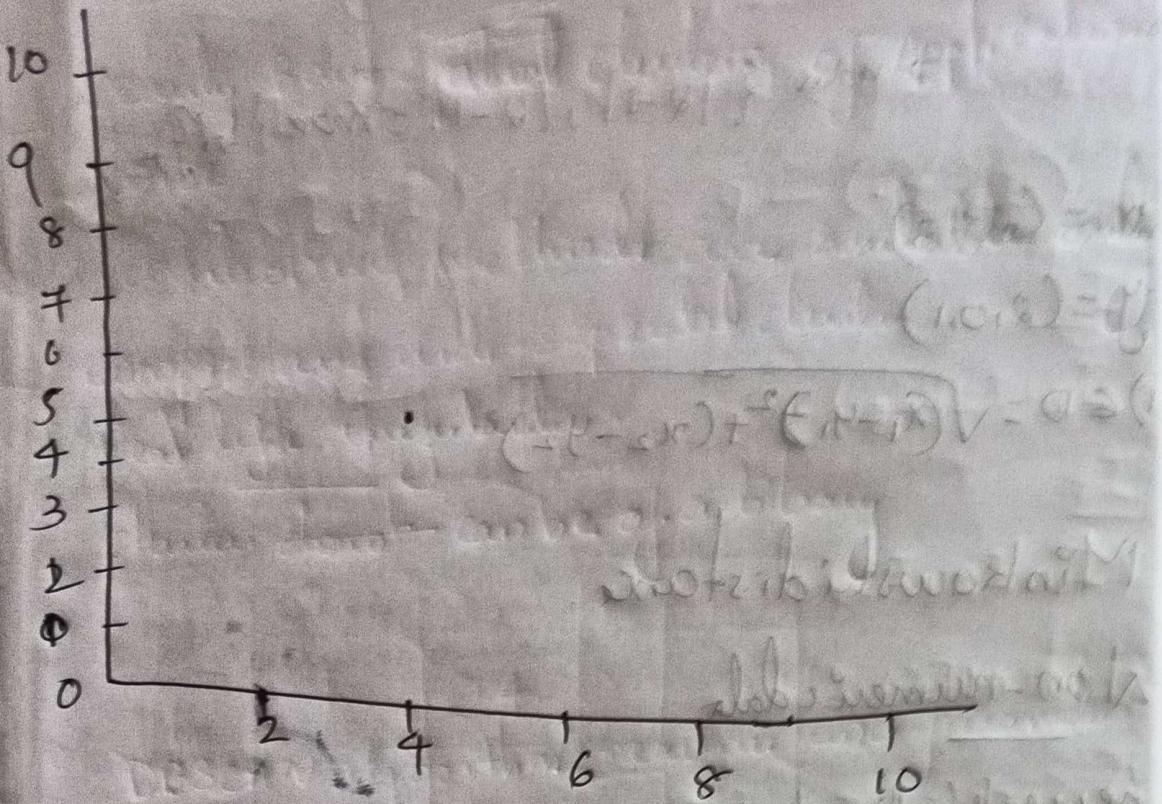
Clustering methods

Hierarchical partition

PAM - Partition Around Medoid Algorithm

	X	Y	Z
0	8	7	8
1	3	7	4
2	4	9	5
3	9	6	7
4	8	5	4
5	5	8	3
6	7	3	4
7	8	4	5
8	7	5	4
9	4	6	6

$$\begin{aligned}
 & |x_1 - y_1| + |x_2 - y_2| \\
 & |8 - 7| + |8 - 4| \\
 & 2 + 3 = 5
 \end{aligned}$$



$C_1(4,5)$ ,  $C_2(8,5)$ .  
calculating cost.

	X	Y	Dissimilarity from $C_1$	Distance from $C_2$
$C_2$	8	7	6	2
$C_1$	3	7	3	1
	4	9	4	8
			6	✓ 2.
			4	6
			5	3
			5	1
			- 38	

Manhattan

$(8,7)$  from  $(4,5)$

$$D = |8-4| + |7-5| =$$

$$\cancel{4} + 2 = 6$$

$(8,7)$  from  $(8,5)$

$$C_1 \rightarrow 2+2+4 = 8$$

$$C_2 \rightarrow 3+\cancel{2} + 1 + 1 + \cancel{4} =$$

$$= \overline{\overline{2}}$$

$$\text{Total} = 20 -$$

new medond ( $C_{8,4}$ )

## DBSCAN

density Based spatial clustering of Applications with noise.

Partitioning methods are suitable only for compact and well-separated clusters.

Real life irregular

clusters ~~can~~ can be of arbitrary.

DBSCAN algorithm requires two parameters.

1.  $\text{eps}(\epsilon)$ : it defines

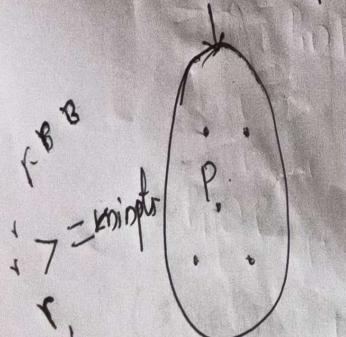
2.  $\text{Minpts}$

3 types of datapoints.

Core point : A point is a core point if it has more than  $\text{Minpts}$  points with  $\epsilon$  dist.

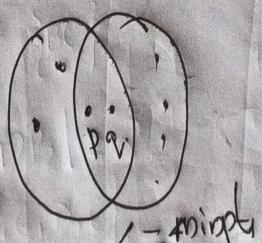
Borderpoint : L.

None point



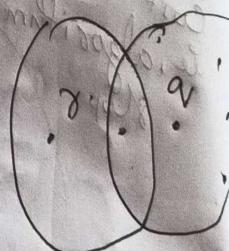
het  $\text{Minpts} = 4$

P - Corepoint

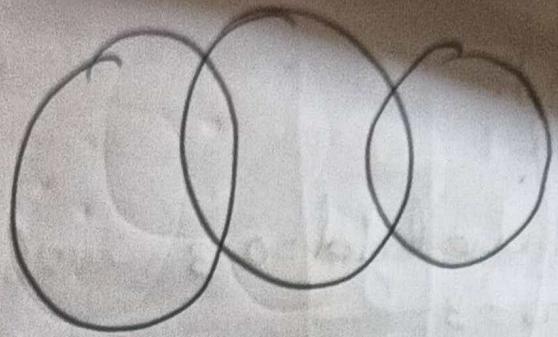


het  $\text{Minpts} = 4$   
P - borderpoint

Q - Corepoint



neither Core nor  
border point  
 $\epsilon$  - noise point



Dates  
Palm  
Puddles

density based - Dbscan  
Partitions - palms

## Categorical clustering - Rock

Robert clustering using links.

Rock belongs to the class of Agglomerative hierarchical clustering Algorithms.

Rock works for categorical attributes

Data  $\rightarrow$  draw random sample  $\rightarrow$  cluster with links.  
neighbors

$$\text{sim}(P_i, P_j) \geq \theta$$

Jaccard coefficient for  $\text{sim}(T_1, T_2)$ ,

$$\text{sim}(T_1, T_2) = \frac{|T_1 \cap T_2|}{|T_1 \cup T_2|} \rightarrow \text{coefficient.}$$

links  $(P_i, P_j)$  no: of common neighbors.

Good new mean

$$g(C_i, C_j) = \frac{\text{link}(C_i, C_j)}{(n_i + n_j) + 2f(\theta) - n_i(1 + 2f(\theta)) - n_j(1 + 2f(\theta))} \cdot \frac{e_{H_2} + e_{H_1}}{e_{H_2} + e_{H_1}}$$

$$P_1 = \{A, B, C, D\}$$

$$P_2 = \{E, B, C\}$$

$$P_3 = \{D, E, B\}$$

$$P_4 = \{E, C, F\}$$

Similarity threshold = 0.3

No. of clusters = 3.

i) Similarity table.

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	1	0.4	0.4	0.17
$P_2$		1	0.5	0.8
$P_3$			1	0.2
$P_4$				1

$$\text{Sim}(P_i, P_j) = \frac{|P_i \cap P_j|}{|P_i \cup P_j|}$$

$$P_1, P_2$$

$$P_1 \cap P_2 = \{B, C\}$$

$$|P_1 \cap P_2| = 2$$

$$P_1 \cup P_2 = \{A, B, C, D, E\}$$

$$|P_1 \cup P_2| = 5$$

$$\text{Sim}(P_i, P_j) = \frac{2}{5}$$

## Adjacency Table

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	1	1	1	0
$P_2$		1	1	1
$P_3$			1	0
$P_4$				1

$$\text{Sim}(P_1, P_2) = 0.4$$

6.470.3

$$\text{Sim}(P_1, P_4) = 0.17$$

0.17 < 0.3  
= 0.

If  $\delta_m(p_i, p_j) \geq \text{threshold} = 1$ .

$$11 \quad 11 \quad \leq \quad 11 \quad = 0.$$

$$NP_1 = 1$$

$$n p_2 = 1$$

3. No. of links / common neighbours.

	$P_1$	$P_2$	$P_3$	$P_4$
$P_1$	$\phi -$	3	3	1
$P_2$		+	3	2
$P_3$			+	<del>4.1</del>
$P_4$				-

No: of links ( $p_i, p_j$ )

- Actual x Adjusted

$$\begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & -1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$= \quad \boxed{2}$$

4) Croodness Measure -

$$g(c_i, c_j) = \text{link}[c_i, c_j]$$

$$g(p_1, p_2) = \frac{\text{link}[p_1, p_2]}{\frac{(p_1 + p_2)^{1+2f(\theta)}}{1+2f(\theta)} - p_1} - \frac{1+2f(\theta)}{1+2f(\theta)} - p_2$$

$$\begin{aligned} f(\theta) &= \frac{1-\theta}{1+\theta} = \frac{1-0.3}{1+0.3} = \underline{\underline{0.57}} \\ &= \frac{3}{2^{1+2 \times 0.57} - 1} = \frac{1}{2^{1+2 \times 0.57} - 1} = \underline{\underline{1.35}} \end{aligned}$$

$$= \underline{\underline{1.35}}$$

Pair	CM
$P_1, P_2$	1.35 -
$P_1, P_3$	1.35 -
$P_1, P_4$	0.45
$P_2, P_3$	1.35 -
$P_2, P_4$	0.90
$P_3, P_4$	0.45 .1

highest causal value for 3 pairs.

cluster  $P_1, P_2$

new hink table.

old hink table  $\rightarrow$  new hink.

	$P_1, P_2$	$P_3$	$P_4$
$P_1, P_2$	-	6	3
$P_3$		-	1
$P_4$			-

$$((P_1, P_2), P_3) \rightarrow (P_1, P_2) (P_2, P_3)$$

$$\text{hink}[(P_1, P_2), P_3] = \text{hink}[P_1, P_2] + \text{hink}[P_2, P_3]$$

$$= \underline{\underline{3+3=6}}$$

$n_i$  = no. of points in  $[P_1, P_2]$

$$= 2.$$

$n_j$  = no. of points in  $[P_3] = \underline{\underline{1}}$