

(iliski kuralları  $\Rightarrow$ )

## \* What is Association Rule mining?

Given a set of transactions, find rules that will predict occurrence of an item based on occurrences of other items in a transaction.

TID	Items
1	Bread, Peanuts, Milk, Fruit, Jam
2	Bread, Jam, Soda, Chips, Milk, Fruit
3	Steak, Jam, Soda, Chips, <u>Bread</u>
4	Jam, Soda, Peanuts, Milk, Fruit
5	Jam, Soda, Chips, Milk, <u>Bread</u>
6	Fruit, Soda, Chips, Milk
7	Fruit, Soda, Peanuts, Milk
8	Fruit, Peanuts, Cheese, Yogurt

→ ex of ass. rules;

$$\{ \text{Bread} \} \xrightarrow{x} \{ \text{Milk} \}$$

$\star \{ \text{Bread} \}$  and  $\{ \text{Milk} \}$  are itemsets

$$\{ \text{Soda} \} \xrightarrow{y} \{ \text{Chips} \}$$

## \* Two evaluation Metrics;

① **Support:** fraction of transactions that contain both X and Y

$$S = \frac{s(\{\text{Bread}, \text{Milk}\})}{\# \text{ of transactions}} = \frac{3}{8}$$

$\star$  Sırası önemli değil  $\text{Milk} \Rightarrow \text{Bread}$  de olur.

② **Confidence:** Measures how often items in Y appear in transaction that contain X

$$C = \frac{s(\{\text{Bread}, \text{Milk}\})}{s(\{\text{Bread}\})} = \frac{3}{4}$$

## \* Association Rule Mining

### \* What are Frequent Itemset?

- An itemset whose support  $\geq$  minsup threshold

\* Given a set of transactions, goal of association rule mining

is to find all rules having; support  $\geq$  minsup and confidence  $\geq$  minconf

## \* Brute-Force Algorithm → computationally prohibitive = (

- List all possible association rules

- Compute support and confidence for each rule

- Prune (bypass) rules that fail the minsup and minconf thresholds.

- Pseudocode:

## Apriori Algorithm ↗ level-wise approach

- If an itemset is frequent, all its subsets must be frequent also
- $\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$  → anti-monotone property

Given the following database and a min support of 3, generate all the frequent itemsets

D	A	B	C	D	E
1	1	1	0	1	1
2	0	1	1	0	1
3	1	1	0	1	1
4	1	1	1	0	1
5	1	1	1	1	1
6	0	1	1	1	0

(a) Binary database

t	i(t)
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

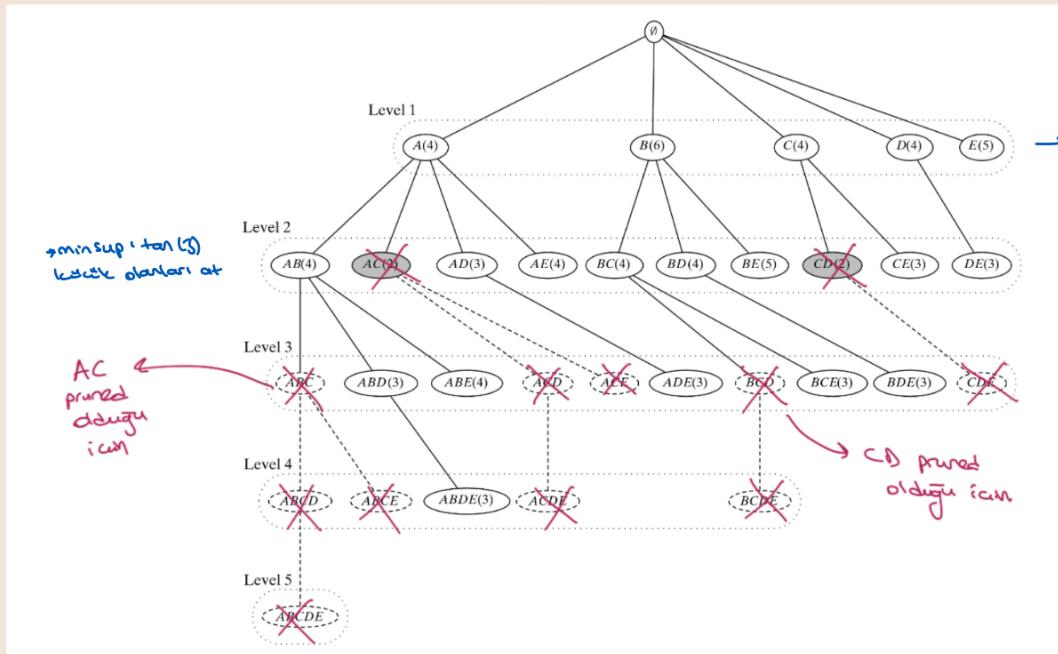
(b) Transaction database

x	A	B	C	D	E
1	1	1	2	1	1
2	3	2	4	3	2
3	4	3	5	3	3
4	5	4	6	6	4
5	5	5	6	5	5
6	6	6	6	6	6

(c) Vertical database

↗ 3 farklı gösterim

→ her lorsunda  
kaçar  
adet  
olduğu  
yazıyor



- Pseudo code:

APRIORI(D, I, minsup):

```

F ← ∅ // set of freq. itemsets
C' ← {∅} // initial prefix tree with single item
for each i ∈ I do
    Add i as child of ∅ in C' with sup(i) ← 0
k ← 1
while C^k ≠ ∅ do
    COMPUTESUPPORT(C^k, D)
    for each leaf X ∈ C^k do
        if sup(X) > minsup then F ← F ∪ {X, sup(X)}
        else remove X from C^k
    C^k+1 ← EXTENDPREFIXTREE(C^k)
    k ← k+1
return F
  
```

$k = i(t)$ 'nın uzunluğu  
; ( $i(t) = 4$ )

COMPUTESUPPORT(X, D):

```

for each  $\langle t, i(t) \rangle \in D$  do
    for each k-subset  $X \subseteq i(t)$  do
        if  $X \in C^k$  then  $\text{sup}(X) \leftarrow \text{sup}(X) + 1$ 
  
```

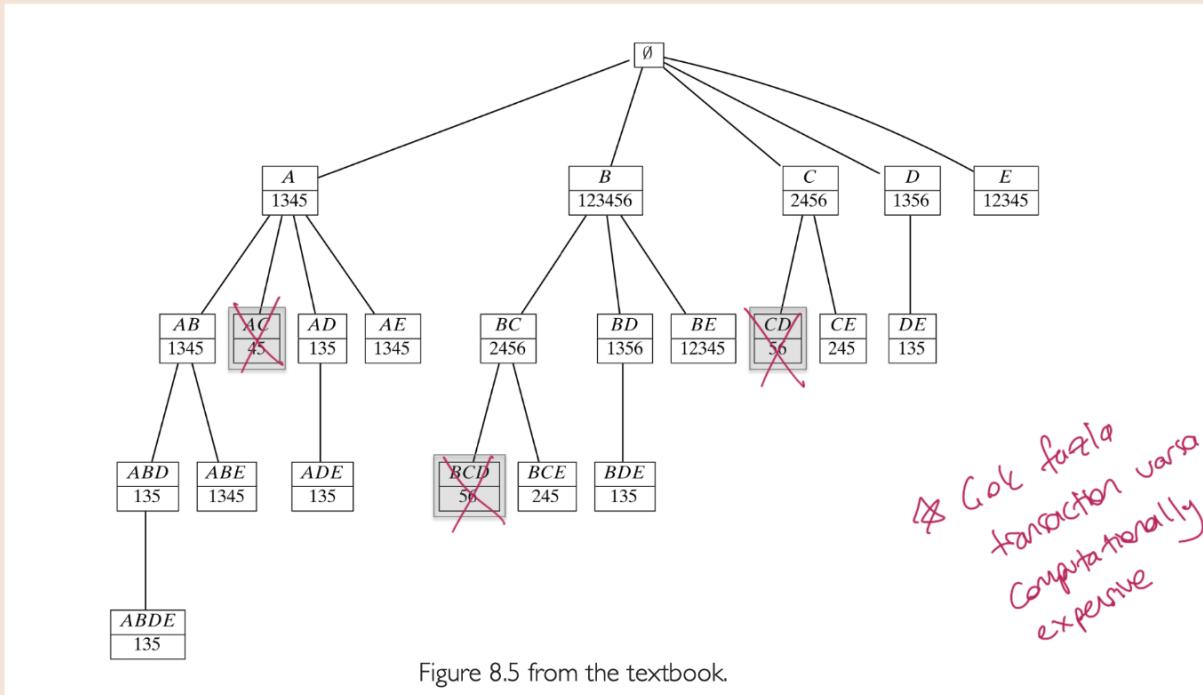
EXTENDPREFIXTREE(C^k):

```

for each leaf  $X_a \in C^k$  do
    for each leaf  $X_b \in \text{SIBLING}(X_a)$  such that  $b > a$  do
         $X_{ab} \leftarrow X_a \cup X_b$ 
        if  $X_j \in C^k$  for all  $X_j \subseteq X_{ab}$ , such that  $|X_j| = |X_{ab}| - 1$  then
            add  $X_{ab}$  as child of  $X_a$  with  $\text{sup}(X_{ab}) \leftarrow 0$ 
        if no extensions from  $X_a$  then
            remove  $X_a$ , and all various ancestors of  $X_a$  with no extensions, from  $C^k$ 
  
```

## Eclat Algorithm:

- Leverages the tidssets (subset of itemset) directly for support computation
- Itemset support is computed by intersecting tidssets of suitably chosen subsets.
- Given  $t(x)$  and  $t(y)$  (for frequent itemsets  $X$  and  $Y$ )  $\Rightarrow t(XY) = t(X) \wedge t(Y)$ ,  $sup(XY) = |t(XY)|$



### - Pseudocode:

#### Algorithm 8.3: Algorithm ECLAT

```

// Initial Call:  $\mathcal{F} \leftarrow \emptyset, P \leftarrow \{(i, t(i)) \mid i \in \mathcal{I}, |t(i)| \geq \text{minsup}\}$ 
ECLAT( $P, \text{minsup}, \mathcal{F}$ ):
1 foreach  $\langle X_a, t(X_a) \rangle \in P$  do
2    $\mathcal{F} \leftarrow \mathcal{F} \cup \{(X_a, \text{sup}(X_a))\}$ 
3    $P_a \leftarrow \emptyset$ 
4   foreach  $\langle X_b, t(X_b) \rangle \in P$ , with  $X_b > X_a$  do
5      $X_{ab} = X_a \cup X_b$ 
6      $t(X_{ab}) = t(X_a) \cap t(X_b)$ 
7     if  $\text{sup}(X_{ab}) \geq \text{minsup}$  then
8        $P_a \leftarrow P_a \cup \{(X_{ab}, t(X_{ab}))\}$ 
9   if  $P_a \neq \emptyset$  then  $\text{ECLAT}(P_a, \text{minsup}, \mathcal{F})$ 

```

## ⇒ FP Growth Algorithm (Frequent Patterns Tree)

### - Core of Apriori Algorithm

- use frequent  $(k-1)$ -itemsets to generate candidate freq.  $k$ -itemsets
- use database scan and pattern matching to collect counts for candidate itemsets

(burada  $k=$ level, yani  $10^4$  frequent 1-itemset,  $10^7$ tane candidate 2-itemset generate edereli)

### - Compress a large database into a compact, FP-tree structure

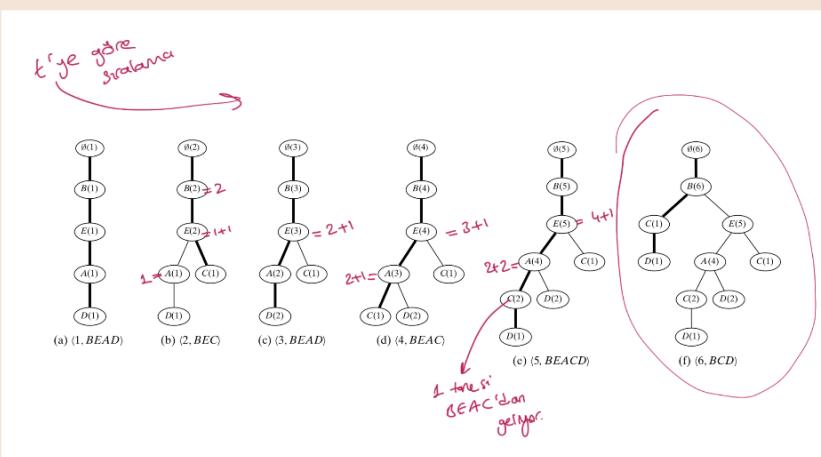
### - Major Steps to mine FP-tree

- Construct the FP-tree
- For each frequent item  $i$ , compute projected FP-tree
- Recursively mine conditional FP-trees and grow frequent patterns obtained
- If the conditional FP-tree contains a single path, enumerate all patterns

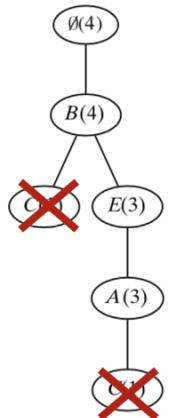
$t$	$I(t)$
1	ABDE
2	BCE
3	ABDE
4	ABCE
5	ABCDE
6	BCD

Transactions  
BEAD  
BEC  
BEAD  
BEAC  
BEACD  
BCD

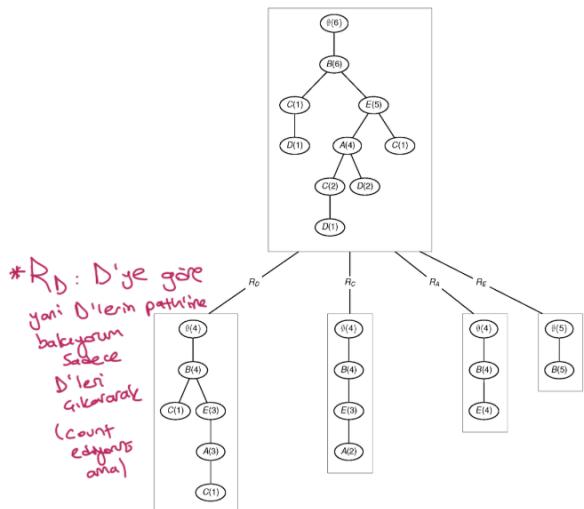
→ sort all transactions based on their frequency (B en çok bolunur)



After removing the infrequent item C (which has support 2), the resulting FP-Tree has only one path B(4)-E(3)-A(3)



We enumerate all the subsets of the path and prefix them with D to obtain the frequent itemsets DB(4), DE(4), DA(3), DBE(3), DBA(3), DEA(3), DBEA(3)



All frequent pattern tree projections. Figure 8.9 of the textbook.

it may need to generate huge candidate sets  
solution

## - Pseudo code:

```

Algorithm 8.5: Algorithm FPGROWTH
=====
// Initial Call: R ← FP-tree(D), P ← ∅, F ← ∅
FPGROWTH(R, P, F, minsup):
1 Remove infrequent items from R
2 if ISPATH(R) then // insert subsets of R into F
3   foreach Y ⊆ R do
4     | X ← P ∪ Y
5     | sup(X) ← minx ∈ Y{cnt(x)}
6     | F ← F ∪ {(X, sup(X))}
7 else // process projected FP-trees for each frequent item i
8   foreach i ∈ R in increasing order of sup(i) do
9     | X ← P ∪ {i}
10    | sup(X) ← sup(i) // sum of cnt(i) for all nodes labeled i
11    | F ← F ∪ {(X, sup(X))}
12    | RX ← ∅ // projected FP-tree for X
13    | foreach path ∈ PATHFROMROOT(i) do
14      | | cnt(i) ← count of i in path
15      | | Insert path, excluding i, into FP-tree RX with count cnt(i)
16    | if RX ≠ ∅ then FPGROWTH(RX, X, F, minsup)

```

## - Benefits of FP - Structure:

- Preserve complete info for frequent pattern mining
- Reduce irrelevant info; infrequent items are gone
- More frequently occurring items → more likely to be shared
- Never be larger than original database
- No candidate generation, no candidate test, no repeated scan of entire database.

## \* Rule Generation

- Given a frequent itemset L, find all non-empty subsets f ⊂ L such that L \ f satisfies the minimum confidence requirement

ex: {A, B, C, D} is a frequent itemset, candidate rules:

$$k = |L| = 4$$

$$\begin{array}{l} ABC \rightarrow D, AC \rightarrow B, BC \rightarrow A, ABD \rightarrow C, AD \rightarrow BC, BD \rightarrow AC, AB \rightarrow CD, \\ \times 2(\text{tersi}) \quad \times 2(\text{tersi}) \end{array} ] 2^4 - 2 = 14 \quad \begin{matrix} \text{association} \\ \text{rules} \end{matrix}$$

or  $2^k - 2$

- Efficient rule generation: Confidence have anti-monotone property  
if rules are generated from some itemset ]  $c(ABC \rightarrow D) \geq c(AB \rightarrow D) \geq c(A \rightarrow BCD)$   
hepsi L' in içinde

## \* Rule Assessment Measures

• How to set minsup?

- if minsup too high; we could miss itemsets containing rare items
- if minsup too low; # of itemsets very large; Computationally expensive!
- a single minsup threshold may not be effective

• Lift: ratio of the observed joint probability of X and Y to the expected joint prob. (if they were statistically independent)

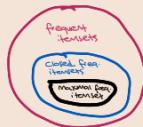
$$\text{lift}(X \Rightarrow Y) = \frac{\text{sup}(X \cup Y)}{\text{sup}(X) \cdot \text{sup}(Y)} = \frac{\text{conf}(X \Rightarrow Y)}{\text{sup}(Y)} \rightarrow \begin{matrix} \text{if } = 1 \text{ then} \\ X \text{ and } Y \text{ statistically independent} \end{matrix}$$

- Lift is measure of deviation

- Lift measures the surprise of rule (if ~1; support of rule is expected)

- Should be much larger or much smaller than 1

## \* Summarizing Itemsets



\* NOT: AB, A' is "super set"

- **Frequent Itemset:** An itemset, whose support is greater than some user-specific minimum support
- **Closed Frequent Itemset:** An itemset  $X$  is closed if all supersets of  $X$  have strictly less support  
 $\text{sup}(X) > \text{sup}(Y)$  for all  $X \subset Y$
- **Maximal Frequent Itemset:** An itemset  $X$  is maximal if has no frequent supersets  $\rightarrow$  en attables  
 $M = \{X \mid X \in F \text{ and } \exists Y \supseteq X \text{ such that } Y \notin F\}$
- **Minimal Generators:** A frequent itemset is a minimal generator if it has no subsets with same support

## \* Mining Frequent Sequences

- **Terminology:**  $s = ACTGAAACG$ 
  - $r_1 = CGAACG$  subsequence of  $s$
  - $r_2 = CTGA$  consecutive subseq. of  $s$
  - $r_3 = ACT$  prefix of  $s$
  - $r_4 = GAAACG$  suffix of  $s$

- Support of a sequence;  $\text{sup}(r) = |\{s_i \in D \mid r \subseteq s_i\}|$

- Relative support;  $r_{\text{sup}}(r) = \text{sup}(r) / N$   $\rightarrow$  % of seq.

-  $r$  is frequent if  $\text{sup}(r) \geq \text{minsup}$

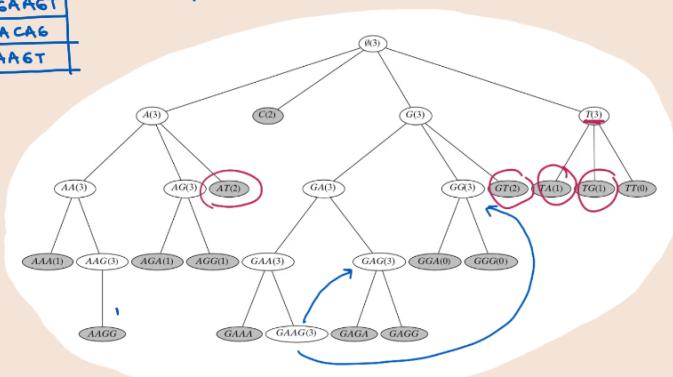
## → GSP Algorithm

- It searches the sequence prefix tree using a level-wise

ex:

k	Sequence
$s_1$	CAGAAGT
$s_2$	TGACAG
$s_3$	GAAGT

$\text{minsup} = 3$



⇒ Given set of freq. sequences at level  $k$ , generate candidates for level  $k+1$ , and compute support, prune not frequent ones.

⇒ For each leaf  $r_a$ , sequence "extended" with last symbol of other leaf  $r_b$  that shares same prefix. So  $r_{ab} = r_a + r_b [k]$

## - Pseudocode:

### Algorithm 10.1: Algorithm GSP

**GSP ( $D, \Sigma, \text{minsup}$ ):**

- $\mathcal{F} \leftarrow \emptyset$
- $\mathcal{C}^{(1)} \leftarrow \{\emptyset\}$  // Initial prefix tree with single symbols
- foreach  $s \in \Sigma$  do Add  $s$  as child of  $\emptyset$  in  $\mathcal{C}^{(1)}$  with  $\text{sup}(s) \leftarrow 0$
- $k \leftarrow 1$  //  $k$  denotes the level
- while  $\mathcal{C}^{(k)} \neq \emptyset$  do
  - COMPUTESUPPORT ( $\mathcal{C}^{(k)}, D$ )
  - foreach leaf  $r \in \mathcal{C}^{(k)}$  do
    - if  $\text{sup}(r) \geq \text{minsup}$  then  $\mathcal{F} \leftarrow \mathcal{F} \cup \{(r, \text{sup}(r))\}$
    - else remove  $r$  from  $\mathcal{C}^{(k)}$
  - $\mathcal{C}^{(k+1)} \leftarrow \text{EXTENDPREFIXTREE} (\mathcal{C}^{(k)})$
  - $k \leftarrow k + 1$
- return  $\mathcal{F}^{(k)}$

**COMPUTESUPPORT ( $\mathcal{C}^{(k)}, D$ ):**

- foreach  $s_i \in D$  do
  - foreach  $r \in \mathcal{C}^{(k)}$  do
    - if  $r \subseteq s_i$  then  $\text{sup}(r) \leftarrow \text{sup}(r) + 1$

**EXTENDPREFIXTREE ( $\mathcal{C}^{(k)}$ ):**

- foreach leaf  $r_a \in \mathcal{C}^{(k)}$  do
  - foreach leaf  $r_b \in \text{CHILDREN}(\text{PARENT}(r_a))$  do
    - $r_{ab} \leftarrow r_a + r_b [k]$  // extend  $r_a$  with last item of  $r_b$ 
      - // prune if there are any infrequent subsequences
      - if  $r_c \in \mathcal{C}^{(k)}$ , for all  $r_c \subset r_{ab}$ , such that  $|r_c| = |r_{ab}| - 1$  then
        - Add  $r_{ab}$  as child of  $r_a$  with  $\text{sup}(r_{ab}) \leftarrow 0$
      - if no extensions from  $r_a$  then
        - remove  $r_a$ , and all ancestors of  $r_a$  with no extensions, from  $\mathcal{C}^{(k)}$
- return  $\mathcal{C}^{(k)}$

## → Sequential Pattern Mining

- In several application we deal with sequence of itemsets (not items)
- Association rules do not consider order of transaction but in many application, ordering is important. (Buying bed  $\Rightarrow$  buying bed sheets)
- In Sequential pattern Mining, a sequence is an ordered list of itemsets  
And it consist of categorical data. When numerical  $\rightarrow$  apply time series

Ex:

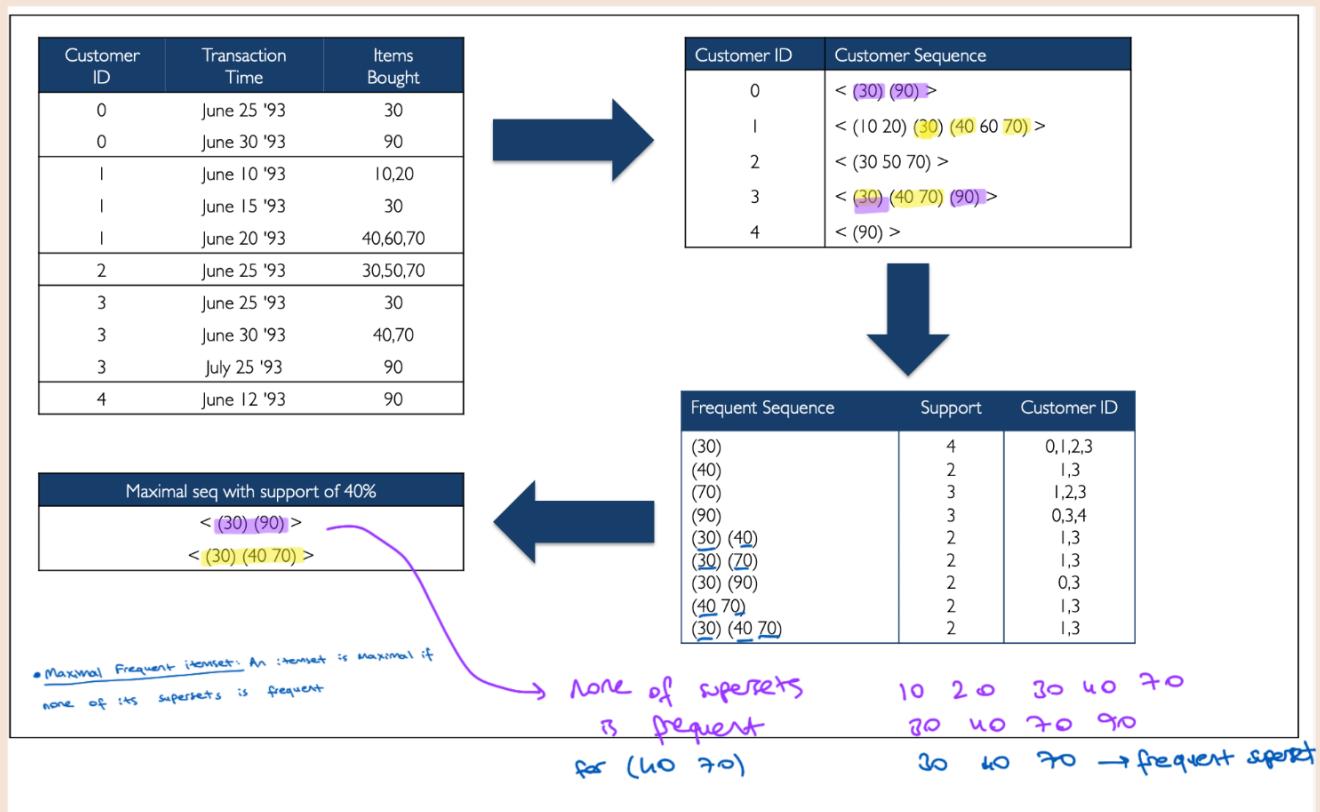
$$I = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

sequence  $\langle \{3\}, \{4, 5\}, \{8\} \rangle$  is subseq. of  $\langle \{6\}, \{3, 7\}, \{9\}, \{4, 5, 8\}, \{3, 8\} \rangle$

$\langle \{3\}, \{8\} \rangle$  is not contained in  $\langle \{3, 8\} \rangle$  or vice-versa

Size of sequence  $\langle \{3\}, \{4, 5\}, \{8\} \rangle$  is 3, length is 4

- Goal: is to find all the sequences that have user-specified min. support (frequent sequences)
- $\text{Sup}(r) = |\{s_i : r \text{ is a subsequence of } s_i\}| \quad \rightarrow \text{rsup}(r) = \text{sup}(r) / N$



- Level-by-level:
  - 1) Generate candidate sequences
  - 2) Scan database to collect support counts
  - 3) Use Apriori property to prune candidates

- Limitations:
  - It can generate massive amount of candidates
  - It requires many scans of the database