# Association Rule Mining

Market basket analysis : to find associations between different itemsets that customers place in shopping basket

Applications : Cross marketing, Catalog/floor design, design attractive packs, web log analysis for e-commerce,

Rule generation :

$$\text{Antecedent} \Rightarrow \text{Consequent (support, Confidence)}$$

Data structure :

$$T = \{t_1, \ldots, t_n\} \text{ set of transcations}$$

Each $t_K$ is an itemset

$$I = \{i_1, \ldots, i_m\}$$

Typical data

| Cid | A1 | A2 | A3 | . | . | . | . | . | . | . | AK |
|-----|----|----|----|---|---|---|---|---|---|---|----|
| $C_1$ | 0 | 0 | 1 | 0 | 1 | 1 | - | - | - | - | 0 |
| $C_2$ | 1 | 0 | 0 | 1 | 1 | 0 | - | - | - | - | 1 |
| ⋮ | | | | | | | | | | | |
| $C_n$ | - | - | - | - | - | - | - | - | | | |

Aim : • Find frequent patterns, i.e. associations among sets of items in T

• Represent these relationships as association rules of the form

$$X \Rightarrow Y \text{ (support, Confidence)}$$

Important definitions:

Support count : # of occurrences of an itemset in the database

$$\sigma(\{itemset\})$$

Support : Fraction of transactions containing the itemset

$$S(itemset) = \frac{\sigma(itemset)}{|T|}$$

Frequent itemset : An itemset whose support
$$\geq a \text{ threshold, minsup}$$

Confidence :. A measure of how often B appears in
[Rule: $(A \Rightarrow B)$]
transactions containing A.

% of transactions containing A which also contains B

$$C(A \Rightarrow B) = \frac{S(A, B)}{S(A)} \quad (\text{est of conditional prob})_{B|A}$$

Example:

| Tid | items |
|-----|-------|
| T1 | bread, egg, peanut-butter |
| T2 | bread, peanut-butter |
| T3 | bread, milk, peanut-butter |
| T4 | beer, bread |
| T5 | beer, milk |

| Rule | Support | Confidence |
|------|---------|-----------|
| bread $\Rightarrow$ peanut-butter | $\frac{3}{5} = 0.6$ | $\frac{3}{4} = 0.75$ |
| peanut-butter $\Rightarrow$ bread | $\frac{3}{5} = 0.6$ | $\frac{3}{3} = 1.0$ |
| beer $\Rightarrow$ bread | $\frac{1}{5} = 0.2$ | $\frac{1}{2} = 0.5$ |
| peanut-butter $\Rightarrow$ egg | $\frac{1}{5} = 0.2$ | $\frac{1}{3} = 0.33$ |
| egg $\Rightarrow$ peanut-butter | $\frac{1}{5} = 0.2$ | $\frac{1}{1} = 1$ |
| egg $\Rightarrow$ milk | 0 | 0 |

**ARM task:** Given a set of transactions, the goal of ARM is to find all rules $\Rightarrow$

    (i) support $\geq$ min sup

    & (ii) confidence $\geq$ min conf

(min sup, min conf) : fixed apriori

Brute force approach
- List all possible association rules
- compute support & confidence for each rule
- prune as per threshold

Approach is computationally prohibitive

## The Apriori Algorithm

Agrawal & Srikant : "Fast algorithms for mining association rules in large databases", Int Conf on VLDB, 1994.

## 2-step ARM of Apriori algorithm

S1: Generate all frequent itemsets with support $\geq$ min sup

S2: Generate association rules using these frequent itemsets

# Anti-monotonicity property of Apriori algorithm

## Downward closure property

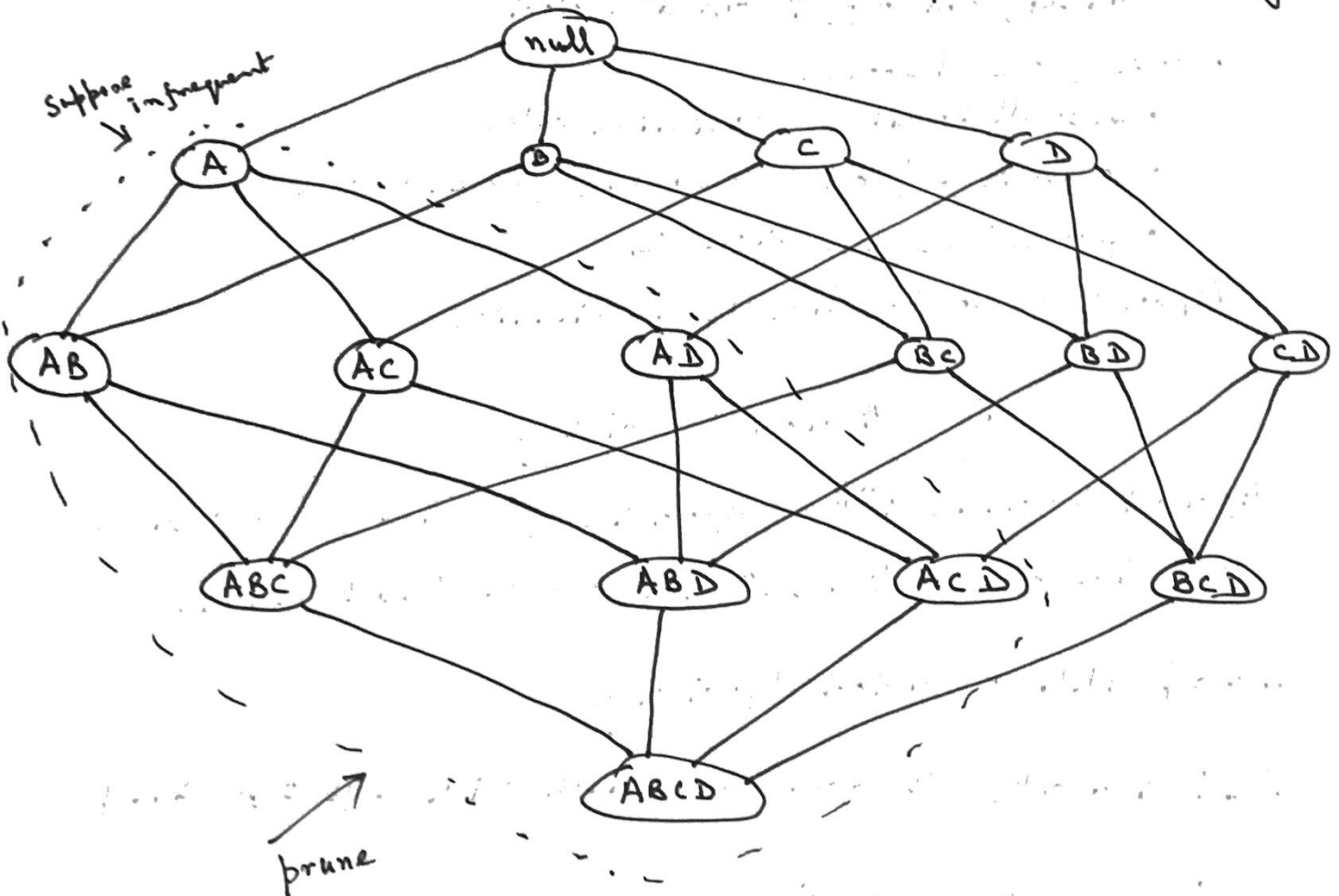- Any subset of a frequent itemset is frequent

$$\text{i.e. } \forall \; X \subseteq Y$$

$$S(X) \geq S(Y) \quad \forall \; X, Y$$

⟹ If an itemset is not frequent, none of it's supersets can be frequent (⟹ prune all such sets)

⟹ If an itemset is not frequent, there is no need to explore its supersets.

<u>Example</u> : Itemset lattice → apriori principle of pruning

- Suppose B is found to be in frequent (support < min sup)

  ⇒ all itemsets with B will also be in frequent

  ⇒ prune all such branches. (i.e. all it's supersets)

    i.e. exclude Itemset AB, BC, BD, ABC, ABD,
    
    BCD, ABCD

- suppose AB is found to be in frequent

  ⇒ all Itemsets with AB will also be in frequent

  ⇒ prune all such branches (i.e. all supersets)

    i.e. exclude Itemsets, ABC, ABD, ABCD.

## Apriori Step-1 in pseudo codes

- K = 1
- Generate frequent Itemsets of length 1
- Prune itemsets of higher orders (i.e. supersets), if necessary
- Generate Itemsets of length K+1 from frequent Itemsets of length K
- Compute the support of new candidate Itemsets w.r.t. min sup. Prune if necessary.
- K = K+1
- Repeat until no frequent itemsets are found.

## Generation (step) of Itemsets for next level

Let $L_K$ denote the frequent itemsets at level $k$
and $C_K$ denote the set of all candidates at level $k$

- Items in $L_{K-1}$ are listed in an order

Step 1: Self joining $L_{K-1} * L_{K-1}$

   i.e. joining of 2 items from $L_{K-1}$

$$\{p.item_1, p.item_2 \dots p.item_{k-1}\}$$

$$\& \quad \{q.item_1 \ q.item_2 \dots q.item_{k-1}\} \quad \substack{\text{under the} \\ \text{given order}}$$

under the given order ($<$)

the itemset

insert into $C_K$, $p.item_1 \ p.item_2 \dots p.item_{k-2} \ p.item_{k-1} \ q.item_{k-1}$

Step 2: Prunning of $C_K$ set

   ∀ itemsets $c$ in $C_K$ and

   ∀ $(k-1)$ order subsets $s$ of $c$   (under the given order)

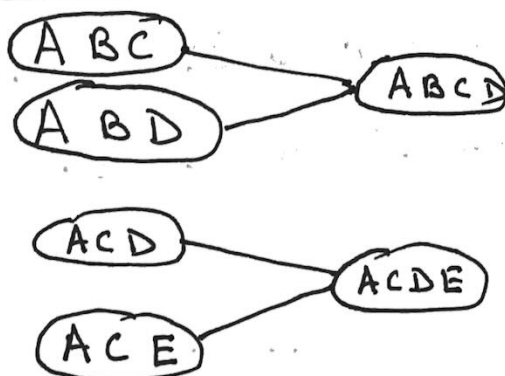   If ($s$ is not in $L_{K-1}$) delete $c$ from $C_K$

   Justification is anti-monotone property

## Example

$$L_3 = \{ABC, ABD, ACD, ACE, BCD\}$$

Self joining step: $L_3 * L_3$

**Prunning step:**

Consider $ACDE \rightarrow$ subsets $ACD, ACE, CDE, ADE$

$CDE/ADE$ is not in $L_3$

$\Rightarrow$ Prune $ACDE$ from $C_4$

Consider $ABCD \rightarrow$ subsets $ABC, ABD, ACD, BCD$

all subsets in $L_3$

$\Rightarrow$ Prunning not req'd for $ABCD$
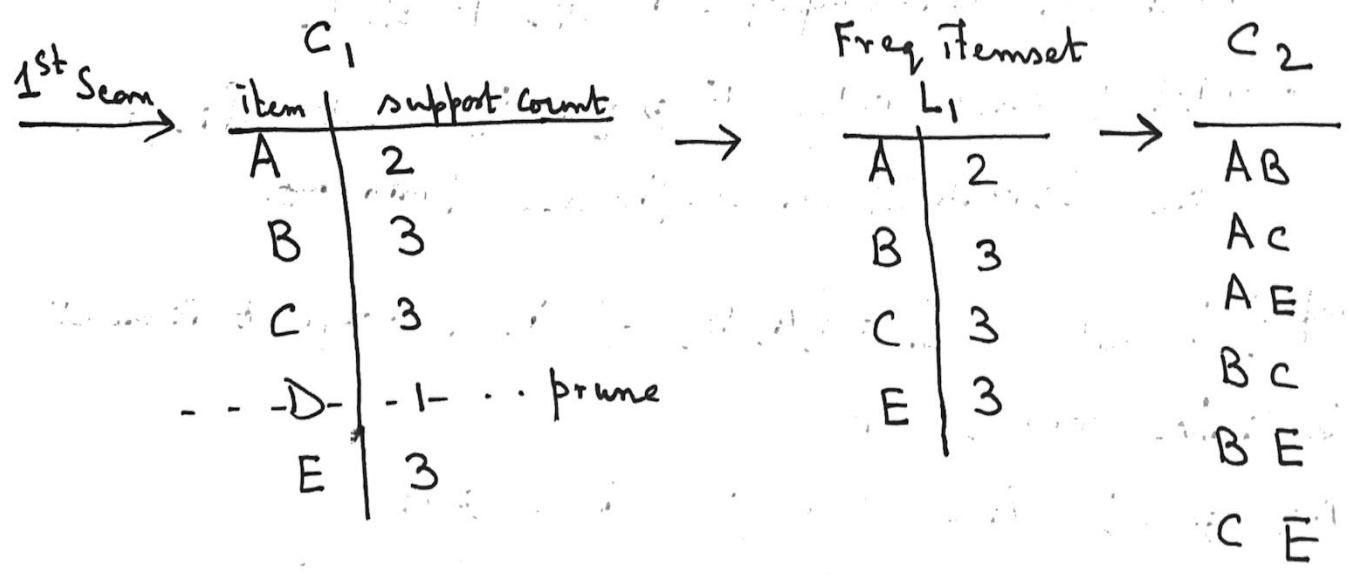
$\Rightarrow$ Pass $ABCD$ to $C_4$

**Example**

$T1 : \{A, C, D\}$

$T2 : \{B, C, E\}$

$T3 : \{A, B, C, E\}$

$T4 : \{B, E\}$

min sup $= .5$

1st Scan $\rightarrow$

$C_1$

| item | support count |
|------|---------------|
| A | 2 |
| B | 3 |
| C | 3 |
| -D- | -1- · · prune |
| E | 3 |

$\rightarrow$

Freq. itemset

$L_1$

| | |
|---|---|
| A | 2 |
| B | 3 |
| C | 3 |
| E | 3 |

$\rightarrow$

$C_2$

| |
|---|
| AB |
| AC |
| AE |
| BC |
| BE |
| CE |

2nd scan →

prune

| Itemset | Support Count |
|---------|---------------|
| ~~A B~~ | ~~1~~ |
| A C | 2 |
| ~~A E~~ | ~~1~~ |
| B C | 2 |
| B E | 3 |
| C E | 2 |

$\underline{L_2}$

| | |
|----|---|
| A C | 2 |
| B C | 2 |
| B E | 3 |
| C E | 2 |

→

$\underline{C_3}$ , $L_3$

BCE

all other joins prunned.

$S(BCE) = 2$

Under $C_3$, the only Itemset BCE has support count 2, i.e $\geqslant$ minsup

## Step 2 of apriori algorithm

Generate association rules using frequent itemsets

Given any frequent itemset $L$;

- find all non-empty subsets F of L
- output each rule $F \Rightarrow \{L-F\}$ that satisfies the threshold on confidence (min conf)

Example : Let $L = \{A, B, C\}$ is a frequent itemset

Candidate rules are

$AB \Rightarrow C$ ; $\quad AC \Rightarrow B$ ; $\quad BC \Rightarrow A$

$A \Rightarrow BC$ ; $\quad B \Rightarrow AC$ ; $\quad C \Rightarrow AB$

In general, there are $2^{|L|} - 2$ candidate rules.

**Remark:** Efficiency in rule generation

Confidence of rules generated from the same Itemset have anti-monotone property

$$C(ABC \Rightarrow D) \geq C(AB \Rightarrow CD) \geq C(A \Rightarrow BCD)$$

Sly $C(ABC \Rightarrow D) \geq C(AC \Rightarrow BD) \geq C(C \Rightarrow ABD)$

$$\left[ \begin{array}{l} \text{Consider, e.g., } ABC \Rightarrow D \quad \& \quad AB \Rightarrow CD \\ \qquad AB.C \; ABC \\ \qquad S(AB) \geq S(ABC) \\ \Rightarrow \dfrac{1}{S(ABC)} \not\geq \dfrac{1}{S(AB)} \\ \Ri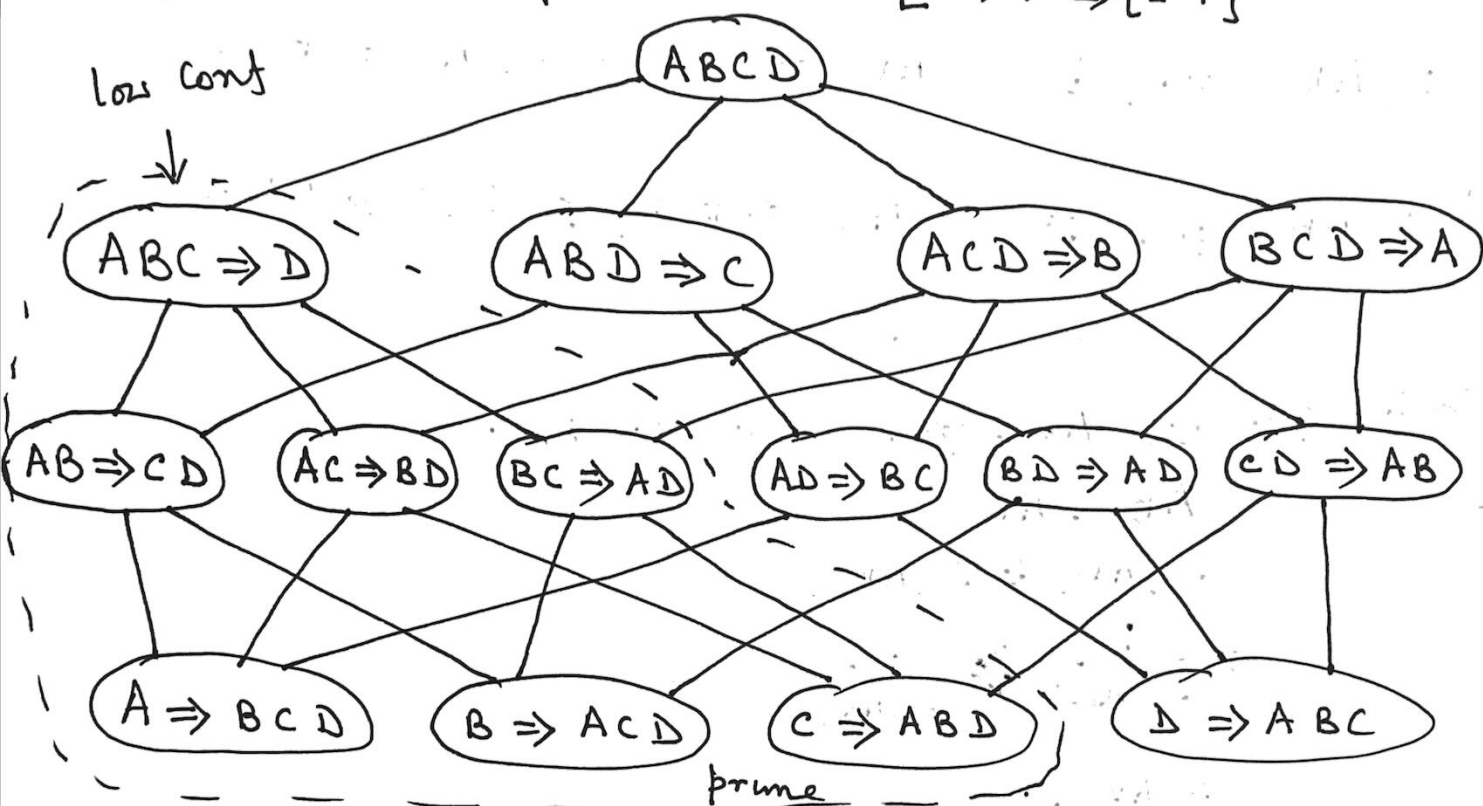ghtarrow \dfrac{S(ABCD)}{S(ABC)} \geq \dfrac{S(ABCD)}{S(AB)} \\ \qquad \text{i.e. } C(ABC \Rightarrow D) \geq C(AB \Rightarrow CD) \end{array} \right]$$

We apply this property to prune rule generation sequence.

# Example of pruning using anti-monotone property

Suppose ABCD is a frequent itemset and look at possible rules for ARM

$$L \to F \Rightarrow \{L - F\}$$

low conf
↓

```
                          ABCD
         ┌──────────┬──────────┬──────────┐
    ABC ⇒ D     ABD ⇒ C     ACD ⇒ B     BCD ⇒ A

 AB⇒CD   AC⇒BD   BC⇒AD   AD⇒BC   BD⇒AD   CD⇒AB

   A ⇒ BCD     B ⇒ ACD     C ⇒ ABD     D ⇒ ABC
```

prune

Suppose the rule ABC ⇒ D is low conf

i.e. $C(ABC \Rightarrow D) < min\ conf$

Then by anti-monotone property, all sub rules below that in the lattice will be $< min\ conf$ and

can be prunned, i.e.

AB ⇒ CD, AC ⇒ BD, BC ⇒ AD

A ⇒ BCD, B ⇒ ACD & C ⇒ ABD

are prunned.

# Example

$T1 : \{A, C, D\}$

$T2 : \{B, C, E\}$

$T3 : \{A, B, C, E\}$

$T4 : \{B, E\}$

Frequent Itemsets with min sup $\geq 0.5$ are

$S(\cdot) = 2$

$\underset{S(\cdot)=2}{BCE}$, $\underset{S(\cdot)=2}{AC}$, $BC$, $\underset{S(\cdot)=3}{BE}$, $\underset{S(\cdot)=2}{CE}$

Consider rules for $BCE$ Itemset

$$BCE$$

$$BC \Rightarrow E \qquad BE \Rightarrow C \qquad CE \Rightarrow B$$

$$B \Rightarrow CE \qquad C \Rightarrow BE \qquad E \Rightarrow BC$$

$$C(BC \Rightarrow E) = \frac{S(BCE)}{S(BC)} = \frac{2}{2} = 1$$

$$\Rightarrow BC \Rightarrow E (.5, 1)$$

$$C(BE \Rightarrow C) = \frac{S(BCE)}{S(BE)} = \frac{2}{3}$$

if min conf $= 0.75$, then $BE \Rightarrow C$

along with $B \Rightarrow CE$ & $E \Rightarrow BC$

are prunned.

$$C(CE \Rightarrow B) = \frac{S(BCE)}{S(CE)} = \frac{2}{2} = 1$$

$$C(C \Rightarrow BE) = \frac{S(BCE)}{S(C)} = \frac{2}{3} \quad \text{prune at min conf } 0.75$$