

DMML, 9 Jan 2020

Market Basket Analysis

Items $I = \{i_1, i_2, \dots, i_n\}$

Transactions $T = \{t_1, t_2, \dots, t_m\}$, each $t_i \subseteq I$

Identify "association rules" $X \rightarrow Y$, $X \cap Y = \emptyset$
 $X, Y \subseteq I$

$X \subseteq I \rightarrow$ support of X is the number of transactions containing X

$$\text{sup}(X) = \left| \{t_i \in T \mid X \subseteq t_i\} \right|$$

Significance of a rule $X \rightarrow Y$

Confidence : $\frac{\text{sup}(X \cup Y)}{\text{sup}(X)}$ Not $X \cap Y$

Two thresholds

- minimum support - what is $\text{sup}(X \rightarrow Y)$?
= $\text{sup}(X \cup Y) \geq \text{min-sup}$

- minimum confidence $\text{conf}(X \rightarrow Y) \geq \text{min-conf}$

Note: There is a single "correct" answer once thresholds are fixed.

How to find frequent itemsets?

Find X s.t. $\text{sup}(X) \geq \text{min-sup}$

If we do this, for each X that is frequent,

for each partition $Y \rightarrow Z$ of X ($Y \cup Z = X$
 $Y \cap Z = \emptyset$)

check if $\text{conf}(Y \rightarrow Z) \geq \text{min-conf}$

Use a priori principle

A Priori

If X is frequent, any $Y \subseteq X$ is also frequent
every

If some subset of X is not frequent, X cannot
be frequent

Level wise search

- Enumerate frequent items = F_1
- Next is F_2 - frequent pairs
 $F_2 \subseteq F_1 \times F_1$

(Naive) A Priori

Frequent sets of size j , F_j

→ Candidate frequent sets of size $j+1$, C_{j+1}

→ Frequent sets of size $j+1$, F_{j+1}

$$F_1 \rightarrow C_2 = (F_1 \times F_1) \setminus \text{Identity} \xrightarrow{\text{count}} F_2$$

$$F_2 \rightarrow C_3 = \{(x, y, z) \mid (x, y), (x, z), (y, z) \in F_2\}$$

$$\rightarrow F_3$$

$$F_j \rightarrow C_{j+1} = \{ (x_1, \dots, x_{j+1}) \mid \text{Every } j\text{-size subset} \\ \text{is in } F_j \}$$

$\rightarrow F_{j+1}$

Enumerate all $j+1$ size
subsets of \mathcal{I} & then prune
 $\binom{N}{j+1}$

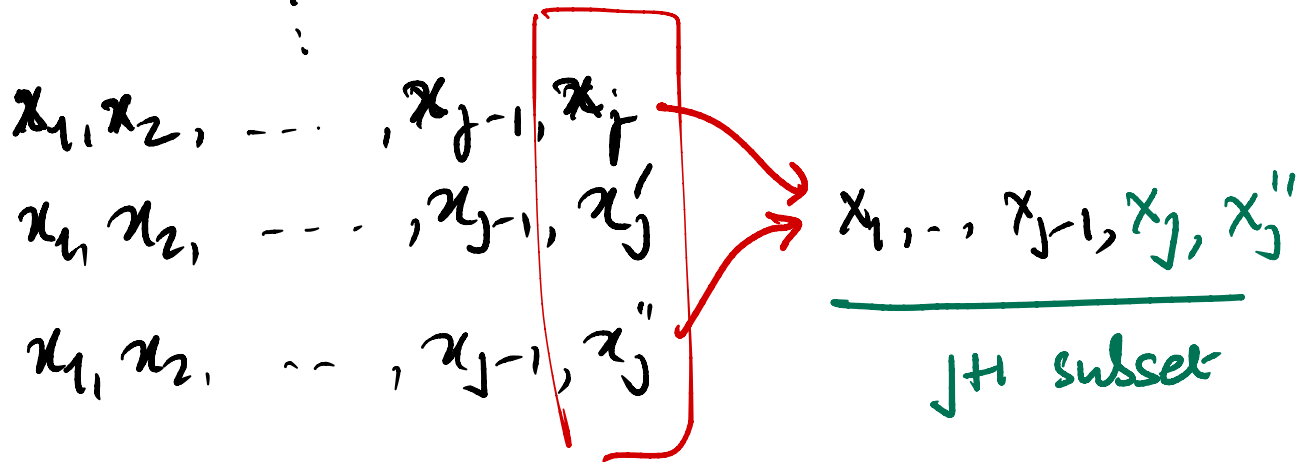
Can we fix this?

Idea: Can we find a $\hat{C}_{j+1} \supseteq C_{j+1}$ more
efficiently s.t. \hat{C}_{j+1} is also "small"

Assume I is ordered $i_1 < i_2 < \dots < i_N$

Any $X \subseteq I$ is listed in ascending order

F_j can be ordered lexicographically (dictionary order)



\hat{C}_{j+1} is constructed by this fusion operation

Is $\hat{C}_{j+1} \supseteq C_{j+1}$?

Suppose $\{y_1, y_2, \dots, y_{j+1}\} \in C_{j+1}$ $y_1 < y_2 < \dots < y_{j+1}$

Every j -subset must be in F_j

y_1, \dots, y_{j-1}, y_j is a j -subset

$y_2, \dots, y_{j-1}, y_{j+1}$ also

\supseteq are in F_j

\downarrow face

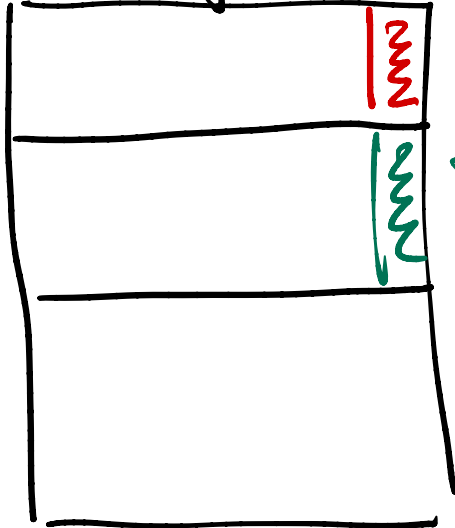
this set would
be in \hat{C}_{j+1}

$$\{x_1, x_2, \dots, x_{j-1}, x_j, x_j''\} \in \hat{C}_{j+1}$$

NOT in C_{j+1}

↓
may not be in F_j

F_j



⇒ all pairwise unions

⇒ all pairwise unions

↓
dictionary
order

Level by level

$F_1 \rightarrow C_1 \rightarrow F_2 \rightarrow C_2 \rightarrow \dots$

When to stop?

- ① $f > \text{max transaction size}$
- ② \hat{C}_{j+1} is empty

As f increases, support must decrease

Part 1 Find frequent itemsets

Part 2 Find rules

For each $X \in \mathcal{F}_1 \cup \mathcal{F}_2 \cup \dots \cup \mathcal{F}_g = \mathcal{F}$

Check if $Y \rightarrow Z$ meets min-conf, $Y \cup Z = X$
 $Y \cap Z = \phi$

$$\begin{array}{c} \Downarrow \\ \frac{\text{sup}(Y \cup Z)}{\text{sup}(Y)} \in \mathcal{F} \\ \in \mathcal{F} \text{ by a-priori} \end{array}$$

values have already been computed

Bottleneck:
All possible decompositions of X

A priori again!

$$X = \{a, b, c, d\}$$

$$Y_1 = \{a, b, c\}$$

$$Z_1 = \{d\}$$

$$Y_2 = \{a, b\}$$

$$Z_2 = \{c, d\}$$

$$Y_1 \rightarrow Z_1$$

$$\frac{\sup(Y_1 \cup Z_1)}{\sup(Y_1)} = \frac{\sup(X)}{\sup(Y_1)}$$

$$Y_2 \subseteq Y_1$$

$$\sup(Y_2) \geq \sup(Y_1)$$

$$Y_2 \rightarrow Z_2$$

$$\frac{\sup(Y_2 \cup Z_2)}{\sup(Y_2)} \stackrel{?}{=} \frac{\sup(X)}{\sup(Y_2)}$$

$$\text{conf}(Y_1 \rightarrow Z_1)$$

$$\geq \text{conf}(Y_2 \rightarrow Z_2)$$

If $R_1 X \setminus \{x_1, x_2\} \rightarrow \{x_1, x_2\}$ meets min-conf

then $R_2 X \setminus \{x_1\} \rightarrow \{x_1\}$ also meet min-conf

$R_3 X \setminus \{x_2\} \rightarrow \{x_2\}$

$$\text{conf}(R_1) \leq \text{conf}(R_2)$$

$$\text{conf}(R_1) \leq \text{conf}(R_3)$$

\therefore if R_2 or R_3
fails min-conf,
 R_1 must fail
min-conf

Check all rules of the form

$$X \setminus \{x\} \rightarrow \{x\}$$

Work up to larger RHS using a-priori

Go on back to stopping

Large rules are difficult to interpret