

ASSociation Rule mining

ASSociation Rules: Relations between Items

↳ mining nearl: sets of method to find this relations

↳ Also called: market Basket Analysis

~~Some definitions~~

Item $\rightarrow I = \{i_1, i_2, \dots, i_m\}$

transaction $\rightarrow t \subseteq I$

↳ Baskets $\quad \quad \quad \downarrow$ subset

Imagine: a shop

$t_1: \{\text{bread, cheese, milk}\}$

$t_2: \{\text{apple, eggs, salt, yogurt}\}$

$t_n: \{\text{biscuit, eggs, milk}\}$

Subsets of Items: X

$$\uparrow X \subseteq t \downarrow \text{transaction}$$

ASSociation rule:

$X \& Y$ are subset of all items $\rightarrow X \& Y$ don't share

$X \rightarrow Y$, where $X, Y \subseteq I$, $X \cap Y = \emptyset$ shared items

↳ Subset of items X is related to subset of items Y

k -itemset \rightsquigarrow itemset that has k items

{milk, bread, cereal} is a 3-Itemset

Assume: Bread \rightarrow Milk {Sup=5%, Conf=100%}]

Support: chance of buying $X \& Y$ together

$$\text{Sup} = \Pr(X \& Y)$$

Confidence: chance that if x is in then y will be too

$$\text{Conf} = \Pr(Y|X)$$

$$\text{support} = \frac{\text{frequency}(X, Y)}{N}$$

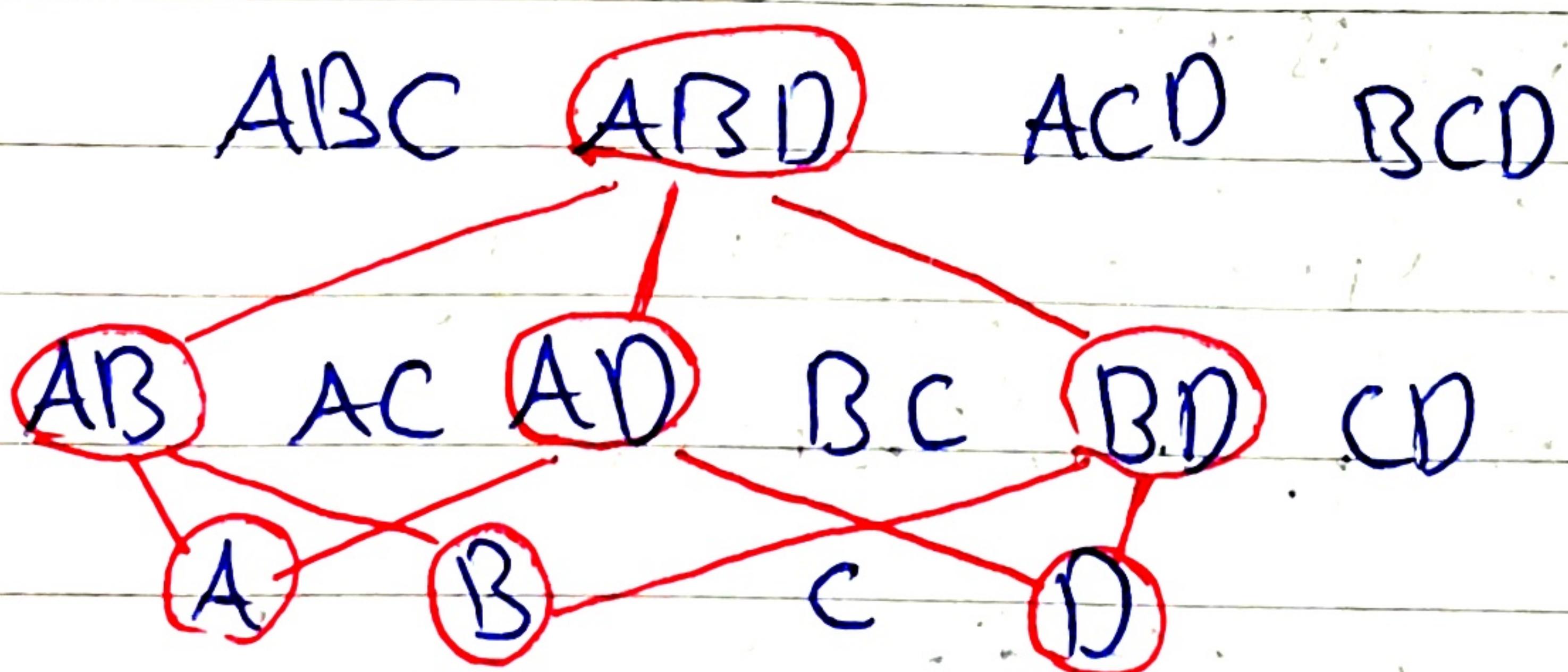
Rule: $X \Rightarrow Y$

$$\text{confidence} = \frac{\text{freq}(X, Y)}{\text{freq}(X)}$$

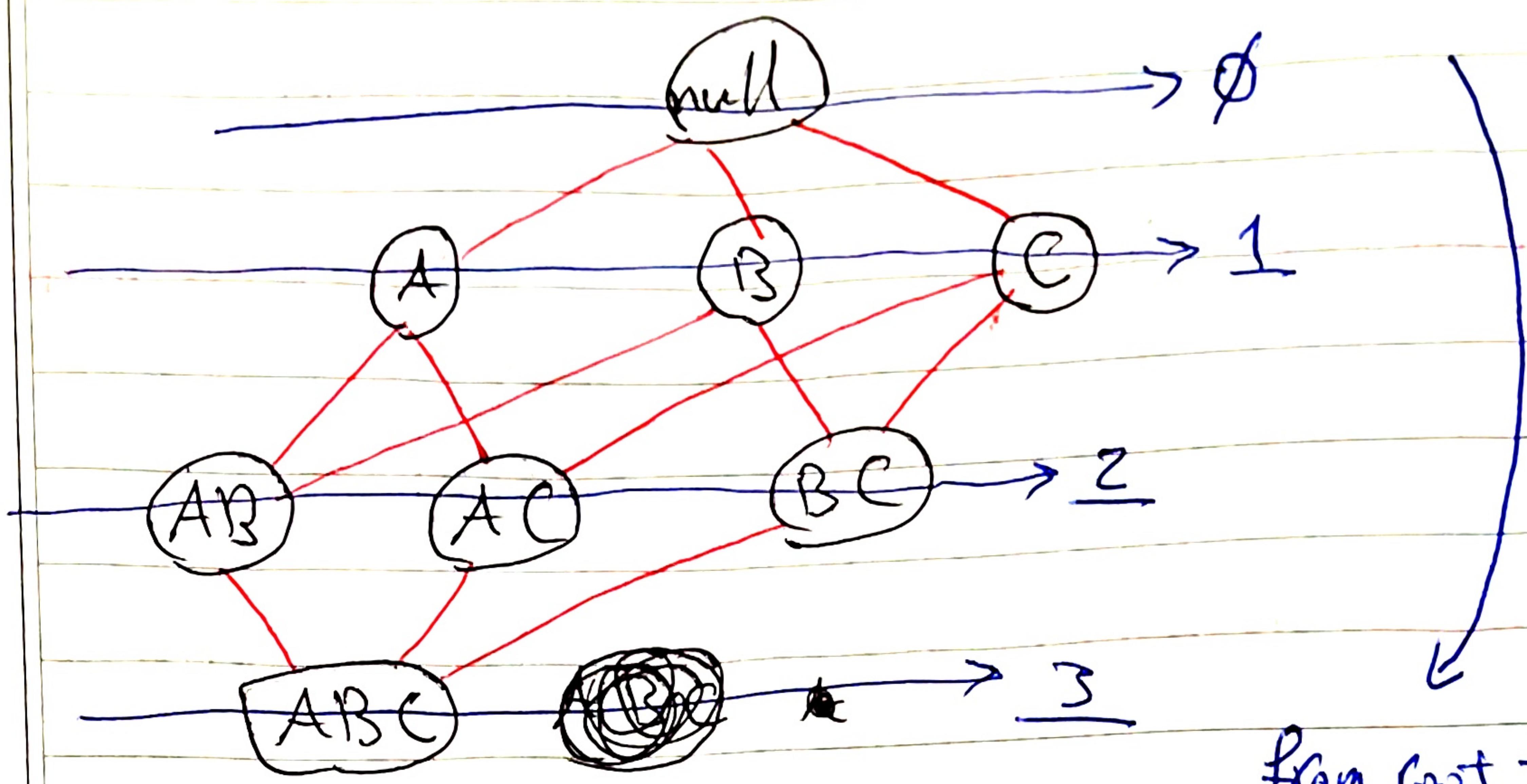
Frequent itemset: set of items that their support is bigger than minSup

APriori Algorithm

each subset of frequent itemset is a frequent set itself.



↳ each of these red circled must be frequent themselves
So ABD can be frequent itemset.



if A is infrequent ~~then~~

then AB, AC, ABC will be cancelled and ignored.

Algorithm APRIORI(+)

$C_1 \leftarrow \text{init-Pass}(+);$

$F_1 \leftarrow \{f \mid f \in C_1, f.\text{count} \geq \text{minSup}\};$

for ($k=2; F_{k-1} \neq \emptyset; k++$) do

$C_k \leftarrow \text{candidate-gen}(F_{k-1});$

for each transaction $t \in T$ do

for each candidate $c \in C_k$ do

if c is contained in t then

$c.\text{count}++;$

$F_k \leftarrow \{c \in C_k \mid c.\text{count} \geq \text{minSup}\}$

return $F \leftarrow \cup_i F_k$

Candidate-gen function

Function candidate-gen(F_{t-1})

$$C_t \leftarrow \emptyset;$$

forall $f_1, f_2 \in F_{t-1}$

$$\text{with } f_1 = \{i_1, \dots, i_{t-2}, i_{t-1}\}$$

$$\text{and } f_2 = \{i_1, \dots, i_{t-2}, i'_{t-1}\}$$

$$\text{and } i_{t-1} < i'_{t-1} \text{ do}$$

$$\textcircled{a} \quad C \leftarrow \{i_1, \dots, i_{t-1}, i'_{t-1}\}; \quad \} \text{Join } f_1, f_2$$

$$C_t \leftarrow C_t \cup \{C\};$$

for each $(t-1)$ -subset S of C do

if ($S \in F_{t-1}$) then

 |
 | delete C from C_t ;

 |
 | end

 |
 | delete non related
 | combinations;

 |
 | end

return C_t ;

$$\text{eg. } F_3 = \{\{1, 2, 3\}, \{1, 2, 4\}, \{1, 3, 4\}, \\ \{1, 3, 5\}, \{2, 3, 4\}\}$$

$$\text{After Join: } C_4 = \{\{1, 2, 3, 4\}, \{1, 3, 4, 5\}\}$$

$$\text{After Pruning: } C_4 = \{\{1, 2, 3, 4\}\} \quad \begin{matrix} \hookrightarrow 1, 4, 5 \text{ not} \\ \text{in } F_3 \end{matrix}$$

generating rules from frequent itemset

* for each frequent itemset X ,
 for each proper nonempty subset A of X

Let $B \subset X - A$

$A \rightarrow B$ is associative rule if:

* $\text{Conf}(A \rightarrow B) \geq \text{minConf}$

* $\text{Sup}(A \rightarrow B) = \text{Sup}(A \cup B) = \text{Sup}(X)$

* $\text{Conf}(A \rightarrow B) = \frac{\text{Sup}(A \cup B)}{\text{Sup}(A)}$

end

end

end

eg. DB

minSup
 $= 2$

~~A, C, D~~

B, C, E

A, B, C, E

B, E

C₁

{A}, 2

{B}, 3

{C}, 3

~~{D}, 1~~

{E}, 3

4

{A}, 2

{B}, 3

{C}, 3

{E}, 3

~~{A, B}, 1~~

C₂ {A, C}, 2

2nd ~~{A, E}, 1~~

{B, C}, 2

{B, E}, 3

{C, E}, 2

→ {B, C, E}, 2

eg. generating rules

if $\text{conf}(A \rightarrow B)$

$> \text{minConf}$

$$\{A, B, C\} \rightarrow \{A\} \rightarrow \{B, C\}$$

it's rule

$$\times \quad \{B\} \rightarrow \{A, C\}$$

$$\{C\} \rightarrow \{A, B\}$$

:

$\hookrightarrow \text{sup}(A \rightarrow B) \text{ is sup}(X)$

$$\text{conf}(A \rightarrow B) \text{ is } \frac{\text{sup}(X)}{\text{sup}(A)}$$

eg. X

$\{2, 3, 4\}$ is freq, $\text{sup} = 50\%$

$$\hookrightarrow \{2\} \xrightarrow{A} \{3, 4\} \xrightarrow{B} \text{conf} = \frac{50\%}{75\%} = 67\%$$

$$\{3\} \rightarrow \{2, 4\} \rightarrow \text{conf} = 67\%$$

$$\{4\} \rightarrow \{2, 3\} \rightarrow \text{conf} = 67\%$$

$$\{2, 3\} \rightarrow \{4\} \rightarrow \text{conf} = 100\%$$

$$\{2, 4\} \rightarrow \{3\} \rightarrow 100\%$$

$$\{3, 4\} \rightarrow \{2\} \rightarrow 67\%$$

$$2 \rightarrow \text{sup} = 75\%$$

$$\{2, 3\} \rightarrow 50\%$$

$$3 \rightarrow \text{sup} = 75\%$$

$$\{2, 4\} \rightarrow 50\%$$

$$4 \rightarrow \text{sup} = 75\%$$

$$\{3, 4\} \rightarrow 75\%$$