# Clustering Association Rules Mining Problem: a Symbolic Approach

Pietro Sala
Department of Computer Science
University of Verona, Italy
pietro.sala@univr.it

Given a universe of item $\texttt{It}$, a <u>linear transaction</u> is a function $\texttt{ltr} : \texttt{It} \to \mathbb{N}$. The the set of all possible transactions is $\mathbb{LTr} = \mathbb{N}^{|\texttt{It}|}$. A <u>Linear Transaction Dataset</u> (LTD) is a **finite** multiset $\texttt{LTd} : \mathbb{LTr} \to \mathbb{N}$, moreover let $|\texttt{LTd}| = \sum_{\texttt{ltr} \in \mathbb{LTr}} \texttt{LTd}(\texttt{ltr})$.

Let $\mathbb{I}_\mathbb{N}$ the set $\mathbb{I}_\mathbb{N} = \{[n, n'] \in \mathbb{N}^2 : n \le n'\}$ the set of all the possible intervals on $\mathbb{N}$, and let $\sqsubseteq$ the relation on $\mathbb{I}_\mathbb{N}^2$ such that $[n, n'] \sqsubseteq [\overline{n}, \overline{n}']$ if and only if $n \le \overline{n} \le \overline{n}' \le n'$. Moreover, given $\overline{n} \in \mathbb{N}$ we say that $\overline{n} \in [n, n']$ if and only if $n \le \overline{n} \le n'$.

A <u>Clustering Itemset</u> is a function $\texttt{X}_\mathcal{C} : \texttt{It} \to \mathbb{I}_\mathbb{N} \cup \{[0, +\infty)\}$, moreover for each $\texttt{it} \in \texttt{It}$, moreover given a $\texttt{X}_\mathcal{C}$ and an item $\texttt{it} \in \texttt{It}$ we say that $\texttt{it}$ is <u>unconstrained</u> by $\texttt{X}_\mathcal{C}$ if and only if $\texttt{X}_\mathcal{C}(\texttt{it}) = [0, +\infty)$ (if $\texttt{X}_\mathcal{C}(\texttt{it}) \neq = [0, +\infty)$ we will say that $\texttt{it}$ is unconstrained ).

Given an LTD $\texttt{LTd}$ and a $\texttt{X}_\mathcal{C}$ we define the support (on $\texttt{LTd}$) of $\texttt{X}_\mathcal{C}$ as:

$$\mathcal{S}up(\texttt{X}_\mathcal{C}) = \frac{\sum_{\substack{\texttt{ltr} \in \mathbb{LTr} \text{ such that:} \\ \text{for each } \texttt{it} \in \texttt{It} \\ \texttt{ltr}(\texttt{it}) \in \texttt{X}_\mathcal{C}(\texttt{it})}} \texttt{LTd}(\texttt{ltr})}{|\texttt{LTd}|}$$

Let us notice that all the concepts of <u>frequent itemsets</u> as well as any measures like <u>confidence</u> may be rephrased in terms of the newly defined support for clustering items.

A <u>Clustering Association Rule</u> is a rule of the form:

$$\texttt{X}_\mathcal{C} \longrightarrow \texttt{Y}_\mathcal{C}$$

where:
- there exist $\texttt{it}, \texttt{it}' \in \texttt{It}$ such that $\texttt{X}_\mathcal{C}(\texttt{it}) \neq [0, +\infty)$ and $\texttt{Y}_\mathcal{C}(\texttt{it}') \neq [0, +\infty)$;
- for each $\texttt{it} \in \texttt{It}$ we have that $\texttt{X}_\mathcal{C}(\texttt{it}) \neq [0, +\infty)$ implies $\texttt{Y}_\mathcal{C}(\texttt{it}) = [0, +\infty)$, and $\texttt{Y}_\mathcal{C}(\texttt{it}) \neq [0, +\infty)$ implies $\texttt{X}_\mathcal{C}(\texttt{it}) = [0, +\infty)$.

Analogously of what has been done for standard association rules we can define the set of <u>frequent clustering itemsets</u> $\mathbb{Fi}_\mathcal{C}$ according to a given threshold $0 \le \epsilon \le 1$.

The <u>Clustering Association Rules Mining Problem</u> (CARM) consists of determining, given an LTD $\texttt{LTd}$ and two thresholds $0 \le \epsilon, \delta \le 1$, the set of all and only the rules $\texttt{X}_\mathcal{C} \longrightarrow \texttt{Y}_\mathcal{C}$ such that $\mathcal{S}up((\texttt{X} \cup \texttt{Y})_\mathcal{C}) \ge \epsilon$ and $\mathcal{C}onf(\texttt{X}_\mathcal{C} \longrightarrow \texttt{Y}_\mathcal{C}) \ge \delta$.

**Assignment:**

Implement an algorithm that solves the CARM problem, i.e., given support/confidence thresholds $0 \le \epsilon, \delta \le 1$ and a dataset $D$ (inputs) enumerates a complete set of Clustering Association Rules on $D$ which holds with support at least $\epsilon$ and confidence at least $\delta$. Test the resulting implementation by extracting the rules on a dataset of your choice such as AirQuality on the UCI machine learning dataset repository.