

# TP DATA MINING

## Résumé

Etant donnée le nombre croissant de grandes bases de données et de la complexité des algorithmes mis en œuvre pour leurs exploitation, la question d'optimisation préoccupe de plus en plus les chercheurs du domaine de la fouille de données.

L'une des techniques les plus utilisées pour extraire ces connaissances est la méthode des règles d'association.

Toutefois, la plupart des algorithmes d'extraction des itemsets fermés fréquents voient leurs performances se dégrader lorsque la taille des données augmente. Pour maintenir les performances de ces algorithmes, l'utilisation de méthodes et outils distribués apparaît comme une solution naturelle.

Le projet demandé, vise à améliorer les performances de l'un des meilleurs algorithmes d'extraction des itemsets fermés fréquents a savoir l'algorithme Close.

Il s'agit de mettre en place une version de l'algorithme Close, dans le cas d'un contexte où les données seront partitionnées. Ceci permettra d'extraire des résultats locaux relatifs à chaque partition. Par la suite une seconde étape permettra de retrouver les résultats finaux à partir de ceux déjà extraits.

Dans notre cas, on considère ainsi  $n$  sources de données  $S_1, S_2, \dots, S_n$ . Soit  $D$  la taille de l'ensemble de nos données et  $d_i$  la taille de chaque source. Pour un itemset  $X$  donné, soit  $\text{Sup}(X)$  et  $\text{Sup}_i(X)$  le support de  $X$  dans  $S_i$  respectivement. Dans ce cas,  $\text{Sup}(X)$  est appelé le support « global » et  $\text{Sup}_i(X)$  le support « local » de l'itemset  $X$  dans  $S_i$ . Pour un support minimum donné  $\text{minsup}$ ,  $X$  est globalement fréquent si:  $\text{Sup}(X) \geq \text{minsup} * D$ .

De même,  $X$  est localement fréquent dans  $S_i$ : si  $\text{Sup}_i(X) \geq \text{minsup} * d_i$ .

Ainsi, le problème d'extraction des règles d'association se réduit à trouver tous les itemsets globalement fréquents et de générer ensuite les règles d'associations correspondantes.

## Objectifs:

Il est demandé d'implémenter l'algorithme Close dans sa nouvelle version et de rédiger un rapport expliquant votre démarche et les différents choix faits.