

Manipulation et fusion de données multidimensionnelles

Franck Ravat, Olivier Teste, Gilles Zurfluh
Institut de Recherche en Informatique de Toulouse / Equipe SIG-ED
118, Route de Narbonne 31062 TOULOUSE cedex 04
mél : {ravat, teste, zurfluh}@irit.fr

Résumé. Cet article définit une algèbre permettant de manipuler des tables dimensionnelles extraites d'une base de données multidimensionnelles. L'algèbre intègre un noyau minimum d'opérateurs unaires permettant d'effectuer les analyses décisionnelles par combinaison d'opérateurs. Cette algèbre intègre un opérateur binaire permettant la fusion de tables dimensionnelles facilitant les corrélations des sujets analysés.

1 Introduction

Nos travaux se situent dans le cadre des systèmes décisionnels intégrant des bases de données multidimensionnelles (BDM). Conceptuellement, ces BDM organisent les données en sujets appelés faits et axes d'analyses appelés dimensions (Kimball, 1996).

1.1 Contexte : notre modèle conceptuel

Définition : Un fait F_j est défini par $(N_{F_j}, M_{F_j}, I_{F_j}, IStar_{F_j})$ où

- N_{F_j} est le nom du fait,
- $M_{F_j} = \{m_1, m_2, \dots, m_w\}$ est un ensemble de mesures (ou indicateurs d'analyse),
- $I_{F_j} = \{I_{F_1}, I_{F_2}, \dots\}$ est l'ensemble des instances de F_j ,
- $IStar_{F_j}$ est une fonction associant chaque instance de I_{F_j} à une instance de chaque dimension liée au fait.

Définition : Une dimension D_i est définie par $(N_{D_i}, A_{D_i}, H_{D_i}, I_{D_i})$ où

- N_{D_i} est le nom de la dimension,
- $A_{D_i} = \{a_{D_i_1}, a_{D_i_2}, \dots, a_{D_i_w}\}$ est un ensemble d'attributs,
- $H_{D_i} = \{h_{D_i_1}, h_{D_i_2}, \dots, h_{D_i_y}\}$ est un ensemble de hiérarchies,
- $I_{D_i} = \{I_{D_i_1}, I_{D_i_2}, \dots\}$ est l'ensemble des instances de D_i .

Définition : Une hiérarchie représente une perspective d'analyse précisant les niveaux de granularité auxquels peuvent être manipulés les indicateurs d'analyse. Une hiérarchie $h_{D_i_x}$ définie sur la dimension D_i est un chemin élémentaire acyclique débutant par l'attribut de plus faible granularité et se terminant par un attribut de plus forte granularité. Elle est définie par $(N_{D_i_x}, Param_{D_i_x}, Suppl_{D_i_x})$ où

- $N_{D_i_x}$ est le nom de la hiérarchie,
- $Param_{D_i_x} = \langle a_{D_i_k}, a_{D_i_1}, \dots, a_{D_i_z} \rangle$ est un ensemble ordonné décrivant la hiérarchie des attributs (chaque attribut est appelé paramètre de la hiérarchie et correspond à un niveau de granularité d'analyse),
- $Suppl_{D_i_x} : Param_{D_i_x} \rightarrow 2^{(A_m - Param_{D_i_x})}$ est une application spécifiant les attributs faibles qui complètent la sémantique des paramètres (chaque paramètre est associé à un ensemble d'attributs faibles).