

Entrepôt de Données Spatiales basé sur GML: Politique de Gestion de Cache

Lionel Savary , Georges Gardarin, Karine Zeitouni

Laboratoire PRiSM, Université de Versailles, 45 Avenue des Etats-Unis - 78035 Versailles

{Lionel.Savary, Georges.Gardarin, Karine.Zeitouni}@prism.uvsq.fr

Motivation : Dans les entrepôts de données, la manipulation de gros volumes de données requière souvent un temps d'exécution important. En particulier, si les requêtes portent sur des données spatiales contenues dans des documents semi-structurés, les temps de réponse deviennent prohibitifs. Afin de réduire le temps de traitement imposé par l'utilisation d'opérateurs spatiaux dans ce type de document, nous proposons une politique de remplacement de cache adaptée aux documents GML. Cette politique prend en compte les données spatiales et non-spatiales, ainsi que le nombre d'opérateurs spatiaux présents dans les requêtes utilisateurs.

Politique de remplacement de cache : Soit D_q le nouveau document de taille T_q à insérer dans le cache. On désigne par T_i la taille du document i ($1 \leq i \leq n$) du cache et par C_{GMLi} le coût d'accès au document du cache. Soient $(X_i)_{i=1..n} \in \{0; 1\}^n$ tel que $X_i = 1$ si le document i est conservé dans le cache, 0 s'il est supprimé. Notons de plus D_{GMLj} le coût d'accès au document j sur disque (lorsqu'il n'est pas en cache). On recherche alors les documents i du cache à supprimer tels que la somme des coûts d'accès soit la plus petite possible:

$$\text{Minimiser Coût d'accès} = \sum_{j=1}^n X_j * C_{GMLj} + \sum_{j=1}^n (1-X_j) * D_{GMLj}$$

Une contrainte est que la somme des tailles des documents éliminés du cache soit supérieure ou égale à T_q : $\sum_{i=1}^n X_i * T_i \geq T_q$

C_{GMLi} représente le coût d'accès à un document GML i en cache, calculé selon la formule d'Arlitt [1], soit $C_{GMLi} = L + F_i * C_i / S_i$, avec : L une constante ; S_i la taille du document ; F_i la fréquence d'accès au document et C_i le coût pour une requête sur des données géographiques. Notre calcul de C_i tient compte du coût sur les données non-spatiales, du coût sur les données spatiales, ainsi que du nombre d'opérateurs spatiaux présents dans la requête. Le problème est de déterminer les X_i qui optimisent le coût total. Nous proposons une adaptation d'algorithmes classiques de recherche opérationnelle pour déterminer les documents à conserver.

Références

1. M. Arlitt, R. Friedrich L. Cherkasova, J. Dille, and T. Jin. Evaluating content management techniques for web proxy caches. In HP Tec. report, Palo Alto, Apr. 1999.