

Sélection simultanée d'index et de vues matérialisées

Nora Maiz, Kamel Aouiche et Jérôme Darmont

Laboratoire ERIC, Université Lumière Lyon 2
5 avenue Pierre Mendès-France
69676 Bron Cedex
{nora.maiz, kamel.aouiche, jerome.darmont}@eric.univ-lyon2.fr
<http://eric.univ-lyon2.fr>

Résumé. Les index et les vues matérialisées sont des structures physiques qui accélèrent l'accès aux données d'un entrepôt. Ces structures engendrent cependant une surcharge de maintenance. Par ailleurs, elles partagent le même espace disque. Les travaux existants dans le domaine de la sélection d'index et de vues matérialisées traitent ces deux structures de manière isolée. Dans cet article, nous couplons au contraire la sélection d'index et de vues matérialisées de façon à prendre en compte les interactions entre ces structures de données et à permettre un partage efficace de l'espace de stockage commun qui leur est alloué. Pour cela, nous avons développé des modèles de coût qui évaluent le bénéfice de la matérialisation de vue et de l'indexation. Ces modèles de coût nous permettent, grâce à un algorithme glouton, de sélectionner une configuration pertinente d'index et de vues matérialisées. Nos expérimentations montrent que notre stratégie se révèle meilleure que celles qui opèrent une sélection isolée des index et des vues matérialisées.

1 Introduction

Les entrepôts de données sont généralement modélisés selon un schéma en étoile contenant une table de faits centrale volumineuse et un certain nombre de tables dimensions représentant les descripteurs des faits Inmon (2002); Kimball et Ross (2002). La table de faits contient des clés étrangères vers les clés primaires des tables dimensions, ainsi que des mesures numériques. Avec ce type de modèle, une requête décisionnelle nécessite une ou plusieurs jointures entre la table de faits et les tables dimensions. De plus, le schéma de l'entrepôt peut comporter des hiérarchies au niveau des dimensions (schéma en flocon de neige), ce qui entraîne des jointures additionnelles. Ces jointures sont très coûteuses en terme de temps de calcul. Ce coût devient prohibitif lorsque les jointures opèrent sur de très grands volumes de données. Il est alors crucial de le réduire.

Les vues matérialisées et les index sont des structures physiques qui permettent de réduire le temps d'exécution des requêtes en précalculant les jointures et en offrant un accès direct aux données. Cependant, lors du rafraîchissement de l'entrepôt de données, ces structures doivent également être mises à jour, ce qui engendre une surcharge pour le système. Par ailleurs, index

et vues matérialisées partagent le même espace de stockage. Il est donc judicieux de ne créer que les plus pertinents.

Les travaux existants dans le domaine de la sélection d'index et/ou de vues matérialisées traitent ces deux structures de manière isolée ou séquentielle. Dans cet article, nous proposons une nouvelle stratégie qui opère une sélection simultanée des vues matérialisées et des index afin de prendre en compte les interactions entre eux.

Dans un premier temps, nous exploitons des stratégies de sélection isolée des vues matérialisées et des index, qui nous fournissent un ensemble d'index et de vues candidats. Nous calculons ensuite grâce à des modèles de coût le bénéfice potentiel de chaque index et de chaque vue matérialisée, en prenant en compte les interactions possibles entre ces deux types de structures. Finalement, un algorithme glouton nous permet de sélectionner simultanément les vues matérialisées et les index les plus pertinents.

Cet article est organisé comme suit. La Section 2 est consacrée à l'état de l'art de ce domaine de recherche. La Section 3 présente globalement notre stratégie de sélection simultanée de vues matérialisées et d'index. Nous détaillons ensuite nos modèles de coût dans la Section 4, puis le calcul des bénéfices d'indexation et de matérialisation de vues dans la Section 5. Nous présentons dans la Section 6 notre algorithme glouton de sélection simultanée de vues matérialisées et d'index. Les premières expérimentations que nous avons menées pour valider notre approche sont présentées dans la Section 7. Nous concluons finalement cet article et évoquons nos perspectives de recherche dans la Section 8.

2 État de l'art

Le problème de sélection d'index et/ou de vues matérialisées consiste à construire une configuration d'index et/ou de vues matérialisées optimisant le coût d'exécution d'une charge donnée, supposée représentative. Cette optimisation peut être réalisée sous certaines contraintes, comme l'espace de stockage alloué aux index et aux vues à sélectionner.

Plus formellement, si $O = \{o_1, \dots, o_n\}$ est un ensemble d'objets (index candidats, vues matérialisées candidates ou index sur les vues), $Q = \{q_1, \dots, q_m\}$ l'ensemble des requêtes de la charge et S la taille de l'espace disque alloué par l'administrateur pour stocker les objets à sélectionner, alors il faut trouver une configuration d'index et de vues matérialisées *Config* tel que :

- le coût d'exécution C des requêtes de la charge soit minimal, c'est-à-dire :

$$C_{/Config}(Q) = \text{Min} (C_{/O}(Q)) ;$$

- l'espace de stockage des index et des vues de *Config* ne dépasse pas S , c'est-à-dire :

$$\sum_{o_i \in Config} \text{taille}(o_i) \leq S.$$

Les problèmes de sélection d'index et de vues matérialisées sont connus pour être NP-complets Comer (1978); Gupta (1999). De ce fait, il n'existe pas d'algorithme qui propose une solution optimale en un temps fini. Plusieurs travaux de recherche proposent des solutions proches de la solution optimale en utilisant des heuristiques réduisant la complexité du problème.

Les travaux traitant la sélection d'index Frank *et al.* (1992); Choenni *et al.* (1993a,b); Valentin *et al.* (2000); Golfarelli *et al.* (2002); Dageville *et al.* (2004) et la sélection de vues Gupta (1999); Gupta et Mumick (1999); Baril et Bellahsene (2003); Smith *et al.* (2004); Gupta et Mumick (2005) s'orientent dans leur majorité vers une sélection séquentielle des vues et des index sur les vues, ou vers une sélection isolée des index ou des vues matérialisées. Cependant, les index et les vues matérialisées sont fondamentalement des structures physiques similaires Agrawal *et al.* (2000). En effet, les deux structures sont redondantes, accélèrent le temps d'exécution des requêtes, partagent la même ressource de stockage et impliquent une surcharge de maintenance pour le système suite aux mises à jour des données. Les vues et les index peuvent alors être en interaction. La présence d'un index sur une vue matérialisée peut en effet rendre celle-ci plus "attractive" et *vice versa*.

Peu de travaux se sont portés sur la sélection simultanée des vues matérialisées et des index. Agrawal *et al.* ont proposé trois alternatives pour l'énumération conjointe de l'espace des index et des vues matérialisées Agrawal *et al.* (2000). La première alternative, dénotée MVFIRST, tend à sélectionner les vues matérialisées en premier, puis les index pour une charge donnée en présence des vues préalablement sélectionnées. La deuxième alternative, dénotée INDFIRST, sélectionne en premier les index, puis les vues. La troisième alternative, dénotée *joint enumeration*, traite la sélection des index, des vues matérialisées et des index sur ces vues en une seule itération. Les auteurs affirment qu'elle est plus efficace que les deux premières.

Bellatreche *et al.* ont traité le problème de distribution de l'espace de stockage entre les vues matérialisées et les index de manière itérative afin de minimiser le coût total d'exécution des requêtes d'une charge donnée Bellatreche *et al.* (2000). Un ensemble de vues et d'index est désigné comme une solution initiale au problème de sélection d'index et de vues. L'approche reconsidère itérativement la solution initiale dans le but de réduire davantage le coût d'exécution des requêtes en redistribuant l'espace de stockage entre les vues et les index. Elle s'appuie sur une compétition perpétuelle entre deux agents, l'espion des index et l'espion des vues. L'espion des index (respectivement, des vues) vole de l'espace réservé pour stocker les vues (respectivement, les index). L'espace ainsi récupéré est utilisé pour créer d'autres index à la place des vues élaguées, suivant des politiques de remplacement. L'opération est validée si le coût d'exécution des requêtes est réduit. La sélection d'index et de vues commence par appliquer l'espion qui réduit le plus le coût des requêtes. Le processus de sélection s'arrête lorsqu'il n'y a plus de réduction du coût des requêtes.

Finalement, Rizzi et Saltarelli ont proposé une approche qui détermine *a priori* un compromis entre l'espace de stockage alloué aux index et aux vues matérialisées en se basant sur les requêtes de la charge Rizzi et Saltarelli (2003). L'idée de Rizzi et Saltarelli est que le facteur clé dans l'optimisation des performances des requêtes est leur niveau d'agrégation, défini par la liste des attributs de la clause **Group by**, et la sélectivité des attributs présents dans les clauses **Having** et **Where**. En effet, la matérialisation offre un grand bénéfice aux requêtes comportant des agrégations de granularité grossière (nombre faible d'attributs dans la clause **Group by**) car elles génèrent peu de groupes dans un grand nombre de n-uplets et, par conséquent, l'accès à une petite vue est moins coûteux que l'accès aux tables de base. D'autre part, les index donnent leur meilleur bénéfice avec des requêtes contenant des attributs dont la sélectivité est élevée car elles sélectionnent peu de n-uplets et, par conséquent, l'accès à un nombre élevé de n-uplets inutiles est évité. Les requêtes avec des agrégations fines et de fortes sélectivités encouragent l'indexation. En revanche, les requêtes avec des agrégations grossières

et de faibles sélectivités encouragent la matérialisation.

Le défaut que nous identifions dans les travaux de Bellatreche *et al.* et de Rizzi et Saltarelli est qu'ils ne prennent pas en compte les interactions entre les vues matérialisées et les index. En effet, les deux types de structures y sont sélectionnés de manière concurrente et non conjointe. Or, des index sur les vues matérialisées déjà sélectionnées peuvent s'avérer des options très intéressantes. Notre approche s'apparente donc à celle d'Agrawal *et al.* (*joint enumeration*). Cependant, ses auteurs ne donnent malheureusement aucun détail sur son fonctionnement, ce qui rend toute comparaison (y compris expérimentale) impossible.

3 Sélection simultanée d'index et de vues matérialisées

Le principe général de notre stratégie de sélection simultanée d'index et de vues matérialisées est représenté à la Figure 1. Rappelons qu'une vue matérialisée est une requête nommée dont les données sont stockées sur disque sous la forme d'une table. La sélection d'index peut donc se faire sur les tables de base ainsi que sur les vues matérialisées. Nous procédons comme suit pour proposer une configuration d'index et de vues matérialisées pertinents :

- extraction d'une charge des requêtes représentative,
- construction de l'ensemble des vues matérialisées candidates à partir de la charge,
- construction de l'ensemble des index candidats à partir de la charge et des vues matérialisées candidates,
- sélection simultanée d'index et de vues matérialisées,
- construction de la configuration finale d'index et de vues matérialisées.

L'extraction de la charge en entrée de notre approche s'effectue à partir du journal des requêtes exécutées sur les données de l'entrepôt. La charge que nous considérons est un ensemble de requêtes de projection, sélection et jointure. De telles requêtes sont composées d'opérations de jointures, de prédicats de sélection et d'opérations d'agrégation. Nous appliquons ensuite une stratégie de sélection isolée de vues matérialisées que nous avons développée et qui est basée sur la classification non supervisée des requêtes Aouiche (2005). Cela permet de construire un ensemble de vues candidates pertinent pour la charge. La classification des requêtes peut être vue comme une sorte de compression de la charge Chaudhuri *et al.* (2002). Cela permet d'assurer la scalabilité de notre approche. En effet, les charges ont tendance à être volumineuses et leur coût de traitement est par conséquent important. Au lieu de réaliser l'optimisation directement à partir de la charge, il est plus judicieux de le faire à partir d'une charge compressée qui conserve les relations existant entre les requêtes de la charge initiale.

Nous avons également développé une stratégie de sélection isolée d'index basée sur la recherche de motifs fréquents fermés Aouiche *et al.* (2003, 2005), qui permet de construire un ensemble d'index pertinent pour les requêtes de la charge et les vues matérialisées candidates générées à l'étape précédente. Notons que, notre approche étant modulaire, nous pourrions utiliser toute autre méthode de sélection isolée d'index ou de vues matérialisées. Finalement, nous appliquons notre stratégie de sélection simultanée de vues matérialisées et d'index, que nous détaillons dans les sections suivantes.

Afin d'illustrer notre propos, nous nous basons sur la charge de requêtes de la Figure 2 et les vues matérialisées et les index candidats qui en découlent (Figures 3 et 4, respectivement).

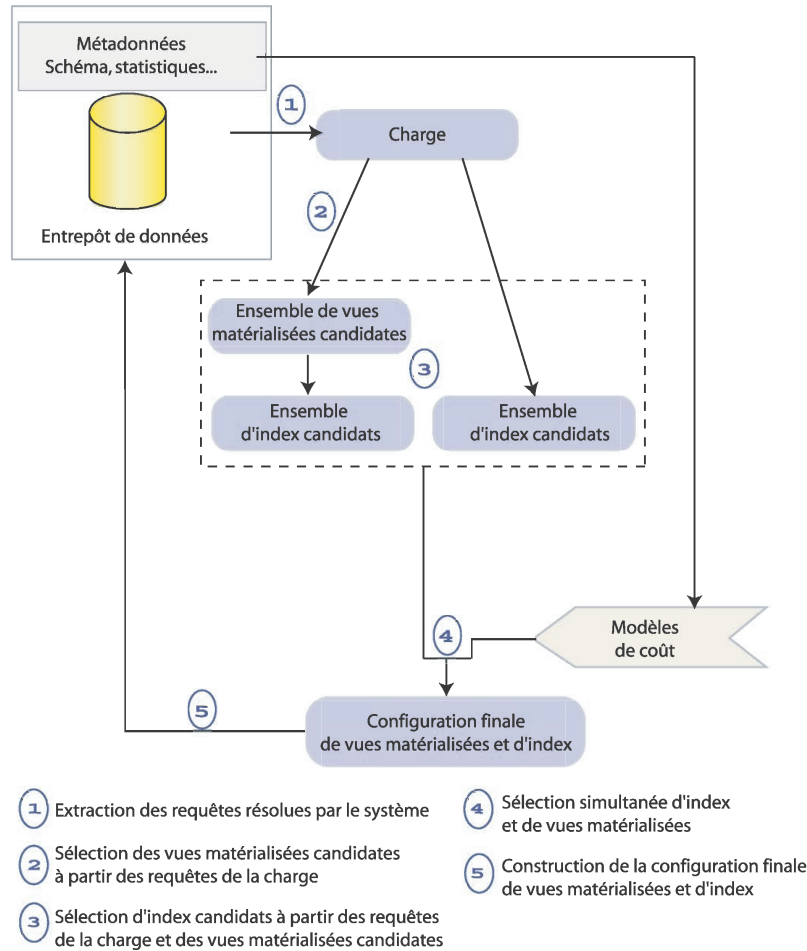


FIG. 1 – Principe de notre sélection simultanée d'index et de vues matérialisées

Nous modélisons les relations existant entre les requêtes, les vues matérialisées et les index candidats à l'aide de trois matrices : *requêtes-vues*, *requêtes-index* et *vues-index*, que nous décrivons dans les sections suivantes.

3.1 Matrice requêtes-vues

La matrice requêtes-vues modélise les relations entre les requêtes de la charge et les vues matérialisées qui en sont extraites, c'est-à-dire les vues exploitées par au moins une requête de la charge. Cette matrice peut être vue comme le résultat de la réécriture des requêtes de la charge en fonction des vues matérialisées. Les lignes et les colonnes de cette matrice sont les requêtes de la charge et les vues matérialisées recommandées par notre stratégie de sélection de vues, respectivement. Le terme général de la matrice est égal à un si une requête donnée

Sélection simultanée d'index et de vues matérialisées

exploite une vue et à zéro sinon. Le Tableau 1 illustre un exemple de matrice requêtes-vues composée de huit requêtes et de neuf vues matérialisées recommandées pour ces requêtes.

<i>q</i> ₁	select sales.time_id, sum (amount_sold) from sales, times where sales.time_id = times.time_id and times.time_fiscal_year = 2000 group by sales.time_id	<i>q</i> ₅	select promotions.promo_name, sum (amount_sold) from sales, promotions where sales.promo_id = promotions.promo_id and promotions.promo_begin_date='30/01/2000' and promotions.promo_end_date='30/03/2000' group by promotions.promo_name
<i>q</i> ₂	select sales.prod_id, sum (amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id and promotions.promo_category = 'news paper' group by sales.prod_id	<i>q</i> ₆	select customers.cust_marital_status, sum (quantity_sold) from sales, customers, products where sales.cust_id = customers.cust_id and sales.prod_id = products.prod_id and customers.cust_gender = 'woman' and products.prod_name = 'shampooing' group by customers.cust_first_name
<i>q</i> ₃	select customers.cust_gender, sum (amount_sold) from sales, customers, products, where sales.cust_id = customers.cust_id and sales.prod_id = products.prod_id and customers.cust_marital_status = 'single' and products.prod_category = 'women' group by customers.cust_gender	<i>q</i> ₇	select products.prod_name, sum (amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id and products.prod_category='tee shirt' and promotions.promo_end_date='30/04/2000' group by products.prod_name
<i>q</i> ₄	select products.prod_name, sum (amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id and promotions.promo_category = 'TV' group by products.prod_name	<i>q</i> ₈	select channels.channel_desc, sum (quantity_sold) from sales, channels where sales.channel_id = channels.channel_id and channels.channel_class = 'Internet' group by channels.channel_desc

FIG. 2 – Exemple de charge

<i>v</i> ₁	<pre> create materialized view v₁ as select sales.time_id, times.time_fiscal_year, sum(amount_sold) from sales, times where sales.time_id = times.time_id group by sales.time_id, times.time_fiscal_year </pre>	<i>v</i> ₅	<pre> create materialized view v₅ as select sales.prod_id, products.prod_category, promotions.promo_category, sum(amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id group by sales.prod_id, products.prod_category, promotions.promo_category </pre>
<i>v</i> ₂	<pre> create materialized view v₂ as select sales.prod_id, sales.cust_id, channels.channel_desc, channels.channel_class, sum(quantity_sold) from sales, channels, products, customers where sales.prod_id = products.prod_id and sales.channel_id = channels.channel_id and sales.cust_id = customers.cust_id group by sales.prod_id, sales.cust_id, channels.channel_desc, channels.channel_class </pre>	<i>v</i> ₆	<pre> create materialized view v₆ as select channels.channel_class, products.prod_name, channels.channel_desc, products.prod_category, sum(sales.quantity_sold), sum(sales.amount_sold) from sales, channels, products where sales.prod_id = products.prod_id and sales.channel_id = channels.channel_id group by channels.channel_class, products.prod_name, products.prod_category, channels.channel_desc </pre>
<i>v</i> ₃	<pre> create materialized view v₃ as select customers.cust_first_name, products.prod_name, products.prod_category, customers.cust_gender, customers.cust_marital_status, sum(sales.quantity_sold) from sales, customers, products where sales.cust_id = customers.cust_id and sales.prod_id = products.prod_id group by customers.cust_first_name, products.prod_name, products.prod_category, customers.cust_gender, customers.cust_marital_status </pre>	<i>v</i> ₇	<pre> create materialized view v₇ as select sales.prod_id, products.prod_category, channels.channel_desc, promotions.promo_name, promotions.promo_begin_date, promotions.promo_end_date, products.prod_name, sum(sales.quantity_sold), sum(sales.amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id and sales.channel_id = channels.channel_id group by sales.prod_id, products.prod_category, channels.channel_desc, promotions.promo_name, promotions.promo_begin_date, promotions.promo_end_date, products.prod_name </pre>
<i>v</i> ₄	<pre> create materialized view v₄ as select products.prod_name, products.prod_category, promotions.promo_category, sum(amount_sold) from sales, products, promotions where sales.prod_id = products.prod_id and sales.promo_id = promotions.promo_id group by products.prod_name, products.prod_category, promotions.promo_category </pre>		

FIG. 3 – Vues matérialisées candidates

Sélection simultanée d'index et de vues matérialisées

index	attributs indexés
i_1	promotions.promo_category
i_2	channels.channel_desc
i_3	channels.channel_class
i_4	customers.cust_marital_status
i_5	customers.cust_gender
i_6	times.time_begin_date
i_7	times.time_end_date
i_8	times.fiscal_year
i_9	products.prod_name
i_{10}	products.prod_category
i_{11}	promotions.promo_name
i_{12}	customers.cust_first_name

FIG. 4 – *Index candidats*

	v_1	v_2	v_3	v_4	v_5	v_6	v_7
q_1	1	0	0	0	0	0	0
q_2	0	0	0	1	0	0	0
q_3	0	0	1	0	0	0	0
q_4	0	0	0	1	0	0	0
q_5	0	0	0	0	0	0	1
q_6	0	0	1	0	0	0	0
q_7	0	0	0	0	0	0	1
q_8	0	1	0	0	0	1	0

TAB. 1 – *Matrice requêtes-vues*

3.2 Matrice requêtes-index

La matrice requêtes-index permet d'identifier les index construits sur les tables de la base. La matrice requêtes-index peut être vue comme la réécriture des requêtes de la charge en fonction des index recommandés par un algorithme de sélection d'index. Les lignes et les colonnes de cette matrice sont respectivement les requêtes de la charge et les index sélectionnés à partir de cette charge. Le terme général de la matrice est égal à un si une requête donnée exploite un index et à zéro sinon. Le Tableau 2 illustre un exemple de matrice requêtes-index composée de huit requêtes et de douze index recommandés pour ces requêtes.

3.3 Matrice vues-index

La matrice vues-index identifie les index construits sur les vues matérialisées recommandées par l'algorithme de sélection de vues. Les lignes et les colonnes de cette matrice sont respectivement les vues matérialisées candidates préalablement sélectionnées et les index candidats recommandés pour ces vues. Le terme général de la matrice est égal à un si une vue donnée exploite un index et à zéro sinon. Le Tableau 3 illustre un exemple de matrice vues-index composée de sept vues et de douze index recommandés pour ces vues.

	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}	i_{11}	i_{12}
q_1	0	0	0	0	0	0	0	1	0	0	0	0
q_2	1	0	0	0	0	0	0	0	0	0	0	0
q_3	0	0	0	1	1	0	0	0	0	1	0	0
q_4	1	0	0	0	0	0	0	0	1	0	0	0
q_5	0	0	0	0	1	1	0	0	0	0	1	0
q_6	0	0	0	0	1	0	0	0	1	0	0	1
q_7	0	0	0	0	0	0	1	0	1	1	0	0
q_8	0	1	1	0	0	0	0	0	0	0	0	0

TAB. 2 – *Matrice requêtes-index*

	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}	i_{11}	i_{12}
v_1	0	0	0	0	0	0	0	1	0	0	0	0
v_2	0	1	0	0	0	0	0	0	0	0	0	0
v_3	0	0	0	1	1	0	0	0	1	1	0	1
v_4	1	0	0	0	0	0	0	0	1	1	0	0
v_5	1	0	0	0	0	0	0	0	0	1	0	0
v_6	0	1	1	0	0	0	0	0	1	1	0	0
v_7	0	1	0	0	0	1	1	0	1	1	1	0

TAB. 3 – *Matrice vues-index*

4 Modèles de coût

Généralement, le nombre d'index et de vues candidats est d'autant plus important que la charge en entrée est volumineuse. La création de tous ces index et vues peut ne pas être réalisable en pratique à cause de la contrainte définie sur l'espace de stockage alloué aux index et vues. Pour pallier ces limitations, nous exploitons des modèles de coût permettant de ne conserver que les index et les vues les plus avantageux. Ces modèles estiment l'espace en octets occupé par les index et les vues, les coûts d'accès aux données à travers ces index et/ou ces vues et le coût de leur maintenance en terme de nombre d'entrées/sorties.

Nous avons développé des modèles qui estiment le coût d'accès aux données à travers des index *bitmap* de jointure, ainsi que les coûts de maintenance et de stockage de ces index Aouiche *et al.* (2005). Nous avons également présenté des modèles qui estiment le coût d'accès aux données à travers des vues matérialisées, ainsi que les coûts de maintenance et de stockage de ces vues. Dans la suite de la section, nous ne développons donc que les nouveaux modèles de coût développés pour ce travail relatifs aux index en B-arbre.

5 Calcul du bénéfice de matérialisation et d'indexation

Le bénéfice apporté par la sélection d'un objet (index, vue matérialisée ou vue matérialisée avec index) est défini comme la différence entre le coût des requêtes de la charge à un moment donné et le coût de ces mêmes requêtes suite à l'ajout de cet objet.

Sélection simultanée d'index et de vues matérialisées

Soient Q une charge de requêtes et $Config$ une configuration composée de vues matérialisées et d'index construits sur les tables de base ou les vues. QI , QV et VI sont respectivement les matrices requêtes-index, requêtes-vues et vues-index. Nous développons dans les sections suivantes le calcul du bénéfice pour un index ou une vue matérialisée.

5.1 Bénéfice apporté par un index

L'ajout d'un index à la configuration $Config$ peut conduire à plusieurs alternatives résumées dans le Tableau 4. En effet, l'ajout d'un index donné à la configuration $Config$ peut améliorer de façon directe le coût des requêtes de la charge ou indirectement à travers des vues auxquelles cet index est associé.

	$VI[v, i] = 1$	$VI[v, i] = 0$
$v \in Config$	min (bénéfice de matérialisation, bénéfice d'indexation de v)	bénéfice d'indexation
$v \notin Config$	—	bénéfice d'indexation

TAB. 4 – Bénéfice apporté par l'ajout d'index

Le bénéfice d'indexation apporté par l'ajout d'un index i est calculé comme suit :

$$benefice(Q, Config \cup \{i\}) = \begin{cases} \frac{C(Q, Config) - C(Q, Config \cup \{i\})}{taille(i)} & \text{si } \forall v \in V, VI[v, i] = 0, \\ \frac{C(Q, Config) - C(Q, Config \cup \{i\} \cup V')}{taille(\{i\}) + \sum_{v' \in V'} taille(\{v'\})} & \text{si } V' = \{v \in Config, VI[v, i] = 1\} \neq \emptyset \\ 0 & \text{sinon.} \end{cases}$$

5.2 Bénéfice apporté par une vue matérialisée

L'ajout d'une vue matérialisée à la configuration $Config$ peut conduire aux alternatives résumées dans le Tableau 5. En effet, l'ajout d'une vue donnée à la configuration $Config$ peut améliorer de façon directe le coût des requêtes de la charge ou de façon collaborative avec les index associés à cette vue.

	$VI[v, i] = 1$	$VI[v, i] = 0$
$i \in Config$	—	bénéfice d'indexation
$i \notin Config$	min (bénéfice d'indexation, bénéfice de matérialisation)	bénéfice de matérialisation

TAB. 5 – Bénéfice apporté par l'ajout de vue

Le bénéfice de matérialisation apporté par la vue matérialisée v est calculé comme suit :

$$benefice(Q, Config \cup \{v\}) = \begin{cases} \frac{C(Q, Config) - C(Q, Config \cup \{v\})}{taille(v)} & \text{si } \forall i \in I, VI[v, i] = 0, \\ \frac{C(Q, Config) - C(Q, Config \cup \{v\} \cup I')}{taille(\{v\}) + \sum_{i' \in I'} taille(\{i'\})} & \text{si } I' = \{i \in Config, VI[v, i] = 1\} \neq \emptyset \\ 0 & \text{sinon.} \end{cases}$$

6 Algorithme de sélection simultanée d'index et de vues matérialisées

Notre algorithme de sélection d'index et de vues matérialisée (Algorithme 1) est basé sur une recherche gloutonne dans l'ensemble O des objets obtenus en réalisant l'union de l'ensemble d'index candidats I et de l'ensemble de vues candidates V ($O = I \cup V$). Soit la fonction objectif F définie comme suit :

$$F_{/Config}(\{o_i\}) = \text{bénéfice}(Q, Config \cup (\{o_i\})) - \beta C_{maintenance}(\{o_i\})$$

où o_i peut être un index ou une vue et $\beta = |Q| p(o_i)$ estime le nombre de mises à jour de o_i . La probabilité de mise à jour $p(o_i)$ est égale à $\frac{1}{\text{nombre d'éléments de } O} \frac{\% \text{rafraîchissement}}{\% \text{interrogation}}$, où le ratio $\frac{\% \text{rafraîchissement}}{\% \text{interrogation}}$ représente la proportion de rafraîchissement par rapport à la proportion d'interrogation de l'entrepôt de données.

Algorithme 1 Sélection simultanée de vues matérialisées et d'index

```

1:  $Config \leftarrow \emptyset$ 
2:  $O \leftarrow I \cup V$ 
3: répéter
4:    $o_{max} \leftarrow \emptyset$ 
5:    $taille(o_{max}) \leftarrow 0$ 
6:    $benefice_{max} \leftarrow 0$ 
7:   pour tout  $o_i \in O - Config$  faire
8:     si  $F_{Config}(\{o_i\}) > benefice_{max}$  alors
9:        $benefice_{max} \leftarrow F_{Config}(\{o_i\})$ 
10:       $o_{max} \leftarrow \arg \max(F_{Config}(\{o_i\}))$ 
11:      {l'ensemble  $o_{max}$  peut contenir des index et des vues matérialisées}
12:   fin si
13:   fin pour
14:   si  $F_{Config}(o_{max}) > 0$  alors
15:     si  $index(o_{max}) = \text{vrai}$  alors
16:        $taille(o_{max}) \leftarrow taille_{index}(o_{max})$ 
17:     sinon
18:       si  $vue(o_{max}) = \text{vrai}$  alors
19:          $taille(o_{max}) \leftarrow taille_{vue}(o_{max})$ 
20:       sinon
21:          $o_{max} = i_{max} \cup v_{max}$ 
22:          $taille(o_{max}) \leftarrow taille(o_{max}) + taille_{index}(i_{max}) + taille_{vue}(v_{max})$ 
23:       fin si
24:     fin si
25:   fin si
26:    $S \leftarrow S - taille(o_{max})$ 
27:    $Config \leftarrow Config \cup (o_{max})$ 
28: jusqu'à ( $F_{Config}(o_{max}) \leq 0$  ou  $O - Config = \emptyset$  ou  $S \leq 0$ )

```

Sélection simultanée d'index et de vues matérialisées

Soit S l'espace disque alloué par l'administrateur de l'entrepôt de données pour stocker les vues matérialisées et les index. À la première itération de notre algorithme, les valeurs de la fonction objectif F sont calculées pour chaque index ou vue de l'ensemble O . Le coût d'exécution de toutes les requêtes de la charge Q est égal au coût d'exécution de ces requêtes à partir des tables de base sans indexation ni vue. L'ensemble o_{max} de vues et/ou d'index qui maximise F , s'il existe ($F_{Config}(\{o_{max}\}) > 0$), est alors ajouté à l'ensemble $Config$ si l'espace de stockage S n'est pas atteint. Si l'ajout est fructueux, l'espace S est diminué de l'espace de stockage occupé par o_{max} . L'espace occupé par o_{max} dépend de son contenu (index et/ou vue). Nous utilisons les fonctions booléennes $index(o_{max})$ et $vue(o_{max})$ qui renvoient vrai si o_{max} est un index ou une vue (lignes 15 à 25 de l'Algorithme 1), respectivement.

Les valeurs de la fonction objectif F sont ensuite recalculées pour chaque élément restant dans $O - Config$, car elles dépendent des vues et des index sélectionnés présents dans $Config$. C'est cela qui permet de prendre en compte les interactions qui peuvent exister entre les index et les vues matérialisées. Rappelons que cette interaction est implicitement incluse dans le calcul du bénéfice, qui exploite les matrices requêtes-index QI , requêtes-vues QV et vues-index VI .

Nous répétons ces itérations jusqu'à ce qu'il n'y ait plus d'amélioration de la fonction objectif ($F_{Config}(o_{max}) \leq 0$), que tous les index et vues aient été sélectionnés ($O - Config = \emptyset$) ou que la limite d'espace de stockage soit atteinte ($S \leq 0$).

7 Expérimentations

Afin de valider notre stratégie de sélection simultanée d'index et de vues matérialisées, nous l'avons expérimentée sur un entrepôt de données test implanté au sein du SGBD Oracle 9i. Nos expérimentations ont été réalisées sur un PC sous Windows XP Pro doté d'un processeur Pentium 4 à 2.4 GHz, d'une mémoire centrale de 512 Mo et d'un disque dur IDE de 120 Go.

Notre entrepôt de données test est composé d'une table de faits **Sales** et de cinq tables dimensions **Customers**, **Products**, **Promotions**, **Times** et **Channels**. Le Tableau 6 détaille le nombre de n-uplets et la taille en Mo de chacune des tables de cet entrepôt.

Table	Nombre de n-uplets	Taille (Mo)
Sales	16 260 336	372,17
Customers	50 000	6,67
Products	10 000	2,28
Times	1 461	0,20
Promotions	501	0,04
Channels	5	0,000 1

TAB. 6 – Caractéristiques de l'entrepôt de données test

Nous avons mesuré le temps d'exécution des requêtes de la charge de la Figure 2 dans les cas suivants : sans index ni vue matérialisée, avec vues matérialisées, et avec index et vues matérialisées. La Figure 5 représente la variation de ce temps de réponse en fonction du pourcentage d'espace de stockage utilisé. Ce pourcentage est calculé par rapport à l'espace total

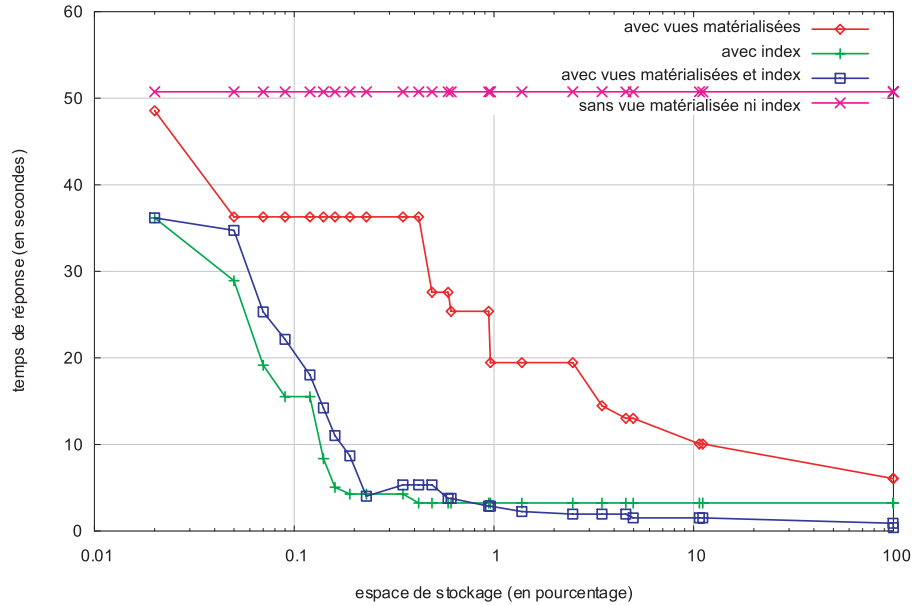


FIG. 5 – Résultats expérimentaux

occupé par tous les index et les vues obtenus en appliquant notre stratégie sans définir aucune contrainte d'espace. Afin de mieux visualiser les résultats, nous avons utilisé une échelle logarithmique sur l'axe des abscisses.

La Figure 5 montre que pour les valeurs élevées de l'espace de stockage, la sélection simultanée d'index et de vues est meilleure que la sélection isolée des index et vues. En revanche, nous constatons que pour les petites valeurs de l'espace de stockage, il peut arriver que la sélection d'index soit plus performante que la sélection simultanée d'index et de vues. Cela peut être expliqué par le fait qu'en général, la taille des index est significativement moins importante que celle des vues. Dans ce cas, on peut mettre dans le même espace (de petite taille) plus d'index que de vues et ainsi améliorer davantage le temps de réponse. En effet, intuitivement, plus on dispose d'index, plus on améliore les performances. Or, ce point n'avait pas été mis en lumière par les travaux antérieurs aux nôtres.

8 Conclusion et perspectives

Nous avons proposé dans cet article une nouvelle démarche de sélection simultanée d'index et de vues matérialisées dans les entrepôts de données. Notre approche prend réellement en compte l'interaction qui peut exister entre les index et les vues matérialisées et les traite simultanément afin de réduire le coût d'exécution de requêtes. Cela donne lieu à une sélection optimale des vues et des index en fonction de l'espace de stockage qui leur est alloué. En effet, nos résultats expérimentaux montrent que la sélection simultanée de vues matérialisées et d'in-

dex est plus performante que la sélection isolée de ces structures lorsque l'espace de stockage est raisonnablement grand.

Les perspectives ouvertes par ces travaux sont de plusieurs ordres. Tout d'abord, il est nécessaire de poursuivre nos expérimentations afin de confronter notre approche à l'existant. Malheureusement, la proposition *joint enumeration* d'Agrawal *et al.*, qui est la plus proche de la nôtre, n'est pas suffisamment documentée pour que nous puissions mener une telle étude à bien. En revanche, comparer nos travaux à ceux de Bellatreche *et al.*, à la fois en termes de gains de performance et de surcharge pour le système, pourrait nous permettre d'établir définitivement que la sélection conjointe d'index et de vues matérialisées est plus intéressante que leur sélection concurrente.

Par ailleurs, cette stratégie d'optimisation des performances s'applique dans un cas statique. La charge sur laquelle est effectuée l'optimisation peut devenir obsolète au bout d'un temps donné. Lorsque cela arrive, il faut resélectionner les index et les vues matérialisées. Les travaux traitant de la détection de sessions basés sur le calcul d'entropie Yao *et al.* (2005) pourraient s'appliquer pour déterminer le moment où il faut lancer la resélection des index et des vues.

Dans ces travaux, nous nous sommes également positionnés dans le cas où l'administrateur optimise les requêtes adressées au système par tous les utilisateurs confondus. Or, les besoins définis par différents profils d'utilisateurs (dans le cas des systèmes multi-utilisateurs) sont différents. Il serait donc plus pertinent d'appliquer nos stratégies sur des groupes de requêtes définis par les utilisateurs identifiés dans chaque profil.

Finalement, nous nous sommes restreints à utiliser seulement les index et les vues matérialisées comme mécanismes d'optimisation des performances. Or, d'autres structures de données peuvent être intégrées ou couplées avec les index et les vues, comme par exemple la gestion de cache, le regroupement et le partitionnement (fragmentation) Agrawal *et al.* (2004); Zilio *et al.* (2004); Bellatreche *et al.* (2005).

Références

- S. Agrawal, S. Chaudhuri, et V.R. Narasayya. Automated selection of materialized views and indexes in SQL databases. In *26th International Conference on Very Large Data Bases (VLDB 2000)*, Cairo, Egypt, pages 496–505, 2000.
- S. Agrawal, S. Chaudhuri, L. Kollár, A.P. Marathe, V.R. Narasayya, et M. Syamala. Database Tuning Advisor for Microsoft SQL Server 2005. In *30th International Conference on Very Large Data Bases (VLDB 2004)*, Toronto, Canada, pages 1110–1121, August–September 2004.
- K. Aouiche, J. Darmont, et L. Gruenwald. Frequent itemsets mining for database auto-administration. In *7th International Database Engineering and Application Symposium (IDEAS 2003)*, Hong Kong, China, pages 98–103, 2003.
- K. Aouiche, J. Darmont, O. Boussaïd, et F. Bentayeb. Automatic selection of bitmap join indexes in data warehouses. In *7th International Conference on Data Warehousing and Knowledge Discovery (DaWaK 2005)*, Copenhagen, Denmark, volume 3589 of LNCS, pages 64–73, August 2005.

- K. Aouiche. *Techniques de fouille de données pour l'optimisation automatique des performances des entrepôts de données*. PhD thesis, Université Lumière Lyon 2, 2005.
- X. Baril et Z. Bellahsene. Selection of materialized views : a cost-based approach. In *15th International Conference on Advanced Information Systems Engineering (CAiSE 2003)*, Klagenfurt, Austria, pages 665–680, 2003.
- L. Bellatreche, K. Karlapalem, et M. Schneider. On efficient storage space distribution among materialized views and indices in data warehousing environments. In *9th International Conference on Information and Knowledge Management (CIKM 2000)*, McLean, USA, pages 397–404, 2000.
- L. Bellatreche, Kamel Boukhalfa, et Mukesh Mohania. An evolutionary approach to schema partitioning selection in a data warehouse environmen. In *7th International Conference on Data Warehousing and Knowledge Discovery (DaWaK 2005)*, Copenhagen, Denmark, volume 3589 of *LNCs*, August 2005.
- S. Chaudhuri, A.K. Gupta, et V.R. Narasayya. Compressing SQL workloads. In *2002 ACM SIGMOD International Conference on Management of Data (SIGMOD 2002)*, Madison, Wisconsin, pages 488–499, June 2002.
- S. Choenni, H. M. Blanken, et T. Chang. Index selection in relational databases. In *5th International Conference on Computing and Information (ICCI 1993)*, Ontario, Canada, pages 491–496, 1993.
- S. Choenni, H. M. Blanken, et T. Chang. On the selection of secondary indices in relational databases. *Data Knowledge Engineering*, 11(3) :207–238, 1993.
- D. Comer. The difficulty of optimum index selection. *ACM Transactions on Database Systems (TODS)*, 3(4) :440–445, 1978.
- B. Dageville, D. Das, K. Dias, K. Yagoub, M. Zaït, et M. Ziauddin. Automatic SQL Tuning in Oracle 10g. In *30th International Conference on Very Large Data Bases (VLDB 2004)*, Toronto, Canada, pages 1098–1109, August-September 2004.
- M. R. Frank, E. Omiecinski, et S. B. Navathe. Adaptive and Automated Index Selection in RDBMS. In *3rd International Conference on Extending Database Technology (EDBT 1992)*, Vienna, Austria, volume 580 of *LNCs*, pages 277–292, 1992.
- M. Golfarelli, S. Rizzi, et E. Saltarelli. Index selection for data warehousing. In *4th International Workshop on Design and Management of Data Warehouses (DMDW 2002)*, Toronto, Canada, pages 33–42, 2002.
- H. Gupta et I.S. Mumick. Selection of views to materialize under a maintenance cost constraint. In *7th International Conference on Database Theory (ICDT 1999)*, Jerusalem, Israel, pages 453–470, 1999.
- H. Gupta et I.S. Mumick. Selection of views to materialize in a data warehouse. *IEEE Transactions on Knowledge and Data Engineering*, 17(1) :24–43, 2005.
- H. Gupta. *Selection and Maintenance of Views in a Data Warehouse*. PhD thesis, Stanford University, 1999.
- W.H. Inmon. *Building the Data Warehouse*. John Wiley & Sons, third edition, 2002.
- R. Kimball et M. Ross. *The Data Warehouse Toolkit : The Complete Guide to Dimensional Modeling*. John Wiley & Sons, second edition, 2002.

- S. Rizzi et E. Saltarelli. View materialization vs. indexing : Balancing space constraints in data warehouse design. In *15th International Conference on Advanced Information Systems Engineering (CAiSE 2003)*, Klagenfurt, Austria, pages 502–519, 2003.
- J.R. Smith, C.S. Li, et A. Jhingran. A wavelet framework for adapting data cube views for OLAP. *IEEE Transactions on Knowledge and Data Engineering*, 16(5) :552–565, 2004.
- G. Valentin, M. Zuliani, D. Zilio, G. Lohman, et A. Skelley. DB2 advisor : An optimizer smart enough to recommend its own indexes. In *16th International Conference on Data Engineering, (ICDE 2000)*, California, USA, pages 101–110, 2000.
- Q. Yao, J. Huang, et A. An. Machine learning approach to identify database sessions using unlabeled data. In *7th International Conference on Data Warehousing and Knowledge Discovery (DaWaK 2005)*, Copenhagen, Denmark, volume 3589 of *LNCS*, pages 254–255, August 2005.
- D.C. Zilio, J. Rao, S. Lightstone, G.M. Lohman, A. Storm, C. Garcia-Arellano, et S. Fadden. DB2 Design Advisor : Integrated Automatic Physical Database Design. In *30th International Conference on Very Large Data Bases (VLDB 2004)*, Toronto, Canada, pages 1087–1097, August-September 2004.

Summary

Indices and materialized views are physical structures that accelerate data access in data warehouses. However, these data structures generate some maintenance overhead. They also share the same storage space. The existing studies about index and materialized view selection consider these structures separately. In this paper, we adopt the opposite stance and couple index and materialized view selection to take into account the interactions between them and achieve an efficient storage space sharing. We develop cost models that evaluate the respective benefit of indexing and view materialization. These cost models are then exploited by a greedy algorithm to select a relevant configuration of indices and materialized views. Experimental results show that our strategy performs better than the independent selection of indices and materialized views.