

Bases de données

Pierre Senellart, Michaël Thomazo, Paul Boniol, Serge Abiteboul
Notes par Antoine Groudiev

Version du 15 février 2024

Table des matières

1	Gestion de données, modèle relationnel, algèbre et SQL	1
1.1	Introduction	1
1.1.1	Gestion de données	1
1.1.2	Types de SGBD	2
1.2	Modèle relationnel	3
1.3	Langages de requêtes	4
1.3.1	Algèbre relationnelle	4
1.3.2	SQL	5
2	Calcul relationnel, théorème de Codd, décidabilité, complexité	6
2.1	Calcul relationnel	6
2.1.1	Syntaxe	6
2.1.2	Sémantique	6
2.1.3	Sous-classes de requêtes	6
2.1.4	Théorème de Codd	7
2.1.5	Sous-classes de FO	7
2.2	Complexité de l'évaluation d'une requête	7
2.3	Analyse statique de requêtes	7

Introduction

Ce document est l'ensemble non officiel des notes du cours *Bases de données* du Département Informatique de l'ENS Ulm. Elles sont librement inspirées des notes de cours sous forme de présentation rédigées par Pierre Senellart.

1 Gestion de données, modèle relationnel, algèbre et SQL

1.1 Introduction

1.1.1 Gestion de données

Les bases de données ont de nombreuses applications (logiciels seuls, sites Web, ...), et servent à gérer des données. En particulier, elles sont utilisées pour :

- structurer des données utiles aux applications
- les conserver de manière persistante, même quand l'application ne tourne pas

- lancer des requêtes efficaces sur un large volume de données
- mettre à jour des données sans violer les contraintes structurelles de la base
- autoriser plusieurs utilisateurs à accéder et modifier les données, possiblement simultanément.

Il est par ailleurs souvent désirable d'accéder aux mêmes données depuis des applications distinctes, sur des ordinateurs différents.

On peut considérer comme exemple le système d'information d'un hôtel. Il est accessible depuis un logiciel interne (bureau d'accueil), un site web (hôtes), et une suite de logiciels de comptabilité. Il nécessite des données structurées pour les chambres, clients, réservations, ... Les données doivent rester stockées de manière pérenne, même en cas de coupure de courant. Les requêtes doivent être quasi instantanées, même avec un long historique de l'hôtel.

L'implémentation naïve nécessite un format de fichier pour stocker les données sur un disque, et un mécanisme de synchronisation, et de récupération en cas d'interruption de la synchronisation. Elle nécessite beaucoup de compétences différentes, de la programmation orientée objet, à des notions de réseau, de maîtrise de structures de données et algorithmes efficaces, de gestion en parallèle, ... Réimplémenter une telle structure à chaque fois qu'un tel outil est nécessaire est une perte de temps.

On utilise donc souvent un *Système de Gestion de Bases de Données* (SGBD), qui est un logiciel simplifiant le design d'une application manipulant des données, en fournissant un accès unifié aux fonctionnalités nécessaires pour la gestion des données, et ce peu importe l'application. Une *base de donnée* désigne une collection de données, gérée par un SGBD.

Un SGBD contient plusieurs fonctionnalités et caractéristiques :

- Indépendance physique
- Indépendance logique
- Facilité de l'accès aux données
- Optimisation des requêtes
- Intégrité logique
- Intégrité physique
- Partage de données
- Standardisation

1.1.2 Types de SGBD

De nombreux SGBD sont largement utilisés et ils ne fournissent pas tous toutes ces fonctionnalités. Les SGBDS peuvent être distingués selon le modèle de données utilisé, les choix techniques pris, la facilité d'utilisation, la possibilité de mise à l'échelle, ou encore l'architecture interne.

Les types de SGBD les plus courants sont les suivants :

- Relationnel : les données sont stockées sous forme de table ; ce modèle permet des requêtes complexes, utilisant le langage SQL
- XML : les données sont stockées sous forme d'arbre
- Graphes/Triplets
- Objets, inspiré de la POO
- Clé-valeur : très proche du fonctionnement d'une table de hachage
- Orienté colonne : modèle de donnée entre le Clé-valeur et Relationnel ; permet de parcourir très rapidement les données à l'intérieur de ce qui correspond à une colonne dans le modèle relationnel

Les SGBD ne suivant pas le modèle relationnel sont appelés **NoSQL**.

Dans le modèle relationnel, les données sont décomposées sous forme de *relation* (i.e. de tables), et utilisent un langage plus ou moins standardisé, le SQL. Les données sont stockées sur le disque, dans des tables de relation elles-mêmes stockées ligne par ligne. C'est un système centralisé, avec des propriétés de distributions limitées. Parmi les systèmes les plus répandus, on pourra citer Oracle, Microsoft SQL Server, MySQL, SQLite, et PostgreSQL. En TP, on utilisera la convention PostgreSQL.

Les SGBD relationnels classiques possèdent de nombreuses forces. Ils introduisent une indépendance entre le modèle de données et les structures de stockage, et les requêtes déclaratives et la façon dont elles sont exécutées. Ils sont capables d'exécuter des requêtes optimisées, finement optimisées, permettant de traiter jusqu'à des To de données. C'est une technologie mature, stable, et efficace, riche en fonctionnalités et interfaces. Le modèle relationnel permet de maintenir des contraintes d'intégrité et des transactions qui garantissent le contrôle de la concurrence des processus, l'isolation entre les utilisateurs, et la récupération en cas d'échec.

Les SGBD relationnels satisfont les propriétés suivantes, dites **ACID** :

- *Atomisme* : une transaction est soit exécutée en entier, soit pas exécutée du tout
- *Consistance* : les contraintes d'intégrité sont respectées
- *Isolation* : deux transactions concurrentes sont exécutées en série
- *Durabilité* : les données peuvent être récupérées même après un échec système

Ces propriétés viennent néanmoins également avec des faiblesses, comme leur incapacité à gérer des volumes de données extrêmement larges, de certains types, et impliquent un coût élevé des requêtes.

Les modèles NoSQL forment un écosystème très divers, abandonnant généralement les propriétés **ACID**, et parfois les requêtes complexes. Des nouveaux systèmes, dits **NewSQL**, essaient d'établir un compromis entre propriétés **ACID** et performances.

1.2 Modèle relationnel

On se fixe les ensembles infinis dénombrables suivants :

- \mathcal{L} , les étiquettes
- \mathcal{V} , les valeurs
- \mathcal{T} , les types, tel que $\forall \tau \in \mathcal{T}, \tau \subseteq \mathcal{V}$

Définition (Schéma relationnel d'arité n). Un schéma relationnel d'arité n est un n -uplet (A_1, \dots, A_n) où chaque A_i est un attribut, c'est-à-dire une paire $(L_i, \tau_i) \in \mathcal{L} \times \mathcal{T}$ où les L_i sont tous distincts.

Définition (Schéma de base de données). Un schéma de base de données est défini par un ensemble fini d'étiquettes $L \subseteq \mathcal{L}$ (les noms des relations), où chaque étiquette L est associée à un schéma de relation.

Exemple. On peut par exemple considérer que \mathcal{L} est l'ensemble des chaînes de caractères alphanumériques commençant par une lettre. \mathcal{V} est l'ensemble des suites de bits, et \mathcal{T} est formé des types *INTEGER*, *REAL*, *TEXT*, et *DATA*, ... Un schéma de base de données est formé de deux noms de relations *Client* et *Reservation*. Un *Client* est de la forme :

$$((id, \text{INTEGER}), (name, \text{TEXT}), (email, \text{TEXT}))$$

et une réservation est de la forme :

$$((id, \text{INTEGER}), (guest, \text{INTEGER}), (room, \text{INTEGER}), (arrival, \text{DATE}), (nights, \text{INT.}))$$

Définition (Instance d'un schéma de relation). Une *instance d'un schéma de relation* (auss appelé une *relation sur ce schéma*) $((L_1, \tau_1), \dots, (L_n, \tau_n))$ est un ensemble fini $\{t_1, \dots, t_k\}$ de tuples de la forme $t_j = (v_{j1}, \dots, v_{jn})$, où $\forall j, \forall i, v_{ij} \in \tau_i$.

Définition (Instance d'un schéma de base de données). Une *instance d'un schéma de base de données* (ou plus simplement *base de donnée sur ce schéma*) est une fonction qui associe à chaque nom de relation une instance du schéma de relation correspondant.

Remarque. Le terme *relation* est utilisé de manière ambiguë pour parler aussi bien d'un schéma de relation ou d'une instance d'un schéma de relation.

Si $A = (L, \tau)$ est le i -ème attribut d'une relation R , et que t est le n -ième tuple d'une instance de R , on note $t[A]$ (ou $t[L]$) la valeur de la i -ème composante de t . De manière similaire, si \mathcal{A} est un k -uplet d'attributs parmi les n attributs de R , $t[\mathcal{A}]$ est le k -uplet formé à partir de t en concaténant les $t[A]$ pour $A \in \mathcal{A}$.

On peut ajouter à un schéma relationnel des *contraintes d'intégrité* de natures différentes, pour définir une notion de *validité* d'une instance.

Définition (Clé). Un uplet d'attributs \mathcal{A} sur un schéma de relation R est une clé s'il ne peut pas exister deux uplets distincts t_1 et t_2 dans une instance de R tels que $t_1[\mathcal{A}] = t_2[\mathcal{A}]$.

Définition (Clé étrangère). Un k -uplet d'attributs \mathcal{A} d'un schéma de relation R est une *clé étrangère* référençant un k -uplet d'attributs \mathcal{B} d'un schéma de relation S si pour toutes les instances I^R et I^S de R et S , et pour tout uplet t de I^R , il existe un *unique* tuple t' de I^S avec $t[\mathcal{A}] = t'[\mathcal{B}]$.

On peut également définir des contraintes vérifiées, c'est-à-dire une condition arbitraire sur les valeurs des attributs d'une relation (en appliquant à chaque tuple de l'instance cette relation).

Exemple. En reprenant l'exemple précédent, on peut donner les contraintes suivantes :

1.3 Langages de requêtes

1.3.1 Algèbre relationnelle

On utilise des *langages algébriques* pour exprimer des requêtes. Une expression d'algèbre relationnelle produit une nouvelle relation à partir des relations de la base de données. Chaque opérateur prend 0, 1 ou 2 sous-expressions. Les opérateurs principaux sont les suivants :

Opérateur	Arité	Description	Condition
R	0	Nom de relation	$R \in \mathcal{L}$
$\rho_{A \rightarrow B}$	1	Renommage	$A, B \in \mathcal{L}$
$\Pi_{A_1 \dots A_n}$	1	Projection	$A_1 \dots A_n \in \mathcal{L}$
σ_φ	1	Sélection	φ est une formule
\times	2	Produit en croix	
\cup	2	Union	
\setminus	2	Différence	
\bowtie_φ	2	Jointure	φ est une formule

FIGURE 1 – Opérateurs principaux

Nom de relation R Retourne la relation (table) spécifiée en argument.

Renommage $\rho_{A \rightarrow B}$ Prend une sous-expression, et renomme une colonne.

Projection $\Pi_{A_1 \dots A_n}$ Prend en argument une séquence d'attributs, et conserve en sortie simplement les attributs spécifiés de la relation donnée.

Sélection σ_φ Prend une formule de sélection pouvant être n'importe quelle combinaison booléenne de comparaisons d'attributs et de constantes, et conserve uniquement les tuples où la condition booléenne est satisfaite.

Produit cartésien \times Retourne le produit cartésien ensembliste des deux ensembles obtenus par les sous-expressions.

Union \cup Retourne l'union des ensembles obtenus par les sous-expressions.

Différence \setminus Retourne l'ensemble obtenu par la première sous-expression, privé de l'ensemble obtenu par la seconde sous-expression. Cet opérateur ajoute de l'expressivité aux requêtes, par exemple avec des requêtes non monotones.

Jointure \bowtie_φ Prends deux tables et une condition de jointure, et retourne les tuples où la condition de jointure est vérifiée. C'est un opérateur crucial dans la perception et l'usage du modèle relationnel, mais est en fait redondant. Il peut être vu comme une combinaison des opérateurs de renommage, produit cartésien, sélection et projection. Si R et S ont pour attributs \mathcal{A} et \mathcal{B} , on note $R \bowtie S$ la *jointure naturelle* de R et S , où la formule de jointure est $\bigwedge_{A \in \mathcal{A} \cap \mathcal{B}} A = A$.

Toutes les expressions de l'algèbre relationnelle ne sont pas valides, et la validité d'une requête dépend en général du schéma de la base de données. Par exemple, si un nom de relation n'existe pas, si l'on essaye d'unifier deux relations avec des attributs différents, ou encore si l'on essaye de produire une table par produit, jointure, ou renommage qui contient deux attributs avec le même nom. L'implémentation de la gestion de ces erreurs par les SGBD est très variable.

En sémantique "sac", qui est utilisée en pratique par les SGBD, toutes les opérations retournent des multiensembles, avec répétitions ; en particulier, la projection et l'union peuvent introduire des multiensembles quand des relations initiales sont définies.

De multiples extensions ont été proposées à l'algèbre relationnelle en ajoutant des fonctionnalités supplémentaires. En particulier, l'extension *agrégation et groupement* des résultats (Klug, 1982 et Libkin, 2004) est beaucoup utilisée en pratique (voir `GROUP BY` en SQL).

1.3.2 SQL

Le *Structured Query Language* (SQL), est un langage plus ou moins standardisé pour interagir avec un SGBD. Seules quelques différences syntaxiques sont remarquables d'une version à une autre.

C'est un langage *déclaratif*, c'est-à-dire qu'il n'est pas nécessaire de penser à la complexité de l'opération, simplement à ce que l'on veut. C'est le SGBD qui prend en charge l'optimisation de la requête.

Les mots-clés et identifiants sont généralement insensibles à la casse. En pratique, les mots clés sont écrits en majuscule, et les identifiants commencent par une majuscule. Les commentaires sont introduits par `--`. Dans certains contextes, on termine les commandes par `;`, mais le point-virgule ne fait pas partie de la commande, il est simplement demandé par l'environnement (par exemple, pour un client en ligne de commande).

NULL La valeur spéciale NULL signifie l'absence de valeur, et introduit une logique tri-valuée étrange : **True**, **False**, et **NULL**. La comparaison normale avec NULL renvoie toujours NULL. On peut vérifier la présence de valeur avec **IS NULL** ou **IS NOT NULL**.

Langage de définition de données

Contraintes

Mises à jour

Requêtes

2 Calcul relationnel, théorème de Codd, décidabilité, complexité

2.1 Calcul relationnel

Le calcul relationnel est un langage logique utilisé pour exprimer des requêtes. Il est constitué de formules logiques du premier ordre, sans symboles de fonctions, mais avec des symboles de relation sur le schéma de base de donnée (plus des prédicats de comparaison). Il approche les requêtes avec une perspective non-nommée et non typée.

2.1.1 Syntaxe

Soient \mathcal{X} un ensemble de variables, un ensemble \mathcal{V} de valeurs et un schéma de base de données S .

- Pour chaque relation $R \in S$ d'arité n , et pour tout $(\alpha_1, \dots, \alpha_n) \in (\mathcal{X} \cup \mathcal{V})^n$, alors $R(\alpha_1, \dots, \alpha_n) \in \text{FO}$.
- Prédicat d'égalité : $\alpha = \alpha' \in \text{FO}$, pour $(\alpha, \alpha') \in (\mathcal{X} \cup \mathcal{V})^2$
- Pour toutes $(\varphi_1, \varphi_2) \in \text{FO}$:
 - $\varphi_1 \wedge \varphi_2 \in \text{FO}$
 - $\varphi_1 \vee \varphi_2 \in \text{FO}$
 - $\neg \varphi_1 \in \text{FO}$
 - $\forall x \varphi_1 \in \text{FO}$
 - $\exists x \varphi_1 \in \text{FO}$

Les variables libres de $\varphi \in \text{FO}$ sont les variables x apparaissant dans φ et qui ne sont pas quantifiées par $\forall x$ ou $\exists x$.

2.1.2 Sémantique

2.1.3 Sous-classes de requêtes

- Requête conjonctive (CQ, *Conjunctive query*) : requête de calcul relationnel sans \vee, \neg .
- Requête positive (PQ, *Positive query*) : requête sans \neg . Elle ne conserve que les requêtes monotones, c'est-à-dire que lorsque la base de donnée augmente en taille, la requête ne peut retourner que plus d'éléments.

- Union de requêtes conjonctives (UCQ, *Union of conjunctive queries*) : cas spécifique des requêtes positives lorsque les \vee et \wedge forment une DNF (Forme Normale Disjonctive).
- Requête booléenne : requête sans variable libre.

2.1.4 Théorème de Codd

Théorème (Codd, 1972). *L'algèbre relationnelle et le calcul relationnel sont équivalents :*

- *Pour toute requête d'algèbre relationnelle q sur un schéma S , il existe une requête de calcul relationnel Q sur S telle que pour toute base de donnée D sur S , $q(D) = Q(D)$.*
- *Pour toute requête de calcul relationnel Q sur un schéma S , il existe une requête d'algèbre relationnelle q sur S telle que pour toute base de donnée D sur S , $q(D) = Q(D)$.*

De plus, la conversion d'un formalisme à un autre peut être fait en temps polynomial.

Ce théorème est très important, car il permet d'utiliser un formalisme déclaratif pour spécifier des requêtes logiques ou en SQL. Ces requêtes sont ensuite compilées grâce à la transformation de Codd vers un formalisme algébrique. Ces requêtes algébriques sont ensuite optimisées, en utilisant les propriétés de l'algèbre relationnelle. Les requêtes optimisées peuvent ensuite être évaluées en exploitant le fait que chaque opérateur de l'algèbre relationnelle peut facilement être implémenté.

2.1.5 Sous-classes de FO

Les requêtes conjonctives (CQ) sont équivalents aux requêtes d'algèbre relationnelle sans \cup et \setminus , et dans lesquelles σ n'a pas de disjonction.

Les unions de requêtes conjonctives (UCQ) sont équivalentes aux requêtes positives (PQ), avec une sur-croissance exponentielle, et sont équivalentes aux requêtes algébriques sans \setminus .

2.2 Complexité de l'évaluation d'une requête

2.3 Analyse statique de requêtes

Conclusion

Le calcul relationnel a pour but de développer un langage de requête logique *équivalent à l'algèbre relationnelle*, ou au cœur de SQL. Grâce à celui-ci, l'évaluation de requêtes est très efficace ; cependant, il ne traite pas toutes les requêtes en temps polynomial, d'où la nécessité des langages de requêtes récurrents. L'analyse statique, souvent indécidable, est traitée en restreignant l'ensemble des requêtes aux requêtes conjonctives, à des fins d'optimisation.