

GGMD - Exercices applicatifs 2

Évaluation de requêtes réparties

UCBL - Département Informatique de Lyon 1 – 2023

L'objectif de ces exercices est de vous sensibiliser au calcul de coût de transfert de données d'une requête répartie.

On considère toujours la base de données *insee_deces* qui a été fragmenté et dont le schéma global est le suivant :

- *personne* (idp, nom, prenom, datenaiss, #lieunaiss, datedeces, #lieudeces, age)
- *region* (reg, nom, #cheflieu, zone)
- *departement* (dep, nom, #cheflieu, #reg)
- *commune* (com, nom, #dep)
- *mairie* (#codeInsee, cp, nomOrga, nomCom, email, tel, url, adresse, latitude, longitude, dateMaj)

reposant sur des fragments définis comme suit :

$$\begin{aligned}
 region_i &= \sigma_{zone=i}(region), i \in \{1; 2; 3\} \\
 departement_i &= departement \bowtie_{reg} region_i, i \in \{1; 2; 3\} \\
 commune_i &= commune \bowtie_{dep} departement_i, i \in \{1; 2; 3\} \\
 personne_age &= \Pi_{idp,age}(personnes) \\
 personne_info_naiss &= \Pi_{idp,nom,prenom,datenaiss,lieunaiss}(personnes) \\
 personne_info_deces &= \Pi_{idp,nom,prenom,datedeces,lieudeces}(personnes) \\
 pers_naiss_i &= personne_info_naiss \bowtie_{lieunaiss=com} commune_i, i \in \{1; 2; 3\} \\
 pers_deces_i &= personne_info_deces \bowtie_{lieudeces=com} commune_i, i \in \{1; 2; 3\} \\
 mairie_i &= mairie
 \end{aligned}$$

A noter, que l'on considère ici la base contenant des données 'propres' avec les contraintes d'intégrités valides.

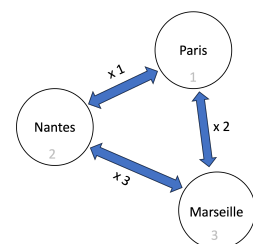
Pour rappel, tous les fragments hébergés sur le site de Paris (P) sont indexés par 1, ceux hébergés sur le site de Nantes (N) sont indexés par 2 et ceux hébergés sur le site de Marseille (M) sont indexés par 3. Pour les besoins de l'exercice, la table *mairie* n'a pas été répliquée mais migrée sur le site de Paris.

On suppose que chaque attribut vaut 0.5 Ko. (hypothèse simplificatrice pour les calculs)

On spécifie les coûts de transfert par la fonction $tr(T, X, Y)$ qui retourne le coût pour transférer d'un ensemble de tuples T depuis le site X vers le site Y.

Pour l'exercice, on suppose que :

$$\begin{aligned}
 tr(T, P, N) &= \text{card}(T) \times \text{taille}(T) \\
 tr(T, P, M) &= 2 \times \text{card}(T) \times \text{taille}(T) \\
 tr(T, N, M) &= 3 \times \text{card}(T) \times \text{taille}(T) \\
 tr(T, X, Y) &= tr(T, Y, X)
 \end{aligned}$$



A noter que des informations statistiques sur les données et les opérateurs sont précisées en annexe.

Exercice 1 : Jointures inter-site

On considère la requête Q1b émise sur le site de Nantes : "Donner les personnes (idp, nom, prénoms) nées dans la région 'Hauts-de-France' et décédées dans la région 'Occitanie'".

NB : En guise de révision, vous pourrez refaire ce TD, avec la requête Q1 pour les personnes nées en 'Auvergne-Rhône-Alpes' et décédées dans la région 'Pays de la Loire'.

Vous allez considérer deux plans d'exécution différents. Le plan P1 consiste à transférer l'ensemble des fragments sur le site de Nantes pour effectuer les calculs sur le seul site de Nantes. Le plan P2 consiste à effectuer le maximum des calculs sur les sites stockant les données.

Plan P1

- Transfert de $region_1$, $departement_1$, $commune_1$ et $pers_naiss_1$ de Paris à Nantes.
- Transfert de $region_3$, $departement_3$, $commune_3$ et $pers_naiss_3$ de Marseille à Nantes.
- Sur le site de Nantes :
 - Calcul de la sélection $\sigma_{nom='Hauts-de-France'}(region_1)$, on obtient T11.
 - Calcul de la jointure entre T11 et $departement_1$, on obtient T12.
 - Calcul de la jointure entre T12 et $commune_1$, on obtient T13.
 - Calcul de la jointure entre T13 et $pers_naiss_1$, on obtient T14.
 - Calcul de la projection $\Pi_{idp,nom,prenoms}(T14)$, on obtient T1.
- Sur le site de Nantes :
 - Calcul de la sélection $\sigma_{nom='Occitanie'}(region_3)$, on obtient T21.
 - Calcul de la jointure entre T21 et $departement_3$, on obtient T32.
 - Calcul de la jointure entre T32 et $commune_3$, on obtient T33.
 - Calcul de la jointure entre T33 et $pers_naiss_3$, on obtient T34.
 - Calcul de la projection $\Pi_{idp,nom,prenoms}(T34)$, on obtient T3.
- Sur le site de Nantes :
- Calcul de l'intersection entre T1 et T3, on obtient le résultat.

Plan P2

- Sur le site de Paris :
 - Calcul de la sélection $\sigma_{nom='Hauts-de-France'}(region_1)$, on obtient T11.
 - Calcul de la jointure entre T11 et $departement_1$, on obtient T12.
 - Calcul de la jointure entre T12 et $commune_1$, on obtient T13.
 - Calcul de la jointure entre T13 et $pers_naiss_1$, on obtient T14.
 - Calcul de la projection $\Pi_{idp,nom,prenoms}(T14)$, on obtient T1.
- Transfert de T1 de Paris à Nantes.
- Sur le site de Marseille :
 - Calcul de la sélection $\sigma_{nom='Occitanie'}(region_3)$, on obtient T31.
 - Calcul de la jointure entre T31 et $departement_3$, on obtient T32.
 - Calcul de la jointure entre T32 et $commune_3$, on obtient T33.
 - Calcul de la jointure entre T33 et $pers_naiss_3$, on obtient T34.
 - Calcul de la projection $\Pi_{idp,nom,prenoms}(T34)$, on obtient T3.
- Transfert de T3 de Marseille à Nantes.
- Sur le site de Nantes :
- Calcul de l'intersection entre T1 et T3, on obtient le résultat.

1. A partir des statistiques proposées en Annexe, calculer les coûts des transferts de données du plan P1.

2. A partir des statistiques proposées en Annexe, calculer les coûts des transferts de données du plan P2.
3. Quelle modification proposeriez-vous au plan P2 pour réduire son coût de transfert ? Proposer un plan P3 et calculer son coût.

Exercice 2 : Algorithme des semi-jointures

On considère la requête Q'1 émise depuis Paris : "Donner les personnes (idp, nom, prénoms, date-naiss, lieu-naiss, date-deces, lieu-deces) nées dans la région 'Hauts-de-France' et décédées dans la région 'Occitanie'"

On considère le plan d'exécution P4 suivant :

- Sur le site de Marseille :
 - Calcul de la sélection $\sigma_{nom='Occitanie'}(region_3)$, on obtient T31.
 - Calcul de la jointure entre T31 et $departement_3$, on obtient T32.
 - Calcul de la jointure entre T32 et $commune_3$, on obtient T33.
 - Calcul de la jointure entre T33 et $pers_deces_3$, on obtient T34.
 - Calcul de la projection $\Pi_{idp,lieudeces,datedeces}(T34)$, on obtient T3.
 - Transfert de T3 de Marseille à Paris.
 - Sur le site de Paris :
 - Calcul de la jointure entre T1 et T3, on obtient le résultat.
1. Vérifier que le coût de transfert du plan P4 est de 6 629 577 Ko
 2. Proposer P5 un plan d'exécution de Q'1 intégrant une gestion de la jointure selon l'algorithme des semi-jointures.
 3. Calculer le coût de P5.

Annexe

On suppose disposer des informations suivantes :

- H1 : $Card(region) = 18$;
- H2 : $Card(region_1) = 4$;
- H3 : $Card(region_2) = 9$;
- H4 : $Card(region_3) = 5$;
- H5 : $Card(departement) = 101$;
- H6 : $Card(departement_1) = 28$;
- H7 : $Card(departement_2) = 28$;
- H8 : $Card(departement_3) = 45$;
- H9 : $Card(commune) = 37\ 600$
- H10 : $Card(commune_1) = 12\ 849$;
- H11 : $Card(commune_2) = 8\ 028$;
- H12 : $Card(commune_3) = 14\ 123$;
- H13 : $Card(personnes) = 24\ 811\ 296$
- H14 : $Card(pers_naiss_1) = 8\ 634\ 189$;
- H15 : $Card(pers_naiss_2) = 5\ 196\ 351$;
- H14 : $Card(pers_naiss_3) = 7\ 444\ 568$;
- H15 : $Card(Q1) = 54\ 357$ tuples
- H16 : $Card(Q2) = 509$ tuples
- H17 : $Card(Q3) = 3\ 036\ 725$ tuples

