

UNIVERSITÉ DE MONTPELLIER
M2 GÉNIE LOGICIEL

Mini projet : Un moteur d'évaluation de
requêtes en étoile

HAI914I : GESTION DES DONNEES AU DELA DE SQL (NoSQL)

Étudiants :

Denis BEAUGET

Hayaat HEBIRET

Encadrant :

Scharffe FRANÇOIS

Rodriguez OLIVIER

Année : 2021 – 2022



UNIVERSITÉ
DE MONTPELLIER



Table des matières

1	GIT	2
2	Création du dictionnaire	2
3	Création des index	3
3.1	Version naïve :	3
3.2	Version optimisée :	4
4	Récupération des requêtes	5
5	Exécution des requêtes	6
6	Résultat obtenu	7
7	Bilan	8

1 GIT

Notre projet est disponible sur le liens suivant : https://github.com/Beauget/HAI914I_Projet

2 Création du dictionnaire

Pour la création du dictionnaire nous avons utilisé une structure HashMap :

```
1 //Structure du dictionnaire
2 public class DictionaryHashMap extends Dictionary {
3
4     private HashMap<Integer,String> dictionary= new HashMap<Integer,String>();
5
6     ....
7 }
```

Nous avons décidé de **mettre à jour le dictionnaire et les index durant la lecture du fichier rdf**. A chaque données rdf nous vérifions si l'objet, le prédicat et le sujet existe déjà dans le dictionnaire. Par exemple :

- L'objet existe : boolean **object** deviens "True" et int **objectIndex** récupère sa position dans le dictionnaire.
- L'objet n'existe pas : boolean **object** reste "False" et int **objectIndex** est défini au moment ou il est ajouté au dictionnaire.

```
1 public Integer[] updateDictionary(Statement st) {
2     double start = System.currentTimeMillis();
3     boolean subject = false;
4     boolean predicate = false;
5     boolean object = false;
6     Integer [] toAdd = {0,0,0};
7     int subjectIndex = 0 , predicateIndex= 0, objectIndex = 0;
8
9     Iterator it = dictionary.entrySet().iterator();
10
11     while (it.hasNext()) {
12
13         HashMap.Entry dic = (Entry) it.next();
14         String value = (String) ((Entry) dic).getValue();
15         Integer key = (Integer) ((Entry) dic).getKey();
16
17         if(value.equals(st.getSubject().toString())) {
18             subject = true;
19             subjectIndex= key;
20         }
```

```

21         if(value.equals(st.getPredicate().toString())) {
22             predicate = true;
23             predicateIndex = key;
24         }
25
26         if(value.equals(st.getObject().toString())) {
27             object = true;
28             objectIndex = key;
29         }
30     }
31
32
33     if(!subject) {
34         Integer compteurHm = dictionary.size()+1;
35         dictionary.put(compteurHm, st.getSubject().toString());
36         toAdd[0] = compteurHm;
37     }
38
39     ....
40     return toAdd;
41 }

```

Ici on observe une partie de l'algorithme , on remarque bien le traitement **“à la volé”**. De plus, updateDictionary **renvoie l'index SPO de la donnée en rdf**.

Si on on génère et regarde **dictionary.txt** on obtient bien le format désiré avec une clé unique pour chaque valeur du dictionnaire.

```

1 1 : http://db.uwaterloo.ca/~galuc/wsdbm/User0
2 2 : http://schema.org/birthDate
3 3 : "1988-09-24"
4 4 : http://db.uwaterloo.ca/~galuc/wsdbm/userId
5 5 : "9764726"

```

3 Création des index

3.1 Version naïve :

Dans un premier temps nous avons eu une approche "naïve" de la création des index. On n'avait tout simplement ajouté les combinaisons de chaque index sous la forme de simple Triplet. Nous avons posé la question par mail car notre implémentation nous semblait incohérente (cela revenait à parcourir chaque index entièrement dans le pire des cas pour répondre à une requête. M.Rodriguez nous a fait un retour qui a comfirmer nos problèmes. Cette solution ne respectait pas le principe Hexastore et la notion d'arbre binaire.

```

1  package deprecated;
2
3  // Triplet = 1,2,3 order give by index
4  public class Triplet {
5
6      public int[] indexing;
7
8      public Triplet(int s, int p, int o) {
9          indexing = new int[] {s,p,o};
10     }
11
12     public int[] getIn() {
13         return indexing;
14     }
15
16
17     public String toString() {
18         String s = indexing[0] + " " + indexing[1] +
19             " " + indexing[2] + " " + "\n";
20         return s;
21     }
22
23 }

```

L'intégralité de notre première version naïve est disponible au sein du package "deprecated" de notre projet

3.2 Version optimisée :

Avec le retour nous nous sommes reposés la question de comment créer les index de la façon la plus "optimisée" possible. Nous avons décidé de garder la notion de type à part entière sur nos index avec un attribut pour identifier l'ordre.

Pour la structure de données on a donc laissé tombé les Triplet pour des HashMap (pour favoriser l'accès aux branches)

```

1  public class IndexOpti extends Index{
2      private String order;
3      private HashMap<Integer,HashMap<Integer,List<Integer>>> index;
4      ...
5  }

```

L'ajout d'un index deviens alors une peut plus complexe car il faut maintenant s'assurer de ne pas écraser une Hashmap déjà existante lors de sa mise à jour.

```

1  public void add(Integer first,Integer second,Integer third) {
2      double start = System.currentTimeMillis();
3      //verifie si la premier hashmap existe
4      HashMap<Integer,List<Integer>> secondHashMap = index.get(first);
5      //si la premiere hash n'existe pas je l'ajoute
6      if(secondHashMap==null) {
7          index.put(first,new HashMap<Integer,List<Integer>>());
8      }
9      //verifie si la deuxieme hashMap existe
10     List<Integer> currentArray = index.get(first).get(second);
11     if (currentArray==null){
12         secondHashMap = index.get(first);
13         currentArray = new ArrayList<Integer>();
14         currentArray.add(third);
15         secondHashMap.put(second,currentArray);
16         index.put(first, secondHashMap);
17     }
18     else {
19         secondHashMap = index.get(first);
20         currentArray.add(third);
21         secondHashMap.put(second,currentArray);
22         index.put(first, secondHashMap);
23     }
24
25     double end = System.currentTimeMillis();
26
27     execIndex += ((end - start) / 1000);
28
29 }

```

Par exemple, si on on génère et regarde **SPO.txt** on obtient bien le format désiré.

```

1  8198 , 8001 : [8201]
2  8198 , 7995 : [2318]
3  8198 , 12011 : [10788]
4  8198 , 7996 : [8199]
5  8198 , 12012 : [11592, 11608, 11600, 11621, 11788]
6  8198 , 7998 : [8200]
7  8198 , 5535 : [5851]

```

4 Récupération des requêtes

Une fois le dictionnaire et les index crée nous utilisons le **Pattern Processor** qui se charge de faire les requêtes. Dans l'optique de mettre en place **un programme qui peut être maintenue**, le dictionnaire et

l'index étende d'une **superClasse**.

Cela aurait permis de faciliter la mise à jour du processeur si nous avions eu plus de temps pour proposer différentes solutions.

```
1  public class Processor {
2      Dictionary dictionary;
3      ArrayList<Index> indexes;
4      ArrayList<Query> queries;
5      ...
6  }
```

5 Exécution des requêtes

Processor.doQueries() traite toute ses requêtes en utilisant les indexs et le dictionnaire et renvoie un **String** contenant le résultat des requêtes (si celle ci existe).

Si une requête contient plus d'un triplet, nous récupérons le résultat de chaque triplet et faisons une **intersection**.

```
1  public class Processor {
2      ...
3      public String doQueries(){
4          double start = System.currentTimeMillis();
5          StringBuilder builder = new StringBuilder();
6          builder.append("Query , Answer(s) \n");
7          for(Query q : queries) {
8              List<String> answer = doAQuery(q);
9              if(answer!=null) {
10                 if(answer.size()!=0) {
11                     builder.append(q.getRealQuery()+" ,");
12                     for(String s :answer) {
13                         builder.append(s.toString()+" ");
14                     }
15                     builder.append("\n");
16                 }
17             }
18         }
19         double end = System.currentTimeMillis();
20         execQuery += ((end - start) / 1000);
21         return builder.toString();
22     }
23     ...
24     private List<Integer> intersectionAnswers(Query q){
25         List<List<Integer>> allAnswers = doQueriesInIntegers(q);
26         List<Integer> output =null;
27         if( allAnswers !=null){
28             output = new ArrayList<Integer>(allAnswers.get(0));
```

```

29         if(output!=null) {
30             for (int i = 1; i < allAnswers.size(); i++) {
31                 List<Integer> l1 = allAnswers.get(i);
32                 if(l1 !=null) {
33                     List<Integer> l2 = intersection(l1,output);
34                     output = l2;
35                 }
36             }
37             else
38                 return null;
39         }
40     }
41     else
42         return null;
43 }
44 return output;
45 }
46 ...
47 }

```

6 Résultat obtenu

Tout les fichiers sont disponible dans le dossier "output" de notre projet à conditions que la demande d'écriture a été faite.

Extrait du fichier **Answers.txt**

```

1  SELECT ?v0 WHERE {
2  ?v0 <http://schema.org/eligibleRegion>
3  <http://db.uwaterloo.ca/~galuc/wsdbm/Country137> . }
4  http://db.uwaterloo.ca/~galuc/wsdbm/Offer195
5  http://db.uwaterloo.ca/~galuc/wsdbm/Offer311
6  http://db.uwaterloo.ca/~galuc/wsdbm/Offer642
7  http://db.uwaterloo.ca/~galuc/wsdbm/Offer642
8
9  ...
10
11 SELECT ?v0 WHERE {      ?v0 <http://purl.org/dc/terms/Location>
12 <http://db.uwaterloo.ca/~galuc/wsdbm/City0> .
13 ?v0<http://schema.org/nationality>
14 <http://db.uwaterloo.ca/~galuc/wsdbm/Country0> .
15 ?v0 <http://db.uwaterloo.ca/~galuc/wsdbm/gender>
16 <http://db.uwaterloo.ca/~galuc/wsdbm/Gender1> . }
17 http://db.uwaterloo.ca/~galuc/wsdbm/User178
18 http://db.uwaterloo.ca/~galuc/wsdbm/User831

```


Pour le calcul des métriques on peut remarqué que **TempsDeLectureDonnee** est NON_DISPONIBLE. En effet la lecture est intimement lié à la création du dictionnaire et de l'index. De plus comme le fichier contenant les données rdf n'est jamais enregistré dans une structure, il n'y a pas de sens de compter son temps de lecture.

Extrait du fichier **option_1_data_output.csv**

```
1  nomDuFichierDonnee,nomDuFichierRequete,
2  NombreDeTriplet,NombreDeRequete,
3  TempsDeLectureDonnee,TempsDeLectureRequete,
4  TempsDeCreationDico,NombreIndex,
5  TempsCreationIndex,TempsTotalEval,
6  TempsTotal
7  data/100K.nt,data/STAR_ALL_workload.queryset,
8  107338,1200,
9  NON_DISPONIBLE,389.0,
10 37604.0,6,
11 325.0,39649.0,
12 39651.0
```

7 Bilan

```
java -jar qengine-0.0.1-SNAPSHOT-jar-with-dependencies.jar
Veuillez entrer le path du fichier contenant les requêtes sparql(si vous souhaitez tester le code avec STAR_ALL_workload.queryset et 100K.nt écrire "default")
default
--- Bienvenue dans notre moteur de requête RDF ---

Options disponible (taper le chiffre correspondant à l'option) :

1 : Création du dictionnaire, index et réponses aux requêtes et du .csv contenant les temps de calculs et des informations sur les fichier
2 : Création du dictionnaire et temps d'exécution (SANS ECRITURE /output)
3 : Création du dictionnaire et temps d'exécution (AVEC ECRITURE /output)
4 : Création des index et temps d'exécution (SANS ECRITURE /output)
5 : Création des index et temps d'exécution (AVEC ECRITURE /output)
6 : Chargement + exécution des requêtes et temps d'exécution (SANS ECRITURE /output)
7 : Chargement + exécution des requêtes et temps d'exécution (AVEC ECRITURE /output)
8 : Toute les données des options précédentes (SANS ECRITURE)
9 : Toute les données des options précédentes (AVEC ECRITURE)
0 : Quittez l'application
```

FIGURE 1 – Lancement du programme depuis le terminal

Pour ce projet, notamment en raison de la deuxième partie du projet nous avons essayé d'extraire un maximum de métrique de notre code. D'une part pour avoir un retour concret du travail que nous avons effectué (un visuel sous un format textuel) et d'autre part pour connaître les limites de notre projet (quel métrique nous ne pouvions pas obtenir par exemple). Nous avons donc permis l'extraction du dictionnaire, des index (un fichier par index) et des réponses aux requêtes sous un format .txt pour avoir un format brut de tout nos résultats. Pour faciliter l'identification des fonctionnalités nous avons rajoutés un petit client CLI ainsi que le projet au format .jar.

Bien que nous pensons que notre projet soit conforme aux attentes nous aurions aimé parcourir plus en détail les possibilités d'optimisé celui ci.

Nous avons vu en cours qu'il est possible de mettre un place un dictionnaire dont **l'ordre des données**

permettait un parcours plus efficace de celui ci, notamment d'ignorer certaine positions car on est sûr qu'elle n'y est pas (par ordre alphabétique par exemple).

Il aurait également été intéressant **d'identifier la structure des requêtes** afin de ne pas répéter les solutions des triplet et de ne pas essayer d'obtenir une solution en passant par tout les indexes.