



Web Sémantique

Projet MovieRetriever

2 0 1 7

C O N T E N U

Ce rapport est une présentation de notre projet et la liste des spécifications pour celui-ci.

Par : Rémy AUDA et DRAVET Jean-Baptiste

Table des matières

Présentation projet	3
Description	3
Explication	3
Consignes	3
Ontologie.....	4
Les classes	4
Les propriétés d'objets.....	4
Le projet	5
Récupérateur des données	5
Architecture	5
Algorithmes.....	6
Chargement des données	6
Problèmes	6
Initialisation	6
Remplir l'ontologie.....	8
Création	9
Conclusion	10
« Akinator »	10
Architecture	10
Algorithmes.....	11
Idée principale.....	11
Système de poids	11
Conclusion	13

Présentation projet

Le but de notre projet est de réaliser un logiciel qui permet de retrouver un film en répondant par oui ou par non à des questions qui seront posées.

Description

Le principe de notre logiciel est, via une interface graphique java, de poser des questions à l'utilisateur à propos d'un film afin de retrouver le titre du dit film. Pour cela, nous allons créer notre propre ontologie simplifiée.

Il s'agit d'un logiciel dans l'esprit de l'Akinator, à la différence qu'avec un Akinator, nous sommes en connaissance du film et ce serait à l'Akinator de le deviner alors qu'ici c'est l'inverse : la personne n'est pas en connaissance du film et notre logiciel essaye d'aider la personne à trouver le titre du film.

Explication

Nous avons décidé de créer notre propre ontologie car celles existantes nous ont semblé trop complexes avec beaucoup trop d'informations. En effet lorsque l'on recherche un film, on a souvent peu d'informations à son sujet. Poser des questions à propos de la musique ou des costumes nous a donc paru inutile. De plus, les grosses ontologies ont un système complexe d'ID (identifiant) qui lie les différentes classes entre elles, ce qui nous semble trop complexe à implémenter dans le temps imparti. Cependant, nous avons implémenté un loader qui va chercher toutes les données des films présents dans une base du linked data, puis nous formatons ces données pour qu'elles correspondent au format de notre propre ontologie. Cette approche a le double avantage de nous permettre de créer notre propre ontologie et ainsi mieux comprendre tous les aspects de ce domaine, tout en nous permettant d'utiliser la puissance du linked data en allant chercher les informations nécessaires pour l'application directement sur le web. **Notre application nous semble idéale pour démontrer la puissance du linked data et du web sémantique. En effet, si nous essayerions de créer cette application avec des technologies classiques comme les bases de données par exemple, cela voudrait dire qu'à chaque sortie de films, il faudrait qu'un administrateur rentre manuellement toutes les données des films dans sa base de données, ce qui est impossible étant donné le volume énorme de données que cela représente. Au contraire, avec le web sémantique, chaque fois qu'un film sort et qu'il est inscrit sur une page wikipédia par exemple, notre application le prend automatiquement en compte en allant directement le chercher sur le web. La base de connaissance de notre application se développe ainsi en même temps que les informations présentent sur le web sans qu'aucune intervention manuelle ne soit nécessaire.**

Consignes

Notre projet est lié aux consignes suivantes :

- RDF data
- RDFS / OWL / SKOS vocabularies, rules
- SPARQL queries (and show the inferences made, taking into account the semantics of the vocabulary)

Nous avons décidé de réaliser notre ontologie en OWL avec l'aide du logiciel Protege.

Ontologie

Voici une représentation de notre ontologie minimaliste. Cette ontologie est utilisée uniquement dans le but de notre logiciel. Elle se justifie par l'utilisation que nous en faisons.

Les classes

Voici la liste des différentes classes que nous avons créées :

- Actor : classe pour représenter les acteurs d'un film. Chaque acteur sera uniquement décrit par son nom, les films dans lesquels il a joué et le personnage qu'il a joué dans chaque film.
- Award : représente les récompenses pour les films.
- Characters : représente les personnages des films. Ils sont joués par des acteurs.
- Country : représente le pays d'origine d'un film.
- Movie : classe qui représente les films. En relation avec toutes les autres classes de l'ontologie.
- Producer : représente les producteurs pour chaque film. Ils font partie de la Team.
- Scriptwriter : représente les scénaristes pour chaque film. Ils font partie de la Team.
- Team : représente l'équipe humaine d'un film.
- Type : représente le type de film (le genre auquel il appartient).

Les propriétés d'objets

Voici la liste des différentes propriétés que nous avons créées :

- Characters : relation en un personnage et un film. Définit qu'un personnage appartient bien à un film.
- OriginalCountry : relation entre un film et son pays d'origine.
- hasActor : relation entre un film et un de ses acteurs. Inverse à la propriété hasPlayIn.
- hasPlayIn : relation entre un acteur et un des films dans lequel il a joué.
- hasScriptWriter : relation entre un film et son scénariste. Inverse de la propriété wasScriptWriterIn.
- wasScriptWriterIn : relation entre un scénariste et un de ses films.
- isTypeOf : relation entre un film et un de ses types (action, aventures, ...)
- wasPlayedBy : relation entre un personnage et l'acteur qui le joue.
- wasProducerIn : relation entre un producteur et un de ses films. Inverse de la propriété wasProductBy.
- wasProductBy : relation entre un film et son producteur.
- wasRewardWith : relation entre un film et une de ses récompenses.

Nous avons décidé d'ajouter un poids sur chaque propriété et chaque valeur. Ces poids vont nous servir pour réaliser un algorithme de pseudo-génération des questions.

Cette partie sera développée ultérieurement dans ce rapport.

Le projet

Nous avons décidé de réaliser le projet en Java avec Jena pour la réalisation des requêtes SPARQL. Suite à la discussion avec notre intervenant, nous allons devoir diviser notre projet en deux projets bien distincts.

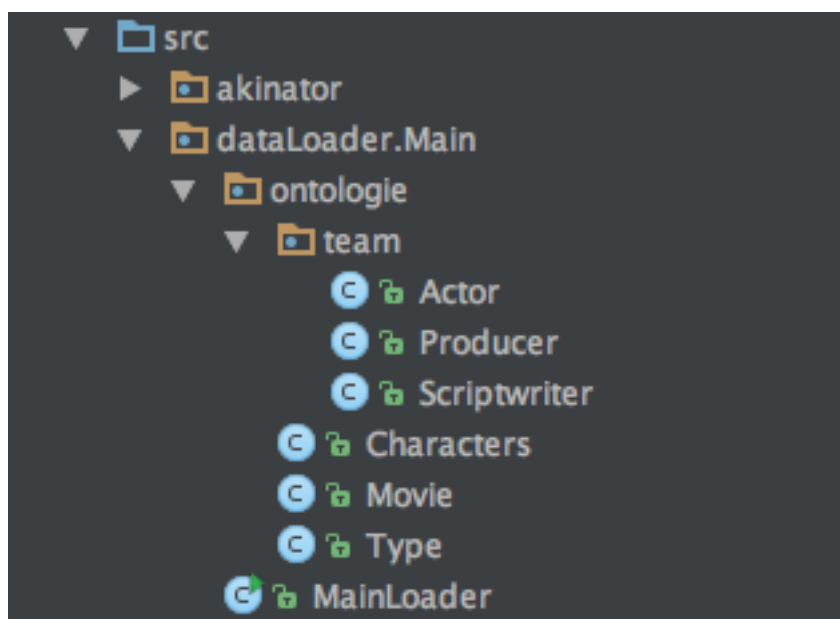
Récupérateur des données

Dans un premier temps, nous allons devoir créer un projet de récupération de données qui va, grâce à un dump d'un dataset du web, récupérer les données des bases de connaissances déjà existantes et remplir notre propre base de données. Pour se faire, il va falloir réaliser un mapping entre l'ontologie de la base de connaissance et la nôtre. Grâce au mapping et à des requêtes SPARQL, nous allons pouvoir remplir une base de données configurée pour notre propre ontologie.

Architecture

Pour réaliser le loader, nous avons réfléchi à une architecture pour répondre à la question de représentation et du mapping. Plus clairement, comment allons-nous faire la transition entre les données sous forme de tuple dans un dataset à des instances de données dans notre ontologie.

Voici un screenshot de l'architecture de notre sous-projet « Loader » :



Chaque classe portant un nom correspondant à une classe de notre ontologie va nous permettre de faire le traitement de la donnée en liaison direct avec la classe. Par exemple, « Actor » s'occupera de représenter les acteurs, d'ajouter les acteurs dans l'objet « Movie » et de rajouter les instances d'acteurs dans notre ontologie.

Donc nous allons représenter chaque instance par un objet (une instance de donnée pour un acteur va être utilisée pour créer l'objet acteur). Ensuite, l'objet Movie sera composé de plusieurs instances de différents objets (plusieurs types, plusieurs acteurs). Enfin, nous allons dans le « MainLoader » charger notre ontologie et la remplir avec tous les objets que nous avons créés.

Algorithmes

Chargement des données

Nous avons tout d'abord voulu récupérer les informations depuis un endpoint sparql de linked database Movie. Nous avons choisi cet endpoint car son ontologie était la plus proche de la nôtre.

Problèmes

Après plusieurs essais de récupération, nous sommes tombés sur un problème majeur : nos requêtes successives se faisaient déconnecter car le site pensait qu'on essayait de le DDOS (flux trop important de requêtes).

Après avoir changé notre façon de faire les requêtes, nous avons réussi à faire passer plus de requêtes mais nous avons aussi réussi à faire tomber le site.

Nous avons alors demandé l'aide de notre intervenant. Il nous a conseillé d'utiliser un dump de linked data movie, ce que nous avons fait. Nous avons enfin pu récupérer toutes les données que nous voulions.

Nous avons encore eu quelques problèmes d'URI mal formé dans le dump, que nous avons corrigé à la main.

Initialisation

Pour récupérer les données, nous avons d'abord dû charger les données. Pour cela nous avons utilisé la classe Model.

```
// initialisation //////////////////////////////////
Model model;
final String inputFileName = "File/RDFS_file/linkedmdb-latest-dump.nt";

// créer un modèle vide
model = ModelFactory.createDefaultModel();

// utiliser le FileManager pour trouver le fichier d'entrée
InputStream in = FileManager.get().open( inputFileName );

// lire le fichier owl
model.read(in,null, "N-TRIPLE");
```

Après avoir chargé le dataset, nous chargeons notre ontologie dans un model différent :

```
//load of our ontologie
Model m = ModelFactory.createDefaultModel();
final String inputFile = "/Users/titanium/Desktop/MoviesInformationRetriever/File/RDFS_file/OntologieMovie.owl";

InputStream test = FileManager.get().open(inputFile);
m.read(test, "RDF/XML");
System.out.println("fini load ontologie");
```

Maintenant, nous allons créer une arraylist de chaque classe pour récupérer les instances qui sont contenu dans le dataset et les transformer en objet :

```
System.out.println("Create Character list");
characteres = Characters.constructListOfCharacter(query, queryexec, resultSet, model);

System.out.println("Create producer list");
producers = Producer.constructListOfProducer(query, queryexec, resultSet, model);

System.out.println("Create type list");
types = Type.constructListOfType(query, queryexec, resultSet, model);

System.out.println("Create writer list");
scriptwriters = Scriptwriter.constructListOfWriter(query, queryexec, resultSet, model);

System.out.println("Create actor list");
actors = Actor.constructListOfActor(query, queryexec, resultSet, model);

System.out.println("Create movie list");
movies = loadMovieFromBigDataset(model);

System.out.println("add all to ontologie");
Movie.fillOntologie(movies, types, actors, characteres, producers, scriptwriters);

System.out.println("ends");
```

Chaque classe possède donc une méthode qui permet de récupérer les informations que nous voulons grâce une requête Sparql.

Voici un exemple de requête pour Acteur, qui ressemble beaucoup aux autres requêtes pour chaque classe :

```
public static String requestforinstance = "prefix mo: <http://data.linkedmdb.org/resource/movie/>" +
    "prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>" +
    "SELECT DISTINCT ?name WHERE" +
    " { ?uri mo:producer ?uriProducer. " +
    "?uri a mo:film. " +
    "?uri mo:initial_release_date ?date." +
    "?uri mo:runtime ?duration." +
    "?uri mo:country ?countrytemp." +
    "?countrytemp rdfs:label ?country." +
    "?uri mo:genre ?genre." +
    "?uri mo:writer ?uriWriter." +
    "?uri mo:actor ?uriactor." +
    "?uriactor rdfs:label ?name." +
    " }";
```

?uri représente l'uri d'un film. Pourquoi une tel requête pour les acteurs ? Car nous souhaitons obtenir uniquement les acteurs des films que nous aurons déjà présélectionné.

Les films présélectionnés devront avoir au minimum un pays d'origine, une date, une durée, et un titre. Tout film n'ayant pas au minimum ces 4 critères ne sera pas récupéré.

Donc pour avoir la liste d'acteur uniquement lié à ces films, nous allons utiliser les paramètres de la requête pour les films comme « filtre » (sans être filter de sparql) et utiliser ceux-ci pour filtrer les acteurs que nous voulons récupérer.

Ensuite nous allons construire une liste d'acteur avec les résultats de la requête :


```

public static ArrayList<Actor> constructListOfActor(Query q, QueryExecution qr, ResultSet r, Model m) {
    String finalRequest = requestForInstance;
    ArrayList<Actor> actors = new ArrayList<>();
    q = create(finalRequest);
    qr = QueryExecutionFactory.create(q, m);
    r = qr.execSelect();
    while (r.hasNext()) {
        QuerySolution binding = r.nextSolution();
        Literal name = binding.getLiteral("name");
        actors.add(new Actor(name.getString().replace("(Actor)", "")));
    }
    return actors;
}

```

Nous pouvons voir dans le code ci-dessous un code exploitant la librairie Jena permettant de manipuler le web sémantique.

Pour Movie, cela sera un peu différent. Nous allons d'abords récupérer chaque film et ensuite construire les différentes listes (acteurs, producteurs, scénaristes et types) liées au film en question :

```

/**
 * load file from dump and return an arrays of movies
 * @param model
 * @return
 */
public static ArrayList<Movie> loadMovieFromBigDataset(Model model){
    Query query = null;
    QueryExecution queryexec = null;
    ResultSet resultSet = null;
    ArrayList<Movie> movies;
    movies = Movie.constructListOfMovie(query, queryexec, resultSet, model);
    movies = Movie.refactorListOfMovie(movies);
    for (Movie movie: movies)
    {
        movie.addListOfType(Type.constructListOfTypeForMovie(movie.getUri(), query, queryexec, resultSet, model));
        movie.addListofScriptwriter(Scriptwriter.constructListofscriptwritersForMovie(movie.getUri(), query, queryexec, resultSet, model));
        movie.addListofProducer(Producer.constructListofProducerForMovie(movie.getUri(), query, queryexec, resultSet, model));
        movie.addListofActors(Actor.constructListofActorForMovie(movie.getUri(), query, queryexec, resultSet, model));
        movie.addListofCharacters(Characters.constructListofCharactersForMovie(movie.getUri(), query, queryexec, resultSet, model));
        // System.out.println(movie.toString());
    }
    return movies;
}

```

Les méthodes « addList. » permettent de récupérer les informations lié à l'Uri du film. Voilà par exemple la requête pour récupérer les types :

```

//permet de récupérer les types pour un film donné
public static String requestpart1 = "prefix mo: <http://data.linkedmdb.org/resource/movie/>" +
    "<http://www.w3.org/2000/01/rdf-schema#" +
    "SELECT DISTINCT ?name WHERE" +
    " { <";
public static String requestpart2 =
    "> mo:genre ?urigenre." +
    "?urigenre rdfs:label ?name." +
    " } ";

```

Remplir l'ontologie

Une fois que nous avons chargé nos fichiers et créé nos objets, nous allons créer les ressources dans notre ontologie. Pour se faire, pour chaque liste d'objet, pour chaque objet, nous allons créer une ressource et les propriétés qui vont avec et écrire dans notre fichier ontologie.


```
for (Type type:types
    ) {
    type.addToOntologie(m);
}

for (Actor actor:actors
    ) {
    actor.addToOntologie(m);
}

for (Characters characters1:characters
    ) {
    characters1.addToOntologie(m);
}

for (Producer producer:producers
    ) {
    producer.addToOntologie(m);
}

for (Scriptwriter scriptwriter:writers
    ) {
    scriptwriter.addToOntologie(m);
}

for (Movie movie:movies
```

Voici un exemple d'ajout de ressource pour un objet type :

```
// ajoute la liste d'instance de type dans l'ontologie
public void addToOntologie(Model m) {
    String prefixemo = "http://www.semanticweb.org/titanium/ontologies/2017/0/untitled-ontology-11#";
    String prefixerdfs = "http://www.w3.org/2000/01/rdf-schema#";
    String prefixerdf = "http://www.w3.org/1999/02/22-rdf-syntax-ns#";

    try {
        Resource resourceType = m.createResource(prefixemo + this.name);
        //type
        Property type = m.createProperty(prefixerdf + "type");
        resourceType.addProperty(type, prefixemo + "Type");
        //add title
        Property label = m.createProperty(prefixerdfs + "label");
        resourceType.addProperty(label, this.name);
    } catch (Exception e) {
    }
}
```

Création

Une fois toutes les ressources créées, nous allons créer le fichier ontologie en utilisant notre modèle qui contient toutes les ressources.

```
try {
    m.write(new FileOutputStream("/Users/titanium/Desktop/MoviesInformationRetriever/File/RDFS_file/OntologieMovie1.owl"));
} catch (FileNotFoundException e) {
    e.printStackTrace();
}
```

Conclusion

Maintenant que nous avons notre ontologie avec une certaine quantité de donnée, nous allons l'utiliser pour notre deuxième partie : l'Akinator.

Cette partie de récupération de données que nous venons de voir est très importante car elle illustre magnifiquement certains des avantages du web sémantique par rapport à d'autres technologies.

En effet, elle illustre le fait que nous pouvons charger des données directement dans le web. Pour le cas de notre application par exemple, à chaque fois que des films seront rajoutés sur le web, ces mêmes films seront donc ajoutés à notre application car nous chargeons les données directement depuis le web. Une application utilisant des bases de données par exemple ne pourrait bénéficier de ces avantages car à chaque fois qu'un nouveau film sortirait sur le web, il faudrait le rajouter à la main dans la base de données pour que l'application puisse prendre en compte ce nouveau film.

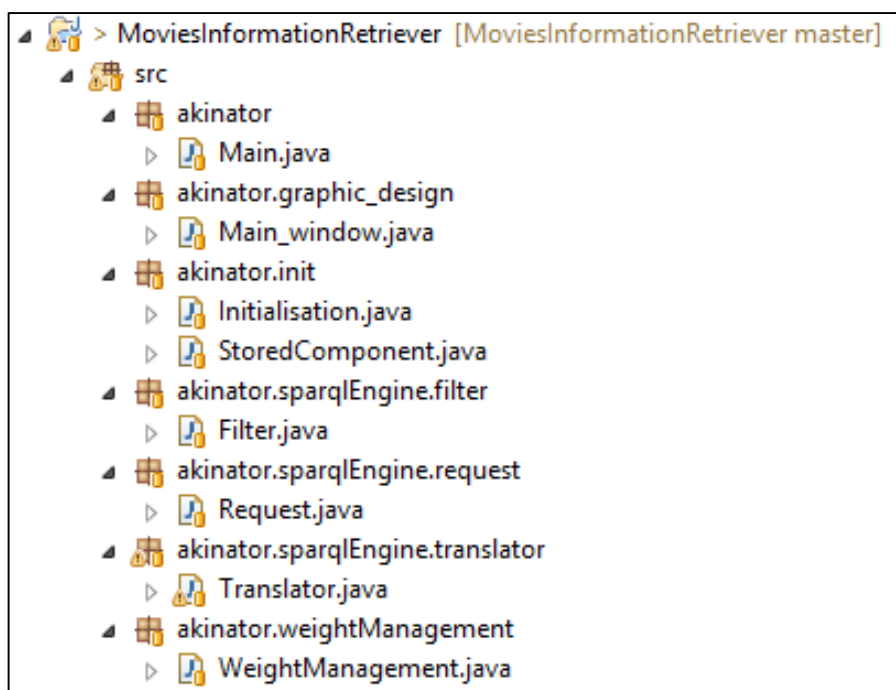
« Akinator »

Une fois les données récupérées, nous avons lancé le second projet. Il s'agit d'un « Akinator » qui va poser des questions à l'utilisateur en utilisant notre ontologie et nos données. Pour cela, nous avons introduit un système de poids. Chaque valeur et chaque propriété ont un poids. Ce poids va nous permettre de définir l'ordre des questions et la sélection des valeurs. Les propriétés qui auront le plus gros poids seront les premières utilisées pour la réalisation des questions. Une fois utilisé, nous allons mettre à jour le poids de la donnée avec des requêtes SPARQL de type DELETE/INSERT pour indiquer à notre algorithme que nous avons déjà utilisé cette propriété.

Cette partie fut également très importante car elle nous permet de comprendre la puissance du web sémantique et de SPARQL notamment, en nous permettant de chercher des informations diverses dans une très grande base de données pour en retirer des connaissances.

Architecture

Voici l'architecture pour le projet Akinator :



Algorithmes

Idée principale

Pour réaliser l'Akinator, nous avons créé notre propre algorithme :

- On commence par récupérer la propriété avec le plus gros poids et la valeur qui a le plus gros poids pour la propriété en question.
- On récupère le label de chacun, on pose la question en utilisant les labels.
- En fonction de la réponse, on diminue le poids de la propriété et on met à 0 le poids de la valeur.
- On crée une requête pour récupérer notre film en fonction de la réponse à la question.
- Si on obtient un seul résultat, on l'affiche.
- Si on obtient plusieurs résultats, on repart à l'étape 1.

Ce système impose de mettre un poids sur chaque propriété et chaque valeur

Système de poids

Le système de poids est basique pour l'instant (on choisit un peu arbitrairement le poids à attribuer sur chaque valeur, avec un générateur de nombre aléatoire) mais dans une prochaine version, on peut imaginer que ce poids serait calculé en fonction de la fréquence de référence pour instance : par exemple, plus un acteur a joué de film, plus son poids sera conséquent.

Au fur et à mesure de notre avancement, nous avons réalisé l'erreur que nous avons faite en ne créant pas de classe « poids » dans notre ontologie. En effet, nous avons mis le score du poids sur la propriété `rdfs:seeAlso` pour les valeurs et `rdfs:isDefinedBy` sur les

Ensuite, lors de la recherche par poids, nous allons faire ça dans l'ordre d'arriver normal, puis pour la décrémentation du poids, nous allons passer de « 1h30 OK » à « 1h30 NONE » pour ne plus le récupérer lorsque nous ferons une recherche pour la durée.

Voici un exemple de requête SPARQL du type DELETE/INSERT que nous avons dû implémenter pour le système de poids.

```
//Decrement the weight for the masterBranch specified by its label.
private void DeleteDataValue(String ValueWithoutOK, String valueWithOK){
    final String queryString = this.prefix+" DELETE { ?uri mo:wasReleasedIn ?value. } INSERT { ?uri mo:wasReleasedIn \"\"+ValueWithoutOK+"NONE\". } where { ?uri mo:wasReleasedIn \"\"+valueWithOK
    //System.out.println(queryString);
    UpdateRequest query = UpdateFactory.create(queryString);
    //System.out.println(this.mainQuery);
    UpdateAction.execute( query, Initialisation.getModel() );

    //this.weight_masterBranch= Integer.toString(max_weight);
}
```

```

t the property, the label and the weight for the masterBranch. Set also these three result in the array Property_Label_Weight_Triplet
private String getDataValue(String Label_MasterBranch){
    final String queryString = this.prefix+"Select distinct ?value where {?uri_value rdfs:label \""+Label_MasterBranch+"\"} .?uri_value rdfs:label ?label.?uri_value rdfs:isDefinedBy ?poids.?ur
    //System.out.print(queryString);
    String getDataValue = new String();
    Query query = QueryFactory.create(queryString);
    //System.out.println(this.mainQuery);
    QueryExecution qexec = QueryExecutionFactory.create(query, Initialisation.getModel());
    try{
        org.apache.jena.query.ResultSet results = qexec.execSelect();
        //System.out.println(results.getRowNumber());

        while (results.hasNext()){
            QuerySolution soln = results.nextSolution();
            Literal valueData = soln.getLiteral("value");
            getDataValue = valueData.toString();
        }
    }
    finally{
        qexec.close();
    }
    return getDataValue;
}
}

```

L'intérêt de l'approche des poids pour poser les questions pertinentes afin trouver le film, outre le fait que cela améliore l'application en elle-même, est qu'elle nous permet de pouvoir expérimenter tous les aspects et la puissance de SPARQL. En effet, cela nous permet d'expérimenter les requêtes SPARQL de type DELETE/INSERT et de type SELECT. Cela nous permet de voir comment faire des recherches dans les données du web et comment faire des modifications dans un graphe afin de modifier les données de notre base de connaissance.

Conclusion

En conclusion, ce projet fut très bénéfique pour comprendre l'intérêt et la puissance du linked data et du web sémantique. Outre le fait d'avoir élaboré une application fonctionnelle qui permet de retrouver les films des utilisateurs, ce projet nous a permis de passer en revue tous les différents aspects du web sémantique. Nous avons pu voir l'intérêt de rechercher les données directement sur le web avec notre implémentation du loader. Cela nous a permis de constater que ce type d'application qui nécessite d'avoir des données constamment mises à jours ne peut pas être implémentée avec des technologies classiques comme des bases de données alors qu'avec le web sémantique, cela devient tout à fait envisageable. Ce projet a été également l'occasion de créer notre propre ontologie owl et d'exploiter toutes ses propriétés. Enfin, nous avons pu mettre à profit nos connaissances en SPARQL afin de rechercher dans les grands volumes de données que nous possédions, les informations nécessaires à l'application et de pouvoir les manipuler selon nos besoins.

Vous pouvez avoir accès à l'intégralité du code source à cette adresse GitHub :

<https://github.com/dj103059/MoviesInformationRetriever>