# CONSTRUCTION D'UN EXTRACTEUR D'INFORMATION GÉOGRAPHIQUE AVEC R ET SPARQL

Juliette Delannoy

Projet encadré par Hadrien Commenges et Thomas Louail

Laboratoire de recherche Géographie-cités

Projet de développement informatique 2018

# Introduction

Années 2000 : début de l'utilisation du web sémantique

⟹ Structuration de la donnée

Interrogation automatique possible

Exemple : Dbpédia = base de données de Wikipédia

⟹ Interrogation avec mon interface

# Plan

- Web sémantique

- Dbpédia

- Langage SPARQL

- Interface
  - Fonctionnalités de base
  - Carte pour résultats spatialisés
  - Informations complémentaires

- Conclusion

# Web sémantique

= Web de la donnée
= Web 3.0

"The Semantic Web provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries"(W3C)

Fournir un **cadre** pour permettre le **partage** et la **réutilisation des données.**

☐ Structurer l'information uniformément pour l'universalité de l'accès et du partage

# Web sémantique

= Web de la donnée (brute)
= Web 3.0

□ Language universel et uniforme :

**RDF**  (Resource Description Framework)

(Cadre de description des ressources)
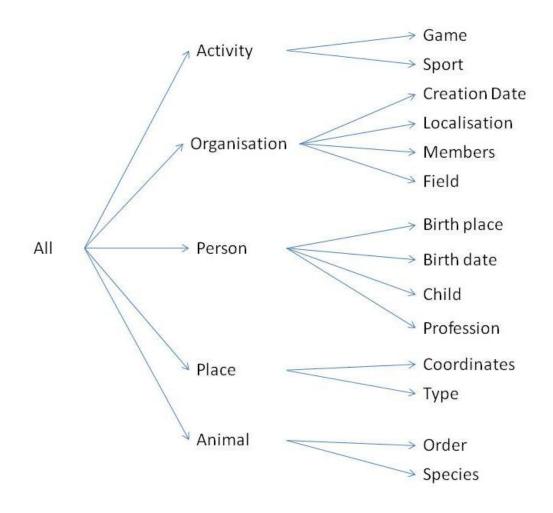
Sujet ⟶ Prédicat ⟶ Objet

Le chien s'appelle Bary.

```
lincolns:Robert
    a gen:Person ;
    gen:birthdate "1843-08-01"^^xsd:date .
lincolns:Mary
    a gen:Woman ;
    gen:child lincolns:Robert .
lincolns:Abraham
    a [ a owl:Class ;
      owl:intersectionOf (gen:Man gen:US_President)
      ] ;
    gen:child lincolns:Robert ;
    gen:married lincolns:Mary .
```

# Web sémantique

= Web de la donnée (brute)
= Web 3.0

- Ontologies

# Web sémantique

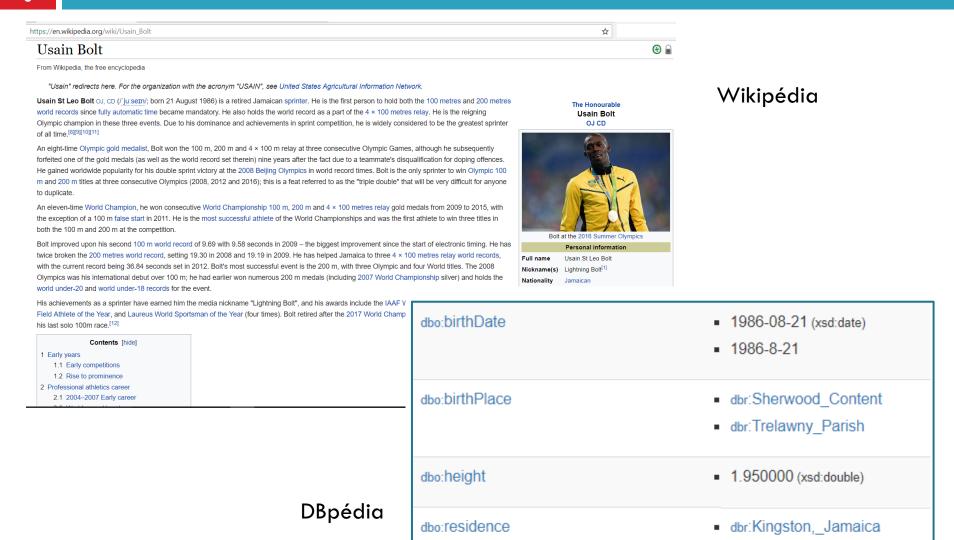But : Pouvoir interroger tout le web automatiquement

Web scraping inutile

Réalité : Structuration partielle

Ex : Comparateur de billets d'avion

Dbpédia : Base de données structurée

# Dbpédia et Wikipédia

Wikipédia

DBpédia

# Dbpédia et Wikipédia

- Dbpedia = Fiche signalétique de Wikipédia

- Permet des questions combinées :

Ex : Tous les athlètes nés entre 1810 et 1880.

Remplissage base de données automatique

Problèmes qui en résultent :

- Produit des erreurs
- Mauvaise structuration possible

# Interroger Dbpédia nécessite :

- ❑ De comprendre ce qu'est le **web sémantique** (notamment pour les ontologies et préfixes)

- ❑ D'être familiariser avec **l'organisation de Dbpédia**

- ❑ De faire des explorations préalables pour savoir quels **prédicats de Dbpédia** sont généralement utilisés dans le cas souhaité

- ❑ Apprendre le **langage SPARQL**

- ❑ De savoir quel **point de terminaison** utiliser (et savoir ce que c'est !)

# Intérêt de mon interface

- ❑ Pas de connaissance du web sémantique nécessaire

- ❑ Pas de connaissance de Dbpédia

- ❑ Les ontologies sont déjà sélectionnées

- ❑ Pas de langage à connaitre

- ❑ Le point de terminaison est déjà fixé

# Schéma de l'organisation

Wikipédia

# Schéma de l'organisation

DBpédia

Wikipédia

# Schéma de l'organisation

Endpoint LiveDBpedia

DBpédia

Wikipédia

# Schéma de l'organisation

Bibliothèque SPARQL dans R

↑

Endpoint LiveDBpedia

↑

DBpédia

↑

Wikipédia

**SPARQL Protocol and RDF Query Language**

# Schéma de l'organisation

Interface

Bibliothèque SPARQL dans R

Endpoint LiveDBpedia

DBpédia

Wikipédia

# SPARQL

- langage universel verbeux

- Permet d'interroger un triplestore

Triplestore = base de données au format RDF

# SPARQL

Chercher toutes les personnes,
Qui sont des athlètes.
Limiter à 300

```
SELECT distinct *
WHERE {
?person a dbo:Person .
?person a dbo:Athlete .
}
LIMIT 300
```

# SPARQL

Chercher le lieu de naissance de Teddy Riner.

```
SELECT distinct *
WHERE {
?person a dbo:Person .
?person rdfs:label « Teddy
Riner »@en .
.?person dbo:birthPlace
?place .
}
```

# Interface

☐ Présentation

# Conclusion

- Facile d'utilisation

- Complète

- Disponible en ligne

- <u>Perspectives</u> :
    - Complexifier le schéma des prédicats
    - Choisir lieu sur la carte
    - Carte de la galaxie

# Apports

- **Découverte du Web sémantique et de ses concepts**
- **Projet à portée très large**
- **Pas de limites d'utilisation**

# Apports

- □ !!! est un groupe de rock indépendant américain

- □ !!! est le nom d'un de leur album

- □ Will Smith est aussi chanteur de rap

- □ Et beaucoup d'autres choses !

# Merci de votre attention !

## Des questions ?

# CSV

| Activity | Organisation | FictionalCharacter | Person | Event | Place | CelestialBody | Species | Work |
|---|---|---|---|---|---|---|---|---|
| All | All | All | All | All | All | All | All | All |
| Game | Company | | Athlete | | ArchitecturalStructure | | | Artwork |
| Sport | EducationalInstitution | | Artist | | NaturalPlace | | | Film |
| | Group | | Politician | | PopulatedPlace | | | MusicalWork |
| | SportsLeague | | Writer | | | | | WrittenWork |
| | SportsTeam | | | | | | | |

# CSV

**Liste des prédicats - Prédicats possibles selon le type du sujet**

| Person | Writer | Politician | Artist | Athlete | Place | PopulatedPlace | NaturalPlace |
|---|---|---|---|---|---|---|---|
| no | no | no | no | no | no | no | no |
| birth name | birth name | birth name | birth name | birth name | type of place | type of place | type of place |
| birth place | birth place | birth place | birth place | birth place | place, localisation | place, localisation | place, localisation |
| birth date | birth date | birth date | birth date | birth date | population | population | population |
| death place | death place | death place | death place | death place | is birth Place of | is birth Place of | is birth Place of |
| death date | death date | death date | death date | death date | is Part of | is Part of | is Part of |
| nationality | nationality | nationality | nationality | nationality | is the city of | is the city of | total area (km²) |
| spouse | spouse | spouse | spouse | spouse | total area (km²) | total area (km²) | |
| children | children | children | children | children | | | |
| parents | parents | parents | parents | parents | | | |
| type of job | title | term start | type of job | type of job | | | |
| team | greatest epoch | term end | title | team | | | |
| term start | is author of | party | greatest epoch | term start | | | |
| term end | style, genre | successor | starred in | term end | | | |
| party | has influenced | predecessor | is author of | title | | | |
| successor | was influenced by | title | style, genre | greatest epoch | | | |
| predecessor | | greatest epoch | has influenced | was influenced by | | | |
| title | | is author of | was influenced by | | | | |
| greatest epoch | | has influenced | | | | | |
| starred in | | was influenced by | | | | | |
| is author of | | | | | | | |
| style, genre | | | | | | | |
| has influenced | | | | | | | |
| was influenced by | | | | | | | |

# CSV

| subtitle | original | direct | unprecise_place | place |
|---|---|---|---|---|
| birth place | dbo:birthPlace | TRUE | TRUE | TRUE |
| starred in | dbo:starring | FALSE | FALSE | FALSE |
| is author of | dbo:author | FALSE | FALSE | FALSE |
| type of place | dbo:type\|dbp:status | TRUE | TRUE | TRUE |
| place, localisation | dbo:city\|dbo:location\|dbo:country\|dbp:country | TRUE | TRUE | TRUE |
| birth date | dbp:birthDate\|dbo:birthDate | TRUE | FALSE | FALSE |
| type of job | dct:description | TRUE | FALSE | FALSE |
| team | dbo:team | FALSE | FALSE | FALSE |
| term start | dbo:activeYearsStartDate\|dbp:termStart\|dbo:activeYearsStartYear | TRUE | FALSE | FALSE |
| term end | dbo:activeYearsEndDate\|dbo:activeYearsEndYear | TRUE | FALSE | FALSE |
| party | dbo:party | TRUE | FALSE | FALSE |
| successor | dbo:predecessor | FALSE | FALSE | FALSE |
| predecessor | dbo:predecessor | TRUE | FALSE | FALSE |
| title | dbp:title | TRUE | FALSE | FALSE |
| greatest epoch | dbp:years | TRUE | FALSE | FALSE |
| population | dbo:populationTotal | TRUE | FALSE | FALSE |
| is birth Place of | dbo:birthPlace | FALSE | FALSE | FALSE |
| is Part of | dbo:isPartOf | TRUE | FALSE | TRUE |
| is the city of | dbo:city | FALSE | FALSE | FALSE |
| total area (km²) | dbo:areaTotal | TRUE | FALSE | FALSE |
| gravity | dbp:gravity | TRUE | FALSE | FALSE |
| period | dbp:period | TRUE | FALSE | FALSE |
| radius | dbp:radius | TRUE | FALSE | FALSE |
| constellation | dbp:constell | TRUE | FALSE | FALSE |
| temperature | dbp:temperature | TRUE | FALSE | FALSE |
| absolute magnitude | dbp:absmagV | TRUE | FALSE | FALSE |
| apparent magnitude | dbp:appmagV | TRUE | FALSE | FALSE |
| family | dbo:family | TRUE | FALSE | FALSE |
| author | dbo:author | TRUE | FALSE | FALSE |
| release date | dbp:date\|dbp:year\|dbo:releaseDate\|dbp:recorded | TRUE | FALSE | FALSE |
| owner | dbp:owner | TRUE | FALSE | FALSE |
| country | dbo:country\|dbp:country | TRUE | FALSE | TRUE |

# CSV

## Résultats de requête au format PDF

| Person | birth place | birth date |
| --- | --- | --- |
| Putter Smith | Bell, California | 19/01/1941 |
| Maggie Smith | Essex | 28/12/1934 |
| Robert Smith (musician) | Lancashire | 21/04/1959 |
| Amery Smith | Los Angeles | 03/08/1964 |
| Richard Smith (English guitarist) | Beckenham | 12/12/1971 |
| Robert Smith (musician) | Blackpool | 21/04/1959 |
| Maggie Smith | Ilford | 28/12/1934 |
| Jabbo Smith | Pembroke, Georgia | 24/12/1908 |
| Broderick Smith | Hertfordshire, England | 17/02/1948 |
| Frederic Marlett Bell-Smith | England | 1846-09-26 |
| Frederic Marlett Bell-Smith | London | 1846-09-26 |
| Major Bill Smith | Checotah, Oklahoma | 21/01/1922 |
| Nick Smith (milliner) | Liverpool | 09/11/1979 |
| Spencer Smith (musician) | Denver | 02/09/1987 |
| Jason Barry-Smith | Queensland | 12/12/1969 |
| O. C. Smith | Mansfield, Louisiana | 21/06/1932 |
| George Logie-Smith | Australia | 02/12/1914 |
| Johnny Smith | Birmingham, Alabama | 25/06/1922 |
| Jason Barry-Smith | Brisbane | 12/12/1969 |
| Dallas Smith (singer) | Canada | 04/12/1977 |
| Tommy Smith (saxophonist) | Edinburgh | 27/04/1967 |
| Nick Glennie-Smith | England | 03/10/1951 |
| Nick Smith (milliner) | England | 09/11/1979 |
| Paul Smith (fashion designer) | England | 05/07/1946 |
| T. V. Smith | England | 05/04/1956 |
| Michael Smith (performance artist) | Illinois | 08/03/1951 |
| Kyla-Rose Smith | Johannesburg | 10/09/1982 |
| Nick Glennie-Smith | London | 03/10/1951 |
| George Logie-Smith | Melbourne | 02/12/1914 |
| Kiki Smith | Nuremberg | 18/01/1954 |
| Sandy Smith | Scotland | 02/08/1983 |
| Jack Smith (artist) | United Kingdom | 18/06/1928 |