

Linked Open Data

Using SPARQL to Retrieve Non-Trivial Information Related to Movies

Primary Topic: SEMI, Secondary Topic: TS¹ & IENLP²

Course: 2021-1B – Group: 126 – Submission Date: 2022-02-06

Ajay Mathew Joseph¹

University of Twente

The Netherlands

a.m.joseph@student.utwente.nl

Nishchit Mahajan²

University of Twente

The Netherlands

n.mahajan@student.utwente.nl

ABSTRACT

Every day new data is being produced and it is often stored in multiple places. Increasing data is an ongoing issue and every day new approaches are being suggested to make use of it. Everyone has some information about a concept but often nobody has the complete information. If the data could be combined from multiple sources, the existing datasets could be enriched and could provide added value. Linked Open Data is one of the interesting solutions to such a problem. In this project, we understand the concept of Linked Open Data and how it could be useful to link existing datasets with it.

KEYWORDS

Linked Open Data, RDF, SPARQL, DBPedia, Data

1 INTRODUCTION

Data is everywhere and everything is data. Everyday data is being produced at an exponential speed and handling such a huge amount of data is becoming troublesome. Storing such huge amounts of data can make data pipelines slow and can be very expensive. As a major resort, the concept of Linked Data can be very helpful. Coined in 2006 by Tim Berners Lee [6], the concept of linking data can be referenced as far as back 1999. This concept of linking helps to share, reuse, and make data machines more readable. This technique helps us to create a small database and link the same topics to multiple different databases. A recent tool to work on such problems is Resource Description Framework (RDF). This framework is a standard model for data interchanges on the Web. This platform forms a directed graph and interlinks the data using named links [4]. This concept has great potential to develop recommender systems and also to enrich existing datasets [5]. In this project, we make use of the concept of RDF in Linked data and try to connect a small movie dataset to DBPedia, a publicly available free and open database, and run some complex queries to test if interlinked data can provide desired results.

The goal of this project is to understand the concept of Linked Open Data by linking an existing movie dataset with DBPedia, an established Linked Open Dataset and thus enriching the initial dataset. We also demonstrate the added value by executing five complex queries and retrieving non-trivial information.

This report has been organized in the following way. Section 2 talks about some related work by other people in similar areas and discusses the dataset used. Section 3 discusses the implementation of the project, including the functions and SPARQL queries used in this project. Section 4 talks about the results obtained from this

project. Section 5 discusses the results and also, lists the limitations of this approach. Section 6 provides a conclusion to the project report which is followed by the references. Finally, appendix A provides the link to the code repository of the complete project.

2 BACKGROUND

This project essentially deals with three concepts which are RDF, DBPedia, and SPARQL. In the following paragraphs, we would describe what they are.

As mentioned in the introduction, RDF is a model for data interchange on the web. In RDF datasets, the data is often stored in the format of triples. The format of a triple is subject-> predicate-> object. Here, subject is the primary entity and it is connected to the object which is the secondary entity with a property/verb which is the predicate [4]. We used the Python library `rdflib` to construct the RDF dataset from an existing dataset.

DBPedia is a project which extracts structured content from the Wikipedia project. Almost all Wikipedia pages contain structured information, which is displayed on the top right of the Wikipedia pages as 'infoboxes'. The DBPedia project extracts this information and a dataset is created out of it which can be queried. The entries in the RDF dataset would be linked to the corresponding DBPedia entry using the owl : sameAs predicate. The links would be retrieved using the DBPedia Spotlight API which is a tool for automatically annotating DBPedia resources [2].

SPARQL is the standard query language and protocol for Linked Open Data and RDF Databases. Information from any data source that can be mapped to and RDF can be queried using SPARQL [1]. The added advantage of SPARQL is that it can be used to query from multiple sources. In this project, we use SPARQL to query information from DBPedia.

3 APPROACH

There are mainly two steps in this project, namely, (a) Constructing an RDF dataset from a base dataset, and (b) running SPARQL queries and retrieving the results for the purpose of demonstration. These will be explained in the following subsections.

The various technologies used in this project include:

- Python programming language
- pandas library
- DBPedia Spotlight API
- rdflib library
- SPARQLWrapper library

3.1 Constructing RDF dataset

In the first step, we obtained the base dataset from Kaggle [3]. The `movies_metadata.csv` file in the dataset contains all the details required for this project. A Python Jupyter notebook was used to construct an RDF dataset from this CSV file. The data was initially loaded into a pandas data frame. It was observed that the dataset consists of 45,466 records with 24 attributes. Some of the attributes were `original_language`, `original_title`, `release_date`, etc.. Further, a function was defined to make use of the DBPedia Spotlight API to find the DBPedia link of a movie given its title. For every movie, its DBPedia link was retrieved from the movie title and along with the movie's `original_title` and `release_date`, an RDF dataset was constructed using the `rdflib` library. The `Graph()` function of `rdflib` was used for this. Other attributes were not stored since this project is for demonstrative purposes and only the movie title is relevant in this context. It was also done so as to not make the dataset larger.

3.2 Tool for executing SPARQL queries

In the second step, we developed a command-line menu interface that could execute five complex SPARQL queries on this dataset and retrieved their results. The user will be shown a five-choice menu and will be asked to enter his choice. Depending on the choice, the corresponding query will be executed using the corresponding functions. These were done using the `SPARQLWrapper` library and Python language. The `rdflib` library was used to parse the `movie.owl` file. The endpoint for the query was set as <http://dbpedia.org/sparql>.

The queries executed were:

- (1) Retrieve related movies
- (2) Get average budget of all movies of a particular movie's director
- (3) Get other movies where the co-actors of a particular movie acted together
- (4) Get youngest main crew member of a movie
- (5) Get the longest movie of a movie's top actor

We also defined three common functions which is useful in all the queries. They are:

- (1) `get_dbpedia_link()`: Given a movie's name, retrieve it's DBPedia link
- (2) `get_actor_name()`: Given an actor's URI, retrieve his/her name
- (3) `get_movie_name()`: Given a movie's URI, retrieve it's name

The logic for each of the queries will be explained in the following subsections.

3.2.1 Related movies. For implementing this query, the input would be a movie's name. The movie's DBPedia link was found using the `get_dbpedia_link()` function and then it was used to find related objects using the `dbo:wikiPageWikiLink` predicate which links to the other Wikipedia pages within this Wikipedia page. The objects found were passed to another function to filter the movies from them and these were returned. The filtering was done by checking the `rdf:type` predicate of the entity.

```
SELECT ?remote_value
WHERE {
    <https://dbpedia.org/page/The_Matrix>
        dbo:wikiPageWikiLink ?remote_value
}
```

Listing 1: Retrieve `dbo:wikiPageWikiLink` objects

3.2.2 Average Budget. This query reads a particular movie's title and returns the average budget of all the movies that this particular movie's director has directed. We find the director from the movie's data and then pass it to a function that finds all movies that he/she has directed and retrieves the budget of all those movies using the `dbo:budget` predicate. The average of the budget was returned as a result using another function that reads the list of budgets as input.

```
SELECT ?director
WHERE {
    <https://dbpedia.org/page/Before_Sunrise>
        dbp:director ?director
}
```

Listing 2: Retrieve Director's name

```
SELECT ?budget
WHERE {
    ?movie dbo:director
        <https://dbpedia.org/page/Richard_Linklater>
    ?movie dbo:budget ?budget
}
```

Listing 3: Retrieve budgets of a director's movies

3.2.3 Other movies of co-actors. In this query, we read a particular movie's title and return all movies in which the top two actors of this movie acted together. The top two actors were retrieved from the movie's page using the `dbo:starring` predicate and then was passed to another function which finds all movies in which these two actors have acted together, using the same predicate.

```
SELECT ?actor
WHERE {
    <https://dbpedia.org/page/The_Matrix>
        dbo:starring ?actor .
}
```

Listing 4: Retrieve actors of a movie

```
SELECT ?film
WHERE {
    ?film dbo:starring
        <https://dbpedia.org/page/Keanu_Reeves> .
    ?film dbo:starring
        <https://dbpedia.org/page/Hugo_Weaving> .
}
```

Listing 5: Retrieve movies in which two actors acted together

3.2.4 Youngest main crew member. For this query, we retrieve the main crew members' URI (director, actors, cinematographer, composer, music, producer) from a movie's data (dbo:director, dbo:starring, dbo:cinematography, dbo: musicComposer, dbp:music, dbo:producer) and then pass them to a function which retrieves the birth dates of all of them in descending order. The birth date's were retrieved using the dbo:birthDate predicate. The first result was returned as an output.

```
SELECT ?dir ?dop ?comp ?music ?screen
      ?prod group_concat(?star; separator=",")
WHERE {
    <https://dbpedia.org/page/The_Matrix>
      dbo:director ?dir .
    <https://dbpedia.org/page/The_Matrix>
      dbo:cinematography ?dop .
    <https://dbpedia.org/page/The_Matrix>
      dbo:musicComposer ?comp .
    <https://dbpedia.org/page/The_Matrix>
      dbp:music ?music .
    <https://dbpedia.org/page/The_Matrix>
      dbo:producer ?prod .
    <https://dbpedia.org/page/The_Matrix>
      dbo:starring ?star .
}
```

Listing 6: Retrieve main crew members of a movie

```
SELECT ?bd ?name WHERE
{
  {SELECT ?bd ?name
  WHERE {
    <https://dbpedia.org/page/Keanu_Reeves>
      dbo:birthDate ?bd .
    <https://dbpedia.org/page/Keanu_Reeves>
      dbp:name ?name .
  }}
  UNION
  {SELECT ?bd ?name
  WHERE {
    <https://dbpedia.org/page/Hugo_Weaving>
      dbo:birthDate ?bd .
    <https://dbpedia.org/page/Hugo_Weaving>
      dbp:name ?name .
  }}
}
ORDER BY DESC(?bd)
```

Listing 7: Retrieve crew members in order of age

3.2.5 Longest movie of an actor. In this query, we take as input a movie's name and then find its lead actor's URI. This was then used to find all movies in which he/she has acted and the runtime values were retrieved using the dbo:runtime predicate. The results were sorted (using ORDER BY) in descending order. The top result was returned as the longest movie of the particular actor.

```
<?xml version="1.0" encoding="utf-8"?>
<rdf:RDF
  xmlns:movie="https://www.imdb.com/"
  xmlns:owl="http://www.w3.org/2002/07/owl#"
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
>
  <rdf:Description rdf:about="@t0061581">
    <movie:title>Divorce American Style</movie:title>
    <movie:releaseDate rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1967-06-21</movie:releaseDate>
    <owl:sameAs rdf:resource="http://dbpedia.org/resource/Divorce_American_Style"/>
  </rdf:Description>
  <rdf:Description rdf:about="@t0110913">
    <movie:title>Pumpkinhead II: Blood Wings</movie:title>
    <movie:releaseDate rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1994-03-16</movie:releaseDate>
    <owl:sameAs rdf:resource="http://dbpedia.org/resource/Jack_Pumpkinhead"/>
  </rdf:Description>
```

Figure 1: movie.owl , RDF Dataset

```
MENU
1. Get related movies
2. Get average budget of all movies of a particular movie's director
3. Get other movies where the coactors of a particular movie acted together
4. Get youngest main crew member of a movie
5. Get the longest movie of a movie's top actor

Enter your choice: █
```

Figure 2: Menu Interface

```
SELECT ?runtime ?name
WHERE
{
  ?movie dbo:starring
    <https://dbpedia.org/page/Robin_Williams> .
  ?movie dbo:runtime ?runtime .
  ?movie dbp:name ?name .
}
ORDER BY DESC(?runtime)
```

Listing 8: Retrieve movies of an actor in decreasing order of runtime

4 RESULTS

4.1 RDF Dataset

The RDF Dataset was constructed with 10,000 movie records out of the 45,466 records since the main objective of the project is to demonstrate the added value of Semantic Web Technology with Linked Open Data. While retrieving the DBpedia links, only 6,905 were successful, and hence, the final dataset contained only those many records. The file was saved as movie.owl. A snippet of the file can be seen in Figure 1.

4.2 SPARQL Queries

As mentioned in the previous section, a menu interface was shown to the user. The interface can be seen in Figure 2.

The results for each of the queries will be shown in the following subsections.

4.2.1 Related Movies. This query took as input a movie's name and printed its related movies. An example is shown in Figure 3.

4.2.2 Average Budget of a Director. This query took as input a movie's name and printed the average budget of its director's movies. An example is shown in Figure 4.

```

Enter your choice: 1
Enter the movie name:The Machinist

Other movies related to The Machinist are:

Kirill Lavrov
The Brothers Karamazov
A Beautiful Mind
Insomnia
Batman Begins

```

Figure 3: Related Movies

```

Enter your choice: 2
Enter the movie name:Before Sunrise
The average budget of Richard Linklater's movies is: Dollars 9943937.5

```

Figure 4: Average Budget of a Director

```

Enter your choice: 3
Enter the movie name:The Matrix

The top 2 actors are:
Keanu Reeves
Hugo Weaving

The movies where they acted together are:
The Matrix Reloaded
The Matrix Revolutions
The Matrix

```

Figure 5: Movies of coactors

```

Enter your choice: 4
Enter the movie name:The Matrix
The youngest main crew member is Keanu Reeves born on 1964-09-02

```

Figure 6: Youngest Crew Member

4.2.3 Other movies of two actors. This query took as input a movie's name and printed all movies in which this particular movie's top two actors acted together. An example is shown in Figure 5.

4.2.4 Youngest main crew member. This query took as input a movie's name and prints its youngest crew member's name and birthdate. An example is shown in Figure 6.

4.2.5 Longest movie of an actor. This query took as input a movie's name and prints the name of its actor's longest movie. An example is shown in Figure 7.

5 DISCUSSION

It has been demonstrated that by linking existing datasets with Linked Open Data, the original dataset could be enriched with new and non-trivial information. By linking related entities, data can be retrieved in unforeseen ways. For example, queries like finding a director's average budget are something that cannot be done with a simple Google search. It has been possible only because all data

```

Enter your choice: 5
Enter the movie name:Jumanji
The longest movie of Robin Williams is Hamlet running 4.033333333333333 hours

```

Figure 7: Longest movie

regarding the director and his movies are linked along with the movies' budget.

Another example is retrieving related movies from a movie's Wikipedia page. Wikipedia pages usually contain links to the Wikipedia pages of other entities related to the original page. This has huge potential for being very relevant entities because Wikipedia pages often contain almost all information about the particular entity and its background. This could be used to build recommender systems. One point to be noted is that the retrieval would be deterministic since the information is structured and not merely predicting anything.

Other queries like retrieving the youngest main crew member, the longest movie of an actor, the movies where two actors acted together, etc. are all non-trivial and would be nearly impossible to be retrieved using machine learning models in a deterministic manner.

Search engines like Google typically use machine learning algorithms and page-rank algorithms to predict the relevance of a search result. However, if there are enough links between the entities of a particular domain, data could be retrieved better and faster using Linked Open Data. Machine learning algorithms could be incorporated along with Linked Open Data tools to retrieve better results.

However, there were some limitations in our approach to this project. First of all, the details regarding an actor or a director cannot be queried by their name alone. The system needs the name of the movie in which they starred or directed. This is a shortcoming since the movie name is irrelevant in the context of this query. Secondly, the DBpedia Spotlight API returned irrelevant results for a few ambiguous movie names, for example, the result for the movie title "JFK" was the DBpedia link of the former American president John F. Kennedy. The results of the API are ranked according to a score. This limitation could be overcome by checking if the result is of the type "movie" or not. Thirdly, some DBpedia entries have missing, redundant, or inconsistent keys. There is also another concern that Wikipedia entries could be edited by anyone and hence there is a challenge in determining the credibility of the information.

6 SUMMARY

It can be seen that with the power of Semantic Web Technology, non-trivial queries can be successfully executed. All the results retrieved in this project cannot be obtained with a normal Google search unless there is a specific website or an article dedicated to this specific kind of information. Also, an existing dataset was enriched with more information by linking the entities with their corresponding entries in the Linked Open Data (here, DBpedia). It was also observed that Semantic Web technology could be used to greatly empower recommender systems and by incorporating them with machine learning models, better (deterministic) results could be obtained. Thus, it can be concluded that Linked Open Data

provides added value by combining data from multiple sources and thus gaining better insights.

A CODE AND DATASET

The source code for this project can be found here: https://github.com/ajktym94/SPARQL_LOD_MOVIES

The RDF dataset was constructed using the Python notebook named `create_movie_rdf.ipynb` and all the SPARQL queries are in the file `main.py`

The RDF dataset constructed is also available in the repository in the file named `movie.owl`. A sample of this file is available as `moviesample.owl`

REFERENCES

- [1] 2013. SPARQL 1.1 Query Language. Retrieved January 31, 2022 from <https://www.w3.org/TR/sparql11-query/>
- [2] 2017. SPOTLIGHT API. Retrieved January 31, 2022 from <https://www.dbpedia-spotlight.org/api>
- [3] 2018. The Movies Dataset. Retrieved January 02, 2022 from https://www.kaggle.com/rounakbanik/the-movies-dataset?select=movies_metadata.csv
- [4] 2021. What is an RDF Triplestore? Retrieved January 31, 2022 from <https://www.ontotext.com/knowledgehub/fundamentals/what-is-rdf-triplestore/>
- [5] Tommaso Di Noia, Roberto Mirizzi, Vito Ostuni, Davide Romito, and Markus Zanker. 2012. Linked open data to support content-based recommender systems. *ACM International Conference Proceeding Series* (09 2012). <https://doi.org/10.1145/2362499.2362501>
- [6] Tim Berners Lee. 2006. Linked Data - Design Issues. <https://www.w3.org/DesignIssues/LinkedData.html>