

Task 1: Dataset Catalogues

Use the DataHub catalogue on <https://old.datahub.io/dataset> to find a dataset under the Creative Commons Attribution license about weather that has a SPARQL endpoint and is organized by the Linked Open Data Cloud.

What is the dataset's name, when was it created and how many triples does it currently contain?

Task 2: Dataset Description

Given is the following description of the YAGO knowledge graph in natural language^{1 2 3}:

YAGO is a huge semantic knowledge base, derived from Wikipedia WordNet and GeoNames. Currently, YAGO has knowledge of more than 10 million entities (like persons, organizations, cities, etc.) and contains more than 120 million facts about these entities.

YAGO is special in several ways:

1. The accuracy of YAGO has been manually evaluated, proving a confirmed accuracy of 95%. Every relation is annotated with its confidence value.
2. YAGO combines the clean taxonomy of WordNet with the richness of the Wikipedia category system, assigning the entities to more than 350,000 classes.
3. YAGO is an ontology that is anchored in time and space. YAGO attaches a temporal dimension and a spatial dimension to many of its facts and entities.
4. In addition to a taxonomy, YAGO has thematic domains such as "music" or "science" from WordNet Domains.
5. YAGO extracts and combines entities and facts from 10 Wikipedias in different languages.

YAGO is developed jointly with the DBWeb group at Télécom ParisTech University.

The YAGO Ontology is licensed under a Creative Commons Attribution 3.0 License by the YAGO team of the Max-Planck Institute for Informatics.

YAGO was first created on October 27, 2009 and last updated on April 28, 2018.

2.1 Give three example dataset profile features that you can find in this text. Also name their corresponding categories (general, qualitative, provenance, links, licensing, statistical, dynamics).

2.2 Write down a formal description of this dataset using the VOID vocabulary⁴.

¹ <https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago>

² <https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/downloads/>

³ <https://old.datahub.io/dataset/yago>

⁴ <https://www.w3.org/TR/void/>

Task 3: Statistical Features of DBpedia

3.1 Write a query to count the total numbers of triples in the English DBpedia (<http://dbpedia.org/sparql/>).

3.2 What does the following query return?

```
SELECT COUNT(*) WHERE {  
    SELECT DISTINCT(?b) WHERE {  
        ?a ?b ?c .  
    }  
}
```

3.3 What does the following query return?

```
SELECT ?b ?c {  
    dbo: ?b ?c .  
}
```

3.4a Write a SPARQL query that, for each type (?type), says how many instances it has.

3.4b Change the query from Task 3.4a, such that the types are returned by decreasing frequency.

3.5 What happens if we include the line “?type rdfs:subClassOf* dbo:Person .” into the query from Task 3.4b?

Task 4: A Priori Type Probabilities

Infer the type a priori probabilities table from the **French (!)** DBpedia (<http://fr.dbpedia.org/sparql>).

4.1 Write a SPARQL query to count the total number of instances in DBpedia with at least one type.

4.2 Use the results from Task 4.1 and Task 3.4b to compute the type distribution in the English DBpedia (a priori probabilities) for the three most frequent types.

Type t	Frequency	P(t) in %

Task 5: Type Distributions for a Given Type

Infer the type distribution table of the property `dbo:location` from the **French** DBpedia.

5.1 Write a SPARQL query that returns the total number of triples that have the property `dbo:location`.

5.2 Write a SPARQL query to find out how many subjects (objects) of the `dbo:location` triples are of a specific type (e.g., 100 subjects in triples with `dbo:location` are having the type `dbo:Person`). Limit the types to those from the `dbo` namespace and order the output by decreasing type count.

5.3 For the three most frequent subject and object types: Compute their type distribution using the outputs of Task 5.1 and 5.2.

Type	Subject Count	Subject (%)	Object Count	Object (%)

Task 6: Type Inference Algorithm

Use the SD-Type Algorithm to infer the type of `dbr:Eiffel_Tower`. Compute the type confidence scores for `dbo:Place`, `dbo:Building` and `dbo:Person`.

Given are the following triples with the entity Eiffel Tower:

```
dbr:Eiffel_Tower dbo:location dbr:Paris
dbr:Eiffel_Tower dbo:architect dbr:Stephen_Sauvestre
dbr:France_Télévisions dbo:location dbr:Eiffel_Tower
```

Additionally, the following property distributions and the a priori type distributions are given:

`dbo:location`

Type	Subject	Object
<code>dbo:Place</code>	0.80	0.75
<code>dbo:Building</code>	0.17	0.01
<code>dbo:Organisation</code>	0.16	0.02
<code>dbo:City</code>	0.03	0.21

`dbo:architect`

Type	Subject	Object
<code>dbo:Place</code>	1.00	0.0
<code>dbo:Building</code>	0.88	0.0
<code>dbo:Person</code>	0.0	0.81
<code>dbo:Architect</code>	0.0	0.36

A priori type distribution:

Type t	P(t)
<code>dbo:Place</code>	0.43
<code>dbo:Building</code>	0.12
<code>dbo:Organisation</code>	0.21

Type t	P(t)
<code>dbo:City</code>	0.08
<code>dbo:Person</code>	0.57
<code>dbo:Architect</code>	0.03