



# **Explorando consultas SPARQL na Wikidata com Python: Tipificação de Metadados e Reconciliação de Dados**

**Luis Felipe Rosa de Oliveira**

**Dalton Lopes Martins**



Programa de Pós-graduação  
Ciência da Informação - PPGCINF  
FACULDADE DE CIÊNCIA DA INFORMAÇÃO

# INTRODUÇÃO

- Contexto:
  - Laboratório de Inteligência de Redes
  - Projeto Tainacan (MARTINS; CARVALHO JÚNIOR; GERMANI, 2019, p.60)
- Temática:
  - Web Semântica
  - Dados Abertos Ligados
- Reconciliação de Dados:
  - Processo permite **reconhecer entidades**, seja no formato textual ou no formato de objeto digital em uma ou mais **bases de conhecimento on-line**, permitindo assim **formar conexões entre bases** por meio de identificadores únicos. (SANDERSON, 2016)

# INTRODUÇÃO

A Wikidata (Base de Conhecimento)

“é uma **base de conhecimento livre** e aberta que pode ser lida e editada por humanos e máquinas” WIKIDATA (2019)

Como ligar meus dados com a Wikidata?

Como desenvolver processos semiautomáticos de reconciliação de dados com a Wikidata?



WIKIDATA

Explorando consultas SPARQL na Wikidata com Python: Tipificação de Metadados e Reconciliação de Dados

label — Douglas Adams (Q42) — item identifier

description — English writer and humorist — aliases  
Douglas Noël Adams | Douglas Noel Adams  
► In more languages

Statements

property — educated at — value — St John's College — value

rank —

statement group —

qualifiers —

opened references —

collapsed reference —

Statement	Value
educated at	St John's College
end time	1974
academic major	English literature
academic degree	Bachelor of Arts
start time	1971

Reference	Value
stated in	Encyclopædia Britannica Online
reference URL	http://www.britannica.com/people/731.000023662/
original language of work	English
retrieved	7 December 2013
publisher	HNDB
title	Douglas Adams (English)

Statement	Value
Brentwood School	Brentwood School
end time	1970
start time	1959

III WIDaT

# INTRODUÇÃO

## Proposta da Pesquisa

- Transformar objetos textuais em objetos digitais (ZENG, 2019 )
- *Linked Open Data*

## Desenvolvimento de Dois Scripts:

- Tipificação de metadados, com o objetivo de reconhecer a qual instância da Wikidata determinado metadado pertence a partir de seu conjunto de valores.
- Reconciliação de dados, em que determinados valores textuais são buscados na Wikidata em busca de identificar seus possíveis objetos representativos.

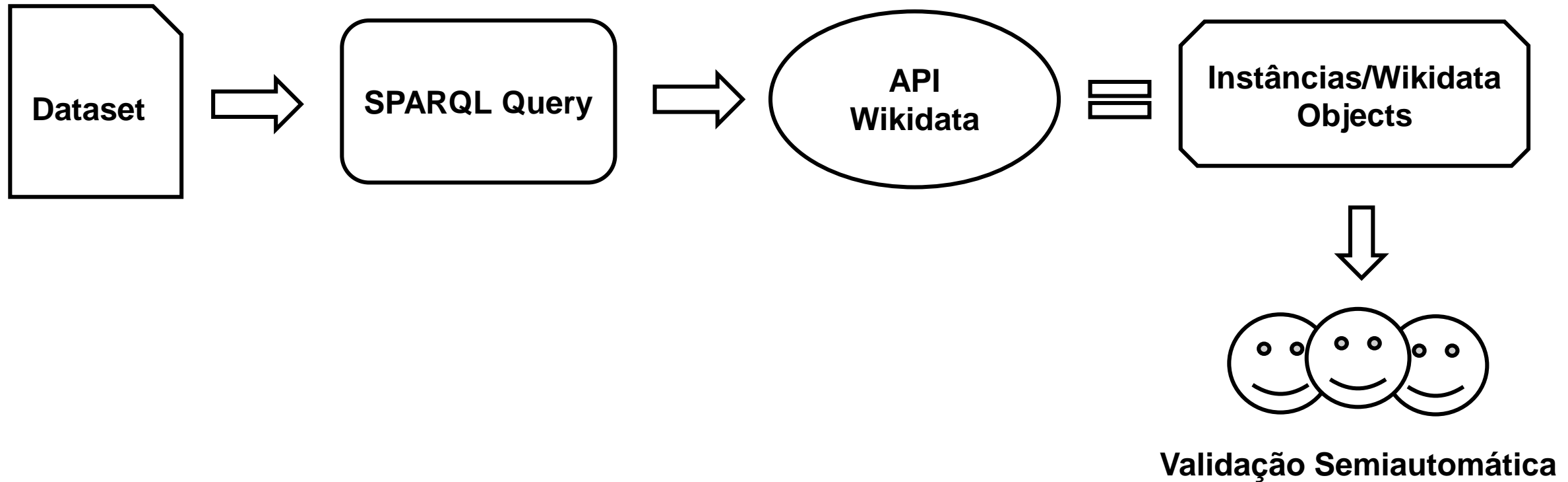
# INTRODUÇÃO

O que se espera:

- Compreender melhor **como dados de diferentes acervos podem ser ligados de forma semiautomática**, complementando propostas de formas mais diretas de se **reconhecer entidades em bases de conhecimento**, e **enriquecer bases de dados locais**.

# Método

Processo de Reconciliação



# Método

## Python

“é uma linguagem de programação que te permite trabalhar rapidamente e integrar sistemas de forma mais efetiva.”

(PYTHON 2019)



### Bibliotecas:

- Requests
- Pandas
- Time
- Fuzzywuzzy

# Método

## Dataset

Cidade	País	Marcas
Goiânia	Brasil	Coca-Cola
São Paulo	Chile	Arno
Brumadinho	França	Epson
Planaltina	Estados U	Dell
Curitiba	Butão	C&A
São Carlos	Argentina	CaBas
Abaetetuba	Marrocos	Calvin Klein



# Método (SPARQL Query – Buscando Instâncias)

Wikidata Query Service

Exemplos Ajuda Mais ferramentas

```
1 SELECT DISTINCT ?sujeito ?instancia_de_que ?instancia_de_queLabel WHERE {
2
3     { ?sujeito ?label "Karajá". }
4     UNION
5     { ?sujeito ?label "Karajá"@en. }
6     UNION
7     { ?sujeito ?label "Karajá"@pt-br. }
8
9     ?sujeito wdt:P31 ?instancia_de_que.
10
11     FILTER(?instancia_de_que NOT IN(wd:Q4167836, wd:Q4167410))
12
13     SERVICE wikibase:label { bd:serviceParam wikibase:language "pt-br", "pt", "en". }
14 }
```

3 resultados em 303 ms

sujeito	instancia_de_que	instancia_de_queLabel
<a href="#">wd:Q10322066</a>	<a href="#">wd:Q34770</a>	linguagem
<a href="#">wd:Q1728924</a>	<a href="#">wd:Q41710</a>	grupo étnico
<a href="#">wd:Q10322066</a>	<a href="#">wd:Q1288568</a>	idioma vivo

# Método (SPARQL Query – Buscando Objetos)

Wikidata Query Service

Exemplos Ajuda Mais ferramentas

```
1 SELECT DISTINCT ?sujeito ?sujeitoLabel ?sujeitoAltLabel WHERE {
2
3     { ?sujeito ?label "Karajá". }
4     UNION
5     { ?sujeito ?label "Karajá"@en. }
6     UNION
7     { ?sujeito ?label "Karajá"@pt-br. }
8
9     ?sujeito wdt:P31 ?instancia_de_que.
10
11     FILTER(?instancia_de_que NOT IN(wd:Q4167836, wd:Q4167410))
12
13     SERVICE wikibase:label { bd:serviceParam wikibase:language "pt-br", "pt", "en". }
14 }
```

2 resultados em 736 ms

sujeito	sujeitoLabel	sujeitoAltLabel
<a href="#">Q10322066</a>	Língua carajá	Língua karajá
<a href="#">Q1728924</a>	Carajás (tribo)	Carajá, Karajá

# Método (Script – Tipificação de Metadados)

```
# Script to map metadata to wikidata entities (instance of)
#Required Packages
import requests
import pandas as pd
import time

#API Endpoint
api_url = 'https://query.wikidata.org/sparql'
#Values to reconcile database
tabela_df = pd.read_csv('planilha_teste.csv', encoding='utf-8')
```

# Método (Script – Tipificação de Metadados)

```
#For each metadata/column on values to reconcile database
for column in database.columns:
    print("Trabalhando no metadado ", column)

#For each value of a column/metadata
for value in tabela_df[column]:
    print("Procurando Resultados para ", value)

#SPARQL query. Returns subject QID, instance QID and instance label, searching the value for subjects
#with matching labels.
query = """ SELECT DISTINCT ?sujeito ?instancia_de_que ?instancia_de_queLabel WHERE {
    { ?sujeito ?label "%s". }
    UNION
    { ?sujeito ?label "%s"@en. }
    UNION
    { ?sujeito ?label "%s"@pt-br. }
    ?sujeito wdt:P31 ?instancia_de_que.
    FILTER(?instancia_de_que NOT IN(wd:Q4167836, wd:Q4167410))
    SERVICE wikibase:label { bd:serviceParam wikibase:language "pt-br", "pt", "en". }
}""" %(value, value, value)
```

```
#Request the query to wikidata api and outputs a json.
r = requests.get(api_url, params = {'format': 'json', 'query': query})
data = r.json()
```

```
#Request the query to wikidata api and outputs a json.
r = requests.get(api_url, params = {'format': 'json', 'query': query})
data = r.json()
```

# Método (Script – Tipificação de Metadados)

```
#Dataframe to store results
result_df = pd.DataFrame(columns = ['column', 'value', 'qid', 'qid_instance', 'instance_label'])

#For each result inset data on result_df
for item in data['results']['bindings']:
    result_df = result_df.append({'column':column, 'value':value,
                                  'qid':item['sujeito']['value'],
                                  'qid_instancia':item['instancia']['value'],
                                  'instance_label':item['instanciaLabel']['value'] },
                                  ignore_index=True)

time.sleep(1)

print(reconcile_database(tabela_df).to_csv("resultado.csv", encoding='utf-8'))
```

# Resultado (Script - Tipificação de Metadados)

column	value	qid	qid_instance	instance_label
Cidade	Goiânia	<a href="http://www.wikidata.org/entity/Q83189">http://www.wikidata.org/entity/Q83189</a>	<a href="http://www.wikidata.org/entity/Q515">http://www.wikidata.org/entity/Q515</a>	cidade
Cidade	Goiânia	<a href="http://www.wikidata.org/entity/Q83189">http://www.wikidata.org/entity/Q83189</a>	<a href="http://www.wikidata.org/entity/Q5119">http://www.wikidata.org/entity/Q5119</a>	capital
País	Brasil	<a href="http://www.wikidata.org/entity/Q155">http://www.wikidata.org/entity/Q155</a>	<a href="http://www.wikidata.org/entity/Q859563">http://www.wikidata.org/entity/Q859563</a>	Estado secular
País	Brasil	<a href="http://www.wikidata.org/entity/Q10344998">http://www.wikidata.org/entity/Q10344998</a>	<a href="http://www.wikidata.org/entity/Q1761072">http://www.wikidata.org/entity/Q1761072</a>	parque estadual
País	Brasil	<a href="http://www.wikidata.org/entity/Q4957796">http://www.wikidata.org/entity/Q4957796</a>	<a href="http://www.wikidata.org/entity/Q7366">http://www.wikidata.org/entity/Q7366</a>	canção
País	Brasil	<a href="http://www.wikidata.org/entity/Q1799348">http://www.wikidata.org/entity/Q1799348</a>	<a href="http://www.wikidata.org/entity/Q3184121">http://www.wikidata.org/entity/Q3184121</a>	município do Brasil
País	Brasil	<a href="http://www.wikidata.org/entity/Q155">http://www.wikidata.org/entity/Q155</a>	<a href="http://www.wikidata.org/entity/Q3624078">http://www.wikidata.org/entity/Q3624078</a>	estado soberano
País	Brasil	<a href="http://www.wikidata.org/entity/Q20274905">http://www.wikidata.org/entity/Q20274905</a>	<a href="http://www.wikidata.org/entity/Q486972">http://www.wikidata.org/entity/Q486972</a>	povoado
Marcas	Epson	<a href="http://www.wikidata.org/entity/Q159614">http://www.wikidata.org/entity/Q159614</a>	<a href="http://www.wikidata.org/entity/Q4830453">http://www.wikidata.org/entity/Q4830453</a>	empresa
Marcas	Epson	<a href="http://www.wikidata.org/entity/Q159614">http://www.wikidata.org/entity/Q159614</a>	<a href="http://www.wikidata.org/entity/Q6881511">http://www.wikidata.org/entity/Q6881511</a>	empresa
Marcas	Epson	<a href="http://www.wikidata.org/entity/Q37054258">http://www.wikidata.org/entity/Q37054258</a>	<a href="http://www.wikidata.org/entity/Q101352">http://www.wikidata.org/entity/Q101352</a>	sobrenome

# Método (Script – Reconciliação de Dados)

```
#Required Packages
import requests
import pandas as pd
import time
from collections import defaultdict
from fuzzywuzzy import fuzz
labels_dict = defaultdict(list)

#API Endpoint
api_url = 'https://query.wikidata.org/sparql'

#Values to reconcile database
tabela_df = pd.read_csv('planilha_teste.csv', encoding='utf-8')

value = 'Karajá'
```

# Método (Script – Reconciliação de Dados)

```
query = """ SELECT DISTINCT ?sujeito ?sujeitoLabel ?sujeitoAltLabel WHERE {  
  
    { ?sujeito ?label "%s". }  
    UNION  
    { ?sujeito ?label "%s"@en. }  
    UNION  
    { ?sujeito ?label "%s"@pt-br. }  
  
    ?sujeito wdt:P31 ?instance.  
  
    SERVICE wikibase:label { bd:serviceParam wikibase:language "pt-br", "pt", "en". }  
}""%(value, value, value)  
  
r = requests.get(api_url, params = {'format': 'json', 'query': query})  
data = r.json()
```



# Método (Script – Reconciliação de Dados)

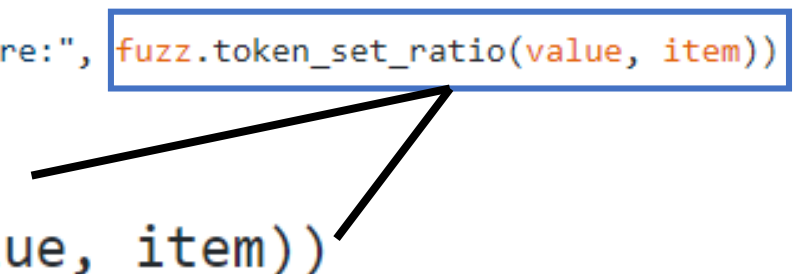
```
#Faz um dicionário com lista das labels para cada resultado retornado
for item in data['results']['bindings']:
    try:
        labels_dict[item['sujeito']['value']] = item['sujeitoAltLabel']['value'].split(',')
        labels_dict[item['sujeito']['value']].append(item['sujeitoLabel']['value'])
    except:
        labels_dict[item['sujeito']['value']].append(item['sujeitoLabel']['value'])

#Limpar Repetições da lista de labels
for key in labels_dict.keys():
    labels_dict[key] = list(set(labels_dict[key]))

#Avaliar o score do match entre o valor procurado com o valor procurado
for key in labels_dict.keys():
    for item in labels_dict[key]:
        print("Valor:", value, "|| Objeto", key, "|| Label:", item, "|| Score:", fuzz.token_set_ratio(value, item))
```

`fuzz.token_set_ratio(value, item)`

`fuzz.token_set_ratio(value, item)`



# Resultado (Script – Reconciliação de Dados)

```
Valor: Karajá || Objeto http://www.wikidata.org/entity/Q10322066 || Label: Língua karajá || Score: 100
Valor: Karajá || Objeto http://www.wikidata.org/entity/Q10322066 || Label: Língua carajá || Score: 50
Valor: Karajá || Objeto http://www.wikidata.org/entity/Q1728924 || Label: Carajás (tribo) || Score: 47
Valor: Karajá || Objeto http://www.wikidata.org/entity/Q1728924 || Label: Carajá || Score: 80
Valor: Karajá || Objeto http://www.wikidata.org/entity/Q1728924 || Label: Karajá || Score: 100
```

# Referências

- MARTINS, D. L.; CARVALHO, J. M. C; GERMANI, L. Projeto Tainacan: Experimentos, Aprendizados e Descobertas da Cultura Digital no Universo dos Acervos das Instituições Memoriais. In: **TIC CULTURA Pesquisa Sobre o Uso das Tecnologias de Informação e Comunicação nos Equipamentos Culturais Brasileiros — 2018**. Comitê Gestor da Internet no Brasil, São Paulo, 2019.
- SANDERSON, R. “**The Linked Data Snowball and Why We Need Reconciliation**”, 2016, Disponível em: <https://www.slideshare.net/azaroth42/linked-data-snowball-or-why-we-need-reconciliation>. Acesso em: 15/08/2019.
- ZENG, M. L. “Semantic enrichment for enhancing LAM data and supporting digital humanities. Review article”. **El profesional de la información**, v. 28, n. 1, 2019.



# OBRIGADO

Luis Felipe Rosa | [luisrosaprof@gmail.com](mailto:luisrosaprof@gmail.com)

Dalton Martins | [dmartins@gmail.com](mailto:dmartins@gmail.com)

**Scripts** – [https://github.com/luisfeliperd/wikidata\\_sparql](https://github.com/luisfeliperd/wikidata_sparql)