



Data mining II Term project

Instructor Dr. NAIT-HAMOUD

Department of Second Cycle - fourth year Applied statistics
National higher School of Mathematics
Algiers, Algeria

Dec 22, 2025

1 Introduction

The aim of this project is to implement an entity linking system. Entity Linking (EL) is an information extraction task that links entity mentions present in a text to knowledge base entries such as Wikipedia, Dbpedia, etc. Figure 2 illustrates this task, the chunks of text highlighted in green represent mentions to entities coherent with the context of the respective texts where they occur. Mainly, the EL task can be broken down into two main sub-tasks.

- First, entity mentions (highlighted in the figure in green) must be located in the text, we refer to this sub-task as "recognition",
- Second, those mentions must be associated to the appropriate Wikipedia page.

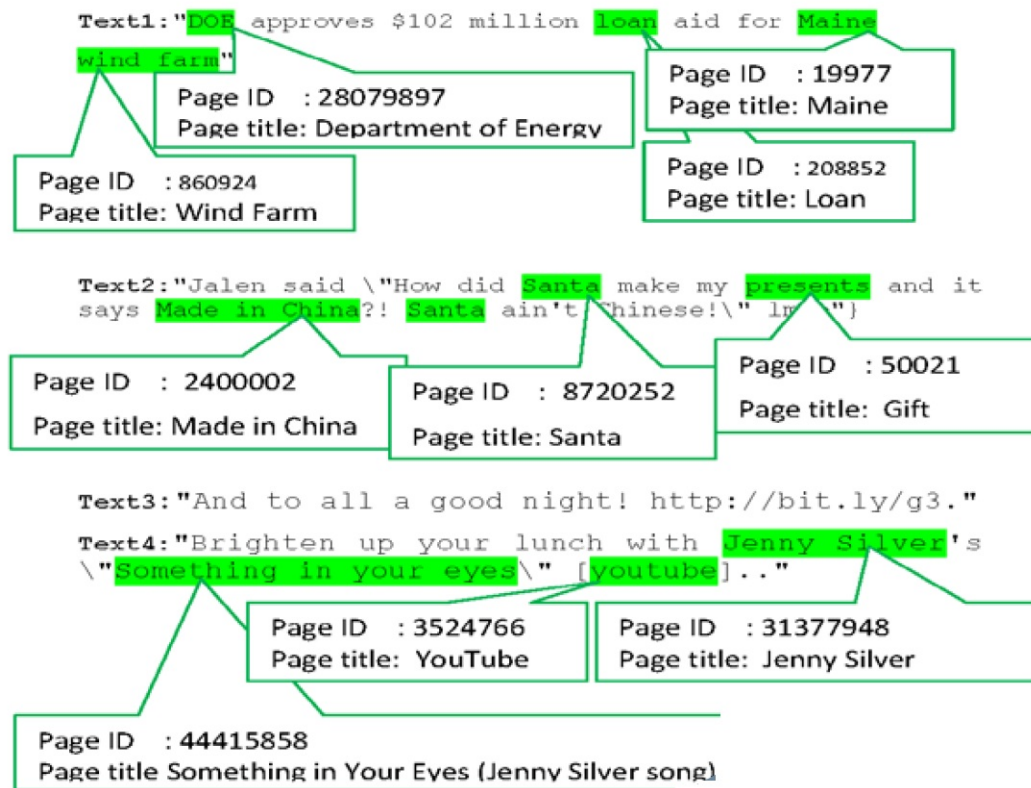


Figure 1: Example of linking mention to entities (Wikipedia page in this case)

Wikipedia knowledge base include three type of pages.

1. Redirect pages: redirects are used to make searching for information in Wikipedia easier,
2. Disambiguation pages (pages that do not cover a unique topic, but help to find the specific article when multiple topics share the same name),
3. Proper Wikipedia pages that cover a unique topic (see definition 2 of section 4 of the provided paper).

One of the methods used by Wikipedia to search for pages is called **direct matching and redirection** that consists in two alternatives.

- If your search query exactly matches an article title, Wikipedia usually navigates you directly to that page instead of showing results (Exact Matches).

- It checks for "redirects"—alternative titles (like "Alan Mathison Turing" redirecting to "Alan Turing") to ensure you find the correct page even with variations (redirects).

For instance, figure 2 depicts the "Alan Turing" Wikipedia page. The information surrounded in red shows that there is a redirect page with the title "Turing" that redirects to the proper page of "Alan Turing". Besides, the information surrounded in green informs that there is a disambiguation page for the word "Turing". Figure 3 shows a snapshot of the disambiguation page corresponding to "Turing".



Figure 2: Example of a Wikipedia page: the text surrounded in blue indicates an example of an anchor (clickable text used to refer to a given wikipedia page)

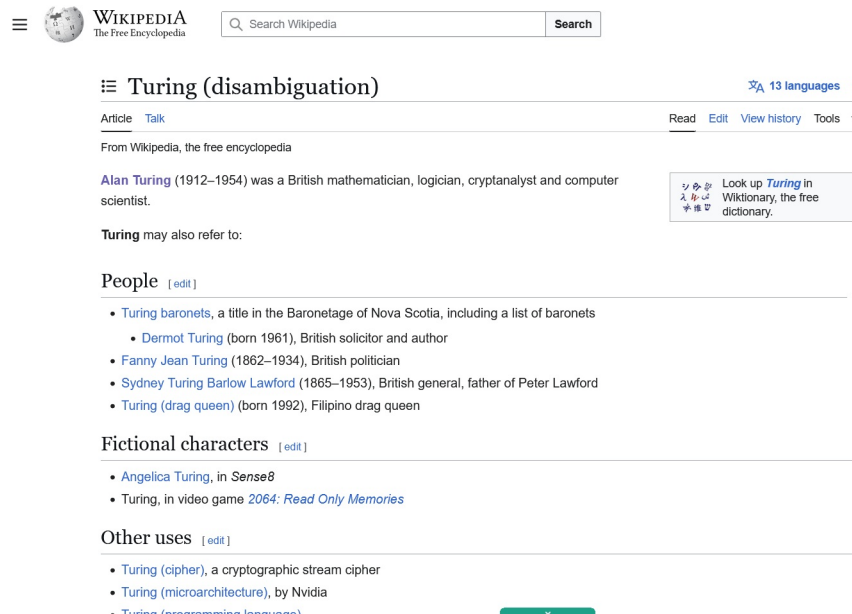


Figure 3: a snapshot of "Turing" Disambiguation page

2 Entity extraction and disambiguation

The EL system TELS (for Tweet Entity Linking System) introduced in the paper of reference [1] is an end-to-end system. This means that the two phases namely: *mention extraction* (or recognition) and *linking* are carried out in one shot. The authors in [1] did not extract the relevant mentions as a first phase, and in a second phase link the relevant extracted mention with the appropriate Wikipedia page.

After a preprocessing phase, described in [1], the authors extract the Ngrams from the cleaned text using a function similar to the one depicted in figure 4. Then, each Ngram is looked for in the inverted index, if it exists; candidate pairs (Ngram,PageID) are generated for each Page referenced by the Ngram. Afterwards, for each pair (Ngram,PageID) a feature vector is build on the basis of the selected features (see the work in [1]). Finally, a Model is trained using the provided revised Meij dataset.

```

1  #-*- coding: utf-8 -*-
2  """
3  Created on Mon Dec 22 23:09:13 2025
4
5  @author: Pc
6  """
7
8  def Get_ngrams(tokens,n):
9
10     ngrams = zip(*[tokens[i:] for i in range(n)])
11     #print(ngrams)
12     tab=[" ".join([elem for elem in ungram]) for ungram in ngrams]
13     #print(tab)
14
15     return(tab)
16
17 text="A step further towards a consensus on linking tweets to Wikipedia"
18 # 1-grams
19 print(Get_ngrams(text.split(),1))
20 print('-----')
21 # 2-grams
22 print(Get_ngrams(text.split(),2))
23 print('-----')
24 # 3-grams
25 print(Get_ngrams(text.split(),3))

```

Usage

Here you can get help of any object by pressing **Ctrl+I** in front of it, either on the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in *Preferences > Help*.

New to Spyder? Read our tutorial

Python 3.10.13 | packaged by Anaconda, Inc. | (main, Sep 11 2023, 13:15:57) [MSC v. 1916 64 bit (AMD64)]

Type "copyright", "credits" or "license()" for more information.

IPython 8.27.0 -- An enhanced Interactive Python.

In [1]: runfile('C:/Users/Pc/Desktop/Travaux Scientifiques/CONFERENCES et papiers/TELS/untitled1.py', wdir='C:/Users/Pc/Desktop/Travaux Scientifiques/CONFERENCES et papiers/TELS')

['A', 'step', 'further', 'towards', 'a', 'consensus', 'on', 'linking', 'tweets', 'to', 'Wikipedia']

['A step', 'step further', 'further towards', 'towards a', 'a consensus', 'consensus on', 'on linking', 'linking tweets', 'tweets to', 'to Wikipedia']

['A step further', 'step further towards', 'further towards a', 'towards a consensus', 'a consensus on', 'consensus on linking', 'on linking tweets', 'linking tweets to', 'tweets to Wikipedia']

Figure 4: A python function to extract Ngrams from a text. If the Ngram position in the text is required the code should be enhanced

The provided material (files) for this project include the paper on TELS system [1] (to read carefully), prebuilt datasets (PostingsLast and PageIdToContext2), the python files to be used for reading records from these datasets, the gold standards, and the Meij dataset (multiple versions).

2.1 Provided prebuilt datasets and manipulation python code

For the needs of the project, the following datasets have been made available for you.

- **Inverted index:** A dictionary of anchors including redirect page titles was created from the English Wikipedia dump of January 2019. The dictionary of anchors is a key-value store that associates to each entry (Wikipedia anchor or redirect page title) the set of Wikipedia pages it refers to, along with a score representing the number of times the entry was used as cross-reference to each page.

The provided files to use for manipulating the inverted index are:

1. The Lmdb **PostingsLast**: To implement the database, LMDB (Lightening-Mapped Database) storage engine was used due to its high read performances. Moreover, to improve serialization, Google Protocol Buffers were used due to their best performances compared with XML, JSON and other existing formats.
2. **InvertedIndexAccess.py**: can be used, through the function call **IndexAccess(ngram,txn)**, to read from the lmdb **PostingsLast**. When you run the python code in InvertedIndexAccess.py,

you are asked to enter a mention (text) **m** or an ngram, if this mention was used as an anchor, or z redirect page, in Wikipedia to refer to specific pages, the result is a list of records. Each record include the page ID of the referenced page, the score and the type of the referenced page.

- **pageId**: The ID of the Wikipedia page
 - **score**: It tells you how many time the text or the mention **m** was used, as an anchor, in the whole Wikipedia to refer to this page.
 - **type**: 0 for anchor (the Ngram was used as an anchor), 1 for redirect (the Ngram was used as a redirect page title) and 2 for both (the Ngram was used as a redirect page title and as an anchor to cross-reference the page)
3. `SerializedListNew_pb2` (it is the defined **Google protocol buffer**), it should be imported in the file `InvertedIndexAccess`)

```

6
7 import SerializedListNew_pb2
8
9 import lmdb
10
11 def Get_Postings(ngram,txn):
12     tt_occrr=0
13     ll=SerializedListNew_pb2.SerializedListNew()
14     if ngram!='':
15         val=txn.get(ngram.encode())
16         my_list=SerializedListNew_pb2.SerializedListNew()
17
18         if val!=None:
19             ll=my_list.FromString(val)
20             #print("ll length:",len(ll.Elements))
21             for ii in range(len(ll.Elements)):
22                 tt_occrr+=ll.Elements[ii].score
23
24     return ll,tt_occrr
25
26 def IndexAccess(ngram,txn):
27     Posting_List=SerializedListNew_pb2.SerializedListNew()
28     Posting_List.TotalOccurr=Get_Postings(ngram,txn)
29     return Posting_List.TotalOccurr
30
31 if __name__=="__main__":
32     ngram=input("Ngrams=")
33     Post_lists=SerializedListNew_pb2.SerializedListNew()
34     Postings=lmdb.open('D:/TELS/PostingsLast',readonly=True)
35     with Postings.begin() as txn:
36         Post_lists.Occrrs=IndexAccess(ngram,txn)
37         #print(Post_lists.Elements)
38         for el in Post_lists.Elements:
39             print("Score= ",el.score)
40             print("PageId= ", el.docId)
41             print("Type= ", el.type)
42             print('-----')
43

```

Usage

Here you can get help of any object by pressing **Ctrl+H** in front of it, either on the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in [Preferences > Help](#).

[New to Spyder? Read our tutorial](#)

Help Variable Explorer Plots Files

Console 3/A X

Python 3.10.13 | packaged by Anaconda, Inc. | (main, Sep 11 2023, 13:15:57) [MSC v. 1916 64 bit (AMD64)]
Type "copyright", "credits" or "license" for more information.

IPython 8.27.0 -- An enhanced Interactive Python.

In [1]: runfile('C:/Users/Pc/Desktop/Travaux Scientifiques/CONFERENCES et papiers/TELS/InvertedIndexAccess.py', wdir='C:/Users/Pc/Desktop/Travaux Scientifiques/CONFERENCES et papiers/TELS')

Ngram=cryptanalyst
Score= 155
PageId= 5715
Type= 2

Score= 1
PageId= 18934432
Type= 0

Total occurrence = 156

Figure 5: Example of a search for the Ngram "cryptanalyst" in the inverted index

- **The dataset of contexts PageIdToContext2**: It is also a key-value store that includes the following information for each Wikipedia page:
 1. Page title,
 2. List of anchors and redirect page titles that link to this page,
 3. Categories of the page included in the Wikipedia categories section,
 4. list of anchors and titles of Wikipedia pages referenced in the abstract,
 5. number of times the page was visited during the year 2019,
 6. Page rank of the page in the Wikipedia knowledge base graph.
- **InterrogatePageIdToContextLmdb.py**: This python code allows for interrogating the database (lmdb) **PageIdToContext2** that represents the dataset of contexts of each Wikipedia page. Figure 6 shows an example of the execution of this python code.

Each Wikipedia page include a "categories" section. This latter is a navigational system used to group related articles. They allow to browse through related topics even if they do not know the exact title of a specific page. Figure 7 shows the "categories" section of the Wikipedia page entitled "Alan Turing".

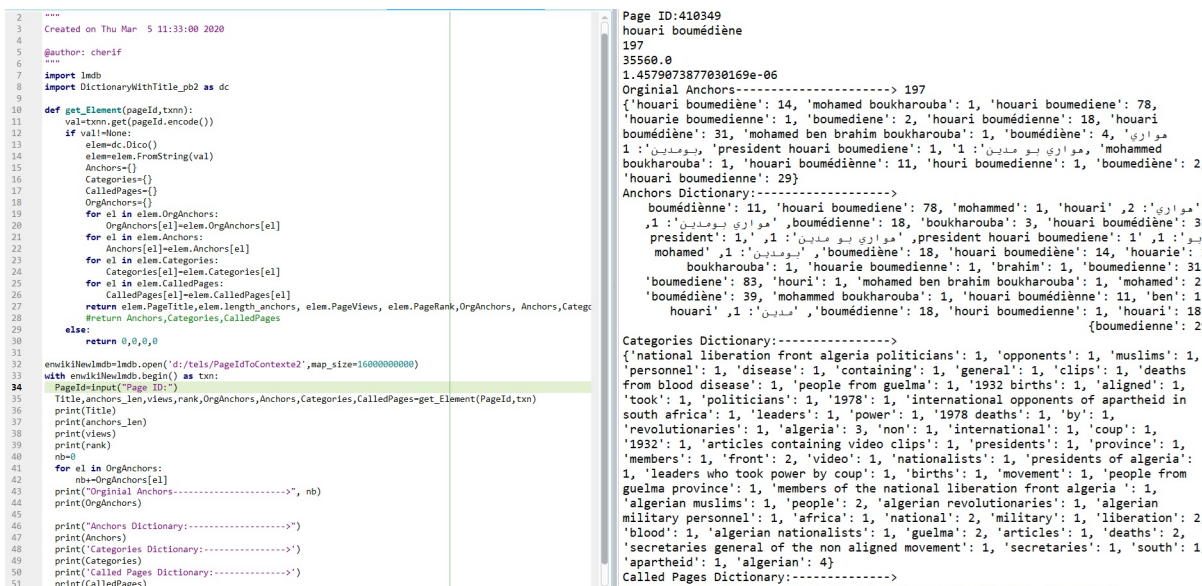


Figure 6: Example of the execution of the python code `InterrogatePageIdToContextLmdb.py`

V·T·E	Alan Turing	[show]
V·T·E	Fellows of the Royal Society elected in 1951	[show]
V·T·E	Timelines of computing	[show]
V·T·E	Early history of video games (1947–1971)	[show]
Authority control databases 		[show]
<div> <div> Portals:</div> <div> <div> Biography</div> <div> England</div> <div> LGBTQ</div> <div> Greater Manchester</div> <div> Mathematics</div> <div> Computer programming</div> </div> </div>		
<div> <div>Categories:</div> <div> <div>Alan Turing</div> <div>1912 births</div> <div>1954 deaths</div> <div>1954 suicides</div> <div>20th-century atheists</div> <div>20th-century English LGBTQ people</div> <div>20th-century English mathematicians</div> <div>20th-century English philosophers</div> <div>Academics of the University of Manchester Institute of Science and Technology</div> <div>Academics of the University of Manchester</div> <div>Alumni of King's College, Cambridge</div> <div>Bayesian statisticians</div> <div>Bletchley Park people</div> <div>British anti-fascists</div> <div>British artificial intelligence researchers</div> <div>British cryptographers</div> <div>British people of World War II</div> <div>Castrated people</div> <div>Computability theorists</div> <div>Computer chess people</div> <div>Computer designers</div> <div>Early history of video games</div> <div>English atheists</div> <div>English computer scientists</div> <div>English gay sportsmen</div> <div>English inventors</div> <div>English LGBTQ scientists</div> <div>English logicians</div> <div>English men long-distance runners</div> <div>British men long-distance runners</div> <div>English people of Irish descent</div> <div>English people of Scottish descent</div> <div>Enigma machine</div> <div>Fellows of King's College, Cambridge</div> <div>Fellows of the Royal Society</div> <div>Foreign Office personnel of World War II</div> <div>Former Protestants</div> <div>Gay academics</div> <div>Gay scientists</div> <div>GCHQ people</div> <div>History of computing in the United Kingdom</div> <div>LGBTQ mathematicians</div> <div>LGBTQ philosophers</div> <div>LGBTQ track and field athletes</div> <div>LGBTQ people who died by suicide</div> <div>Officers of the Order of the British Empire</div> <div>People convicted for homosexuality in the United Kingdom</div> <div>People educated at Sherborne School</div> <div>People from Maida Vale</div> <div>People from Wilmslow</div> <div>People who have received posthumous pardons</div> <div>Princeton University alumni</div> <div>Recipients of British royal pardons</div> <div>Scientists of the National Physical Laboratory (United Kingdom)</div> <div>Suicides by cyanide poisoning</div> <div>Suicides in England</div> <div>Theoretical biologists</div> <div>British theoretical computer scientists</div> <div>People from St Leonards-on-Sea</div> </div> </div>		

Figure 7: Categories of the Wikipedia page corresponding to Alan Turing

2.2 Gold standards and training set

The gold standards and the Meij datasets, used in [1], are also provided for the comparison of your results with the baselines: TagMe, AIDA, DBpedia Spotlight and WAT.

Each gold standard is provided in a separate folder that includes two files:

- a "tsv" file for the tweet texts
 - a "tsv" file for the annotations.
- **Remark:** The columns in the annotation file of the gold standards are not the same.

3 Required tasks

1. Train a model, using Meij revised dataset, with at least two different techniques seen and experienced during lab and lecture sessions (**DNN** is mandatory). your code must be modular and include:
 - a text processing module,
 - a feature extraction module (use at least 7 features defined in [1])
2. Test your models on all the provided gold standards, as in the work in [1])
3. Compare your implemented EL system with TagMe and AIDA.
4. Provide a detailed and well-organized report on your project, including an argumentation for all your choices.

4 Submission deadline

The deadline will be announced later.

References

- [1] Nait-Hamoud, M.C., Lahfa, F. & Ennaji, A. *A step further towards a consensus on linking tweets to Wikipedia*. Evol. Intel. 16, 1825–1840 (2023). <https://doi.org/10.1007/s12065-020-00549-8>