

## **Task : Natural Language Processing**

The following text was extracted from the Wikipedia page about Edward Snowden<sup>1</sup>.

In 2013, Snowden was hired by an NSA contractor, Booz Allen Hamilton, after previous employment with Dell and the CIA. On May 20, 2013, Snowden flew to Hong Kong after leaving his job at an NSA facility in Hawaii, and in early June he revealed thousands of classified NSA documents to journalists Glenn Greenwald, Laura Poitras, and Ewen MacAskill. Snowden came to international attention after stories based on the material appeared in The Guardian and The Washington Post.

1.1 Insert this text into the Stanford Core NLP demo: <http://corenlp.run/> .

1.2 What types of output are returned?

1.3 What is the difference between Named Entity Recognition (NER) and Entity Linking (EL)?

## **Task 2: Relation Extraction**

2.1 Use Ollie to annotate the text from Task 1: <https://github.com/knowitall/ollie> (Section “Local Machine”).

2.2 Describe the output and list the most important advantages and disadvantages of this relation extraction.

## **Task 3: DBpedia Languages**

3.1 What will the following query return when executed on the English DBpedia?

```
SELECT LANG(?label) WHERE {  
    ?sub rdfs:label ?label .  
}  
GROUP BY LANG(?label)
```

3.2 Rewrite the query to add the frequency of labels per language.

3.3 Execute your query from Task 3.2 on <https://dbpedia.org/sparql> (English DBpedia) and <https://fr.dbpedia.org/sparql> (French DBpedia). List the three most frequent languages per DBpedia. What do the results tell us about the use of languages in different DBpedia language versions?

---

<sup>1</sup> [https://en.wikipedia.org/w/index.php?title=Edward\\_Snowden&oldid=836189168](https://en.wikipedia.org/w/index.php?title=Edward_Snowden&oldid=836189168)

#### **Task 4: Language Identification**

Given are the following sentence  $s$  and two language profiles:

- sentence  $s$ : “in may they may go in their therme”
- English: [ the, may, in\_, by\_, and, one, ing, his, ted ]
- German: [ sch, her, das, die, bei, ein, in\_, the, go\_ ]

4.1 Create the document profile of  $s$  as follows:

4.1a First, write down all character-based 3-grams of  $s$ . Ignore 3-grams starting with a blank and those at the start or end of the sentence.

4.1b Now, sort the set of 3-grams by descending frequency and assign ranks to them. The most frequent 3-gram gets rank 1. If more 3-grams appear equally often, they get the same rank.

4.2 Determine the language of  $s$  based on its document profile and the language profiles as follows:

4.2a Compute the Jaccard similarities between the document profile and both language profiles. Which language is identified for  $s$  according to this similarity?

4.2b Compute the “out of place” measures between the document profile and both language profiles (with default penalty  $K = 10$ ). Which language is identified for  $s$  according to this measure?

**Task 5: Text Similarity**

Given are the following three sentences<sup>2 3</sup>:

$s_{E1}$  (English): *Theresa May has announced plans to call a snap general election for the UK on 8 June.*

$s_{E2}$  (English): *Donald Tusk said the 27 other EU states would forge ahead with Brexit, saying the UK election of the parliament would not change their plans.*

$s_{G1}$  (German): *Regierungschefin Theresa May hat angekündigt, dass die Neuwahl des Parlaments in Großbritannien am 8. Juni stattfinden sollen.*

Compute the sentence similarities based on the following three measures:

- **Longest common subsequence:** length of the longest subsequence of characters that appears in both sentences
- **Biword similarity:** Split the text into biwords (word-based 2-grams, ignore commas and dots) and compute their overlap using the **Dice coefficient**.
- **Knowledge-based similarity:** Manually annotate named entities and time expressions. Then compute the knowledge-based similarity as the mean of the Jaccard similarity computed with named entities and the **Jaccard similarity** computed with time expressions.

English Sent.	German Sent.	Longest common subsequence	Biword similarity	Knowledge-based similarity
$s_{E1}$	$s_{G1}$			
$s_{E2}$	$s_{G1}$			

<sup>2</sup> <http://www.bbc.com/news/uk-politics-39629603>

<sup>3</sup> <http://www.spiegel.de/politik/ausland/theresa-may-grossbritannien-soll-am-8-juni-vorzeitig-waehlen-a-1143687.html>