

Natural Language Processing

Aims, tools, consequences

Marco Petolicchio

October 2020

A first overview to NLP

- **NLP** is an acronym for Natural Language Processing
- With the use of computational methods, we are able to retrieve quantitative information about language use, patterns...
- NLP is used either at a morphological, syntactical, phonological, semantical level, and also for comparing different languages

- We start with a text, or a collection of texts, as the input for the machine
- Then we elaborate the input in order to parse some information, e.g. as we want to retrieve a syntactical analysis
- We can make the machine to recognize all the Part-of-Speech for each word (POS-tagging)

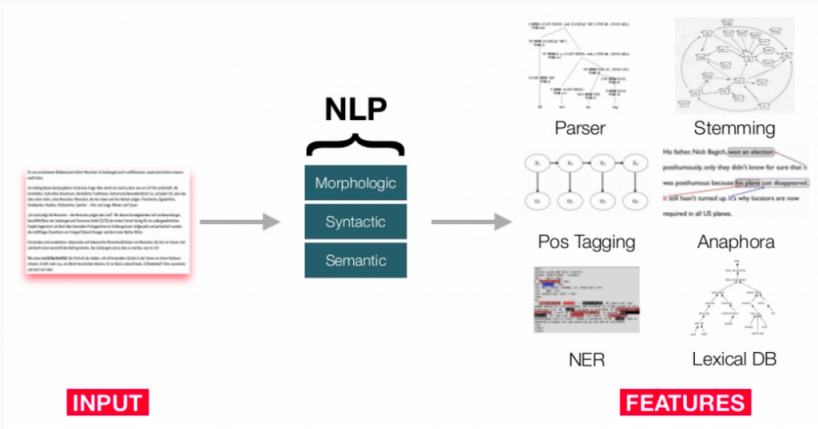


Figure 1: NLP Pipeline (from wordlift.io)

- The use of NLP is a powerful tool in the field of Human-Machine interaction, as in automated translation, summary creation from texts, assistive technologies, speech analysis, automated caption generation, recognition of text from scan and images (OCR), and so on
- But how can a machine work this manner?

Corpora

- The first thing we should be aware of is that the machine works using a set of information as a model
- A text should be compared to other texts in some respects, in order to retrieve quantitative information

- We can define a **corpus** as a collection of texts, with a certain grade of homogeneity, that we can use for a certain kind of generalization about language use and structures

AntConc 3.2.2w (Windows) 2008

File Global Settings Tool Preferences About

Corpus Files

example.txt

Concordance Concordance Plot File View Clusters Collocates Word List Keyword List

H#	KWIC	File
1	1>Quantitative corpus linguistics with R: a practi	example.txt
2	Why another introduction to corpus linguistics? 4 1.2 Outline (example.txt
3	1 Corpora 11 2.1.1 What is a corpus ? 11 2.1.2 What kinds of corp	example.txt
4	pus? 11 2.1.2 What kinds of corpora are there? 13 2.2 Frequency	example.txt
5	uency list of an unannotated corpus 113 4.1.2 A reverse frequen	example.txt
6	uency list of an unannotated corpus 117 4.1.3 A frequency list (example.txt
7	equency list of an annotated corpus 119 4.1.4 A frequency list (example.txt
8	sequences from an annotated corpus 120 4.1.5 A frequency list (example.txt
9	f word pairs from an annotated corpus 123 4.1.6 A frequency list (example.txt
10	equency list of an annotated corpus (with one word per line) 128	example.txt
11	f word pairs of an annotated corpus (with one word per line) 130	example.txt
12	iles of an (SGML POS-tagged) corpus 138 4.2.3 More complex conce	example.txt
13	iles of an (SGML POS-tagged) corpus 143 4.2.4 A lemma-based conc	example.txt
14	a POS-tagged and lemmatized corpus 147 4.3 Collocations 150 4.	example.txt
15	s 1: processing multi-tiered corpora 156 4.5 Excursus 2: Unicode	example.txt
16	Why another introduction to corpus linguistics? This book is a	example.txt

Total No. 1

Files Processed

Reset

Search Term ☒ Words ☐ Case ☒ Regex

(corpus|corpora) Advanced

Concordance Hits 667

Search Window Size 50

Start Stop Sort

Kwic Sort

☒ Level 1 0 ☐ Level 2 0 ☐ Level 3 0

Save Window

Exit

Figure 2: Example of concordances in a corpus

Any natural corpus will be skewed. Some sentences won't occur because they are obvious, others because they are false, still others because they are impolite. The corpus, if natural, will be so wildly skewed that the description [based upon it] would be no more than a mere list.

(Chomsky 1962; 159)

A good question now can arise:

*what we define with the term **text**?*

As we said, a corpus is a collection of texts.

- Texts do not mean that they are solely related to written language, because we can have audio transcript as well
- Texts are collections of words (which are collections of letters (XD) ?)
- Barely said, words are **tokens**

- In order to provide a better input for our software, it could be useful to add some information to the texts. In informatics, we can add some pieces of information to our work, using *metadata* and *markup languages*.
- For example, if we want to annotate a certain text, specifying where is the subject in each sentence, we can add this information and use it later. We can use for doing it a markup language, which has a clear structure and encodes all these information in a formal way. Examples of markup languages are HTML (the language for the majority of the web pages), XML (which is used in RSS feeds), and from XML we have different dialects as XML-TEI (for encoding literary texts, manuscripts, etc..), MML (for encoding MIDI files for music), KML (for geospatial information), etc...

- In computational linguistics the quite-default markup language is the CONLL-U language, which provides syntactic, morphological, dependency information.

ID	FORM	LEMMA	UPOSTAG	XPOSTAG	FEATS	HEAD	DEPREL	DEPS	NER
# newdoc url=http://www.poweredbyosteons.org/2012/01/brief-history-of-bioarchaeological.html									
# newdoc s3 = s3://aws-publicdatasets/common-crawl/crawl-data/CC-MAIN-2016-07/segments...									
...									
# sent_id=http://www.poweredbyosteons.org/2012/01/brief-history-of-bioarchaeological.html#60									
# text = The American Museum of Natural History was established in New York in 1869.									
0	The	the	DT	DT	-	2	det	2:det	O
1	American	American	NNP	NNP	-	2	nn	2:nn	B-Organization
2	Museum	Museum	NNP	NNP	-	7	nsubjpass	7:nsubjpass	I-Organization
3	of	of	IN	IN	-	2	prep	-	I-Organization
4	Natural	Natural	NNP	NNP	-	5	nn	5:nn	I-Organization
5	History	History	NNP	NNP	-	3	pobj	2:prep_of	I-Organization
6	was	be	VBD	VBD	-	7	auxpass	7:auxpass	O
7	established	establish	VBN	VBN	-	7	ROOT	7:ROOT	O
8	in	in	IN	IN	-	7	prep	-	O
9	New	New	NNP	NNP	-	10	nn	10:nn	B-Location
10	York	York	NNP	NNP	-	8	pobj	7:prep_in	I-Location
11	in	in	IN	IN	-	7	prep	-	O
12	1869	1869	CD	CD	-	11	pobj	7:prep_in	O
13	-	7	punct	7:punct	O
...									

Figure 3: Example of CONLL-U

As we have seen, a corpus can be annotated in order to result in a better **model** for the machine, an input that permits to the various algorithms the possibility to perform different analyses on the text that we want to process.

There are different technologies that we can use for doing it ==>

Artificial Intelligence, Machine Learning, and other things

NLP and Computational Linguistics are two faces of the same medal: one is problem-oriented and it's the technology where the theoretical approach grounds on.

With the use of Artificial Intelligence, we can elaborate a bit our understanding of linguistic facts.

Machine learning (and Deep Learning) is the application of the algorithm based on the assumptions made upon the model we have as input.

We continue it later