

NATURAL LANGUAGE PROCESSING (NLP)

Carlos J. Costa





Learning Goals

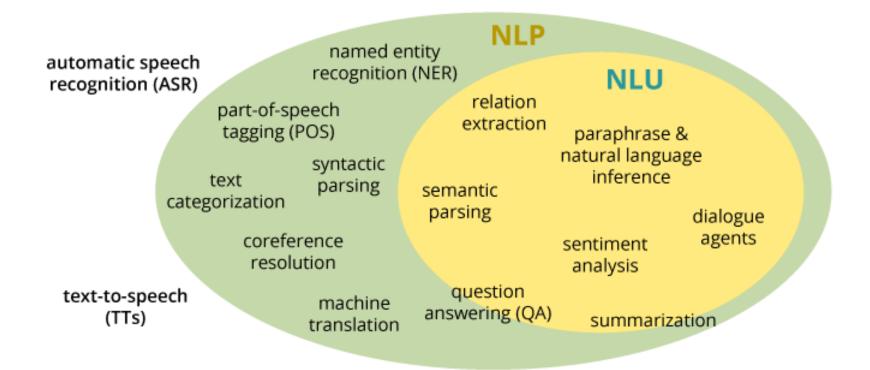
- Understand the context of NLP
- Explain main concepts
- Use main libraries



Table of Contents

- Definition
- Main Concepts
- Libraries







- Natural language processing
- Subfield of artificial intelligence, linguistics, and computer science
- Create software to process and analyze large amounts of natural language data

Information Retrieval



Machine Translation



Sentiment Analysis



Natural Language
Processing
(NLP)

Information Extraction



Question Answering



NLP

- Tokenization
- Stemming
- Lemmatization
- Part-of-Speech (POS) Tagging
- Named Entity Recognition (NER)
- Sentiment Analysis
- Language Modeling
- Machine Translation

- Text Summarization
- Information Extraction
- Question Answering
- Text Classification
- Topic Modeling
- Dependency Parsing
- Discourse Analysis



Tokenization

Breaking down text into smaller units such as words, phrases, or sentences.

```
import nltk
sentence_data = "First, I will explain you how this work. Then, you will do it. "
nltk_tokens = nltk.sent_tokenize(sentence_data)
print (nltk_tokens)

['First, I will explain you how this work.', 'Then, you will do it.']
```

```
text = "First, I will explain you how this work. Then, you will do it."
import nltk
new_text = nltk.word_tokenize(text)
print (new_text)

['First', ',', 'I', 'will', 'explain', 'you', 'how', 'this', 'work',
'.', 'Then', ',', 'you', 'will', 'do', 'it', '.']

from nltk.tokenize import RegexpTokenizer
tokenizer = RegexpTokenizer(r'\w+')
new_text=tokenizer.tokenize(text)
print (new_text)

['First', 'I', 'will', 'explain', 'you', 'how', 'this', 'work', 'Then',
'you', 'will', 'do', 'it']
```

Stemming

Removing suffixes or prefixes from words to obtain their root form.

```
from nltk.stem import PorterStemmer
e_words= ["studies", "studying", "cries", "cry"]
ps =PorterStemmer()
for w in e_words:
    rootWord=ps.stem(w)
    print(rootWord),
```

studi studi cri cri

Lemmatization

Reducing words to their base or dictionary form while still ensuring they are valid words.

```
import nltk
from nltk.stem import --> WordNetLemmatizer
wordnet_lemmatizer = WordNetLemmatizer()
text = "studies studying cries cry"
tokenization = nltk.word_tokenize(text)
for w in tokenization:
    print(wordnet_lemmatizer.lemmatize(w))
```

study studying cry cry

Part-of-Speech (POS) Tagging

Assigning grammatical categories (such as noun, verb, adjective, etc.) to words in a sentence.

Part-of-Speech (POS) Tagging

6 tokens = nltk.word_tokenize(text)
7 pos_tags = nltk.pos_tag(tokens)

text = "I am waiting for the end of the class."

[('I', 'PRP'), ('am', 'VBP'), ('waiting', 'VBG'), ('for', 'IN'), ('the', 'D
T'), ('end', 'NN'), ('of', 'IN'), ('the', 'DT'), ('class', 'NN'), ('.', '.')]

import nltk

print(pos_tags)

like to read books

POS Tagging

PRP VBP TO VB NNS

Named Entity Recognition (NER)

Identifying and classifying named entities (such as person names, locations, organizations, etc.) in text.

```
1 # Named Entity Recognition (NER):
   import nltk
 5 text = "The campus of ISEG is in Lisbon, Portugal. Carlos works there."
 6 tokens = nltk.word tokenize(text)
 7 tags = nltk.pos_tag(tokens)
 8 entities = nltk.chunk.ne_chunk(tags)
 9 print(entities)
(S
 The/DT
 campus/NN
 of/IN
  (ORGANIZATION ISEG/NNP)
 is/VBZ
 in/IN
 (GPE Lisbon/NNP)
 ,/,
 (GPE Portugal/NNP)
 (PERSON Carlos/NNP)
 works/VBZ
 there/RB
 ./.)
```

Sentiment Analysis

Determining the sentiment or opinion expressed in text, typically categorized as positive, negative, or neutral.

```
# Sentiment Analysis:

from nltk.sentiment.vader import SentimentIntensityAnalyzer

analyzer = SentimentIntensityAnalyzer()

text = "I love this movie! It's amazing."

sentiment_scores = analyzer.polarity_scores(text)

print(sentiment_scores)

{'neg': 0.0, 'neu': 0.266, 'pos': 0.734, 'compound': 0.8516}
```

Language Modeling

Building statistical models to predict the next word in a sequence of words, often used in machine translation, autocomplete, and speech recognition.

```
: 1 # Language Modeling:
2 import nltk
4 # nltk.download('gutenberg')
5 #nltk.corpus.gutenberg.fileids()
6 
7 text = nltk.corpus.gutenberg.raw('bible-kjv.txt')
8 #text = nltk.corpus.gutenberg.raw('shakespeare-hamlet.txt')
9 words = nltk.word_tokenize(text)
10 bigrams = nltk.ngrams(words, 2)
11 model = nltk.ConditionalFreqDist(bigrams)
12 print(model["face"].most_common(5)) # Predict next word after "to"
[('of', 109), (',', 79), ('to', 43), ('.', 31), ('from', 20)]
```

Machine Translation

Translating text from one language to another automatically using computational methods.

```
from deep translator import GoogleTranslator
 6
    def translate text(text, dest lang='es'):
 8
        try:
            translated_text = GoogleTranslator(source='auto', target=dest_lang).translate(text)
 9
            return translated text
10
        except Exception as e:
11
12
            return f"An error occurred: {str(e)}"
13
14
    # Example usage
    translated text = translate text("Das funktioniert tatsächlich gut.", dest lang='pt')
    print(translated text)
16
Isso realmente funciona bem.
```

Text Summarization

Generating concise summaries of longer text documents while preserving the most important information.

```
from transformers import pipeline

# Load the summarization pipeline with a specific model
summarizer = pipeline("summarization", model="t5-base", tokenizer="t5-base")

# Example text
text = "We are currently exploring the applications of Artificial Intelligence in management and economics. Our focus exte

# Generate summary
summary = summarizer(text, max_length=50, min_length=10, do_sample=False)[0]['summary_text']

# Print the summary
print(summary)

**The summary of the summary
```

our focus extends beyond studying the use of AI to include the management of artificial intelligence . we have been significant in organizing annual Artificial Intelligence and Management workshops since 2019 .

```
# Information Extraction:

import re

text = "João Duque is the dean of ISEG, University of Lisbon. Luís Carriço is the dean of FCUL, University of Lisbon."

matches = re.findall(r'([A-Za-zÀ-ÖØ-öØ-ÿ]+)\s+([A-Za-zÀ-ÖØ-öØ-ÿ]+)\s+is\s+the\s+dean\s+of\s+([^,.]+)', text)

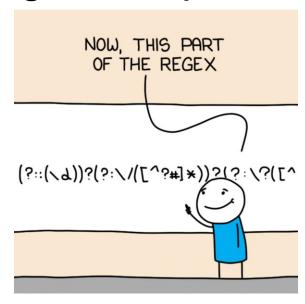
print(matches)
```

[('João', 'Duque', 'ISEG'), ('Luís', 'Carriço', 'FCUL')]

Information Extraction

Automatically extracting structured information from unstructured text, such as extracting entities, relations, and events.

Regular Expression



- Sequence of characters that specifies a search pattern.
- https://docs.python.org/3/howto/regex.html



Question Answering

Building systems that can understand and answer questions posed in natural language.

Natural Language Processing

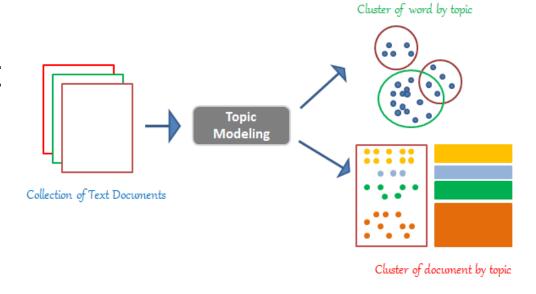
Text Classification

Categorizing text documents into predefined categories or classes, such as spam detection, sentiment classification, topic classification, etc.

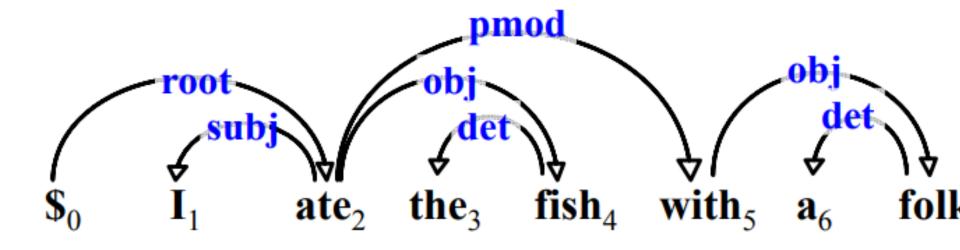
```
# Text Classification:
  from sklearn.feature extraction.text import CountVectorizer
   from sklearn.naive bayes import MultinomialNB
   from sklearn.pipeline import make pipeline
   from sklearn.model selection import train test split
   from sklearn.metrics import accuracy score
   X = ["I love this movie!", "This movie is terrible."]
   y = [1, 0]
11
12 X train, X test, y train, y test = train test split(X, y, test size=0.2, random state=42)
   model = make pipeline(CountVectorizer(), MultinomialNB())
   model.fit(X train, y train)
   predictions = model.predict(X test)
   accuracy = accuracy score(y test, predictions)
   print("Accuracy:", accuracy)
18
19
```

Topic Modeling

Identifying topics or themes present in a collection of documents, often using techniques like Latent Dirichlet Allocation (LDA).







Dependency Parsing Analyzing the grammatical structure of sentences to determine the relationships between words.



Discourse Analysis

Studying the organization and structure of connected texts beyond the sentence level, including coherence, cohesion, and rhetorical relations.

```
# Discourse Analysis
   from nltk.corpus import stopwords
 4 from nltk.tokenize import sent_tokenize, word_tokenize
6 # Example text with multiple sentences
   Natural Language Processing (NLP) is a subfield of artificial intelligence.
   It focuses on the interaction between computers and humans through natural language.
   NLP techniques are used in various applications such as machine translation, sentiment
   analysis, and text summarization.
12
13
14 # Tokenize the text into sentences
   sentences = sent tokenize(text)
17 # Tokenize each sentence into words and remove stopwords
18 | stop words = set(stopwords.words('english'))
19 tokenized sentences = [word tokenize(sentence) for sentence in sentences]
   filtered_sentences = [[word for word in tokens if word.lower() not in stop_words] for tokens in tokenized_sentences]
   # Compute lexical chains based on word similarity
   def compute lexical chains(sentences):
       for i, sentence in enumerate(sentences):
27
           for word in sentence:
               for j, prev_sentence in enumerate(sentences[:i]):
                   if word in prev sentence:
                       chain.append((word, j))
           chains.append(chain)
       return chains
35 # Compute lexical chains for the example
   lexical_chains = compute_lexical_chains(filtered_sentences)
38 # Print the lexical chains for each sentence
39 for i, chain in enumerate(lexical chains):
        print(f"Sentence {i+1} Lexical Chain: {chain}")
41
```

```
Sentence 1 Lexical Chain: []
Sentence 2 Lexical Chain: [('.', 0)]
Sentence 3 Lexical Chain: [('NLP', 0), ('.', 0)]
```

Web Scraping

- . Beautiful Soup is a HTML parser.
- This Python library is designed for screen-scraping projects.
- . Three features make it powerful:
 - navigating, searching, and modifying a parse tree
 - converts incoming documents to Unicode and outgoing documents to UTF-8
 - sits on top of popular Python parsers like lxml and html5lib.
- https://www.crummy.com/software/ BeautifulSoup/



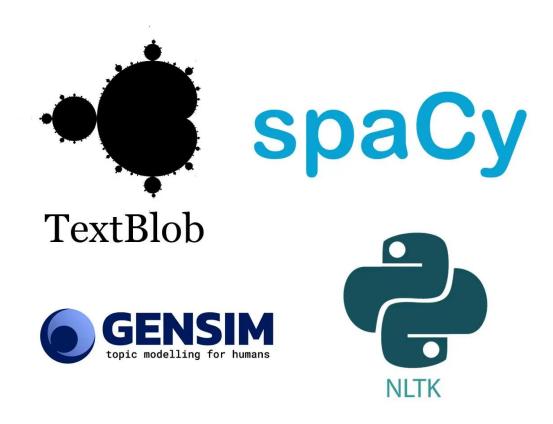


 Removing stop words is an essential step in NLP text processing

Stopword

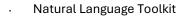
- filtering out high-frequency words that add little or no semantic value to a sentence
- for example to, at, for, is, etc.

Libraries









- A suite of text processing libraries for:
 - Classification
 - Tokenization
 - Stemming
 - . Tagging
 - Parsing
 - . Semantic reasoning
 - http://www.nltk.org/book/



NLTK

Conclusions

- NLP as subset of Al and CS
- Main concepts and Applications
- Main libraries

