



Lisbon School
of Economics
& Management
Universidade de Lisboa

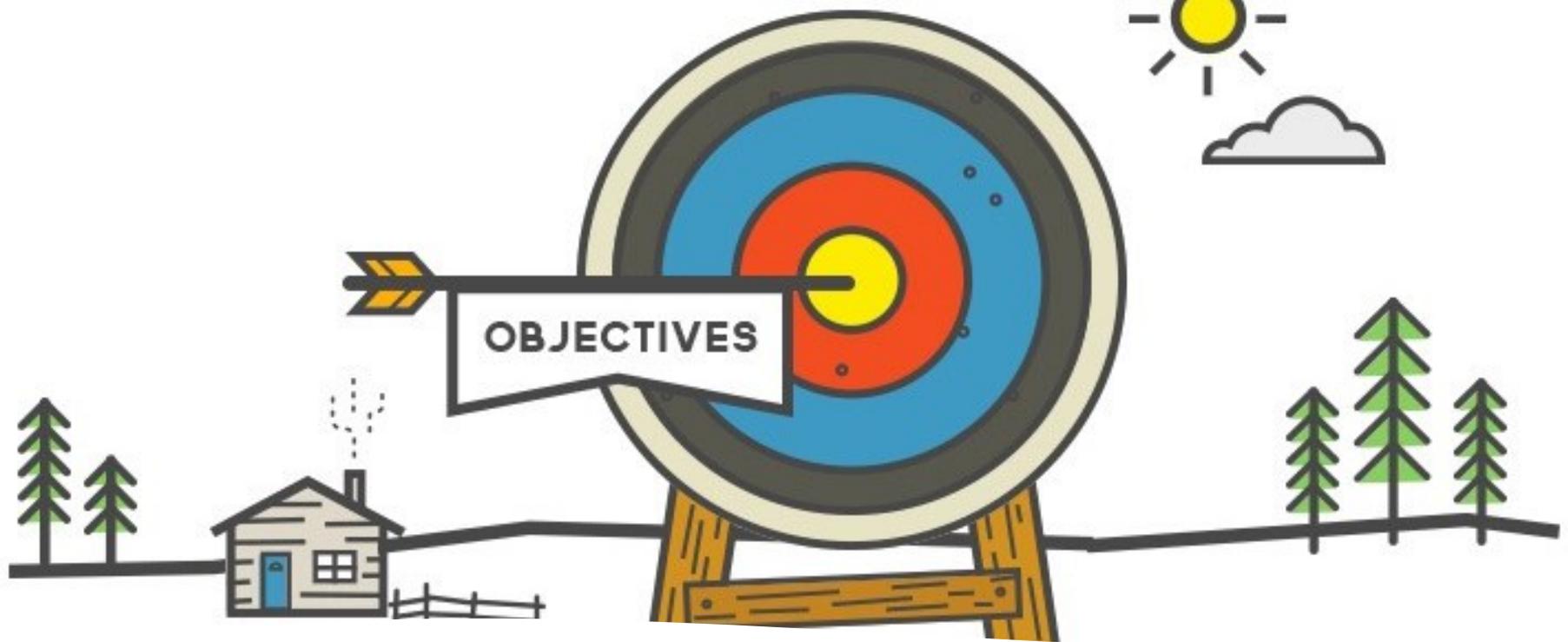


LISBOA

UNIVERSIDADE
DE LISBOA

NATURAL LANGUAGE PROCESSING (NLP)

Carlos J. Costa

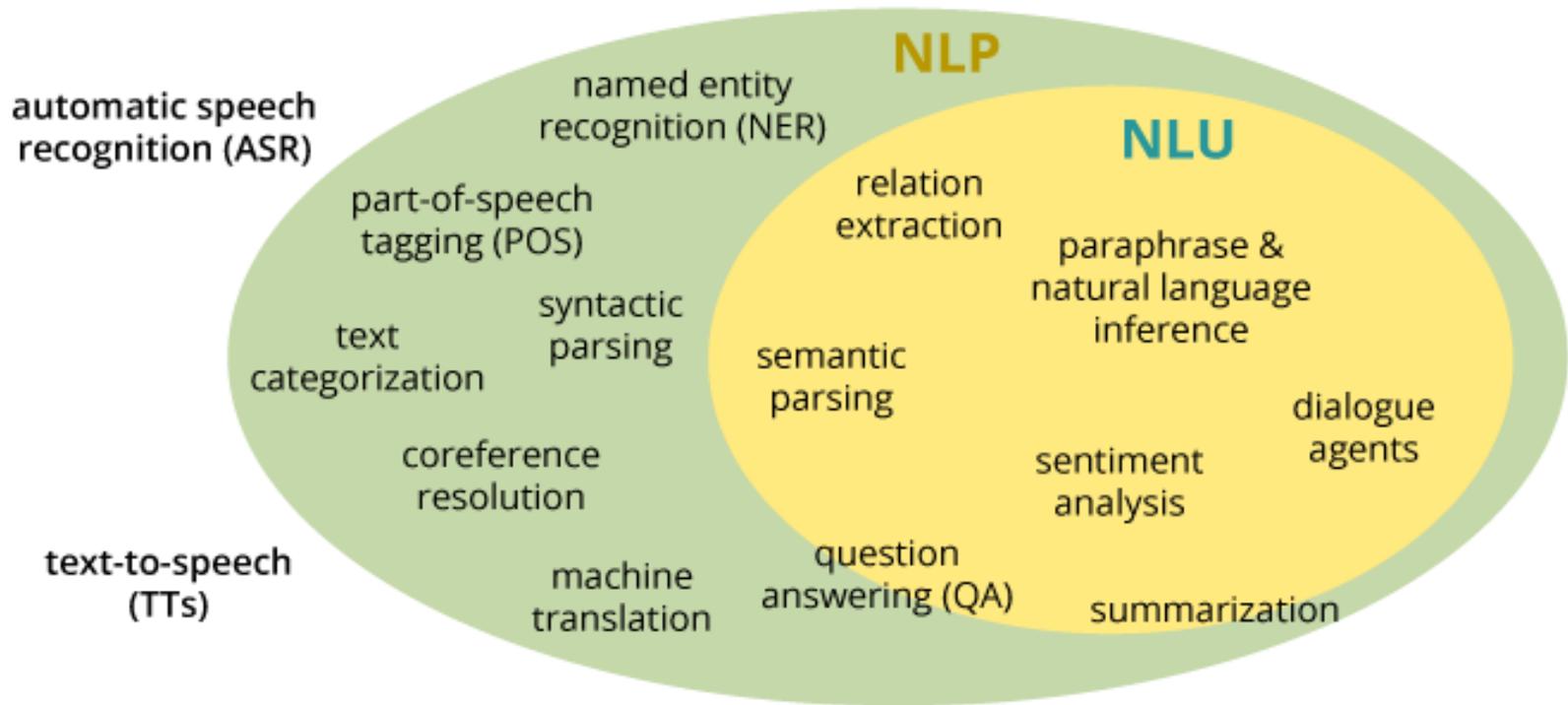


Learning Goals

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

Table of Contents

- Definition
- Main Concepts
- Libraries



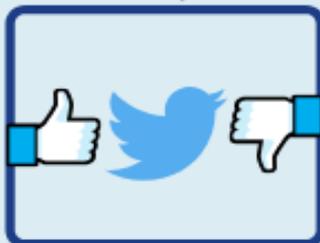
NLP

- 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



Sentiment
Analysis



Information
Extraction



Machine
Translation



Natural Language Processing (NLP)

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.

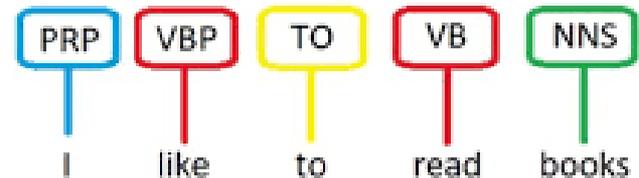
```
import nltk
from nltk.stem import WordNetLemmatizer
wordnet_lemmatizer = WordNetLemmatizer()
text = "studies studying cries cry"
tokenization = nltk.word_tokenize(text)
for w in tokenization:
    print(wordnet_lemmatizer.lemmatize(w))
```

```
study
studying
cry
cry
```

Part-of-Speech (POS) Tagging

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

POS Tagging



```
1 # Part-of-Speech (POS) Tagging
2
3 import nltk
4
5 text = "I am waiting for the end of the class."
6 tokens = nltk.word_tokenize(text)
7 pos_tags = nltk.pos_tag(tokens)
8 print(pos_tags)
9
```

```
[('I', 'PRP'), ('am', 'VBP'), ('waiting', 'VBG'), ('for', 'IN'), ('the', 'DT'), ('end', 'NN'), ('of', 'IN'), ('the', 'DT'), ('class', 'NN'), ('.', '.')]
```

Named Entity Recognition (NER)

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

```
1 # Named Entity Recognition (NER):
2
3 import nltk
4
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.

```
1 # Sentiment Analysis:
2
3 from nltk.sentiment.vader import SentimentIntensityAnalyzer
4 |
5 analyzer = SentimentIntensityAnalyzer()
6 text = "I love this movie! It's amazing."
7 sentiment_scores = analyzer.polarity_scores(text)
8 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
  3 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"
 13 |
```

```
[('of', 109), (',', 79), ('to', 43), ('.', 31), ('from', 20)]
```

Machine Translation

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

```
4
5 from deep_translator import GoogleTranslator
6
7 def translate_text(text, dest_lang='es'):
8     try:
9         translated_text = GoogleTranslator(source='auto', target=dest_lang).translate(text)
10        return translated_text
11    except Exception as e:
12        return f"An error occurred: {str(e)}"
13
14 # Example usage
15 translated_text = translate_text("Das funktioniert tatsächlich gut.", dest_lang='pt')
16 print(translated_text)
```

Isso realmente funciona bem.

Text Summarization

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

```
4 from transformers import pipeline
5
6 # Load the summarization pipeline with a specific model
7 summarizer = pipeline("summarization", model="t5-base", tokenizer="t5-base")
8
9 # Example text
10 text = "We are currently exploring the applications of Artificial Intelligence in management and economics. Our focus exte
11
12 # Generate summary
13 summary = summarizer(text, max_length=50, min_length=10, do_sample=False)[0]['summary_text']
14
15 # Print the summary
16 print(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 .

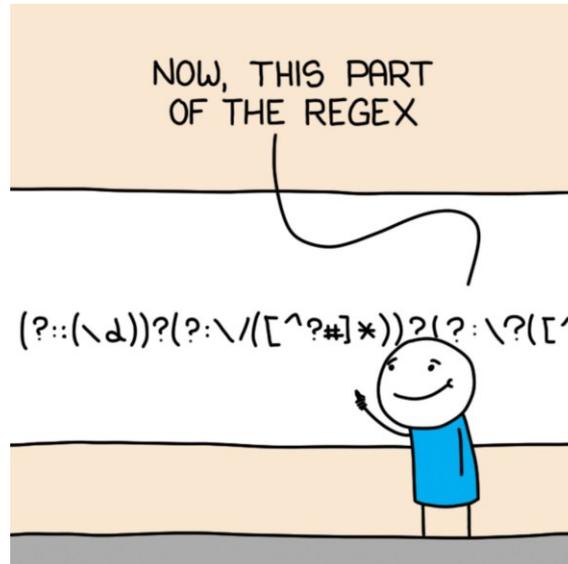
```
1 # Information Extraction:
2
3 import re
4 |
5 text = "João Duque is the dean of ISEG, University of Lisbon. Luís Carriço is the dean of FCUL, University of Lisbon."
6 matches = re.findall(r'([A-Za-zÀ-ÖØ-öø-ÿ]+\s+([A-Za-zÀ-ÖØ-öø-ÿ]+\s+is\s+the\s+dean\s+of\s+([^.]+)', text)
7 print(matches)
8
```

```
[('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.

```
1 #Question Answering:
2
3 from transformers import pipeline
4
5 qa_pipeline = pipeline("question-answering", model="bert-large-uncased-whole-word-masking-finetuned-squad",
6                       tokenizer="bert-large-uncased-whole-word-masking-finetuned-squad")
7 context = """Natural Language Processing (NLP) is a subfield of artificial intelligence that focuses on the
8             interaction between computers and humans through natural language. """
9 question = "What is NLP?"
10 answer = qa_pipeline(question=question, context=context)
11 print(answer['answer'])
12
```

Natural Language Processing

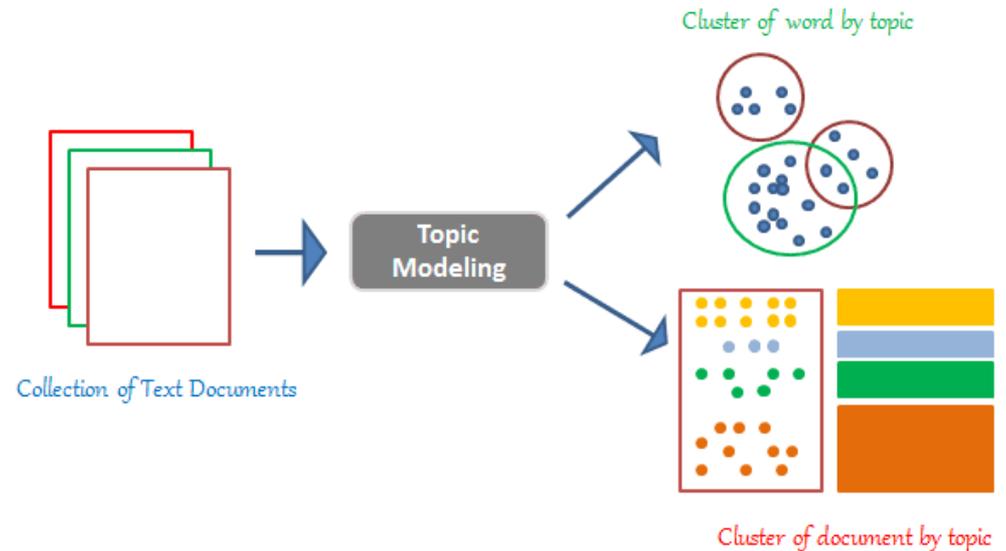
Text Classification

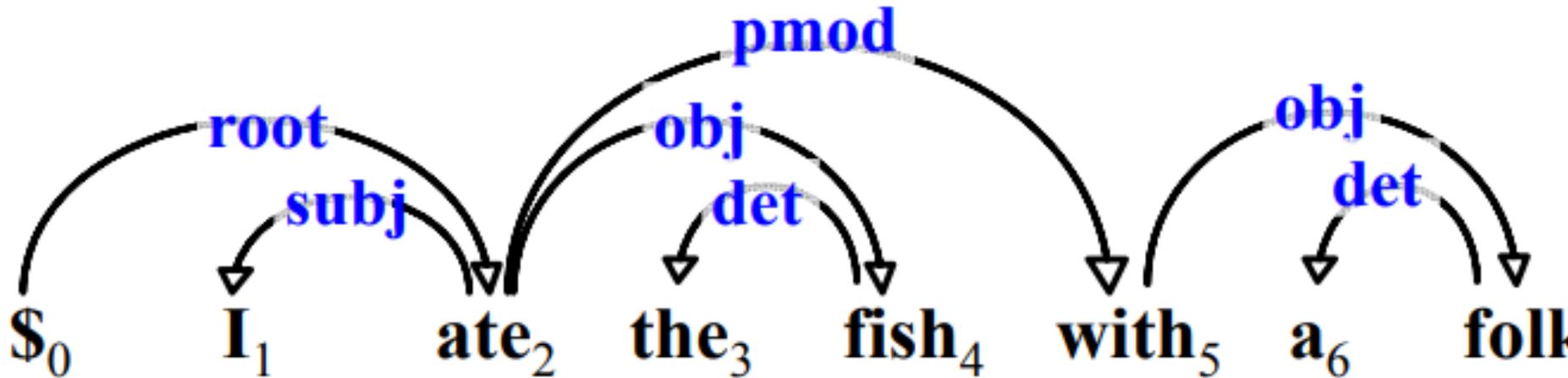
Categorizing text documents into predefined categories or classes, such as spam detection, sentiment classification, topic classification, etc.

```
1 # Text Classification:
2
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.pipeline import make_pipeline
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score
8
9 X = ["I love this movie!", "This movie is terrible."]
10 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)
13 model = make_pipeline(CountVectorizer(), MultinomialNB())
14 model.fit(X_train, y_train)
15 predictions = model.predict(X_test)
16 accuracy = accuracy_score(y_test, predictions)
17 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.

```
1 # Discourse Analysis
2
3 from nltk.corpus import stopwords
4 from nltk.tokenize import sent_tokenize, word_tokenize
5
6 # Example text with multiple sentences
7 text = """
8 Natural Language Processing (NLP) is a subfield of artificial intelligence.
9 It focuses on the interaction between computers and humans through natural language.
10 NLP techniques are used in various applications such as machine translation, sentiment
11 analysis, and text summarization.
12 """
13
14 # Tokenize the text into sentences
15 sentences = sent_tokenize(text)
16
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]
20 filtered_sentences = [[word for word in tokens if word.lower() not in stop_words] for tokens in tokenized_sentences]
21
22 # Compute lexical chains based on word similarity
23 def compute_lexical_chains(sentences):
24     chains = []
25     for i, sentence in enumerate(sentences):
26         chain = []
27         for word in sentence:
28             for j, prev_sentence in enumerate(sentences[:i]):
29                 if word in prev_sentence:
30                     chain.append((word, j))
31                     break
32             chains.append(chain)
33     return chains
34
35 # Compute lexical chains for the example
36 lexical_chains = compute_lexical_chains(filtered_sentences)
37
38 # Print the lexical chains for each sentence
39 for i, chain in enumerate(lexical_chains):
40     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

- BeautifulSoup 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/>

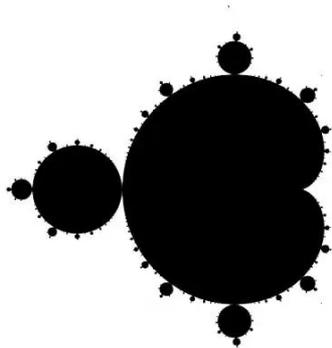




Stopword

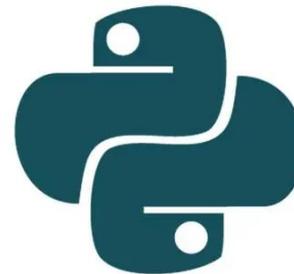
- Removing stop words is an essential step in NLP text processing
- filtering out high-frequency words that add little or no semantic value to a sentence
- for example to, at, for, is, etc.

Libraries



TextBlob

spaCy



NLTK



NLTK Or Natural Language Toolkit

NLTK

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

Conclusions

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