

# Theory: Overview of SpaCy

🕒 17 minutes   0 / 5 problems solved

Skip this topic

Start practicing

65 users solved this topic. Latest completion was about 8 hours ago.

As you may know, NLTK is one of the most popular NLP libraries. It supports tokenization, normalization, and other text processing procedures. In this topic, we will introduce **SpaCy**, another NLP library. While NLTK is mostly used by students and scholars, SpaCy is designed for developers of language-related projects. SpaCy includes fewer algorithms, but they provide more precise results. We will discuss what differs SpaCy from NLTK and highlight the main features of it.

## §1. Installation

You need to install SpaCy first; use *pip* for this:

```
1 | pip install spacy
```

You also need to import it afterward.

```
1 | import spacy
```

SpaCy (as well as NLTK) contains a good of various *models* in different languages: English, Spanish, German, Chinese, French, etc. The whole list is available [on the official SpaCy site](#).

In this topic, we will work with the `en_core_web_sm` model that was trained on blogs, news, and comments. It is used for POS-tagging, named entity recognition, and other tasks. You can install the model from the command line, as shown below:

```
1 | python -m spacy download en_core_web_sm
```

You can use one of the following ways to upload it to your program:

1. Import the model using the integrated loader.

```
1 | en_sm_model = spacy.load("en_core_web_sm")
```

2. Import it as a Python module.

```
1 | import en_core_web_sm
2 |
3 |
4 | en_sm_model = en_core_web_sm.load()
```

The model import process is basic, you can do it for any other models.

SpaCy requires a great amount of free space, and sometimes it also requires other additional programs. You may encounter several problems with the installation. Unfortunately, we can't cover every one of them, but you can find answers on special forums.

## §2. NLTK vs SpaCy

Before we start working with SpaCy, it would be good to compare the basic features of these libraries.

Criteria	NLTK	SpaCy
Multi-Language Support	Yes	Yes
Types of Inputs and Outputs	Strings	Objects
Word Vector Support	No	Yes
Performance	Slow	Fast

Current topic:

[Overview of SpaCy](#) ...

Topic depends on:

✗ [Intro to text representations](#) ...

✗ [Text normalization](#) ...

Topic is required for:

[Introduction to Stanza](#) ...

Table of contents:

[1 Overview of SpaCy](#)

[§1. Installation](#)

[§2. NLTK vs SpaCy](#)

[§3. Getting Started](#)

[§4. Other NLP Procedures](#)

[§5. Word Vectors](#)

[§6. Conclusion](#)

[Feedback & Comments](#)

Target Audience	Researchers	Developers
-----------------	-------------	------------

To sum up, SpaCy is a modern alternative for traditional NLTK.

### §3. Getting Started

Let's overview the basic NLP procedures and how they are implemented. It is essential to remember that SpaCy turns processed text not into strings, but objects.

#### 1. Tokenization

SpaCy automatically divides your document into tokens when you use a for-loop.

```
1 |
doc = en_sm_model("Microsoft News delivers news from the most popular and trusted publishers.")
2 |     for i in doc:
3 |         print(i.text)
4 |
5 |     # Microsoft
6 |     # News
7 |     # delivers
8 |     # news
9 |     # from
10 |    # the
11 |
12 |    # most
13 |
14 |    # popular
15 |
16 |    # and
17 |
18 |    # trusted
19 |
20 |    # publishers
21 |
22 |    # .
```

First, we create a doc-object using the already installed model. Then, with the help of a for-loop, we print all the token objects of the sentence.

#### 2. POS-tagging and Lemmatization

You can establish the lemma for each token as well as its part of speech. Use the `token.lemma_` method for lemmas and the `token.pos_` method for parts of speech.

```
1 |
doc = en_sm_model("Microsoft News delivers news from the most popular and trusted publishers.")
2 |     for i in doc:
3 |         print("{0} - {1} - {2}".format(i.text, i.lemma_, i.pos_))
4 |
5 |     # Microsoft - Microsoft - PROPN
6 |     # News - News - PROPN
7 |     # delivers - deliver - VERB
8 |     # news - news - NOUN
9 |     # from - from - ADP
10 |
11 |    # the - the - DET
12 |
13 |    # most - most - ADV
14 |
15 |    # popular - popular - ADJ
16 |
17 |    # and - and - CCONJ
18 |
19 |    # trusted - trusted - ADJ
20 |
21 |    # publishers - publisher - NOUN
22 |
23 |    # . - . - PUNCT
```

As you can see, we can easily get the necessary information about the tokens. You can find more token attributes in the [Linguistic Features](#) sections of the official SpaCy documentation.

## §4. Other NLP Procedures

SpaCy also provides more ways to process texts for further analysis. Their basic overview is presented below.

### 1. Stopword Removal

Some tasks may require stopwords removal. To implement this feature in SpaCy, you can import the built-in stopwords.

```
1 | from spacy.lang.en.stop_words import STOP_WORDS
```

Let's have a look at them. It will produce a set:

```
1 | print(STOP_WORDS)
2 |
# {'afterwards', 'would', 'others', 'thence', 'itself', 'besides', 'five' ...}
```

So, the example given in the previous section can be processed in the following way:

```
1 |
doc = en_sm_model("Microsoft News delivers news from the most popular and trusted publishers.")
2 | for i in doc:
3 |     if i.lemma_ not in STOP_WORDS:
4 |         print(i.lemma_)
5 |
6 | # Microsoft
7 | # News
8 | # deliver
9 | # news
10 |
11 | # popular
12 |
13 | # trusted
14 |
15 | # publisher
16 |
17 | # .
```

Be careful with lemmatization. If you don't lemmatize your tokens in advance, some stopwords won't be recognized.

All articles, conjunctions, and some frequent adverbs are omitted. Mind the dot; the punctuation marks **are not** included in the stop words.

### 2. Named Entity Recognition

A **named entity** is an object of the real world, for example, a person, an organization, etc. SpaCy can recognize different named entities in a document by forecasting it with the help of a built-in model. The standard way of named entity recognition (NER) is `doc.ents`.

```
1 |
doc = en_sm_model("Microsoft News delivers news from the most popular and trusted publishers.")
2 | for i in doc.ents:
3 |     print(i.text)
4 |
5 | # Microsoft News
```

The output is a line with the company name.

### 3. Dependency Tree

You can also see a syntactic structure of sentences. A dependency tree reflects syntactic relations between words in a sentence. In the example below, we will print a head element and a child (dependent) element of each phrase.

```

1 |
doc = en_sm_model("Microsoft News delivers news from the most popular and trusted
publishers.")
2 |   for i in doc:
3 |       print("{0} --> {1}".format(i.head.text, i.text))
4 |
5 |   # News --> Microsoft
6 |   # delivers --> News
7 |   # delivers --> delivers
8 |   # delivers --> news
9 |   # news --> from
10 |
11 |   # publishers --> the
12 |
13 |   # popular --> most
14 |
15 |   # publishers --> popular
16 |
17 |   # popular --> and
18 |
19 |   # popular --> trusted
20 |
21 |   # from --> publishers
22 |
23 |   # delivers --> .

```

The first element of the output is the head element, the second one is a child. Mind the verb that is the head of the sentence. It doesn't depend on anything. In the tree representation, the head element is dependent on itself, so the line `delivers --> delivers` is correct.

You can also visualize your syntactic structure using the `displacy` library. More information on this library is available in the [displaCy](#) section of the official SpaCy documentation. The result is a picture of visualized relations between words.



In the picture, you can also find different types of dependencies. For example, `amod` stands for an **adjectival modifier**. You can learn more about types of dependencies in the [Annotations Specifications](#) of the official documentation.

## §5. Word Vectors

Word similarity can be determined by comparing **word vectors** that are also known as word embeddings. Some models in SpaCy provide built-in vectors. Unfortunately, `en_core_web_sm` cannot deal with vectors. In the example below, we use `en_core_web_lg`, which is the biggest model for the English language.

```
1 en_lg_model = en_core_web_lg.load()
2 doc = en_lg_model("Cats don't like dogs.")
3 for token1 in doc:
4     for token2 in doc:
5         print(token1.text, token2.text, token1.similarity(token2))
6
7 # Cats Cats 1.0
8 # Cats do 0.34179398
9 # Cats n't 0.3543607
10
11 # Cats like 0.3767576
12
13 # Cats dogs 0.83117634
14
15 # ...
```

The words represented in the sentence are common in English, so we use them to demonstrate how the built-in vectors work in SpaCy. In our case, the model's predictions are to the point. A "*cat*" is very similar to "*dog*", but "*cat*" has a lower percentage of similarity with the word "*like*". The same words are 100% similar to each other.

## \$6. Conclusion

In this topic, we have covered the main aspects of working with SpaCy. So far, we have learned:

- how to install SpaCy and download modules;
- the common and distinct features of NLTK and SpaCy;
- how to implement the basic NLP procedures: tokenization, lemmatization, and POS-tagging;
- how to use the built-in stopwords set for text processing;
- how to implement NER;
- how to create a syntactic tree of a sentence;
- how to work with word vectors.

Of course, SpaCy has more possibilities, the [official documentation](#) can help you with them.

 Report a typo

9 users liked this theory. 0 didn't like it. What about you?



Start practicing

[Comments \(0\)](#)

[Hints \(0\)](#)

[Useful links \(0\)](#)

[Show discussion](#)