



Đăng nhập vào Medium bằng Google



Kiên
htrungkien03012003@gmail.com




Kiên Hoàng
kienhthe170770@fpt.edu.vn

2 tài khoản khác

Two minutes NLP — SpaCy sheet

POS tagging, dependency parsing, NER, and sentence

**Fabio Chiusano** · Follow

Published in NLPlanet · 3 min read · Jan 26, 2022


 85  1



Photo by [Shelbey Fordyce](#) on [Unsplash](#)

SpaCy is a free, open-source library for advanced Natural Language Processing in Python. It is designed specifically for production use and helps build applications that process and understand large volumes of text. It can be used for a multitude of use cases, such as information extraction, natural language understanding systems or to pre-process text for deep learning.

List of spaCy tasks

Here's a list of NLP tasks that spaCy can perform.

NAME	DESCRIPTION
Tokenization	Segmenting text into words, punctuations marks etc.
Part-of-speech (POS) Tagging	Assigning word types to tokens, like verb or noun.
Dependency Parsing	Assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.
Lemmatization	Assigning the base forms of words. For example, the lemma of "was" is "be", and the lemma of "rats" is "rat".
Sentence Boundary Detection (SBD)	Finding and segmenting individual sentences.
Named Entity Recognition (NER)	Labelling named "real-world" objects, like persons, companies or locations.
Entity Linking (EL)	Disambiguating textual entities to unique identifiers in a knowledge base.
Similarity	Comparing words, text spans and documents and how similar they are to each other.
Text Classification	Assigning categories or labels to a whole document, or parts of a document.
Rule-based Matching	Finding sequences of tokens based on their texts and linguistic annotations, similar to regular expressions.
Training	Updating and improving a statistical model's predictions.
Serialization	Saving objects to files or byte strings.

List of tasks that spaCy can perform. Image from <https://spacy.io/usage/spacy-101>.

While some of spaCy's features work independently, others require trained pipelines to be loaded. SpaCy currently offers trained pipelines for a variety of languages, which can be installed as individual Python modules. Here's an example where we download the trained pipeline *en_core_web_sm*.

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

spacy_1.sh hosted with ❤ by GitHub

[view raw](#)

The trained pipeline you choose always depends on your use case and the texts you're working with. For a general-purpose use case, the small, default pipelines (i.e. the ones ending in *_sm*) are always a good start.

Tokenization

Tokenization consists in segmenting text into words, punctuations marks, etc. This is done by applying rules specific to each language.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("The cat is on the table")
5  for token in doc:
6      print(token.text)
7
8  # The
9  # cat
10 # is
11 # on
12 # the
13 # table
```

spacy_2.py hosted with ❤ by GitHub

[view raw](#)

POS Tagging

POS (Part of Speech) Tagging refers to categorizing words in a text in correspondence with a particular part of speech, depending on the definition of the word and its context.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("The cat is on the table")
5  for token in doc:
6      print(f"{token.text} --- POS: {token.pos_}, {token.tag_}")
7
8  # The --- POS: DET, DT
9  # cat --- POS: NOUN, NN
10 # is --- POS: AUX, VBZ
11 # on --- POS: ADP, IN
12 # the --- POS: DET, DT
13 # table --- POS: NOUN, NN
```

spacy_3.py hosted with ❤ by GitHub

[view raw](#)

The `pos_` attribute contains the simple UPOS part-of-speech tag, whereas the `tag_` attribute contains the detailed POS tag.

Dependency Parsing

Dependency Parsing consists in assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("The cat is on the table")
5  for token in doc:
6      print(f"{token.text} --- dependency label: {token.dep_}")
7
8  # The --- dependency label: det
9  # cat --- dependency label: nsubj
10 # is --- dependency label: ROOT
11 # on --- dependency label: prep
12 # the --- dependency label: det
13 # table --- dependency label: pobj
```

spacy_4.py hosted with ❤ by GitHub

[view raw](#)

Stopwords

Stopwords are the most common words of a language, which are often ignored in NLP tasks as they usually carry little meaning to the sentences.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("The cat is on the table")
5  for token in doc:
6      print(f"{token.text} --- is stopwords: {token.is_stop}")
7
8  # The --- is stopwords: True
9  # cat --- is stopwords: False
10 # is --- is stopwords: True
11 # on --- is stopwords: True
12 # the --- is stopwords: True
13 # table --- is stopwords: False
```

spacy_5.py hosted with ❤️ by GitHub

[view raw](#)

Lemmatization

Lemmatization assigns the base forms of words. For example, the lemma of “was” is “be”, and the lemma of “dogs” is “dog”.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("The cat is on the table")
5  for token in doc:
6      print(f"{token.text} --- lemma: {token.lemma_}")
7
8  # The --- lemma: the
9  # cat --- lemma: cat
10 # is --- lemma: be
11 # on --- lemma: on
12 # the --- lemma: the
13 # table --- lemma: table
```

spacy_6.py hosted with ❤️ by GitHub

[view raw](#)

Named Entity Recognition (NER)

Named Entity Recognition refers to labeling named “real-world” objects in texts, like persons, companies, or locations.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_sm")
4  doc = nlp("Elon Musk cofounded the electronic-payment firm PayPal and formed SpaceX.")
5
6  for ent in doc.ents:
7      print(ent.text, ent.start_char, ent.end_char, ent.label_)
8
9  # Elon Musk 0 9 PERSON
10 # PayPal 48 54 ORG
```

spacy_7.py hosted with ❤️ by GitHub

[view raw](#)

Word embeddings

A word embedding is a learned representation (usually a vector of numbers) for text where words that have the same meaning have a similar

representation.

To make them compact and fast, spaCy’s small pipeline packages (all packages that end in `_sm`) don’t ship with word vectors and only include context-sensitive tensors. This means you can still use the `similarity()` methods to compare sentences and words, but the result won’t be as good, and individual tokens won’t have any vectors assigned. So in order to use real word vectors, you need to download a larger pipeline package.

```
1 python -m spacy download en_core_web_md
```

spacy_8.sh hosted with ❤️ by GitHub [view raw](#)

This is how you get word embeddings with spaCy.

```
1 import spacy
2
3 nlp = spacy.load("en_core_web_md")

6 vectors = []
7 for token in tokens:
8     print(token.text, token.has_vector, token.is_oov)
9     vectors.append(token.vector)
10
11 # The True False
12 # cat True False
13 # is True False
14 # on True False
15 # the True False
16 # aofafgag False True
17
18 print(vectors[0])
19
20 # [ 2.7204e-01 -6.2030e-02 -1.8840e-01  2.3225e-02 -1.8158e-02  6.7192e-03
21 # -1.3877e-01  1.7708e-01  1.7709e-01  2.5882e+00 -3.5179e-01 -1.7312e-01
22 #  4.3285e-01 -1.0708e-01  1.5006e-01 -1.9982e-01 -1.9093e-01  1.1871e+00
23 #  ...
```

spacy_9.py hosted with ❤️ by GitHub [view raw](#)

Sentence similarity

With spaCy, you can compute similarities between sentences. This is done by averaging the word embeddings of the words in each sentence and then computing similarity with a similarity measure.

```
1  import spacy
2
3  nlp = spacy.load("en_core_web_md") # make sure to use larger package!
4  doc1 = nlp("I like salty fries and hamburgers.")
5  doc2 = nlp("Fast food tastes very good.")
6  doc3 = nlp("Where is the cat.")
7
8  # Similarity of doc1 and doc2
9  print(doc1.similarity(doc2))
10
11 # 0.7799485853415737
12
13 # Similarity of doc1 and doc3
14 print(doc1.similarity(doc3))
15
16 # 0.6210606690259671
```

spacy_10.py hosted with ❤ by GitHub [view raw](#)

Thank you for reading! If you are interested in learning more about NLP, remember to follow NLPlanet on [Medium](#), [LinkedIn](#), and [Twitter](#)!

NLPlanet related posts

Awesome NLP — 21 popular NLP libraries of 2022 The landscape of NLP libraries medium.com	
Two minutes NLP — A Taxonomy of Tokenization Methods Word-level, Character-level, BPE, WordPiece, and SentencePiece medium.com	
Two minutes NLP — Sentence Transformers cheat sheet Sentence Embeddings, Text Similarity, Semantic Search, and Image Search medium.com	



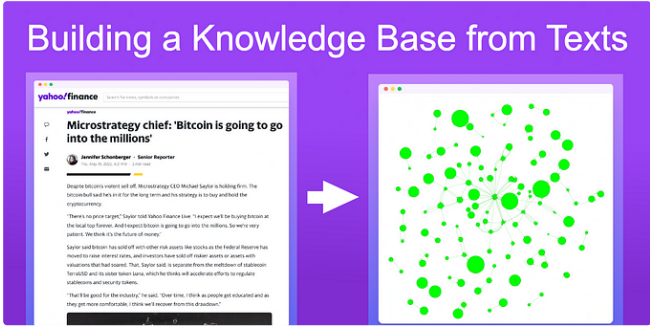
Written by Fabio Chiusano

5.2K Followers · Editor for NLPlanet

Freelance data scientist — Top Medium writer in Artificial Intelligence

Follow

More from Fabio Chiusano and NLPlanet



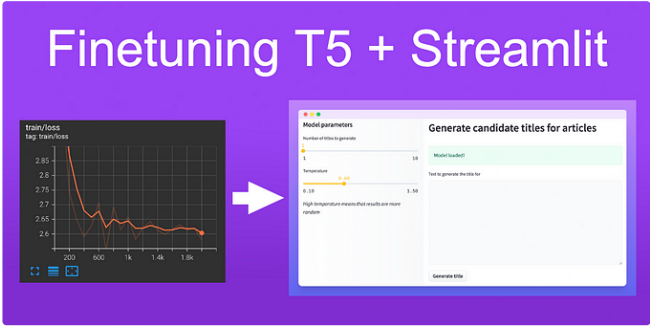
Fabio Chiusano in NLPlanet

Building a Knowledge Base from Texts: a Full Practical Example

Implementing a pipeline for extracting a Knowledge Base from texts or online articles

12 min read · May 25, 2022

1.1K 10



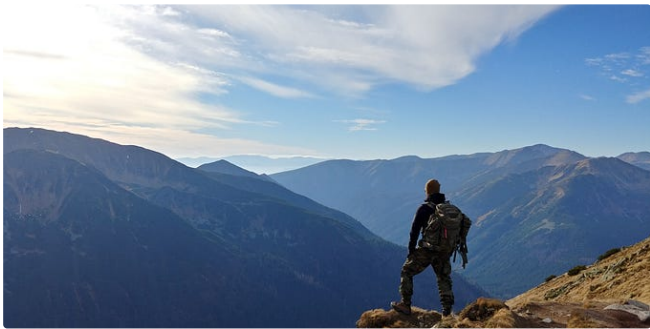
Fabio Chiusano in NLPlanet

A Full Guide to Finetuning T5 for Text2Text and Building a Demo...

All you need to know to build a full demo: Hugging Face Hub, Tensorboard, Streamlit,...

15 min read · May 17, 2022

440 10



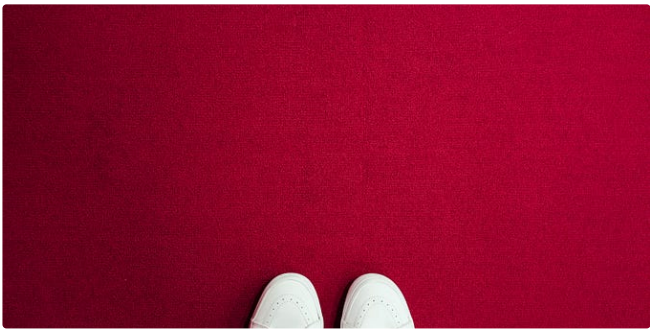
Fabio Chiusano in NLPlanet

Two minutes NLP—Perplexity explained with simple probabilities

Language models, sentence probabilities, entropy

5 min read · Jan 27, 2022

190 3



Fabio Chiusano in NLPlanet

Two minutes NLP—Learn the ROUGE metric by examples

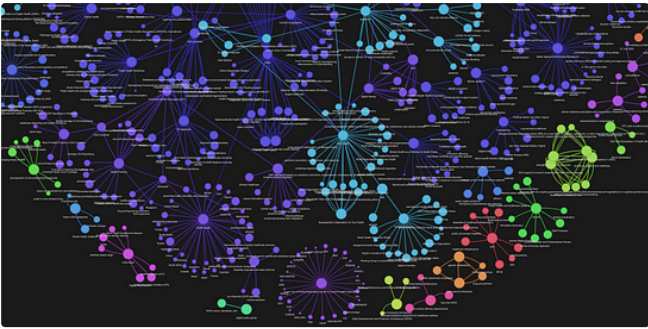
ROUGE-N, ROUGE-L, ROUGE-S, pros and cons, and ROUGE vs BLEU

5 min read · Jan 19, 2022

132 2

See all from Fabio Chiusano See all from NLPlanet

Recommended from Medium



 Rahul Nayak in Towards Data Science

How to Convert Any Text Into a Graph of Concepts

A method to convert any text corpus into a Knowledge Graph using Mistral 7B.

12 min read · Nov 10

 4.4K  42 



 Youssef Hosni in Towards AI

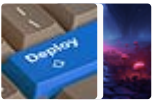
Andrej Karpathy LLM Paper Reading List for LLM Mastery

LLM Paper Recommendation from World Leading AI Researcher

🌟 · 10 min read · Dec 7

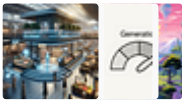
 443  1 

Lists



Predictive Modeling w/ Python

20 stories · 690 saves



Natural Language Processing

986 stories · 475 saves



Practical Guides to Machine Learning


10 stories · 785 saves



data science and AI

38 stories · 4 saves



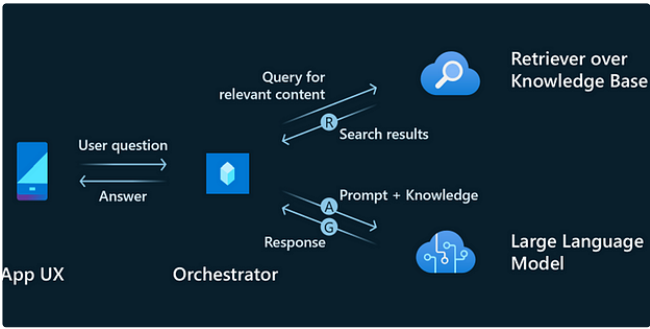
 Gursev Pirge in John Snow Labs

Text Preprocessing: Splitting texts into sentences with Spark NLP

Using Sentence Detection in Spark NLP for Text Preprocessing

8 min read · Jul 21

 2  




 James Nguyen

Forget RAG: Embrace agent design for a more intelligent grounded...

The Retrieval Augmented Generation (RAG) design pattern has been commonly used to...

6 min read · Nov 18

 314  4 



Chetankumar Khadke

Information extraction with Zephyr 7B

Information extraction with LLM uses advanced language models like Zephyr to...

11 min read · Nov 27



122



Dilip Kashyap

Google’s new AI Model Gemini now available in Bard, here is how to use

Google recently introduced Gemini AI, its latest large language model (LLM), built to...

3 min read · Dec 7



598



13



See more recommendations