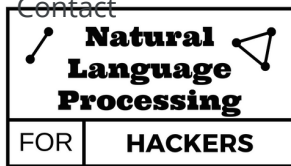
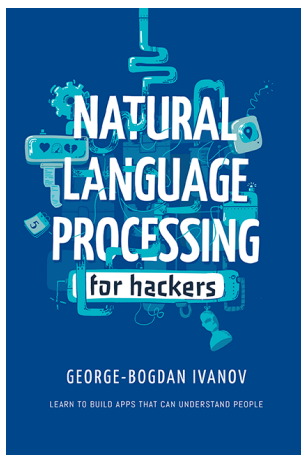


Contact



SEARCH

THE NLP-FOR-
HACKERS BOOK

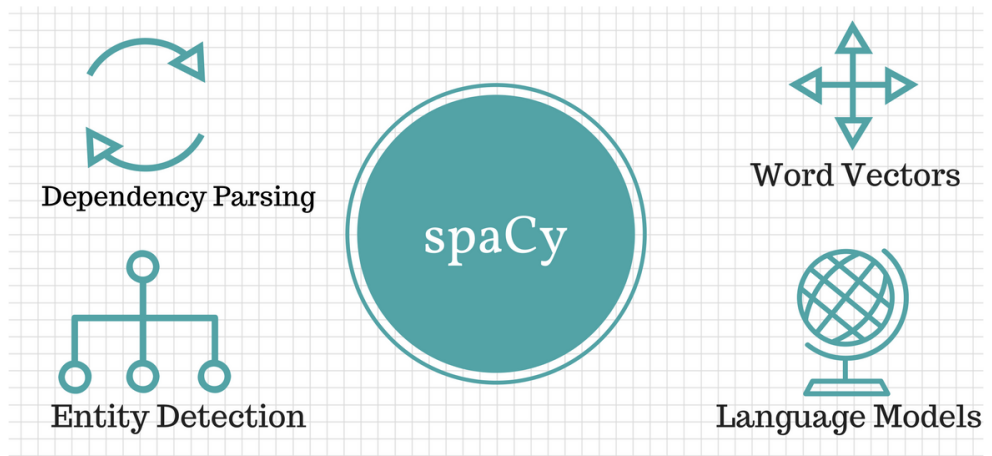


Get the

book

LIKE MY
TUTORIALS?

 Buy me a
coffee



Complete Guide to spaCy

Updates

- **29-Apr-2018** – Fixed import in extension code (Thanks Ruben)

spaCy is a relatively new framework in the Python Natural Language Processing environment but it quickly gains ground and will most likely become the de facto library. There are some really good reasons for its popularity:


➡ It's really **FAST**


Written in Cython, it was specifically designed to be as fast as possible


➡ It's really **ACCURATE**

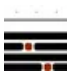
spaCy implementation of its dependency parser is one of the best-performing in the world:


[It Depends: Dependency Parser](#)

 [Complete Guide to Topic Modeling](#)

 [Complete Guide to spaCy](#)

 [Complete guide to build your own Named Entity Recognizer with Python](#)

 [Splitting text into sentences](#)

 [Recipe: Text clustering using NLTK and scikit-learn](#)

TOPICS

[Deep Learning \(4\)](#)

[Deep Natural Language Processing \(4\)](#)

[Machine Learning \(5\)](#)

[NLTK \(11\)](#)

[Resources \(2\)](#)

[Scikit-Learn \(8\)](#)

[spaCy \(2\)](#)

[Text Mining \(4\)](#)

[Uncategorized \(2\)](#)

[Wordnet \(3\)](#)

➔ Batteries included

- **Index preserving tokenization** (details about this later)
- Models for *Part Of Speech tagging*, *Named Entity Recognition* and *Dependency Parsing*
- Supports **8 languages** out of the box
- Easy and **beautiful visualizations**
- Pretrained **word vectors**

➔ Extensible

It plays nicely with all the other already existing tools that you know and love: **Scikit-Learn**, **TensorFlow**, **gensim**

➔ DeepLearning Ready

It also has its own deep learning framework that's especially designed for NLP tasks:
[Thinc](#)

Quickstart

spaCy is easy to install:

```
1 pip install -U spaCy
2 python -m spacy download en
3
```

Notice that the installation doesn't automatically download the English model. We need to do that ourselves.

```
1 import spacy
2 nlp = spacy.load('en')
3 doc = nlp('Hello World!')
4 for token in doc:
5     print("'" + token.text + "'", token.idx)
6
7 # "Hello" 0
8 # " " 6
9 # "World" 10
10 # "!" 15
11
```

```
1 doc = nlp("Next week I'll be in Madrid.")
2 for token in doc:
3     print("{0}\t{1}\t{2}\t{3}\t{4}\t{5}\t{6}\t{7}\t{8}\t{9}\t{10}\t{11}\t{12}\t{13}\t{14}\t{15}\t{16}\t{17}\t{18}\t{19}\t{20}\t{21}\t{22}\t{23}\t{24}\t{25}\t{26}\t{27}\t{28}\t{29}\t{30}\t{31}\t{32}\t{33}\t{34}\t{35}\t{36}\t{37}\t{38}\t{39}\t{40}\t{41}\t{42}\t{43}\t{44}\t{45}\t{46}\t{47}\t{48}\t{49}\t{50}\t{51}\t{52}\t{53}\t{54}\t{55}\t{56}\t{57}\t{58}\t{59}\t{60}\t{61}\t{62}\t{63}\t{64}\t{65}\t{66}\t{67}\t{68}\t{69}\t{70}\t{71}\t{72}\t{73}\t{74}\t{75}\t{76}\t{77}\t{78}\t{79}\t{80}\t{81}\t{82}\t{83}\t{84}\t{85}\t{86}\t{87}\t{88}\t{89}\t{90}\t{91}\t{92}\t{93}\t{94}\t{95}\t{96}\t{97}\t{98}\t{99}\t{100}\t{101}\t{102}\t{103}\t{104}\t{105}\t{106}\t{107}\t{108}\t{109}\t{110}\t{111}\t{112}\t{113}\t{114}\t{115}\t{116}\t{117}\t{118}\t{119}\t{120}\t{121}\t{122}\t{123}\t{124}\t{125}\t{126}\t{127}\t{128}\t{129}\t{130}\t{131}\t{132}\t{133}\t{134}\t{135}\t{136}\t{137}\t{138}\t{139}\t{140}\t{141}\t{142}\t{143}\t{144}\t{145}\t{146}\t{147}\t{148}\t{149}\t{150}\t{151}\t{152}\t{153}\t{154}\t{155}\t{156}\t{157}\t{158}\t{159}\t{160}\t{161}\t{162}\t{163}\t{164}\t{165}\t{166}\t{167}\t{168}\t{169}\t{170}\t{171}\t{172}\t{173}\t{174}\t{175}\t{176}\t{177}\t{178}\t{179}\t{180}\t{181}\t{182}\t{183}\t{184}\t{185}\t{186}\t{187}\t{188}\t{189}\t{190}\t{191}\t{192}\t{193}\t{194}\t{195}\t{196}\t{197}\t{198}\t{199}\t{200}\t{201}\t{202}\t{203}\t{204}\t{205}\t{206}\t{207}\t{208}\t{209}\t{210}\t{211}\t{212}\t{213}\t{214}\t{215}\t{216}\t{217}\t{218}\t{219}\t{220}\t{221}\t{222}\t{223}\t{224}\t{225}\t{226}\t{227}\t{228}\t{229}\t{230}\t{231}\t{232}\t{233}\t{234}\t{235}\t{236}\t{237}\t{238}\t{239}\t{240}\t{241}\t{242}\t{243}\t{244}\t{245}\t{246}\t{247}\t{248}\t{249}\t{250}\t{251}\t{252}\t{253}\t{254}\t{255}\t{256}\t{257}\t{258}\t{259}\t{260}\t{261}\t{262}\t{263}\t{264}\t{265}\t{266}\t{267}\t{268}\t{269}\t{270}\t{271}\t{272}\t{273}\t{274}\t{275}\t{276}\t{277}\t{278}\t{279}\t{280}\t{281}\t{282}\t{283}\t{284}\t{285}\t{286}\t{287}\t{288}\t{289}\t{290}\t{291}\t{292}\t{293}\t{294}\t{295}\t{296}\t{297}\t{298}\t{299}\t{300}\t{301}\t{302}\t{303}\t{304}\t{305}\t{306}\t{307}\t{308}\t{309}\t{310}\t{311}\t{312}\t{313}\t{314}\t{315}\t{316}\t{317}\t{318}\t{319}\t{320}\t{321}\t{322}\t{323}\t{324}\t{325}\t{326}\t{327}\t{328}\t{329}\t{330}\t{331}\t{332}\t{333}\t{334}\t{335}\t{336}\t{337}\t{338}\t{339}\t{340}\t{341}\t{342}\t{343}\t{344}\t{345}\t{346}\t{347}\t{348}\t{349}\t{350}\t{351}\t{352}\t{353}\t{354}\t{355}\t{356}\t{357}\t{358}\t{359}\t{360}\t{361}\t{362}\t{363}\t{364}\t{365}\t{366}\t{367}\t{368}\t{369}\t{370}\t{371}\t{372}\t{373}\t{374}\t{375}\t{376}\t{377}\t{378}\t{379}\t{380}\t{381}\t{382}\t{383}\t{384}\t{385}\t{386}\t{387}\t{388}\t{389}\t{390}\t{391}\t{392}\t{393}\t{394}\t{395}\t{396}\t{397}\t{398}\t{399}\t{400}\t{401}\t{402}\t{403}\t{404}\t{405}\t{406}\t{407}\t{408}\t{409}\t{410}\t{411}\t{412}\t{413}\t{414}\t{415}\t{416}\t{417}\t{418}\t{419}\t{420}\t{421}\t{422}\t{423}\t{424}\t{425}\t{426}\t{427}\t{428}\t{429}\t{430}\t{431}\t{432}\t{433}\t{434}\t{435}\t{436}\t{437}\t{438}\t{439}\t{440}\t{441}\t{442}\t{443}\t{444}\t{445}\t{446}\t{447}\t{448}\t{449}\t{450}\t{451}\t{452}\t{453}\t{454}\t{455}\t{456}\t{457}\t{458}\t{459}\t{460}\t{461}\t{462}\t{463}\t{464}\t{465}\t{466}\t{467}\t{468}\t{469}\t{470}\t{471}\t{472}\t{473}\t{474}\t{475}\t{476}\t{477}\t{478}\t{479}\t{480}\t{481}\t{482}\t{483}\t{484}\t{485}\t{486}\t{487}\t{488}\t{489}\t{490}\t{491}\t{492}\t{493}\t{494}\t{495}\t{496}\t{497}\t{498}\t{499}\t{500}\t{501}\t{502}\t{503}\t{504}\t{505}\t{506}\t{507}\t{508}\t{509}\t{510}\t{511}\t{512}\t{513}\t{514}\t{515}\t{516}\t{517}\t{518}\t{519}\t{520}\t{521}\t{522}\t{523}\t{524}\t{525}\t{526}\t{527}\t{528}\t{529}\t{530}\t{531}\t{532}\t{533}\t{534}\t{535}\t{536}\t{537}\t{538}\t{539}\t{540}\t{541}\t{542}\t{543}\t{544}\t{545}\t{546}\t{547}\t{548}\t{549}\t{550}\t{551}\t{552}\t{553}\t{554}\t{555}\t{556}\t{557}\t{558}\t{559}\t{560}\t{561}\t{562}\t{563}\t{564}\t{565}\t{566}\t{567}\t{568}\t{569}\t{570}\t{571}\t{572}\t{573}\t{574}\t{575}\t{576}\t{577}\t{578}\t{579}\t{580}\t{581}\t{582}\t{583}\t{584}\t{585}\t{586}\t{587}\t{588}\t{589}\t{590}\t{591}\t{592}\t{593}\t{594}\t{595}\t{596}\t{597}\t{598}\t{599}\t{600}\t{601}\t{602}\t{603}\t{604}\t{605}\t{606}\t{607}\t{608}\t{609}\t{610}\t{611}\t{612}\t{613}\t{614}\t{615}\t{616}\t{617}\t{618}\t{619}\t{620}\t{621}\t{622}\t{623}\t{624}\t{625}\t{626}\t{627}\t{628}\t{629}\t{630}\t{631}\t{632}\t{633}\t{634}\t{635}\t{636}\t{637}\t{638}\t{639}\t{640}\t{641}\t{642}\t{643}\t{644}\t{645}\t{646}\t{647}\t{648}\t{649}\t{650}\t{651}\t{652}\t{653}\t{654}\t{655}\t{656}\t{657}\t{658}\t{659}\t{660}\t{661}\t{662}\t{663}\t{664}\t{665}\t{666}\t{667}\t{668}\t{669}\t{670}\t{671}\t{672}\t{673}\t{674}\t{675}\t{676}\t{677}\t{678}\t{679}\t{680}\t{681}\t{682}\t{683}\t{684}\t{685}\t{686}\t{687}\t{688}\t{689}\t{690}\t{691}\t{692}\t{693}\t{694}\t{695}\t{69
```

```
16 # 1 10 -PRUN- False False x PRUN PKP
17 # 'll 11 will False False 'xx VERB
18 # 15 False True SPACE _SP
19 # be 17 be False False xx VERB VB
20 # in 20 in False False xx ADP IN
21 # Madrid 23 madrid False False Xxxxx
22 # . 29 . True False . PUNCT .
23
```

The spaCy toolbox

Let's now explore what are the models bundled up inside spaCy.

Sentence detection

Here's how to achieve one of the most common NLP tasks with spaCy:

```
1 doc = nlp("These are apples. These are oranges.")
2
3 for sent in doc.sents:
4     print(sent)
5
6 # These are apples.
7 # These are oranges.
8
```

Part Of Speech Tagging

We've already seen how this works but let's have another look:

```
1 doc = nlp("Next week I'll be in Madrid.")
2 print([(token.text, token.tag_) for token in doc])
3
4 # [('Next', 'JJ'), ('week', 'NN'), ('I', 'PRP')]
5
```

Named Entity Recognition

Doing NER with spaCy is super easy and the pretrained model performs pretty well:

```
1 doc = nlp("Next week I'll be in Madrid.")
2 for ent in doc.ents:
3     print(ent.text, ent.label_)
4
5 # Next week DATE
```

You can also view the IOB style tagging of the sentence like this:

```
1 from nltk.chunk import conlltags2tree
2
3
4 doc = nlp("Next week I'll be in Madrid.")
5 iob_tagged = [
6     (
7         token.text,
8         token.tag_,
9         "{0}-{1}".format(token.ent_iob_, token.
10     ) for token in doc
11 ]
12
13 print(iob_tagged)
14
15 # In case you like the nltk.Tree format
16 print(conlltags2tree(iob_tagged))
17
18 # [('Next', 'JJ', 'B-DATE'), ('week', 'NN', 'I-
19
20 # (S
21 #   (DATE Next/JJ week/NN)
22 #   I/PRP
23 #   'll/MD
24 #   be/VB
25 #   in/IN
26 #   (GPE Madrid/NNP)
27 #   ./.)
28
```

The spaCy NER also has a healthy variety of entities. You can view the full list here:

Entity Types

```
1 doc = nlp("I just bought 2 shares at 9 a.m. bec
2 for ent in doc.ents:
3     print(ent.text, ent.label_)
4
5 # 2 CARDINAL
6 # 9 a.m. TIME
7 # 30% PERCENT
8 # just 2 days DATE
9 # WSJ ORG
10
```

Let's use `displaCy` to view a beautiful visualization of the Named Entity annotated sentence:

```
1 from spacy import displacy
```

I just bought 2 **CARDINAL** shares at 9 a.m. **TIME** because the stock went up 30% **PERCENT** in just 2 days **DATE** according to the **WSJ ORG**

Chunking

spaCy automatically detects noun-phrases as well:

```
1 doc = nlp("Wall Street Journal just published a  
2 for chunk in doc.noun_chunks:  
3     print(chunk.text, chunk.label_, chunk.root.  
4  
5 # Wall Street Journal NP Journal  
6 # an interesting piece NP piece  
7 # crypto currencies NP currencies  
8
```

Notice how the chunker also computes the root of the phrase, the main word of the phrase.

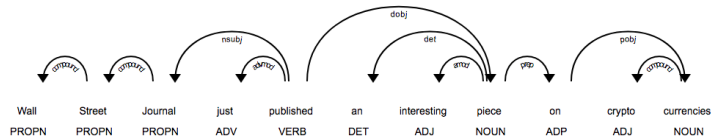
Dependency Parsing

This is what makes spaCy really stand out. Let's see the dependency parser in action:

```
1 doc = nlp('Wall Street Journal just published a  
2  
3 for token in doc:  
4     print("{0}/{1} <-- {2}-- {3}/{4}".format(  
5         token.text, token.tag_, token.dep_, tok  
6  
7 # Wall/NNP <--compound-- Street/NNP  
8 # Street/NNP <--compound-- Journal/NNP  
9 # Journal/NNP <--nsubj-- published/VBD  
10 # just/RB <--advmod-- published/VBD  
11 # published/VBD <--ROOT-- published/VBD  
12 # an/DT <--det-- piece/NN  
13 # interesting/JJ <--amod-- piece/NN  
14 # piece/NN <--dobj-- published/VBD  
15 # on/IN <--prep-- piece/NN  
16 # crypto/JJ <--compound-- currencies/NNS  
17 # currencies/NNS <--pobj-- on/IN  
18
```

depending on how you display some information.

```
1 from spacy import displacy
2
3 doc = nlp('Wall Street Journal just published a
4 displacy.render(doc, style='dep', jupyter=True,
5
```



Word Vectors

spaCy comes shipped with a Word Vector model as well. We'll need to download a larger model for that:

```
1 python -m spacy download en_core_web_lg
2
```

The vectors are attached to spaCy objects: **Token**, **Lexeme** (a sort of unattached token, part of the vocabulary), **Span** and **Doc**. The multi-token objects average its constituent vectors.

Explaining word vectors(aka word embeddings) are not the purpose of this tutorial. Here are a few properties word vectors have:

- If two words are similar, they appear in similar contexts
- Word vectors are computed taking into account the context (surrounding words)

vectors

- Using vectors we can derive relationships between words

Let's see how we can access the embedding of a word in spaCy:

```
1 nlp = spacy.load('en_core_web_lg')
2 print(nlp.vocab['banana'].vector)
3
```

There's a really famous example of word embedding math: "man" - "woman" + "queen" = "king". It sounds pretty crazy to be true, so let's test that out:

```
1 from scipy import spatial
2
3 cosine_similarity = lambda x, y: 1 - spatial.di
4
5 man = nlp.vocab['man'].vector
6 woman = nlp.vocab['woman'].vector
7 queen = nlp.vocab['queen'].vector
8 king = nlp.vocab['king'].vector
9
10 # We now need to find the closest vector in the
11 maybe_king = man - woman + queen
12 computed_similarities = []
13
14 for word in nlp.vocab:
15     # Ignore words without vectors
16     if not word.has_vector:
17         continue
18
19     similarity = cosine_similarity(maybe_king,
20     computed_similarities.append((word, similar
21
22 computed_similarities = sorted(computed_similar
23 print([w[0].text for w in computed_similarities
24
25 # ['Queen', 'QUEEN', 'queen', 'King', 'KING', '
26
```

Surprisingly, the closest word vector in the vocabulary for "man" - "woman" + "queen"

...you're getting everything you need...

Computing Similarity

Based on the word embeddings, spaCy offers a similarity interface for all of its building blocks: **Token**, **Span**, **Doc** and **Lexeme**. Here's how to use that similarity interface:

```
1 banana = nlp.vocab['banana']
2 dog = nlp.vocab['dog']
3 fruit = nlp.vocab['fruit']
4 animal = nlp.vocab['animal']
5
6 print(dog.similarity(animal), dog.similarity(fr
7 print(banana.similarity(fruit), banana.similari
8
```

Let's now use this technique on entire texts:

```
1 target = nlp("Cats are beautiful animals.")
2
3 doc1 = nlp("Dogs are awesome.")
4 doc2 = nlp("Some gorgeous creatures are felines")
5 doc3 = nlp("Dolphins are swimming mammals.")
6
7 print(target.similarity(doc1)) # 0.89017652184
8 print(target.similarity(doc2)) # 0.91158284491
9 print(target.similarity(doc3)) # 0.78229567528
10
```

Extending spaCy

The entire spaCy architecture is built upon three building blocks: **Document** (the big encompassing container), **Token** (most of the time, a word) and **Span** (set of consecutive **Tokens**). The extensions you create can add extra functionality to anyone of these components. There are some examples out there for what you can do. Let's create an extension ourselves.

```
1 import spacy
2 from spacy.tokens import Doc
3 from nltk.sentiment.vader import SentimentIntensityAnalyzer
4
5 sentiment_analyzer = SentimentIntensityAnalyzer
6 def polarity_scores(doc):
7     return sentiment_analyzer.polarity_scores(doc._.text)
8
9 Doc.set_extension('polarity_scores', getter=polarity_scores)
10
11 nlp = spacy.load('en')
12 doc = nlp("Really Whaaat event apple nice! it!")
13 print(doc._.polarity_scores)
14 # {'neg': 0.0, 'neu': 0.596, 'pos': 0.404, 'compound': 0.404}
15
```

One can easily create extensions for every component type. Such extensions only have access to the context of that component.

What happens if you need the tokenized text along with the Part-Of-Speech tags.

Let's now build a custom `pipeline`.

Pipelines are another important abstraction of spaCy. The `nlp` object goes through a list of pipelines and runs them on the document. For example the `tagger` is ran first, then the `parser` and `ner` pipelines are applied on the already POS annotated document. Here's how the `nlp` default pipeline structure looks like:

```
1 nlp = spacy.load('en')
2 print(nlp.pipeline)
3
4 # [('tagger', <spacy.pipeline.Tagger object at 0x100000000>),
5 #  ('parser', <spacy.pipeline.DependencyParser object at 0x100000000>),
6 #  ('ner', <spacy.pipeline.EntityRecognizer object at 0x100000000>)]
7
```

Creating a custom pipeline

Let's build a custom pipeline that needs to be applied after the `tagger` pipeline is ran.

```
1 from nltk.corpus import wordnet as wn
2 from spacy.tokens import Token
3
4
5 def penn_to_wn(tag):
6     if tag.startswith('N'):
7         return 'n'
8
9     if tag.startswith('V'):
10        return 'v'
11
12    if tag.startswith('J'):
13        return 'a'
14
15    if tag.startswith('R'):
16        return 'r'
17
18    return None
19
20
21 class WordnetPipeline(object):
22     def __init__(self, nlp):
23         Token.set_extension('synset', default=None)
24
25     def __call__(self, doc):
26         for token in doc:
27             wn_tag = penn_to_wn(token.tag_)
28             if wn_tag is None:
29                 continue
30
31             ss = wn.synsets(token.text, wn_tag)
32             token._.set('synset', ss)
33
34         return doc
35
36
37 nlp = spacy.load('en')
38 wn_pipeline = WordnetPipeline(nlp)
39 nlp.add_pipe(wn_pipeline, name='wn_synsets')
40 doc = nlp("Paris is the awesome capital of France.")
41
42 for token in doc:
43     print(token.text, "-", token._.synset)
44
45 # Paris - Synset('paris.n.01')
46 # is - Synset('be.v.01')
47 # the - None
48 # awesome - Synset('amazing.s.02')
49 # capital - Synset('capital.n.01')
50 # of - None
51 # France - Synset('france.n.01')
52 # . - None
53
```

Let's see how the pipeline structure looks like:

```
1 print(nlp.pipeline)
2
```

```
b7 | # ('wn_synsets', <__main__.wordnetPipeline obj
```

Conclusions

spaCy is a modern, reliable NLP framework that quickly became the standard for doing NLP with Python. Its main advantages are: speed, accuracy, extensibility. It also comes shipped with useful assets like word embeddings. It can act as the central part of your production NLP pipeline.

Related

[Roundup of Python NLP Libraries](#)
July 17, 2018
In "Resources"

[What is Natural Language Processing?](#)
November 30, 2016
In "general"

[Complete Guide to Word Embeddings](#)
April 6, 2018
In "Deep Learning"

| DATE | AUTHOR | CATEGORY |
|----------------|---------|----------|
| March 28, 2018 | bogdani | spaCy |

| TAG | COMMENTS |
|--|-------------|
| embeddings framework IOB-tagging part-of-speech similarity spaCy tagging text analysis word vectors wordnet | 18 Comments |

TO CONTENTS



Liza

03/29/2018 - 9:30 AM

Hi Bogdan,

I was playing around with spaCy and I have a problem with Matcher, it doesn't work... and I couldn't find any support...do I have to load it separately? Or could it be because of wrong setting?

and another question (sorry I kinda have a lot...) I was reading your blog about training NER (an amazing tutorial thought), but do you have similar articles, if I want to build and train my own one (or maybe some tips:))?

Thanks a lot!

REPLY



bogdani

03/29/2018 - 12:52 PM

Hi Liza,

Can you maybe paste the Matcher code? I don't have much experience with it either.

...ing. What do you need more for
training your NER? Does the 2nd
episode shed some light?
<https://nlpforhackers.io/training-ner-large-dataset/>

Thanks,
Bogdan.

REPLY



ahmed mohammed

04/07/2018 - 3:00 AM

this is clear precise and informative thank
you Bogdani

REPLY



Ruben

04/29/2018 - 9:49 AM

Hi Bogdani, great tutorial!
I've replicated it on a jupyter notebook, I
just found a typo in line 22 in "Creating a
custome pipeline", it should be

```
1 | token.set_extension('synset', default=None)
```

instead of

```
1 | Token.set_extension('synset', default=None)
```

comment on this article

Otherwise thanks for the article, it was very useful!

REPLY



bogdani

04/29/2018 - 9:49 PM

Hi Ruben,

I was missing this import `from spacy.tokens import Token`. Fixed, thanks for reporting this 😊

REPLY

Banu Selin Tosun

04/21/2019 - 5:41 PM

Hi,

I also needed another adjustment in the above mentioned snippet.

1st I needed to download wordnet as

```
nltk.download('wordnet')
```

2nd at the above mentioned line, the token extensions should be set to force=True via

```
def __init__(self, nlp):  
    Token.set_extension('synset',  
                        default=None, force=True)
```

Everything works after that.
I am not sure if this is system

...articles are all up to date.

Maybe a newer issue in the new version.

Thank you,
Selin

REPLY

transhuaman

06/07/2018 - 10:24 AM

simple witing with high information. great!

best npl training

REPLY

dat

06/21/2018 - 9:24 AM

great article, thank for sharing.

REPLY

bogdani

06/21/2018 - 9:35 AM

Thank you, Sir 😊

REPLY

Archana

10/29/2018 - 10:10 AM

REPLY

bogdani

10/29/2018 - 11:32 AM

Unfortunately, no

REPLY

Arpit

12/18/2018 - 8:23 PM

Hello Bogdani,

its one of the best tutorial for SpaCy specially adding the pipeline part.

I was having a doubt relating to the .similarity function in SpaCy. Suppose when comparing two sentences does it consider the POS tagging and parsing pipelines?? I doubt it happens because it uses GloVe vector representations which does not support the POS tagging etc.

Do you have any ideas how can i use the parts of speech and dependency parsing features (like provided by spacy) in word vectors models ??

Thanks

REPLY

bogdani

12/20/2018 - 12:03 PM

cosine similarity between the mean vectors of the 2 sentences. Because embeddings don't take into account the POS, the similarity function won't take that into account either. I have a tutorial on a different similarity function here:
<https://nlpforhackers.io/wordnet-sentence-similarity/>

REPLY

Arpit

12/19/2018 - 8:03 PM

Can you explain how the similarity function of SpaCy works ? Does it use the tagging and dependency parsing information into account when finding the similarity score ?

REPLY

bogdani

12/20/2018 - 12:17 PM

Think it's the cosine similarity between the mean of the word vectors

REPLY

Kalyan

12/20/2018 - 3:13 AM

...spacy or dependency parser using
displacy, I got this error

```
1 | <span class="ansi-red-intense-fg ansi-bold">Typ
```

REPLY

bogdani

12/20/2018 - 12:17 PM

What Python version are you using?

REPLY

cristo

06/03/2019 - 3:21 PM

spacy is only a shit tool...if You need faster
tokenizer go with nltk...spacy
become uselessly slow
and cumbersome..please; ps. package like
tm in R, text2vec, and others works 100
time better than this crock...

REPLY

YOUR EMAIL ADDRESS WILL NOT BE PUBLISHED. REQUIRED FIELDS ARE
MARKED *

Visual

Text

Paragraph ▼

B

I

☰

☰
1
3

“

☰

☰



NAME *

EMAIL *

WEBSITE

Post Comment

☐ NOTIFY ME OF FOLLOW-UP COMMENTS BY EMAIL.

☐ NOTIFY ME OF NEW POSTS BY EMAIL.