## ⌄ IST664 : Introduction to spaCy

Edited by Jeff Stanton and Preeti Jagadev

SpaCy is an open source natural language processing library written by Matthew Honnibal and Ines Montani. Most of spaCy is written natively in Python. Unlike NLTK, which was designed for teaching and research, spaCy was created from the start to support production applications - real world activities that require natural language processing. SpaCy uses a "pipeline" metaphor such that input documents and data go through a variety of typical processing stages where each stage feeds into the next one. Examples of these stages include tokenization, part of speech tagging, named entity recognition, and transformation into word vectors.

Try searching for "spaCy" on Kaggle.com. At this writing there were more than 4600 projects that used spaCy. Part of the appeal is that spaCy makes it easy to get started with a project. SpaCy contains support for dozens of different languages and its integration with word- and sentence-embedding approaches provides access to the advantages of pre-trained deep learning models.

In this lab you will get a more comprehensive view of the architecture and capabilities of spaCy.

Sections of this lab:

- Basics: Getting Started
- Lemmatization
- Token Extracting / Removing / Transforming
- Sentence Segmentation
- Part of Speech Tagging
- Named Entity Recognition
- Dependency Parsing
- Word Vectors
- Sentence Similarity
- Customizing pipeline components

## ⌄ Basics: Getting Started

```
1 # Every spaCy project begins with importing the package and
2 # instantiating a processing object that is initialized with a particular
3 # language model. In this case we will start with a small English pipeline
4 # trained from text harvested from the web.
5 import spacy
6 nlp = spacy.load("en_core_web_sm")
7
8 # That's equivalent to:
9 # import en_core_web_sm
10 # nlp = en_core_web_sm.load()
11
12 type(nlp)
```

```
spacy.lang.en.English
def __call__(text: Union[str, Doc], *, disable: Iterable[str]=SimpleFrozenList(), component_cfg:
Optional[Dict[str, Dict[str, Any]]]=None) -> Doc

/usr/local/lib/python3.11/dist-packages/spacy/lang/en/__init__.py
A text-processing pipeline. Usually you'll load this once per process,
and pass the instance around your application.

Defaults (class): Settings, data and factory methods for creating the `nlp`
    object and processing pipeline.
```

```
1 # There are lots of things this object can do. Let's use
2 # dir to get a list of them:
3
4 [m for m in dir(nlp) if m[0] != "_"]
```

```
['Defaults',
 'add_pipe',
 'analyze_pipes',
 'batch_size',
 'begin_training',
 'component',
```

```
'component_names',
'components',
'config',
'create_optimizer',
'create_pipe',
'create_pipe_from_source',
'default_config',
'default_error_handler',
'disable_pipe',
'disable_pipes',
'disabled',
'enable_pipe',
'evaluate',
'factories',
'factory',
'factory_names',
'from_bytes',
'from_config',
'from_disk',
'get_factory_meta',
'get_factory_name',
'get_pipe',
'get_pipe_config',
'get_pipe_meta',
'has_factory',
'has_pipe',
'initialize',
'lang',
'make_doc',
'max_length',
'memory_zone',
'meta',
'path',
'pipe',
'pipe_factories',
'pipe_labels',
'pipe_names',
'pipeline',
'rehearse',
'remove_pipe',
'rename_pipe',
'replace_listeners',
'replace_pipe',
'resume_training',
'select_pipes',
'set_error_handler',
'set_factory_meta',
'to_bytes',
'to_disk',
'tokenizer',
'update',
'use params',
```

Note that all of the methods shown above pertain to pipelines - a modular sequence of processing steps. In general these follow a standard order:

- Tokenizer
- Part of speech tagger
- Dependency parser (organizes each sentence into its constituent parts
- Named entity recognizer
- Lemmatizer
- Additional elements (including document classification)

You know enough of the essential foundations of NLP to know why these pipeline elements appear in this order. For example, you could not apply part of speech tags to words without first tokenizing the raw text.

```
1 # At the most basic level, and at the beginning of most
2 # NLP pipelines, we tokenize a document:
3 doc = nlp("Hello World!") # This is the most basic way to use the instance
4 type(doc), len(doc) # What is the result?
5 #Output is 3 - Maybe for Hello, World !?
```

⤷  (spacy.tokens.doc.Doc, 3)

```
1 # A spaCy "tokens-doc" behaves like a list, such that
2 # we can use a list comprehension to access the individual
3 # tokens in the document:
4 [token.text for token in doc]
```

```
['Hello', 'World', '!']
```

```python
1 # And because it behaves like a list, we can also use
2 # slicing to get access to the individual tokens.
3 first_token = doc[0] # Slice the first token
4 print(type(first_token)) # What is its type?
5 print(first_token.text) # Show the text of the token
```

```
<class 'spacy.tokens.token.Token'>
Hello
```

```python
1 # In spaCy terminology, a span is any contiguous set of tokens.
2 # Spans are often used to break up a document into sentences. Here
3 # we are just using slicing to create a span with the first two
4 # of our three tokens.
5 span = doc[0:2]
6 [token.text for token in span]
```

```
['Hello', 'World']
```

```python
1 # For this first exercise, tokenize a longer text excerpted from Wikipedia.
2 # Use slicing to show the first five tokens:
3
4 longtext = """A neural network is either a biological neural network or an
5 artificial neural network for solving artificial intelligence (AI) problems.
6 The connections of the biological neuron are modeled as weights. A positive
7 weight reflects an excitatory connection, while negative values mean
8 inhibitory connections."""
9
10 # 5.1: Tokenize longtext
11 doc = nlp(longtext)
12
13 # 5.2: Display the *texts* of tokens in a span consisting of the first 5 tokens
14 print([token.text for token in doc[0:5]])
15
16 # 5.2a: (Challenge) Use Python slicing notation to show the *last* 5 tokens
17 print([token.text for token in doc[-5:]])
18
```

```
['A', 'neural', 'network', 'is', 'either']
['mean', '\n', 'inhibitory', 'connections', '.']
```

```python
1 # SpaCy uses the language model to make better tokenization decisions. Let's
2 # compare spaCy tokenization with the primitive use of split(). Remember that
3 # split() defaults to splitting on spaces.
4 headline = "Rare Bird's Detection Highlights Promise of 'Environmental DNA'"
5
6 splitspacy = nlp(headline) # Use spaCy tokenization
7 splitspace = headline.split() # Use simple splitting on spaces
8
9 print("SpaCy tokens:")
10 print([t.text for t in splitspacy])
11 print(len(splitspacy), "tokens.")
12
13 print("\nSimple splitting:")
14 print([s for s in splitspace])
15 print(len(splitspace), "tokens.")
```

```
SpaCy tokens:
['Rare', 'Bird', ''s', 'Detection', 'Highlights', 'Promise', 'of', '‘', 'Environmental', 'DNA', '’']
11 tokens.

Simple splitting:
['Rare', 'Bird's', 'Detection', 'Highlights', 'Promise', 'of', '‘Environmental', 'DNA’']
8 tokens.
```

```python
1 # Why do you think it might be helpful to tokenize the possessive "Bird's" into
2 # two tokens? Add a comment that explains your reasoning. Then find or write
3
4 #Maybe to show Punctuation or the lack thereof - Wanting to Breakout the Root Words?
5
6 # a new sentence that contains a hyphenated noun phrase. How does spaCy treat that?
7 headline = "I really love Mac-and-Cheese from Chick-Fil-A"
8
```

```
 9 # 5.2b: Use spaCy to tokenize a sentence that contains a hyphenated phrase.
10 splitspacy = nlp(headline)
11
12 print("SpaCy tokens:")
13 print([t.text for t in splitspacy])
14 print(len(splitspacy), "tokens.")
15
```

```
SpaCy tokens:
['I', 'really', 'love', 'Mac', '-', 'and', '-', 'Cheese', 'from', 'Chick', '-', 'Fil', '-', 'A']
14 tokens.
```

## Lemmatization

Lemmatization is the process of reducing inflected forms, sometimes derivationally related forms of a word to a common base form. This reduced form or root word is called a lemma. Lemmas have an advantage over simple stemming: Lemmas are always dictionary words. Lemmatizing can be a valuable data reduction technique because it aggregates various inflective forms of a word down to a single root.

```
1 # Demonstrate spaCy lemmatization with a verb form
2 text = "I am, you are, and he is." # All the verbs are variations on the verb "to be"
3
4 # Note that the underscore following the attribute name in
5 # the expression token.lemma_ provides the human readable form of the attribute.
6 [token.lemma_ for token in nlp(text)]
```

```
['I', 'be', ',', 'you', 'be', ',', 'and', 'he', 'be', '.']
```

```
 1 # Look at the non-text form of the lemma:
 2
 3 # 5.3: use token.lemma instead of token.lemma_
 4 text = "I am, you are, and he is." # All the verbs are variations on the verb "to be"
 5 [token.lemma for token in nlp(text)]
 6
 7 # Write a comment describing what you see. These values are
 8 # ID numbers for spaCy's "StringStore." More information here:
 9 # https://spacy.io/usage/spacy-101#vocab
10 #Looks like the token.lemma somehow converted the Words to Assigned Numbers? IE you could input these numbers and convert them back to t
```

```
[4690420944186131903,
 10382539506755952630,
 2593208677638477497,
 7624161793554793053,
 10382539506755952630,
 2593208677638477497,
 2283656566040971221,
 1655312771067108281,
 10382539506755952630,
 12646065887601541794]
```

```
1 # Here's another example
2 text = "Look! It looks like he looked."
3 doc = nlp(text)
4 for token in doc:
5     print("token:{} -> lemma:{}".format(token.text,token.lemma_ ))
```

```
token:Look -> lemma:look
token:! -> lemma:!
token:It -> lemma:it
token:looks -> lemma:look
token:like -> lemma:like
token:he -> lemma:he
token:looked -> lemma:look
token:. -> lemma:.
```

```
1 # Add your own example, this time using different forms of a noun
2 text = "The Geese and Mice like Mangos"
3 # 5.4: Lemmatize two or more inflective forms of a noun. What about
4 # irregular inflections (such as the plural of mouse)?
5 doc = nlp(text)
6 for token in doc:
7     print("token:{} -> lemma:{}".format(token.text,token.lemma_ ))
```

```
token:The -> lemma:the
token:Geese -> lemma:geese
token:and -> lemma:and
token:Mice -> lemma:mouse
token:like -> lemma:like
token:Mangos -> lemma:mangos
```

## Token Extracting / Removing / Transforming

Here's an overview of all of the bound methods and attributes that a token has. When creating an NLP pipeline, it is helpful not to have to write our own code to find out these things.

| Attribute Name | Type | Description |
|---|---|---|
| lemma | int | Base form of the token, with no inflectional suffixes. |
| lemma_ | unicode | Base form of the token, with no inflectional suffixes. |
| norm | int | The token's norm, i.e. a normalized form of the token text. Usually set in the language's tokenizer exceptions or norm exceptions. |
| norm_ | unicode | The token's norm, i.e. a normalized form of the token text. Usually set in the language's tokenizer exceptions or norm exceptions. |
| lower | int | Lowercase form of the token. |
| lower_ | unicode | Lowercase form of the token text. Equivalent to Token.text.lower(). |
| shape | int | Transform of the tokens's string, to show orthographic features. For example, "Xxxx" or "dd". |
| shape_ | unicode | Transform of the tokens's string, to show orthographic features. For example, "Xxxx" or "dd". |
| prefix | int | Hash value of a length-N substring from the start of the token. Defaults to N=1. |
| prefix_ | unicode | A length-N substring from the start of the token. Defaults to N=1. |
| suffix | int | Hash value of a length-N substring from the end of the token. Defaults to N=3. |
| suffix_ | unicode | Length-N substring from the end of the token. Defaults to N=3. |
| is_alpha | bool | Does the token consist of alphabetic characters? Equivalent to token.text.isalpha(). |
| is_ascii | bool | Does the token consist of ASCII characters? Equivalent to all(ord(c) < 128 for c in token.text). |
| is_digit | bool | Does the token consist of digits? Equivalent to token.text.isdigit(). |
| is_lower | bool | Is the token in lowercase? Equivalent to token.text.islower(). |
| is_upper | bool | Is the token in uppercase? Equivalent to token.text.isupper(). |
| is_title | bool | Is the token in titlecase? Equivalent to token.text.istitle(). |
| is_punct | bool | Is the token punctuation? |
| is_left_punct | bool | Is the token a left punctuation mark, e.g. (? |
| is_right_punct | bool | Is the token a right punctuation mark, e.g. )? |
| is_space | bool | Does the token consist of whitespace characters? Equivalent to token.text.isspace(). |
| is_bracket | bool | Is the token a bracket? |
| is_quote | bool | Is the token a quotation mark? |
| is_currency V2.0.8 | bool | Is the token a currency symbol? |
| like_url | bool | Does the token resemble a URL? |
| like_num | bool | Does the token represent a number? e.g. "10.9", "10", "ten", etc. |
| like_email | bool | Does the token resemble an email address? |

## Extracting

The list of attributes that spaCy makes available on the token object provide a variety of type tests. The is_ attributes allow testing for alpahnumeric, uppercase, title case, left punctuation mark, right punctuation mark, any punctuation mark, a bracket, a quote mark or a currency symbol. These are all very helpful in navigating within a string of tokens: Later we will show a search capability that allows us to include these in pattern matching.

There are also three "like" attributes that show if a token looks like a web address, a numeric string, or an email address.

Let's run some tests on a long and complex sentence from Wikipedia:

```
1 text='''An information retrieval technique using latent semantic structure was
2 patented in 1988 (US Patent 4,839,853, now expired) by Scott Deerwester,
3 Susan Dumais, George Furnas, Richard Harshman, Thomas Landauer, Karen Lochbaum
4 and Lynn Streeter. In the context of its application to information retrieval,
5 it is sometimes called latent semantic indexing (LSI).'''
```

```
1 my_list=[] # Initialize a blank list
2 doc = nlp(text) # Tokenize the text string
3 for token in doc: # Check each token
4     if token.is_punct: # Run the bound method
```

```
5          my_list.append(token) # Append to the list
6
7 for item in my_list: # Review each item in the list
8      print(item) # Print the item
```

```
(
,
)
,
,
,
,
.
,
(
)
.
```

```
1 # Now do something similar but use a list comprehension
2 [tok for tok in doc if tok.is_left_punct]
```

```
[(, (]
```

```
1 # Add a line of code to display right punctuation
2 [tok for tok in doc if tok.is_right_punct]
3 # 5.5: Use the bound method to detect and print right puncutation
4 bound_method = [tok for tok in doc if tok.is_right_punct]
5 print(bound_method)
```

```
[), )]
```

```
1 # Add a line of code to detect tokens that seem like numbers
2 num_tok = [tok for tok in doc if tok.like_num]
3 # 5.6: Use the bound method to detect and display numbers
4 bound_method = [tok for tok in doc if tok.like_num]
5 print(bound_method)
```

```
[1988, 4,839,853]
```

```
1 # For diagnostic purposes, it may be useful to examine these attributes
2 # all together. Here's a code fragment that sets up a pandas df of token
3 # attributes:
4
5 import pandas as pd # Use a pandas DF
6
7 # These will be out column names
8 cols = ("text", "lemma_","is_punct", "is_stop", "is_alpha","is_space","lower_")
9
10 rows = [] # A blank list to hold the rows
11
12 for t in doc: # Iterate through the tokens - will work for any length document
13     # build the next row
14     row = [t.text, t.lemma_,  t.is_punct,  t.is_stop,  t.is_alpha,  t.is_space,  t.lower_]
15     rows.append(row) # Append the row to the existing rows
16
17 # Create the pandas data frame from the column names and the list of rows
18 attri_pdf = pd.DataFrame(rows, columns=cols)
19
20 attri_pdf # Gives a preview, but may not show all rows
```

| | text | lemma_ | is_punct | is_stop | is_alpha | is_space | lower_ |
|---|---|---|---|---|---|---|---|
| 0 | An | an | False | True | True | False | an |
| 1 | information | information | False | False | True | False | information |
| 2 | retrieval | retrieval | False | False | True | False | retrieval |
| 3 | technique | technique | False | False | True | False | technique |
| 4 | using | use | False | True | True | False | using |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 62 | indexing | indexing | False | False | True | False | indexing |
| 63 | ( | ( | True | False | False | False | ( |
| 64 | LSI | LSI | False | False | True | False | lsi |
| 65 | ) | ) | True | False | False | False | ) |
| 66 | . | . | True | False | False | False | . |

67 rows × 7 columns

Next steps: ( Generate code with `attri_pdf` )   ( ◉ View recommended plots )   ( New interactive sheet )

In previous weeks we have considered stop words and why in some cases it makes sense to remove them from the token stream. Let's examine spaCy's stop word list.

```
1 import spacy
2 spacy_stopwords = spacy.lang.en.stop_words.STOP_WORDS
3 len(spacy_stopwords)
```

326

```
1 list(spacy_stopwords)[:8]
```

```
['without',
 'being',
 'and',
 'myself',
 'wherever',
 'nevertheless',
 'take',
 'amount']
```

It is good to know what is on a stop list because sometimes these lists contain words that we do not want to discard. Because spaCy tags each token with an attribute showing whether that token is a stop word, but does not discard the stop words, we have the opportunity to do diagnostics on the results.

Let's process another piece of text from Wikipedia to focus on the stop words.

```
1 text = """In natural language processing, the Latent Dirichlet Allocation (LDA)
2 is a generative statistical model that allows sets of observations to be
3 explained by unobserved groups that explain why some parts of the data are
4 similar. For example, if observations are words collected into documents,
5 it posits that each document is a mixture of a small number of topics and that
6 each word's presence is attributable to one of the document's topics. LDA is
7 an example of a topic model and belongs to the machine learning field and in
8 a wider sense to the artificial intelligence field."""
9
10 doc = nlp(text)
11 type(doc), len(doc)
```

(spacy.tokens.doc.Doc, 113)

```
1 # Display the tokens that are stop words:
2 print([token for token in doc if token.is_stop])
```

[In, the, is, a, that, of, to, be, by, that, why, some, of, the, are, For, if, are, into, it, that, each, is, a, of, a, of, and, that, ε

Note how spaCy has correctly tagged stop words even when they begin with a capital letter.

```
1 # Make a list of the tokens that are not stop-words
2 no_stops = [token for token in doc if not token.is_stop]
3 type(no_stops), len(no_stops)
```

    (list, 64)

```
1 # Use slicing to view the first few non-stop words.
2 no_stops[0:12]
```

    [natural,
     language,
     processing,
     ,,
     Latent,
     Dirichlet,
     Allocation,
     (,
     LDA,
     ),
     ,
     generative]

In some applications we may have uses for the punctuation tokens, but it is also good to know how to remove them. Conveniently, spaCy has also tagged every token with an indicator of whether it is punctuation.

```
1 # Also remove punctuation tokens
2 no_stops_or_punct = [token for token in no_stops if not token.is_punct]
3 type(no_stops_or_punct[0]), len(no_stops_or_punct)
```

    (spacy.tokens.token.Token, 56)

```
1 # Use slicing to view the first few non-stop, non-punct words.
2 no_stops_or_punct[0:10]
```

    [natural,
     language,
     processing,
     Latent,
     Dirichlet,
     Allocation,
     LDA,
     ,
     generative,
     statistical]

```
 1 # Hmm, no_stops_or_punct[7] seems to be blank. Write a few lines of
 2 # code to find out what it is/was. Hint: Try creating a pandas
 3 # data frame of token attributes like the one we made in an earlier
 4 # section of this lab.
 5
 6 # 5.7: Investigate the mysterious token.
 7 no_stops_or_punct[7] #Confirmed Blank
 8 text = """In natural language processing, the Latent Dirichlet Allocation (LDA)
 9 is a generative statistical model that allows sets of observations to be
10 explained by unobserved groups that explain why some parts of the data are
11 similar. For example, if observations are words collected into documents,
12 it posits that each document is a mixture of a small number of topics and that
13 each word's presence is attributable to one of the document's topics. LDA is
14 an example of a topic model and belongs to the machine learning field and in
15 a wider sense to the artificial intelligence field."""
16
17 doc = nlp(text)
18 no_stops = [token for token in doc if not token.is_stop]
19 no_stops_or_punct = [token for token in no_stops if not token.is_punct]
20
21 df = pd.DataFrame(
22     [(i, t.text, t.orth_, repr(t.text), t.is_space, t.is_alpha, t.pos_) for i, t in enumerate(no_stops_or_punct)],
23     columns=["Index", "Text", "Orth_ID", "Repr", "Is_Space", "Is_Alpha", "POS"]
24 )
25
26 print(df)
27 #Used Source - https://www.phind.com/
28
29 #I guess 7 is just a Space?
```

```
     Index         Text        Orth_ID          Repr  Is_Space  Is_Alpha  \
0        0      natural       natural     'natural'     False      True
1        1     language      language    'language'     False      True
2        2   processing    processing  'processing'     False      True
3        3       Latent        Latent      'Latent'     False      True
4        4    Dirichlet     Dirichlet   'Dirichlet'     False      True
5        5   Allocation    Allocation  'Allocation'     False      True
6        6          LDA           LDA         'LDA'     False      True
7        7          \n            \n          '\n'      True     False
8        8   generative    generative  'generative'     False      True
9        9  statistical   statistical 'statistical'     False      True
10      10        model         model       'model'     False      True
11      11       allows        allows      'allows'     False      True
12      12         sets          sets        'sets'     False      True
13      13 observations  observations 'observations'    False      True
14      14          \n            \n          '\n'      True     False
15      15    explained     explained   'explained'     False      True
16      16   unobserved    unobserved  'unobserved'     False      True
17      17       groups        groups      'groups'     False      True
18      18      explain       explain     'explain'     False      True
19      19        parts         parts       'parts'     False      True
20      20         data          data        'data'     False      True
21      21          \n            \n          '\n'      True     False
22      22      similar       similar     'similar'     False      True
23      23      example       example     'example'     False      True
24      24 observations  observations 'observations'    False      True
25      25        words         words       'words'     False      True
26      26    collected     collected   'collected'     False      True
27      27    documents     documents   'documents'     False      True
28      28          \n            \n          '\n'      True     False
29      29       posits        posits      'posits'     False      True
30      30     document      document    'document'     False      True
31      31      mixture       mixture     'mixture'     False      True
32      32        small         small       'small'     False      True
33      33       number        number      'number'     False      True
34      34       topics        topics      'topics'     False      True
35      35          \n            \n          '\n'      True     False
36      36         word          word        'word'     False      True
37      37     presence      presence    'presence'     False      True
38      38  attributable  attributable 'attributable'   False      True
39      39     document      document    'document'     False      True
40      40       topics        topics      'topics'     False      True
41      41          LDA           LDA         'LDA'     False      True
42      42          \n            \n          '\n'      True     False
43      43      example       example     'example'     False      True
44      44        topic         topic       'topic'     False      True
45      45        model         model       'model'     False      True
46      46      belongs       belongs     'belongs'     False      True
47      47      machine       machine     'machine'     False      True
48      48     learning      learning    'learning'     False      True
49      49        field         field       'field'     False      True
50      50          \n            \n          '\n'      True     False
51      51        wider         wider       'wider'     False      True
52      52        sense         sense       'sense'     False      True
53      53   artificial    artificial  'artificial'     False      True
54      54 intelligence  intelligence 'intelligence'    False      True
55      55        field         field       'field'     False      True
```

## Discuss and Collaborate

Task 5.7, just above, may require some ingenuity. Check in with someone else in the lab to get ideas on how to tackle this. When you have uncovered the result, discuss the implications with your partner.

```
1 # Another attribute on a token contains the lowercase version
2 # of the token. Why does this attribute end with an underscore?
3 lowercased = [ token.lower_ for token in no_stops_or_punct]
4 lowercased[0:9]
```

```
['natural',
 'language',
 'processing',
 'latent',
 'dirichlet',
 'allocation',
 'lda',
 '\n',
 'generative']
```

When we called nlp() on the text object and created the tokens, spaCy also automatically guessed at the lemma for each token and stuck that in as an attribute. Knowing what you know about lemmatization, what does this imply about other processing that spaCy may have done to this text?

```
1 # Make an additional list of lemma tokens
2 lemma_list = [token.lemma_ for token in no_stops_or_punct]
3 type(lemma_list[0]), len(set(lemma_list))
```

    (str, 39)

```
1 # The output above suggests that some lemmas appear in the token
2 # list more than one time. Use Counter from the collections package
3 # to count instances of lemmas.
4 from collections import Counter
5
6 # 5.8: Instantiate a counter object with Counter(lemma_list). Assign this
7 # to a new variable such as wc_lemmas
8 wc_lemmas = Counter(lemma_list)
9 print(wc_lemmas)
10
```

    Counter({'\n': 7, 'document': 3, 'topic': 3, 'LDA': 2, 'model': 2, 'observation': 2, 'explain': 2, 'example': 2, 'word': 2, 'field': 2,

```
1 # 5.9: Display the frequency counts of the five most common lemmas. Hint:
2 # a Counter has a bound method called most_common() that takes one
3 # argument called "n"
4 most_common = wc_lemmas.most_common(5)
5 print(most_common)
```

    [('\n', 7), ('document', 3), ('topic', 3), ('LDA', 2), ('model', 2)]

```
1 # 5.10: Grab a new long string from Wikipedia or another source
2 #       and remove stop words and punctuation, then lemmatize and
3 #       count the frequencies of the top five lemmas. Write new code
4 #       below, based on what you did above.
5
6 text = "Software packages on Android, which use the APK format, are generally distributed through a proprietary application store; non-Go
7 doc = nlp(text)
8 filter_token = [token for token in doc if not token.is_stop and not token.is_punct]
9 lemma_list = [token.lemma_ for token in filter_token]
10 wc_lemmas = Counter(lemma_list)
11 most_common = wc_lemmas.most_common(5)
12 print(most_common)
13 #Used Source - https://www.phind.com/
14
```

    [('operating', 3), ('system', 3), ('Android', 2), ('software', 1), ('package', 1)]

## ∨ Sentence Segmentation

The spaCy doc object contains an element called "sents" that records the beginning and ending position (counting by tokens) of each sentence in the document. Finding sentence boundaries requires a substantial amount of algorithmic complexity, because the ending punctuation in strings such as U.S. or etc. may or may not indicate a sentence boundary. There are four strategies for sentence boundary detection in spaCy: dependency parser (default), statistical segmenter, rule-based segmenter, or custom function. Let's tokenize a fragment of Wikipedia text using the default (dependency parser) and then examine the resulting sentences.

```
1 text = """"Sentence boundary disambiguation (SBD), also known as sentence breaking, sentence boundary detection, and sentence segmentatic
2 text[-26:] # Show the end of the string
3
```

    ' computer code, and slang.'

```
1 doc = nlp(text)
2 type(doc.sents)
```

    generator

In Python, a generator is a special kind of iterator function that creates the requested elements on demand and "on the fly." This is a helpful approach when working with large sets of data elements that it would be challenging to represent in memory all at once. It is as easy to create a generator as it is to create a list comprehension:

```
1 (n ** 2 for n in range(50)) # First try it without assigning it to an object
```

    <generator object <genexpr> at 0x78c3a4fffac0>

```
1 sq_gen = (n ** 2 for n in range(20)) # This time, save the generator and check its type
2 type(sq_gen)
```

    generator

```
1 # Try running this cell twice. What happens the second time?
2 # After the second run, what happens if you run the previous cell again?
3 for num in sq_gen:
4     if num < 100:
5         print(num)
6 #First Time Run
7 #0
8 #1
9 #4
10 #9
11 #16
12 #25
13 #36
14 #49
15 #64
16 #81
17
18 #Second Time it Returns Nothing or Blank on OutPut
```

Looking back a couple of code blocks, doc.sents is a generator object, which means we can iterate through it's elements to find what we need.

```
1 for sent in doc.sents:
2     print("start_pos={}, end_pos={}, text:{}".format(sent.start, sent.end, sent.text))
```

    start_pos=0, end_pos=36, text:Sentence boundary disambiguation (SBD), also known as sentence breaking, sentence boundary detection, and
    start_pos=36, end_pos=67, text:Natural language processing tools often require their input to be divided into sentences; however, senten
    start_pos=67, end_pos=103, text:In written English, a period may indicate the end of a sentence, or may denote an abbreviation, a decima
    start_pos=103, end_pos=139, text:About 47% of the periods in the Wall Street Journal corpus denote abbreviations.[1] Question marks and

```
1 # Use the doc.sent generator object to iterate through the sentences. Create
2 # a pandas data frame containing the starting and ending token numbers of each
3 # sentence. Then use the data frame to fetch the first token from each sentence.
4 import pandas as pd # Use a pandas DF
5
6 # These will be our column names
7 cols = ("start", "end")
8 rows = [] # A blank list to hold the rows
9
10 # 5.11: Iterate through the sentences, appending a row of start and end
11 #        positions for each sentence.
12 text = """Sentence boundary disambiguation (SBD), also known as sentence breaking, sentence boundary detection, and sentence segmentatic
13 doc = nlp(text)
14 for sent in doc.sents:
15     rows.append([sent.start, sent.end])
16
17 sent_pdf = pd.DataFrame(rows, columns=cols)
18
19 sent_pdf
20
21
```

|       | start | end |
|-------|-------|-----|
| **0** | 0     | 36  |
| **1** | 36    | 67  |
| **2** | 67    | 103 |
| **3** | 103   | 139 |

Next steps:  ( Generate code with `sent_pdf` )   ( ⊙ View recommended plots )   ( New interactive sheet )

```
1 # 5.12: Create a pandas data frame from the rows and column names
2 df = pd.DataFrame(rows, columns=cols)
```

```
   start  end
0      0   36
1     36   67
2     67  103
3    103  139
```

```
1 # 5.13: Display the pandas data frame
2 print(df)
```

```
   start  end
0      0   36
1     36   67
2     67  103
3    103  139
```

```
1 # 5.14: Iterate through the data frame and print the first token
2 #        from each sentence.
3
4 for index, row in df.iterrows():
5     print(doc[row['start']])
```

```
Sentence
Natural
In
About
```

## ⌄ Part of Speech Tagging

POS tagging is a critical processing step in most pipelines, and much more difficult in some languages (e.g., Chinese) than in others. The standard spaCy pipeline always includes a "tagger" that assigns POS tags based on a statistical analysis of a training corpus.

```
1 doc = nlp("I was reading an article about Berkeley Avenue in Reading, which was closed due to a police investigation.")
2 type(doc), len(doc)
```

```
(spacy.tokens.doc.Doc, 20)
```

```
1 from tabulate import tabulate # To make a neat table
2
3 tabdata = [ (token, token.tag_, token.pos_, spacy.explain(token.tag_)) for token in doc]
4
5 print(tabulate(tabdata,  headers=["Token", "Token Tag", "POS", "Explanation"]))
6
```

```
Token         Token Tag    POS    Explanation
------------  -----------  -----  ----------------------------------------
I             PRP          PRON   pronoun, personal
was           VBD          AUX    verb, past tense
reading       VBG          VERB   verb, gerund or present participle
an            DT           DET    determiner
article       NN           NOUN   noun, singular or mass
about         IN           ADP    conjunction, subordinating or preposition
Berkeley      NNP          PROPN  noun, proper singular
Avenue        NNP          PROPN  noun, proper singular
in            IN           ADP    conjunction, subordinating or preposition
Reading       NNP          PROPN  noun, proper singular
,             ,            PUNCT  punctuation mark, comma
which         WDT          PRON   wh-determiner
was           VBD          AUX    verb, past tense
closed        VBN          VERB   verb, past participle
```

```
due            IN       ADP    conjunction, subordinating or preposition
to             IN       ADP    conjunction, subordinating or preposition
a              DT       DET    determiner
police         NN       NOUN   noun, singular or mass
investigation  NN       NOUN   noun, singular or mass
.              .        PUNCT  punctuation mark, sentence closer
```

```
1 # Make a list of tokens for all of the proper nouns
2 propnlist = [token for token in doc if token.pos_ == "PROPN"]
3
4 [ (token, token.is_ascii, token.is_title) for token in propnlist]
```

```
[(Berkeley, True, True), (Avenue, True, True), (Reading, True, True)]
```

```
1 # Write another sentence that includes a place name. Tokenize it,
2 # display the POS tags, and excerpt the proper noun(s).
3
4 sent = "I like to go to the Movies on Elm Street Called Cinemark"
5 # 5.15: Create a new text object for tokenizing.
6 doc = nlp(sent)
7
8 for token in doc:
9     print(f"{token.text}: {token.pos_}")
10
11 proper_nouns = [token.text for token in doc if token.pos_ == "PROPN"]
12 print(proper_nouns)
13
14 #Used Source - https://www.phind.com/
15
```

```
I: PRON
like: VERB
to: PART
go: VERB
to: ADP
the: DET
Movies: PROPN
on: ADP
Elm: PROPN
Street: PROPN
Called: VERB
Cinemark: PROPN
['Movies', 'Elm', 'Street', 'Cinemark']
```

```
1 # 5.16: Tokenize the text object.
2 doc = nlp(sent)
3 for token in doc:
4     print(token.text)
```

```
I
like
to
go
to
the
Movies
on
Elm
Street
Called
Cinemark
```

```
1 # 5.17: Display the POS tags for all tokens.
2 pos_tags = [token.pos_ for token in doc]
3 print(pos_tags)
```

```
['PRON', 'VERB', 'PART', 'VERB', 'ADP', 'DET', 'PROPN', 'ADP', 'PROPN', 'PROPN', 'VERB', 'PROPN']
```

```
1 # 5.18: Extract the proper nouns and display them.
2 prop_noun = [token.text for token in doc if token.pos_ == "PROPN"]
3 print(prop_noun)
```

```
['Movies', 'Elm', 'Street', 'Cinemark']
```

## Dependency Parsing

```
1 # Let's take a closer look at the dependency structure:
2 from tabulate import tabulate # To make a neat table
3
4 tabdata = [ (token.text, token.tag_, token.dep_, token.head.text, token.head.tag_) for token in doc]
5
6 print(tabulate(tabdata,  headers=["Token", "Token POS", "Dependency", "Head Token", "Head POS"]))
7
```

```
Token      Token POS    Dependency    Head Token    Head POS
--------   -----------  ------------  ------------  ----------
I          PRP          nsubj         like          VBP
like       VBP          ROOT          like          VBP
to         TO           aux           go            VB
go         VB           xcomp         like          VBP
to         IN           prep          go            VB
the        DT           det           Movies        NNPS
Movies     NNPS         pobj          to            IN
on         IN           prep          Movies        NNPS
Elm        NNP          compound      Street        NNP
Street     NNP          pobj          on            IN
Called     VBN          acl           Movies        NNPS
Cinemark   NNP          dobj          like          VBP
```

Take a close look at the output just above. For each token in the sentence, the text of the token is shown along with its part of speech. Then the dependency relation is shown. For example, the first token, "I", is the noun/subject of the sentence and is therefore dependent on the main verb, "reading", which is the gerund form of the verb to read.

Take the time to examine each row of the output and make sure you understand the dependency relation that is being documented. And remember that spaCy's ability to diagram the relations in this way works because of the language model we originally loaded: "en_core_web_sm". Also important: The default sentence segmentation that we examined in a previous block works because spaCy's dependency parser accounts for all of the elements in a sentence, and therefore "knows" when the period character is closing a sentence.

```
1 # Grab another sentence from the web, but this time, cut off the sentence
2 # before the end so that some key grammatical element is missing. Do paste
3 # a period on the end, though, just to see if you can confuse spaCy.
4
5 # 5.19: Cut and paste part of a sentence from the web into a text variable.
6 sent = " After having built a prototype internally known as the Fadden demo predominantly by purchasing licensing agreements for most of
```

```
1 # 5.20: Tokenize the sentence.
2 doc = nlp(sent)
3 doc
```

```
 After having built a prototype internally known as the Fadden demo predominantly by purchasing licensing agreements for most of the
 software components built around a custom JavaScript front-end, the company failed to convince investors, and so in April 2004 they
 pivoted to building an Operating System for Phones at the.
```

```
1 # 5.21: Generate a table showing the dependency relations in the sentence.
2 from tabulate import tabulate
3
4 tabdata = [ (token.text, token.tag_, token.dep_, token.head.text, token.head.tag_) for token in doc]
5
6 print(tabulate(tabdata,  headers=["Token", "Token POS", "Dependency", "Head Token", "Head POS"]))
```

```
Token          Token POS    Dependency    Head Token    Head POS
-------------  -----------  ------------  ------------  ----------
               _SP          dep           After         IN
After          IN           prep          failed        VBD
having         VBG          aux           built         VBN
built          VBN          pcomp         After         IN
a              DT           det           prototype     NN
prototype      NN           dobj          built         VBN
internally     RB           advmod        known         VBN
known          VBN          acl           prototype     NN
as             IN           prep          known         VBN
the            DT           det           demo          NN
Fadden         NNP          compound      demo          NN
demo           NN           pobj          as            IN
predominantly  RB           advmod        by            IN
by             IN           prep          built         VBN
purchasing     VBG          pcomp         by            IN
```

```
licensing       NN         compound      agreements      NNS
agreements      NNS        dobj          purchasing      VBG
for             IN         prep          purchasing      VBG
most            JJS        pobj          for             IN
of              IN         prep          most            JJS
the             DT         det           components      NNS
software        NN         compound      components      NNS
components      NNS        pobj          of              IN
built           VBN        acl           components      NNS
around          IN         prep          built           VBN
a               DT         det           custom          NN
custom          NN         pobj          around          IN
JavaScript      NNP        nmod          end             NN
front           JJ         compound      end             NN
-               HYPH       punct         end             NN
end             NN         appos         custom          NN
,               ,          punct         failed          VBD
the             DT         det           company         NN
company         NN         nsubj         failed          VBD
failed          VBD        ROOT          failed          VBD
to              TO         aux           convince        VB
convince        VB         xcomp         failed          VBD
investors       NNS        dobj          convince        VB
,               ,          punct         failed          VBD
and             CC         cc            failed          VBD
so              RB         advmod        pivoted         VBD
in              IN         prep          pivoted         VBD
April           NNP        pobj          in              IN
2004            CD         nummod        April           NNP
they            PRP        nsubj         pivoted         VBD
pivoted         VBD        conj          failed          VBD
to              IN         prep          pivoted         VBD
building        VBG        pcomp         to              IN
an              DT         det           System          NNP
Operating       VBG        compound      System          NNP
System          NNP        dobj          building        VBG
for             IN         prep          System          NNP
Phones          NNPS       pobj          for             IN
at              IN         prep          building        VBG
the             DT         pobj          at              IN
.               .          punct         pivoted         VBD
```
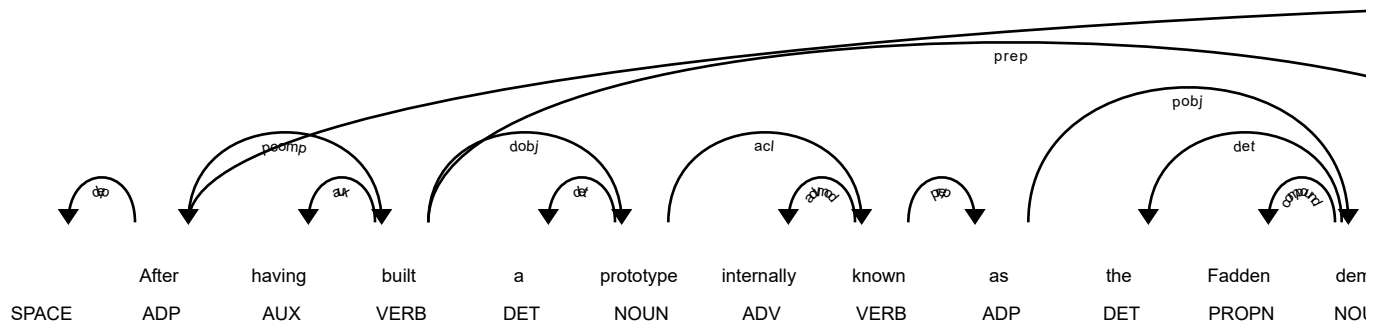
```
1 # 5.22: Add a comment to document any mistakes that spaCy made.
2 #Documentation of potetntial Mistakes
3 #-              HYPH       punct         end             NN
4 #,              ,          punct         failed          VBD
5 #.              .          punct         pivoted         VBD
```

```
1 # As we saw in a previous lab, there is a graphical display module for
2 # spaCy that supports drawing a figure of the dependency relations.
3 from spacy import displacy
4 displacy.render(doc, style='dep', jupyter=True, options={'distance': 90})
```



Now that you can see the dependency relations as a graph, do you notice any problems with the parsing? Did spaCy make any mistakes in connecting the various elements of the sentence?

## ∨ Named Entity Recognition

Whenever spaCy finds a token that looks like a proper noun, it tags it as a predicted named entity. "Predicted," because each spaCy language model has a trained classifier that makes predictions of whether or not a token might be a named entity and also what type of entity it is (e.g., an organization, a country, or something else).

The IOB tagging method is a a straightforward way of notating the status of tokens. Tokens can have one of the following four statuses:

| TAG | ID | DESCRIPTION |
|-----|-----|-------------|
| I | 1 | Token is inside an entity. |
| O | 2 | Token is outside an entity. |
| B | 3 | Token begins an entity. |
|  | 0 | No entity tag is set (missing value). |

```
1 text='''We're bringing the celebration of Syracuse University's 150 years of impact to Chicago'''
2 doc = nlp(text)
3 type(doc.ents)
```

⇥ tuple

```
1 # So the list of entities is a tuple, which means we should be able to slice it.
2 doc.ents[0]
```

⇥ Syracuse University

```
1 # We can iterate through all of the entities in the document.
2 for ent in doc.ents:
3     print("{}, [{},{}), {}".format(ent.text, ent.start_char, ent.end_char, ent.label_))
4
5     # For each entity, we can also access each entity as a span and
6     # iterate through its tokens
7     for token in ent.as_doc():
8         print("    {} {} {}".format(token, token.ent_iob_, token.ent_type_))
```

```
⇥ Syracuse University, [34,53), ORG
        Syracuse B ORG
        University I ORG
    150 years, [56,65), DATE
        150 B DATE
        years I DATE
    Chicago, [79,86), GPE
        Chicago B GPE
```

Notice that for each of the three entities above, there's a numeric expression that shows the start and end of the character span where the entities can be found in the original string.

In the output shown above, ORG and DATE are pretty clear, but what does GPE stand for? Do a web search on "spacy GPE" to find out.

```
1 # The displacy module can also provide a graphical view of the named entities:
2 from spacy import displacy
3 displacy.render(doc, style="ent", jupyter=True)
```

⇥ We're bringing the celebration of  Syracuse University **ORG**  's  150 years **DATE**  of impact to  Chicago **GPE**

```
1 # Now write or find another sentence that contains some named entities.
2 text = "Duke University has had 25 Good Years of Basketball to the Durham Community"
3 # 5.23: Assign a new sentence containing some entities into a text variable.
4
```

```
1 # 5.24: Tokenize the sentence.
2 doc = nlp(text)
3 doc
```

⇥ Duke University has had 25 Good Years of Basketball to the Durham Community

```
1 # 5.24: Generate a displacy graphic with the named entities.
2 from spacy import displacy
3 displacy.render(doc, style="ent", jupyter=True)
```

Duke University **ORG** has had 25 Good Years **MONEY** of Basketball to the Durham Community **ORG**

## Word Vectors

Word vectors are numeric representations of words in multidimensional space. Notable implementations of word vectors include GloVe, a technique that obtains vectors from a large word co-occurence matrix developed from a corpus, and Word2vec, an predictive neural network training process that does repetitive, incremental training over strings of tokens in a corpus. Word vectors have the interesting property that semantically similar words appear close to each other in multidimensional space. Given any two word vectors, we can calculate the similarity (or distance) between them.

SpaCy language models come in several sizes. Small models do not contain word vectors - a space saving strategy. So if we want to use word vectors, we need to load the medium or large version of a language model as shown in the next cell. Note that because of the size of a large language model file, the next code cell may take about three minutes to complete. That's a lot of time and therefore a lot of data. Imagine how this might impact the operation of a deployed system.

```
1 import spacy.cli # Use the command line interface
2 spacy.cli.download("en_core_web_lg") # This imports the large model onto your virtual machines
3 import en_core_web_lg # Now that it is downloaded, we can import it
4 nlp_lg = en_core_web_lg.load() # Create an instance for further use
5 type(nlp_lg)
```

```
✓ Download and installation successful
You can now load the package via spacy.load('en_core_web_lg')
⚠ Restart to reload dependencies
If you are in a Jupyter or Colab notebook, you may need to restart Python in
order to load all the package's dependencies. You can do this by selecting the
'Restart kernel' or 'Restart runtime' option.
```

```
spacy.lang.en.English
def __call__(text: Union[str, Doc], *, disable: Iterable[str]=SimpleFrozenList(), component_cfg:
Optional[Dict[str, Dict[str, Any]]]=None) -> Doc
```

```
/usr/local/lib/python3.11/dist-packages/spacy/lang/en/__init__.py
A text-processing pipeline. Usually you'll load this once per process,
and pass the instance around your application.

Defaults (class): Settings, data and factory methods for creating the `nlp`
    object and processing pipeline.
```

```
1 # Let's get the vectors for three words to examine similarities
2 apple = nlp_lg.vocab["apple"]
3 banana = nlp_lg.vocab["banana"]
4 car = nlp_lg.vocab["car"]
5
6 print(apple.vector) # Take a look at one of them
```

```
[-3.6391e-01  4.3771e-01 -2.0447e-01 -2.2889e-01 -1.4227e-01  2.7396e-01
 -1.1435e-02 -1.8578e-01  3.7361e-01  7.5339e-01 -3.0591e-01  2.3741e-02
 -7.7876e-01 -1.3802e-01  6.6992e-02 -6.4303e-02 -4.0024e-01  1.5309e+00
 -1.3897e-02 -1.5657e-01  2.5366e-01  2.1610e-01 -3.2720e-01  3.4974e-01
 -6.4845e-02 -2.9501e-01 -6.3923e-01 -6.2017e-02  2.4559e-01 -6.9334e-02
 -3.9967e-01  3.0925e-02  4.9033e-01  6.7524e-01  1.9481e-01  5.1488e-01
 -3.1149e-01 -7.9939e-02 -6.2096e-01 -5.3277e-03 -1.1264e-01  8.3528e-02
 -7.6947e-03 -1.0788e-01  1.6628e-01  4.2273e-01 -1.9009e-01 -2.9035e-01
  4.5630e-02  1.0120e-01 -4.0855e-01 -3.5000e-01 -3.6175e-01 -4.1396e-01
  5.9485e-01 -1.1524e+00  3.2424e-02  3.4364e-01 -1.9209e-01  4.3255e-02
  4.9227e-02 -5.4258e-01  9.1275e-01  2.9576e-01  2.3658e-02 -6.8737e-01
 -1.9503e-01 -1.1059e-01 -2.2567e-01  2.4180e-01 -3.1230e-01  4.2700e-01
  8.3952e-02  2.2703e-01  3.0581e-01 -1.7276e-01  3.2536e-01  5.4696e-03
 -3.2745e-01  1.9439e-01  2.2616e-01  7.4742e-02  2.2033e-01 -4.0301e-01
 -3.1594e-01 -2.8910e-02  9.7858e-01  7.1860e-01  1.4995e-01  6.3421e-02
  2.8332e-01 -1.5231e-01  3.9330e-04  1.8076e-01 -4.0199e-01  6.0187e-02
 -2.7543e-01  1.6590e-01 -2.5774e-01  1.6150e-01  3.7247e-01 -3.8273e-01
  2.4012e-01 -4.2617e-02 -6.6785e-01 -9.4437e-01  2.7916e-01  1.0476e-01
  1.3952e+00 -1.4296e-01 -5.5049e-01  5.3982e-02 -7.7524e-01 -2.8255e-01
 -2.3323e-02  2.4801e-01  2.2855e-01 -3.7408e-01  7.6012e-02  2.4031e-01
  1.0746e-01  1.2411e-01 -2.0676e-01 -2.5804e-01 -1.6791e-01  4.3499e-01
  6.1762e-01 -2.9955e-02  1.6196e-01 -2.9001e-01 -3.1159e-01 -8.7262e-01
  4.3167e-01 -1.5071e-01 -4.1420e-01 -5.3730e-01 -1.9910e-01  1.3270e-01
 -1.5018e-01 -4.9335e-01 -2.5127e+00  3.1660e-01  3.6396e-01 -5.9248e-02
  3.1120e-02  4.1071e-02  1.6917e-02  5.8410e-01 -2.0201e-01  7.0238e-02
  8.7547e-01 -2.0114e-01  5.1920e-01  2.6786e-01 -5.5643e-01 -3.1247e-01
 -3.7992e-01  4.2857e-01  4.1780e-01  3.0608e-01 -2.1657e-01  7.2464e-01
  6.1734e-01  5.8085e-02 -6.2708e-01  5.2895e-02 -2.5628e-01 -3.2688e-01
```

```
         -6.1280e-01  6.2609e-01 -1.7965e-01  8.8925e-01  2.1963e-01 -3.4052e-03
         -7.8663e-02  3.4799e-01 -2.6062e-01  8.0410e-03  1.1721e-01 -4.5147e-01
         -1.2178e-01 -5.7030e-01  4.6602e-01  2.5059e-02  5.3986e-02 -7.6693e-01
          1.3173e-01 -2.8776e-02 -4.1915e-01 -2.4415e-01 -4.0295e-01 -4.1520e-01
          3.7643e-02 -1.4843e-01  2.6094e-02  1.5315e-01  3.8310e-01 -5.5825e-01
         -3.3433e-01 -2.7939e-02 -4.3712e-01 -3.1802e-01 -3.1731e-01  9.2891e-02
         -9.9397e-02 -1.8846e-01  5.2270e-02  2.9061e-01  1.0639e+00  9.9584e-02
         -5.6775e-01  2.9446e-01  3.7797e-01 -2.1905e-01 -5.2616e-01 -4.1744e-01
         -6.5951e-01 -4.0820e-01 -6.0945e-01  1.1759e-02 -2.9122e-01 -3.1457e-01
          5.7076e-02  4.1503e-01  3.7345e-01 -4.7119e-02 -7.1996e-02  1.4587e-01
         -3.0763e-01  1.0759e-01 -5.9447e-01 -4.0205e-01  3.0677e-01 -1.9891e-01
         -7.0775e-01 -1.1513e-01  3.0866e-01 -6.9235e-01  2.1219e-01  1.0554e-01
          2.2617e-01 -2.6145e-01 -3.9298e-01 -2.3585e-01  3.0795e-02 -1.0193e-01
          3.2070e-01  3.0505e-01 -5.3470e-01 -7.9272e-02 -1.6817e-01 -2.2115e-01
         -3.5143e-01 -9.2376e-02  1.4686e-01 -1.9859e-01  2.0460e-01  2.0276e-01
          3.6144e-01 -3.5867e-01  4.0095e-01  6.3686e-02 -1.2763e-01 -1.6226e-01
         -3.1763e-01 -5.8732e-01 -5.4009e-01 -4.9035e-01 -4.6035e-01 -1.9794e-01
         -2.5209e-01  2.5706e-01  4.0110e-01  5.2830e-02 -3.2079e-01  3.9563e-01
         -4.4512e-01 -9.1862e-02 -1.9243e-01  1.5397e-01 -2.8923e-01  6.0561e-01
          5.8133e-01  3.2268e-01  6.3892e-02  8.5438e-02  1.4956e-01  3.8134e-01
         -1.1820e-01 -2.3951e-01 -6.7731e-01  2.8090e-01 -5.1770e-01 -4.1098e-01
         -4.1292e-01 -6.7856e-02 -3.3721e-02 -7.2958e-01 -4.7891e-01  7.2956e-01]
```

```python
1 # A spaCy lexeme, with a variety of attribute tests available
2 type(apple), [m for m in dir(apple) if m[0:3] == "is_"]
```

```
(spacy.lexeme.Lexeme,
 ['is_alpha',
  'is_ascii',
  'is_bracket',
  'is_currency',
  'is_digit',
  'is_left_punct',
  'is_lower',
  'is_oov',
  'is_punct',
  'is_quote',
  'is_right_punct',
  'is_space',
  'is_stop',
  'is_title',
  'is_upper'])
```

```python
1 # We can also confirm that a lexeme has a vector representation available
2 banana.has_vector
```

```
True
```

```python
1 # The .similarity() bound method allows us to generate a cosine similarity
2 # score for two vectors.
3 apple.similarity(banana)
```

```
0.5831844806671143
```

```python
1 # The similarity is transitive, and by the calculation spaCy uses, seems
2 # to result in values in the interval 0 to 1. Higher values indicate
3 # closer similarity.
4 banana.similarity(apple)
```
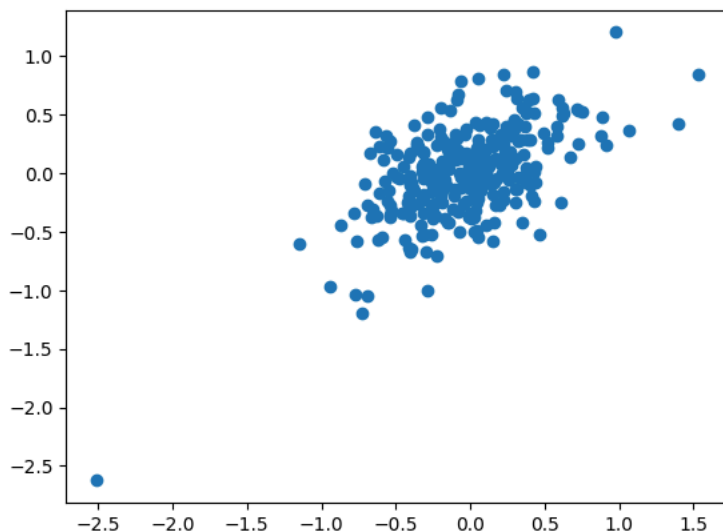
```
0.5831844806671143
```

```python
1 apple.similarity(car)
```
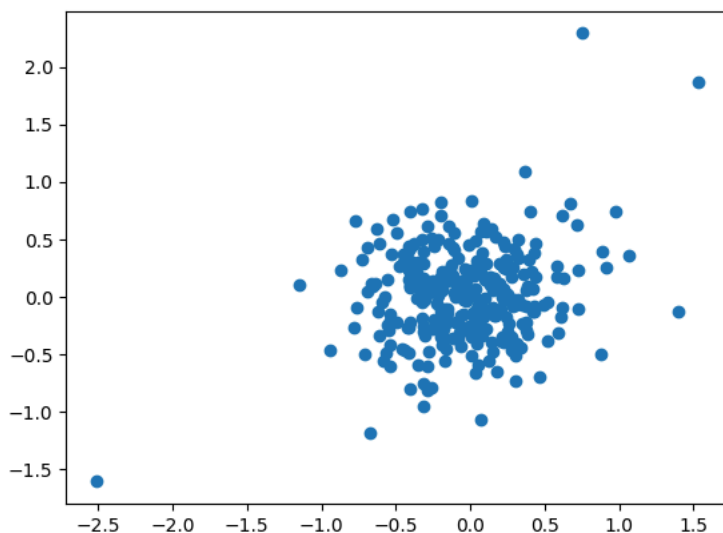
```
0.2174709141254425
```

```python
1 # A scatterplot of two similar vectors should indicate linearity
2 import matplotlib.pyplot as plt # Make a simple plot
3 plt.scatter(apple.vector, banana.vector) # Show the result
```

<matplotlib.collections.PathCollection at 0x78c4981197d0>



```
1 # A scatterplot of two different vectors should be a circular cloud
2 plt.scatter(apple.vector, car.vector)
```

<matplotlib.collections.PathCollection at 0x78c39df13fd0>



```
1 # There are many options for a cosine similarity calculation,
2 # including this one from scipy that we are not using here.
3
4 #from scipy.spatial.distance import cosine
```

```
1 # With some help from numpy, we can calculate our own cosine scores
2 # using a vector dot product, computed with np.dot()
3 import numpy as np
4 def cosine(x,y):
5     return np.dot(x,y) / (np.sqrt(np.dot(x,x)) * np.sqrt(np.dot(y,y)))
```

```
1 # One of the notable aspects of word vectors is the possibility of
2 # doing simple math to create analogies.
3 man = nlp_lg.vocab["man"].vector
4 woman = nlp_lg.vocab["woman"].vector
5 king = nlp_lg.vocab["king"].vector
6 queen = nlp_lg.vocab["queen"].vector
7
8 analogy = king - man + woman # Calculate the analogy with vector math
9 cosine(analogy, queen)
```

np.float32(0.78808445)

```
1 # Let's compare the analogy to another vector
2 cosine(analogy, car.vector) # Note how we needed to extract the vector for car
```

    np.float32(0.13843127)

```
1 # Now you compute an anology and see what kind of results you get
2
3 # 5.25: Obtain the vectors for Paris, France, Berlin, and Germany
4 paris = nlp_lg.vocab["paris"]
5 berlin = nlp_lg.vocab["Berlin"]
6 germany = nlp_lg.vocab["Germany"]
7 france = nlp_lg.vocab["france"]
8
9 print(germany.vector) # Take a look at one of them
```

    [-3.5532e-01  5.9025e-01  1.7082e-01 -2.9313e-03  7.3128e-01 -5.3945e-02
     -4.7846e-02  2.4606e-01 -1.8597e-01  8.0269e-01 -1.2702e+00 -1.2990e-01
     -1.7094e-01 -1.5757e-01  8.6082e-01  4.9421e-03 -2.7626e-01  1.3027e+00
      6.9247e-01  1.3457e-01  7.2636e-01  2.7259e-01  3.3717e-01 -2.1307e-01
     -4.1304e-01 -3.2571e-01 -3.9285e-01  9.8547e-02 -4.2460e-02  7.9824e-01
      2.2305e-01 -5.6391e-02 -1.9127e-03  4.4519e-01  3.8532e-01  2.1436e-01
      3.3141e-01  6.0297e-02  2.6125e-02  1.5253e-03  3.9830e-01 -2.2976e-01
     -1.9558e-01  3.7927e-01 -2.0930e-01 -1.5060e-01  8.1397e-02 -1.0093e+00
     -2.7606e-02  4.2331e-02 -5.1963e-01  1.5031e-01 -3.1286e-01 -8.1624e-01
      7.6164e-02  9.4180e-02  9.1425e-02 -5.0202e-01  3.7440e-01 -3.6910e-01
     -4.2274e-01  2.0528e-02  3.5944e-01 -5.9502e-01 -4.3034e-01  6.5409e-01
     -1.7446e-01  5.4362e-01  1.1969e-01 -2.2330e-01 -1.7697e-01  3.1174e-01
     -6.6029e-02 -5.2093e-01  5.4626e-01 -3.9850e-01 -3.9863e-01  2.2108e-01
      1.3207e-01  6.8904e-01 -3.0168e-01  6.0699e-01  5.7867e-01  3.3999e-01
     -1.9682e-01 -5.4443e-01  1.0600e+00 -4.6676e-01 -5.7628e-02 -1.1523e-01
     -4.0842e-01 -3.0253e-01  1.6459e-01 -2.9100e-01  2.5571e-01 -2.0017e-01
     -1.6753e-01 -1.9953e-01 -1.2312e-01  2.2200e-02  1.8670e-01 -1.5847e-02
      2.0861e-01  5.6133e-01  3.8892e-01 -1.3207e+00  8.6137e-02 -5.7136e-02
     -1.1727e-01 -7.7089e-01  2.8236e-01 -3.2194e-01 -1.4631e-01 -3.8873e-01
      2.4520e-01 -3.8108e-02  9.0389e-01 -7.2336e-02 -8.8831e-01 -2.7025e-01
      9.4136e-02  4.5059e-01  3.1964e-01 -2.5763e-01 -1.7744e-01  5.8418e-01
      4.3581e-01  2.8566e-02  5.5962e-01  2.8713e-01 -1.0694e-03 -5.8154e-01
     -2.2003e-01 -3.0540e-01  3.9144e-01 -6.6091e-01 -2.2653e-01  9.5001e-02
     -1.7457e-01 -3.6791e-02 -1.5287e+00 -2.7696e-01  3.0186e-01 -3.1090e-01
     -1.9702e-01  5.8790e-01  9.8363e-02  5.0177e-01 -1.4607e-02 -1.2669e-01
     -1.8857e-02 -1.9668e-01 -4.9711e-01 -4.1141e-01  1.7557e-01  2.2816e-01
     -1.8279e-01  6.2052e-01  3.8148e-01  8.7166e-01  6.1474e-01 -1.4620e-02
      1.3988e-01 -1.7228e-01  1.6204e-01 -1.5001e-01 -2.4463e-01  2.1547e-01
     -5.0795e-01  2.4293e-01 -5.7434e-01 -2.3912e-01 -5.1797e-01  6.5299e-01
      1.8023e-01  2.3604e-01 -8.8867e-02  4.2646e-01  8.1626e-01  2.2606e-01
      2.3301e-01 -2.3897e-01 -2.3303e-01  2.1172e-01  3.4523e-01 -2.8430e-01
     -6.1216e-03 -1.1932e-01 -3.1599e-01  1.7572e-01 -3.2610e-02  1.5303e-01
      2.4861e-01 -3.8134e-01 -4.8207e-01  6.9897e-01  9.0488e-01 -3.6598e-01
      6.9226e-01 -6.1707e-01  1.8368e-01 -2.8152e-01  6.8885e-01 -1.6421e-01
      1.4285e-01 -2.5463e-01  6.4269e-01 -2.0218e-01  3.7522e-01  3.0399e-02
     -1.2285e-01 -1.2358e-01  4.3093e-01 -4.4876e-01 -2.4861e-01  5.6688e-02
      8.9564e-02  1.2192e-01 -5.2339e-01  6.1985e-01  2.6750e-01 -3.3201e-01
      4.9707e-01 -2.8828e-01 -6.8352e-02  6.5845e-02  3.0167e-01  2.0211e-01
      1.5745e-01  2.7303e-01 -3.4166e-01  4.5752e-01 -6.3572e-01 -4.9968e-01
     -1.9339e-01  1.8362e-01  1.6941e-01 -6.2114e-02 -1.3672e-02 -9.6475e-04
     -3.3685e-01  2.8164e-01  6.8517e-02 -4.5437e-02  3.8221e-01 -7.3306e-02
      1.5074e-01  1.3741e-01  2.4259e-01 -4.2146e-01  7.5622e-01  4.2130e-01
      8.0590e-01 -6.4212e-01  2.9125e-01 -2.9884e-01 -2.5247e-01 -1.8504e-01
     -2.7947e-01 -1.2723e-01  4.0416e-01  1.2838e-01 -4.7400e-01  3.8117e-01
     -4.3376e-01 -1.9804e-01  4.7486e-01 -2.1127e-01 -5.5112e-01  3.2791e-01
     -7.0879e-01 -1.2255e-01 -8.5436e-03  3.0136e-01  4.2968e-01  6.6084e-01
      4.2748e-01 -1.7740e-01  2.3307e-01 -1.9780e-01 -1.3235e-01 -2.5483e-01
      1.1043e-01  8.4103e-02  3.6823e-02  5.7533e-01 -1.3012e-01 -5.6691e-02
     -9.3585e-02  5.4244e-01 -1.4315e-01  3.4506e-01  1.2744e-01  2.1567e-01
      1.5834e-01 -2.1950e-01  5.3764e-01 -2.7061e-01 -2.4100e-01  8.9329e-01]
```

```
1 # 5.26: Compute an analogy using three of these vectors
2
3 paris = nlp_lg.vocab["paris"].vector
4 germany = nlp_lg.vocab["Germany"].vector
5 france = nlp_lg.vocab["france"].vector
6
7 analogy = france - paris + germany
8 cosine(analogy, germany)
```

    np.float32(0.8795712)

```
1 # 5.27: Find the cosine similarity of the analogy with the fourth vector
2
3 paris = nlp_lg.vocab["paris"].vector
```

```
4 germany = nlp_lg.vocab["Germany"].vector
5 france = nlp_lg.vocab["france"].vector
6 berlin = nlp_lg.vocab["berlin"].vector
7
8 analogy = france - paris + germany
9 cosine(analogy, berlin)
```

np.float32(0.52002084)

```
1 # 5.28: Find the cosine similarity of the analogy with banana
2 apple = nlp_lg.vocab["apple"].vector
3 banana = nlp_lg.vocab["banana"].vector
4 car = nlp_lg.vocab["car"].vector
5
6 analogy = apple - banana + car
7 cosine(analogy, car)
```

np.float32(0.7700715)

## ⌄ Sentence Similarity

As you have seen above, word vectors for similar words have similar patterns and though we don't know what each position in a vector signifies, these patterns support simple vector math for analogical reasoning. This suggests the possibility - which research seems to back up - that averaging two or more vectors provides a semantic summary of the vectors. Document and span objects in spaCy have attached vectors that repesents the average of all of the component word vectors.

```
1 doc1 = nlp_lg("The quick brown fox jumps over the lazy dog.")
2 doc2 = nlp_lg("The lazy dog jumps over the quick brown fox.")
3 doc1.similarity(doc1)
```

1.0

```
1 # We can also manually calculate the cosine similarity. It should be transitive
2 # such that cosine(A,B) == cosine(B,A)
3 doc1_vec = doc1.vector
4 doc2_vec = doc2.vector
5 cosine(doc2_vec, doc1_vec), cosine(doc1_vec, doc2_vec)
```

(np.float32(1.0), np.float32(1.0))

```
1 # Let's get the similarity of the first sentence (doc1) with a new sentence.
2 docdiff = nlp_lg("Four score and seven years ago, our fathers brought forth upon this continent a new nation.")
3 cosine(docdiff.vector, doc1_vec)
```

np.float32(0.72331554)

```
1 # The semantics can be quite muddled by the averaging process. These two
2 # sentences express the opposite sentiment
3 doc3 = nlp_lg("I like snow")
4 doc4 = nlp_lg("I hate snow")
5 cosine(doc3.vector, doc4.vector)
```

np.float32(0.94878143)

Think through the implication of the result shown above. By creating an averaged word vector, we have given every word in a sentence the same weight in creating the new summary vector. Another exercise to try would be to take out the word snow and see how that changes the results.

```
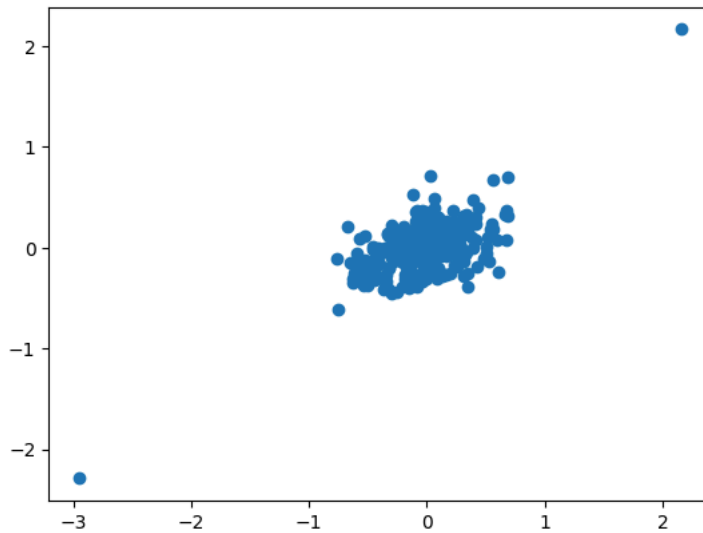1 # 5.29: Recalculate the previous cell after removing the word snow from both sentences.
2 # Feel free to try any other experiments of interest.
3 doc3 = nlp_lg("I like")
4 doc4 = nlp_lg("I hate")
5 cosine(doc3.vector, doc4.vector)
```

np.float32(0.9049743)

```
1 # And here we see part of why this is so.
2 plt.scatter(nlp_lg.vocab["hate"].vector, nlp_lg.vocab["like"].vector)
```

<matplotlib.collections.PathCollection at 0x78c39dca3590>



```
1 # Now you create vector summaries of three sentences and look at the
2 # cosine similarities for each pair.
3
4 # 5.29: Create nlp_lg() objects for three sentences that you write yourself or
5 # copy and paste from the web.
6
7 sent_0 = "I like the Red Balloon on the fence"
8 sent_1 = "I hate that Cat on the Ledge"
9 sent_2 = "I wish I could fly in a Spaceship"
10
11 sent_0_doc = nlp_lg(sent_0)
12 sent_1_doc = nlp_lg(sent_1)
13 sent_2_doc = nlp_lg(sent_2)
14
15
```

```
1 # 5.30: Compute pairwise similarities for their vector summaries
2 sent_0_vec = sent_0_doc.vector
3 sent_1_vec = sent_1_doc.vector
4 sent_2_vec = sent_2_doc.vector
5
6 cosine(sent_0_vec, sent_1_vec)
```

np.float32(0.8767786)

```
1 # 5.31: Add a comment explaining what you see
2 #0.87 Correlation or Similar
```

## ˅ Putting it All Together

Put everything you have learned from the exercises above to use in processing two Wikipedia articles that discuss similar topics (for example, the Super Bowl and the Rose Bowl). Extract the text from two Wikipedia articles. Then process both texts with spaCy. Extract all of the named entities and see which named entities the two articles have in common. Calculate a ratio of matching named entities to all detected entities.

```
1 # 5.32: Extract Named Entities from Two Wikipedia Articles and Find Matches
2 sent_0 = "As the NCAA does not organize or award an official national championship for FBS football (instead merely recognizing the deci
3 sent_1 = "College Football Playoff officials commissioned the trophy for the new playoff system, preferring a new award that was unconne
4
5 doc0 = nlp(sent_0)
6 doc1 = nlp(sent_1)
7
8 ents_0 = set(ent.text for ent in doc0.ents)
9 ents_1 = set(ent.text for ent in doc1.ents)
```

```
10
11 common_ents = ents_0 & ents_1
12 total_unique_ents = ents_0 | ents_1
13
14 ratio = len(common_ents) / len(total_unique_ents)
15 print(common_ents)
16 print(total_unique_ents)
17 print(ratio)
18
```

```
{'FBS'}
{'2014', 'FBS', 'the 2013 season', 'NCAA', 'first', 'the BCS National Championship Game', 'Bowl Championship Series', 'CFP', 'AFCA'}
0.1111111111111111
```

## ⌄ Optional Advanced Topic: Processing pipelines

Throughout this lab, we have been calling a function that we often referred to as nlp(). After loading a spaCy language model such as en_core_web_sm, we instantiate a pipeline object to conduct all of the steps that we will routinely want to accomplish with a document.

The spaCy pipeline can be modified to change the default components or to add new components. Here's a list of the default components from the spaCy documentation:

| NAME | COMPONENT | CREATES | DESCRIPTION |
|------|-----------|---------|-------------|
| tokenizer | Tokenizer | Doc | Segment text into tokens. |
| tagger | Tagger | Doc[i].tag | Assign part-of-speech tags. |
| parser | DependencyParser | Doc[i].head, Doc[i].dep, Doc.sents, Doc.noun_chunks | Assign dependency labels. |
| ner | EntityRecognizer | Doc.ents, Doc[i].ent_iob, Doc[i].ent_type | Detect and label named entities. |
| textcat | TextCategorizer | Doc.cats | Assign document labels. |
| … | custom components | `Doc._.xxx`, `Token._.xxx`, `Span._.xxx` | Assign custom attributes, methods or properties. |

```
1 # We can also examine the pipeline for an instantiated object like this:
2 import spacy
3 nlp = spacy.load("en_core_web_sm")
4
5 nlp.pipeline
```

```
[('tok2vec', <spacy.pipeline.tok2vec.Tok2Vec at 0x78c4d493bb30>),
 ('tagger', <spacy.pipeline.tagger.Tagger at 0x78c3b3bba930>),
 ('parser', <spacy.pipeline.dep_parser.DependencyParser at 0x78c4993699a0>),
 ('attribute_ruler',
  <spacy.pipeline.attributeruler.AttributeRuler at 0x78c3b3c5e290>),
 ('lemmatizer',
  <spacy.lang.en.lemmatizer.EnglishLemmatizer at 0x78c3b3c240d0>),
 ('ner', <spacy.pipeline.ner.EntityRecognizer at 0x78c3b3bd4c80>)]
```

You may notice in the list above that there is no tokenizer. In the spaCy pipeline model, it is assumed that tokenization was accomplished before pipeline processing begins. All pipeline elements receive a doc object, work on it and return a doc object. The tokenizer has a different kind of job becuase it receives a raw character string and returns a list of tokens. Thus, the language object has a different slot where the tokenizer is listed:

```
1 nlp.tokenizer
```

```
<spacy.tokenizer.Tokenizer at 0x78c3fb163e20>
```

The spaCy pipeline was designed to balance simplicity and computational effort. Simplicity is important for getting started quickly with a language processing tasks, so the default pipeline contains all the stuff that most people need to address a realistic task. But if a component is not needed, it can save a lot of compute time to take a task out of the pipeline. Take a look at this example:

```
1 print(spacy.__version__) # Some version dependent stuff below
```

```
3.8.5
```

```
1 # Let's skip the entity recognition and the dependency parsing
2 nlp_simple = spacy.load("en_core_web_sm", exclude=["parser","ner"])
3 # Version 3 also adds facilities for enabling
4 # and disabling pipeline elements on the fly.
5
6 print(nlp_simple.pipeline) # SHow the pipeline
```

[('tok2vec', <spacy.pipeline.tok2vec.Tok2Vec object at 0x78c3b3124170>), ('tagger', <spacy.pipeline.tagger.Tagger object at 0x78c37de22d

```
1 # Groucho Marx is an entity, but this pipeline doesn't detect it, because the
2 # ner is disabled.
3 simple_doc = nlp_simple("Groucho Marx shot an elephant in his underpants.")
4
5 [ent for ent in simple_doc.ents]
```

[]

```
1 # And the pipeline did not do dependency parsing
2 [token.dep_ for token in simple_doc]
```

['', '', '', '', '', '', '', '', '']

```
1 # 5.32: Create a pipeline that contains a dependency parser, but no
2 #        entity recognition. Process a sentence and show that the entity
3 #        recognition is disabled.
4
5 nlp_simple = spacy.load("en_core_web_sm", exclude=["ner"])
6 simple_doc = nlp_simple("Ellie dropped her bone down a rabbit hole into the darkness.")
7 [ent for ent in simple_doc.ents]
8
```

[]

## ⌄ Adding Custom Pipeline Components

A component receives a Doc object and can modify it. By adding a component to the pipeline, you'll get access to the Doc at any point during processing – instead of only being able to modify it afterwards. You can control the position of the new component in the pipeline with the last, first, before, and after arguments.

| ARGUMENT | TYPE | DESCRIPTION |
|---|---|---|
| doc | Doc | The Doc object processed by the previous component. |
| RETURNS | Doc | The Doc object processed by this pipeline component. |

| ARGUMENT | TYPE | DESCRIPTION |
|---|---|---|
| last | bool | If set to True, component is added last in the pipeline (default). |
| first | bool | If set to True, component is added first in the pipeline. |
| before | unicode | String name of component to add the new component before. |
| after | unicode | String name of component to add the new component after. |

```
 1 import spacy
 2 from spacy.tokens import Doc, Span, Token
 3 import json
 4 from spacy.language import Language
 5
 6 # Here's a custom function for removing stopwords, starting with a decorato
 7 @Language.component("remstop_component")
 8 def remove_stopwords(doc):
 9     # A pipeline element would not normally contain this kind of diagnostic
10     # but we just want to show what the method received before processing.
11     print("Before stopwords_removal, this doc is: {}".format(doc))
12     space_list = [t.whitespace_  for t in doc if not t.is_stop]
13     new_doc = Doc(doc.vocab,
14             words=[t.orth_ for t in doc if not t.is_stop],
15             spaces=space_list
16             )
17     return new_doc
18
19 # Instantiate a default pipeline
20 nlp = spacy.load("en_core_web_sm")
```

```
21
22 # Add our stopword finder/remover
23 nlp.add_pipe("remstop_component", name="stopwords_removal", first=True)
24
25 # Show the pipeline
26 print(nlp.pipe_names)  # ['stopwords_removal', 'tagger', 'parser', 'ner']
27
28 # Process a sentence
29 doc = nlp("This is a sentence.")
30
31 # See the result
```

```
['stopwords_removal', 'tok2vec', 'tagger', 'parser', 'attribute_ruler', 'lemmatizer', 'ner']
Before stopwords_removal, this doc is: This is a sentence.
After stopwords_removal, this doc is: sentence.
```

## ⌄ Attribute and Method Extensions

For advanced users, it is possible to add new attributes and methods to spaCy objects. In a complex language processing system, additional attributes and methods could be used to annotate, control, and modify specialized features of a document or other spaCy object.

```
1 # Here we set a new attribute, called "classified" to indicate whether
2 # a document's contents should be kept secret.
3 Doc.set_extension("classified", default=True)
4 assert doc._.classified
5
6 # Note that this block will throw an error if it is run more than once
```

```
1 # Now we can set or retrieve the attribute
2 doc._.classified = False
3
```