

Natural Language Toolkit (Spacy)

Natural Language Processing

Lecture 3



THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

Finding Words, Phrases, Names and Concepts

- Before going further you should install Spacy.
- Check the spacy Webster for installation guide.
<https://spacy.io/usage>

Install spaCy

Operating system: macOS / OSX, **Windows**, Linux

Package manager: **pip**, conda, from source

Hardware: **CPU**, GPU

Configuration: ☐ virtual env ⓘ, ☐ train models ⓘ

Trained pipelines: ☐ Catalan, ☐ Chinese, ☐ Danish, ☐ Dutch, ☒ English, ☐ French, ☐ German, ☐ Greek, ☐ Italian, ☐ Japanese, ☐ Lithuanian, ☐ Macedonian, ☐ Multi-language, ☐ Norwegian Bokmål, ☐ Polish, ☐ Portuguese, ☐ Romanian, ☐ Russian, ☐ Spanish

Select pipeline for: **efficiency ⓘ**, accuracy ⓘ

```
$ pip install -U pip setuptools wheel
$ pip install -U spacy
$ python -m spacy download en_core_web_sm
```

- The nlp object.

```
1 | from spacy.lang.en import English
2 | nlp = English()
```

- At the center of spaCy is the object **containing the processing pipeline**. We usually call this variable "nlp".
- For example, to create an English nlp object, you can import the English language class from **spacy.lang.en** and instantiate it. You can use the **nlp object like a function** to analyze text.
- It contains all the different **components in the pipeline**.
- It also includes **language-specific rules** used for tokenizing the text into words and punctuation.

• The Doc object.

- The Doc lets you access information about the text in a structured way, and no information is lost.

```
1 | doc = nlp("Hello world!")
2 | for token in doc:
3 |     print(token.text)
```

• The Token object.

- To get a token at a specific position, you can index into the doc.

```
1 | token = doc[1]
2 | print(token.text)
```

• The Span object.

- A Span object is a slice of the document consisting of one or more tokens.

```
1 | span = doc[1:3]
2 | print(span.text)
```

● The Lexical object.

```
1 | doc = nlp("It costs $5.")
2 | print("Index:   ", [token.i for token in doc])
3 | print("Text:    ", [token.text for token in doc])
4 | print("is_alpha:", [token.is_alpha for token in doc])
5 | print("is_punct:", [token.is_punct for token in doc])
6 | print("like_num:", [token.like_num for token in doc])
```

- `i` is the index of the token within the parent document.
- `text` returns the token text.
- `is_alpha`, `is_punct` and `like_num` return boolean values indicating whether the token consists of alphabetic characters, whether it's punctuation or whether it resembles a number.
- These attributes are also called lexical attributes

- What are statistical models?
 - It enable spaCy to predict linguistic attributes in contex.
 - Part-of-speech tags.
 - Syntactic dependencies.
 - Named entities.
 - It can be trained.
 - It can ne be fine tuned.

- Some of the most interesting things, you can analyze are **context-specific**.
- For example, whether a word is a verb or whether a span of text is a person name.
- Statistical models enable spaCy to make **predictions in context**.
- This usually includes **part-of speech tags, syntactic dependencies and named entities**.
- Models are **trained** on large datasets of labeled example texts.
- They can be updated with more examples to **fine-tune** their predictions.
- For example, to **perform better** on your specific data.

- spaCy provides a number of pre-trained model packages you can download using the spacy download command.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
```

- For example, the "en_core_web_sm" package is a small English model that supports all core capabilities and is **trained on web text**.
- The spacy.load method loads a **model package by name** and returns an nlp object.
- The package provides the **binary weights** that enable spaCy to make predictions.
- It also includes the **vocabulary, and meta information** to tell spaCy which language class to use and how to configure the processing pipeline.

Predicting Part-of-speech Tags

- Lets load and process a text.

```
1 | import spacy
2 |
3 | nlp = spacy.load("en_core_web_sm")
4 | doc = nlp("She ate the pizza")
5 | for token in doc:
6 |     print(token.text, token.pos_)
```

- In this example, we're using spaCy to predict part-of-speech tags, the word types in context.
- For each token in the doc, we can print the text and the .pos_ attribute, the predicted part-of-speech tag.
- In spaCy, attributes that return strings usually end with an **underscore** – attributes without the underscore return an **integer ID value**.
- Here, the model correctly predicted "ate" as a verb and "pizza" as a noun.

Predicting Syntactic Dependencies

- Lets find more insight about the text.

```
1 | import spacy
2 |
3 | nlp = spacy.load("en_core_web_sm")
4 | doc = nlp("She ate the pizza")
5 | for token in doc:
6 |     print(token.text, token.pos_, token.dep_, token.head.text)
```

- In addition to the part-of-speech tags, we can also predict how the words are related.
- For example, whether a word is the subject of the sentence or an object.
- The **.dep_ attribute** returns the predicted dependency label.
- The **.head** attribute returns the syntactic head token. You can also think of it as the parent token this word is attached to.

Dependency Label Scheme

- To describe syntactic dependencies, spaCy uses a standardized label scheme.
 - The pronoun "She" is a nominal subject attached to the verb – in this case, to "ate".
 - The pronoun "She" is a nominal subject attached to the verb – in this case, to "ate".
 - The pronoun "She" is a nominal subject attached to the verb – in this case, to "ate".

Label	Description	Example
nsubj	nominal subject	She
dobj	direct object	pizza
det	determiner (article)	the

Predicting Named Entities

- Lets find more insight about the text.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
3 | doc = nlp("Apple is looking at buying U.K. startup for $1 billion")
4 |
5 | for ent in doc.ents:
6 |     print(ent.text, ent.label_)
```

- Named entities are "real world objects" that are assigned a name – for example, a person, an organization or a country.
- The **doc.ents** property lets you access the named entities predicted by the model.
- It returns an iterator of Span objects, so we can print the entity text and the entity label using the **.label_ attribute**.
- In this case, the model is correctly predicting "Apple" as an organization, "U.K." as a geopolitical entity and "\$1 billion" as money.

The Spacy Explain Method

- Get quick definitions of the most common tags and labels.

```
1 | import spacy
2 |
3 | print(spacy.explain("GPE"))
4 | print(spacy.explain("NNP"))
5 | print(spacy.explain("dobj"))
```

- To get definitions for the most common tags and labels, you can use the **spacy.explain** helper function.
- For example, "GPE" for geopolitical entity isn't exactly intuitive – but `spacy.explain` can tell you that it refers to countries, cities and states.
- The same works for part-of-speech tags and dependency labels.

- Match on Doc objects, not just strings.
- Match on tokens and token attributes
- It's also more flexible: you can search for texts but also other lexical attributes.
- You can even write rules that use the model's predictions.
- For example, find the word "duck" only if it's a verb, not a noun.

- Match exact token texts.

```
1 | [{"TEXT": "iPhone"}, {"TEXT": "X"}]
```

- Match lexical attributes

```
1 | [{"LOWER": "iphone"}, {"LOWER": "x"}]
```

- Match any token attributes

```
1 | [{"LEMMA": "buy"}, {"POS": "NOUN"}]
```

- Match patterns are lists of dictionaries. Each dictionary describes one token.
- The keys are the names of token attributes, mapped to their expected values.

- In this example, we're looking for **two tokens** with the text "iPhone" and "X".
- Also we can match two tokens whose lowercase forms equal "iphone" and "x".
- We can even write patterns using attributes predicted by the model.
- For example, find the **word "duck" only if it's a verb, not a noun.**
- Here, we're matching a token with the lemma "buy", plus a noun. The lemma is the base form, so this pattern would match phrases like "buying milk" or "bought flowers".

- To use a pattern, we first import the matcher from `spacy.matcher`.

```
1 | import spacy
2 | from spacy.matcher import Matcher
3 |
4 | nlp = spacy.load("en_core_web_sm")
5 |
6 | matcher = Matcher(nlp.vocab)
7 | pattern = [{"TEXT": "iPhone"}, {"TEXT": "X"}]
8 | matcher.add("IPHONE_PATTERN", [pattern])
9 |
10 | doc = nlp("Upcoming iPhone X release date leaked")
11 |
12 | matches = matcher(doc)
```

- The matcher is initialized with the shared vocabulary, `nlp.vocab`.
- The `matcher.add` method lets you add a pattern.
- The first argument is a unique ID to identify which pattern was matched. The second argument is the pattern.

- When you call the matcher on a doc, it returns a list of tuples.

```
1 | import spacy
2 | from spacy.matcher import Matcher
3 |
4 | nlp = spacy.load("en_core_web_sm")
5 | matcher = Matcher(nlp.vocab)
6 | pattern = [{"TEXT": "iPhone"}, {"TEXT": "X"}]
7 | matcher.add("IPHONE_PATTERN", [pattern])
8 |
9 | doc = nlp("Upcoming iPhone X release date leaked")
10 | matches = matcher(doc)
11 |
12 | for match_id, start, end in matches:
13 |     matched_span = doc[start:end]
14 |     print(matched_span.text)
```

- Each tuple consists of three values: the match ID, the start index and the end index of the matched span.
- The `matcher.add` method lets you add a pattern.
- This means we can iterate over the matches and create a `Span` object: a slice of the doc at the start and end index.

Matcher 3 - Lexical Attributes Example

```
1 | import spacy
2 | from spacy.matcher import Matcher
3 | nlp = spacy.load("en_core_web_sm")
4 | matcher = Matcher(nlp.vocab)
5 | pattern = [{"IS_DIGIT": True}, {"LOWER": "fifa"}, {"LOWER": "world"},
6 |           {"LOWER": "cup"}, {"IS_PUNCT": True}]
7 |
8 | matcher.add("FIFA", [pattern])
9 | doc = nlp("2018 FIFA World Cup: France won!")
10 | matches = matcher(doc)
11 |
12 | for match_id, start, end in matches:
13 |     matched_span = doc[start:end]
14 |     print(matched_span.text)
```

- We're looking for five tokens:
- A token consisting of only digits.
- Three case-insensitive tokens for "fifa", "world" and "cup".
- And a token that consists of punctuation.

Matcher 4 - Other Token Attributes

```
1 | import spacy
2 | from spacy.matcher import Matcher
3 |
4 | nlp = spacy.load("en_core_web_sm")
5 | matcher = Matcher(nlp.vocab)
6 |
7 | pattern = [
8 |     {"LEMMA": "love", "POS": "VERB"},
9 |     {"POS": "NOUN"}
10 | ]
11 | matcher.add("Other", [pattern])
12 |
13 | doc = nlp("I loved dogs but now I love cats more.")
14 | matches = matcher(doc)
15 |
16 | for match_id, start, end in matches:
17 |     matched_span = doc[start:end]
18 |     print(matched_span.text)
```

- In this example, we're looking for two tokens:
- A token consisting of only digits.
- A verb with the lemma "love", followed by a noun.

Matcher 5 - Operators and Quantifiers

```
1 | import spacy
2 | from spacy.matcher import Matcher
3 | nlp = spacy.load("en_core_web_sm")
4 | matcher = Matcher(nlp.vocab)
5 |
6 | pattern = [{"LEMMA": "buy"}, {"POS": "DET", "OP": "?"}, {"POS": "NOUN"}]
7 | matcher.add("Other", [pattern])
8 |
9 | doc = nlp("I bought a smartphone. Now I'm buying apps.")
10 | matches = matcher(doc)
11 |
12 | for match_id, start, end in matches:
13 |     matched_span = doc[start:end]
14 |     print(matched_span.text)
```

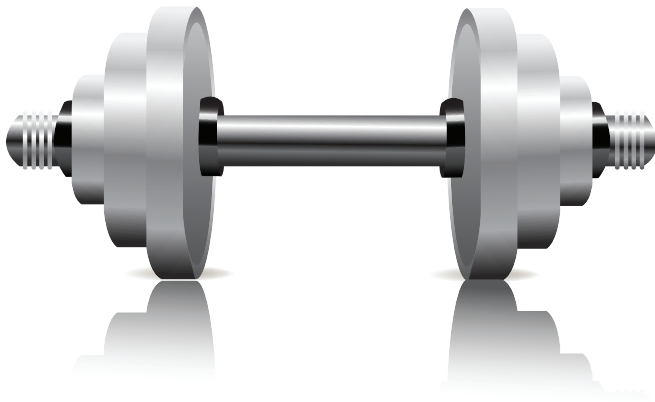
- Operators and quantifiers let you define how often a token should be matched.
- They can be added using the "OP" key.
- Here, the "?" operator makes the determiner token optional, so it will match a token with the lemma "buy", an optional article and a noun.

Using Operators And Quantifiers

Example	Description
{ "OP": "!" }	Negation: match 0 times
{ "OP": "?" }	Optional: match 0 or 1 times
{ "OP": "+" }	Match 1 or more times
{ "OP": "*" }	Match 0 or more times

- An "!" negates the token, so it's matched 0 times.
- A "?" makes the token optional, and matches it 0 or 1 times.
- A "+" matches a token 1 or more times.
- And finally, an "*" matches 0 or more times.
- Operators can make your patterns a lot more powerful, but they also add more complexity

- Exercise 1 to 7
- [Class-Ex-Lecture3.py](#)



Data Structures: Vocab, Lexemes and String Store

Shared Vocabulary and String Store

- spaCy stores all shared data in a vocabulary, the Vocab.
- To save memory, all strings are encoded to hash IDs. If a word occurs more than once, we don't need to save it every time.
- The string store is available as `nlp.vocab.strings`. It's a lookup table that works in both directions.

```
1 | import spacy
2 |
3 | nlp = spacy.load("en_core_web_sm")
4 |
5 | print(nlp.vocab.strings.add("text"))
6 | coffee_hash = nlp.vocab.strings["text"]
7 | coffee_string = nlp.vocab.strings[coffee_hash]
8 | print(coffee_string)
```

- A Lexeme object is an entry in the vocabulary

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
3 |
4 | doc = nlp("I love natural language processing.")
5 | lexeme = nlp.vocab["natural"]
6 | print(lexeme.text, lexeme.orth, lexeme.is_alpha)
```

- You can get a lexeme by looking up a string or a hash ID in the vocab.
- Lexemes expose attributes, just like tokens.
- They hold context-independent information about a word, like the text, or whether the word consists of alphabetic characters.
- Lexemes don't have part-of-speech tags, dependencies or entity labels.

- The Doc is one of the central data structures in spaCy.

```
1 | from spacy.lang.en import English
2 | nlp = English()
3 |
4 | from spacy.tokens import Doc
5 |
6 | words = ["Hello", "world", "!"]
7 | spaces = [True, False, False]
8 |
9 | doc = Doc(nlp.vocab, words=words, spaces=spaces)
```

- It's created automatically when you process a text with the nlp object. But you can also instantiate the class manually.
- Here we're creating a doc from three words. The spaces are a list of boolean values indicating whether the word is followed by a space. Every token includes that information – even the last one!
- The Doc class takes three arguments: the shared vocab, the words and the spaces.

- A Span is a slice of a doc consisting of one or more tokens.

```
1 | from spacy.tokens import Doc, Span
2 | from spacy.lang.en import English
3 | nlp = English()
4 |
5 | words = ["Hello", "world", "!"]
6 | spaces = [True, False, False]
7 |
8 | doc = Doc(nlp.vocab, words=words, spaces=spaces)
9 | span = Span(doc, 0, 2)
10 | span_with_label = Span(doc, 0, 2, label="GREETING")
11 | doc.ents = [span_with_label]; print(doc.ents)
```

- To create a Span manually, we can also import the class from `spacy.tokens`.
- We can then instantiate it with the doc and the span's start and end index, and an optional label argument.
- The `doc.ents` are writable, so we can add entities manually by overwriting it with a list of spans.

Word Vectors And Semantic Similarity

- spaCy can compare two objects and predict similarity
- **Doc.similarity()**, **Span.similarity()** and **Token.similarity()** are the syntaxes.
- It is suggested to use either `en_core_web_md` (medium model) or `en_core_web_lg` (large model).
- The Doc, Token and Span objects have a `.similarity` method that takes another object and returns a floating point number between 0 and 1, indicating how similar they are.
- In order to use similarity, you need a larger spaCy model that has word vectors included.

Similarity examples 1

- Here's an example. Let's say we want to find out whether two documents are similar.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_md")
3 |
4 |
5 | doc1 = nlp("I like fast food")
6 | doc2 = nlp("I like pizza")
7 | print(doc1.similarity(doc2))
8 |
9 | doc = nlp("I like pizza and pasta")
10 | token1 = doc[2]
11 | token2 = doc[4]
12 | print(token1.similarity(token2))
```

- First, we load the medium English model, "en_core_web_md".
- We can then create two doc objects and use the first doc's similarity method to compare it to the second.
- The same works for tokens.

- You can also use the similarity methods to compare different types of objects.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_md")
3 |
4 | doc = nlp("I like pizza")
5 | token = nlp("soap")[0]
6 |
7 | print(doc.similarity(token))
8 | span = nlp("I like pizza and pasta")[2:5]
9 | doc = nlp("McDonalds sells burgers")
10 |
11 | print(span.similarity(doc))
```

- Here, the similarity score is pretty low and the two objects are considered fairly dissimilar.
- Here's another example comparing a span – "pizza and pasta" – to a document about McDonalds.

How does spaCy predict similarity?

- Similarity is determined using word vectors
- Multi-dimensional meaning representations of words
- Generated using an algorithm like Word2Vec and lots of text
- Can be added to spaCy's statistical models
- Cosine similarity, but can be adjusted.
- Vectors can be added to spaCy's statistical models.

- To give you an idea of what those vectors look like, here's an example.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_md")
3 |
4 | doc = nlp("I have a banana")
5 | print(doc[3].vector)
```

- First, we load the medium model again, which ships with word vectors.
- Next, we can process a text and look up a token's vector using the `.vector` attribute.
- The result is a 300-dimensional vector of the word "banana".

Combining Statistical Models and Rules

- Combining statistical models with rule-based systems is one of the most powerful tricks you should have in your NLP toolbox.
- For instance, detecting product or person names usually benefits from a statistical model.
- Instead of providing a list of all person names ever, your application will be able to predict whether a span of tokens is a person name.
- Similarly, you can predict dependency labels to find subject/object relationships.
- To do this, you would use spaCy's entity recognizer, dependency parser or part-of-speech tagger.

Adding Statistical Predictions

- Here's an example of a matcher rule for "golden car".

```
1 | from spacy.matcher import Matcher
2 | import spacy
3 | nlp = spacy.load("en_core_web_sm")
4 | matcher = Matcher(nlp.vocab)
5 |
6 | matcher.add("CAR", [[{"LOWER": "golden"}, {"LOWER": "car"}]])
7 | doc = nlp("I have a Golden Car")
8 | for match_id, start, end in matcher(doc):
9 |     span = doc[start:end]
10 |     print("Matched span:", span.text)
11 |     # Get the span's root token and root head token
12 |     print("Root token:", span.root.text)
13 |     print("Root head token:", span.root.head.text)
14 |     # Get the previous token and its POS tag
15 |     print("Previous token:", doc[start - 1].text, doc[start - 1].pos_)
```

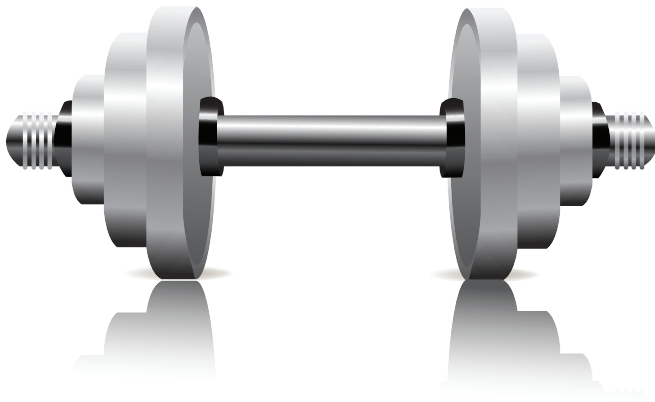
- If we iterate over the matches returned by the matcher, we can get the match ID and the start and end index of the matched span.
- Span objects give us access to the original document and all other token attributes and linguistic features predicted by the model.

- The phrase matcher can be imported from `spacy.matcher` and follows the same API as the regular matcher.

```
1 | import spacy
2 | from spacy.matcher import PhraseMatcher
3 |
4 | nlp = spacy.load("en_core_web_sm")
5 | matcher = PhraseMatcher(nlp.vocab)
6 |
7 | pattern = nlp("Golden Car")
8 | matcher.add("CAR", [pattern])
9 | doc = nlp("I have a Golden Car")
10 |
11 | for match_id, start, end in matcher(doc):
12 |     span = doc[start:end]
13 |     print("Matched span:", span.text)
```

- Instead of a list of dictionaries, we pass in a Doc object as the pattern.
- We can then iterate over the matches in the text, which gives us the match ID, and the start and end of the match.

- Exercise 8 to 11
- [Class-Ex-Lecture3.py](#)



Built-in Pipeline Components

Name	Description	Creates
tagger	Part-of-speech tagger	Token.tag, Token.pos
parser	Dependency parser	Token.dep, Token.head, Doc.sents, Doc.noun_chunks
ner	Named entity recognizer	Doc.ents, Token.ent_job, Token.ent_type
textcat	Text classifier	Doc.cats

- The part-of-speech tagger sets the token.tag and token.pos attributes.
- token.dep and token.head attributes and is also responsible for detecting sentences and base noun phrases.
- The named entity recognizer adds the detected entities to the doc.ents property.
- Finally, the text classifier sets category labels that apply to the whole text.

- To see the names of the pipeline components present in the current `nlp` object, you can use the `nlp.pipe_names` attribute.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
3 | print(nlp.pipe_names)
4 | print(nlp.pipeline)
```

- For a list of component name and component function tuples, you can use the `nlp.pipeline` attribute.
- The component functions are the functions applied to the doc to process it and set attributes – for example, part-of-speech tags or named entities.

Custom Pipeline Components

- Custom pipeline components let you add your own function to the spaCy pipeline that is executed when you call the `nlp` object on a text
- Custom components are executed automatically when you call the `nlp` object on a text.
- They're especially useful for adding your own custom metadata to documents and tokens.
- You can also use them to update built-in attributes, like the named entity spans.

Anatomy of a Component

Argument	Description	Example
last	If True, add last	<code>nlp.add_pipe(component, last=True)</code>
first	If True, add first	<code>nlp.add_pipe(component, first=True)</code>
before	Add before component	<code>nlp.add_pipe(component, before="ner")</code>
after	Add after component	<code>nlp.add_pipe(component, after="tagger")</code>

- Setting `last` to `True` will add the component last in the pipeline. This is the default behavior.
- Setting `first` to `True` will add the component first in the pipeline, right after the tokenizer.
- The `before` and `after` arguments let you define the name of an existing component to add the new component before or after. For example, `before="ner"` will add it before the named entity recognizer.

Example - Simple Component

- Here's an example of a simple pipeline component.

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
3 | from spacy.language import Language
4 |
5 | @Language.component("component1")
6 | def custom_component1(doc):
7 |     print("Doc length:", len(doc))
8 |     return doc
9 | nlp.add_pipe('component1', name="component-info-1", first=True)
10 | print("Pipeline:", nlp.pipe_names)
11 |
12 | @Language.component("component2")
13 | def custom_component2(doc):
14 |     print("Doc length:", len(doc))
15 |     return doc
16 |
17 | nlp.add_pipe('component2', name="component-info-2", first=True)
18 | doc = nlp("Hello world!")
19 | print(doc)
```

Setting Custom Attributes

- Custom attributes let you add any metadata to docs, tokens and spans.
- The data can be added once, or it can be computed dynamically.
- Custom attributes are available via the `._` (dot underscore) property.
- This makes it clear that they were added by the user, and not built into spaCy, like `token.text`.
- The first argument is the attribute name.
- Keyword arguments let you define how the value should be computed.

```
1 | from spacy.tokens import Doc, Token, Span
2 | Doc.set_extension("title", default=None)
3 | Token.set_extension("is_color", default=False)
4 | Span.set_extension("has_color", default=False)
```

- Define a getter and an optional setter function. Getter only called when you retrieve the attribute value.

```
1 | from spacy.tokens import Token
2 | import spacy
3 | nlp = spacy.load("en_core_web_sm")
4 |
5 | def get_is_color(token):
6 |     colors = ["red", "yellow", "blue"]
7 |     return token.text in colors
8 |
9 | Token.set_extension("is_color", getter=get_is_color)
10 |
11 | doc = nlp("The sky is blue.")
12 | print(doc[3]._.is_color, "-", doc[3].text)
```

- The getter function is only called when you retrieve the attribute.
- This lets you compute the value dynamically, and even take other custom attributes into account.

- Method extensions make the extension attribute a callable method.

```
1 | from spacy.tokens import Doc
2 | import spacy
3 | nlp = spacy.load("en_core_web_sm")
4 |
5 | def has_token(doc, token_text):
6 |     in_doc = token_text in [token.text for token in doc]
7 |     return in_doc
8 |
9 | Doc.set_extension("has_token", method=has_token)
10 |
11 | doc = nlp("The sky is blue.")
12 | print(doc._.has_token("blue"), "- blue")
13 | print(doc._.has_token("cloud"), "- cloud")
```

- You can then pass one or more arguments to it, and compute attribute values dynamically
- In this example, the method function checks whether the doc contains a token with a given text. The first argument of the method is always the object itself – in this case, the doc.

- Use `nlp.pipe` method

```
1 | import spacy
2 | nlp = spacy.load("en_core_web_sm")
3 |
4 | # docs = [nlp(text) for text in LOTS_OF_TEXTS]---Slow
5 | # docs = list(nlp.pipe(LOTS_OF_TEXTS))---Fast
6 |
7 | # doc = nlp("Hello world")
8 | # doc = nlp.make_doc("Hello world!")
9 |
10 | text = 'I love performance'
11 | with nlp.disable_pipes("tagger", "parser"):
12 |     doc = nlp(text)
13 |     print(doc.text)
```

- Using only the tokenizer.
- Disabling pipeline components

Training and Updating Models

- Statistical models make predictions based on the examples they were trained on.
- You can usually make the model more accurate by showing it examples from your domain.
- You often also want to predict categories specific to your problem, so the model needs to learn about them.
- This is essential for text classification, very useful for entity recognition and a little less critical for tagging and parsing.

- The training data tells the model what we want it to predict.
- This could be texts and named entities we want to recognize, or tokens and their correct part-of-speech tags.
- To update an existing model, we can start with a few hundred to a few thousand examples.
- To train a new category we may need up to a million.
- spaCy's pre-trained English models for instance were trained on 2 million words labelled with part-of-speech tags, dependencies and named entities.
- Training data is usually created by humans who assign labels to texts.

- Exercise 12 to 14
- [Class-Ex-Lecture3.py](#)

