

drawing



Application of Deep Learning to Text and Image Data

Module 2, Lab 1: Processing Text

In this notebook, you will learn techniques to analyze and process text data. Text processing is known as *natural language processing (NLP)* and is an important topic because of how much information is communicated through text. Knowing how to handle text will help you build models that perform better and are more useful.

You will learn the following:

- What a word cloud is and how to create one
- How to use stemming and lemmatization
- What part-of-speech tagging is and how it impacts text processing
- How to use named entity recognition to sort data

You will be presented with two kinds of exercises throughout the notebook: activities and challenges.

Activity



Challenge



No coding is needed for an activity. You try to understand a concept, answer questions, or run a code cell.

Challenges are where you can practice your coding skills.

Index

- [Word cloud](#)
- [Part-of-speech tagging](#)
- [Stemming and lemmatization](#)
- [Named entity recognition](#)

Initial Setup

First let's put everything in place.

```
In [1]: !pip install -U -q -r requirements.txt
```

```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This
behaviour is the source of the following dependency conflicts.
autovizwidget 0.21.0 requires pandas<2.0.0,>=0.20.1, but you have pandas 2.0.3 which is incompatible.
hdiipyterutils 0.21.0 requires pandas<2.0.0,>=0.17.1, but you have pandas 2.0.3 which is incompatible.
sparkmagic 0.21.0 requires pandas<2.0.0,>=0.17.1, but you have pandas 2.0.3 which is incompatible.
```

Install the [spaCy](#) library. This will be used to perform some NLP tasks in the lab.

```
In [2]: !python -m spacy download en_core_web_sm
```

```

/home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages/torch/cuda/__init__.py:551: UserWarning
: Can't initialize NVML
  warnings.warn("Can't initialize NVML")
Collecting en-core-web-sm==3.7.1
  Downloading https://github.com/explosion/spacy-models/releases/download/en_core_web_sm-3.7.1/en_core_web_sm-3
.7.1-py3-none-any.whl (12.8 MB)
12.8/12.8 MB 100.4 MB/s eta 0:00:0000:010:01
Requirement already satisfied: spacy<3.8.0,>=3.7.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from en-core-web-sm==3.7.1) (3.7.2)
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/py
thon3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.0.12)
Requirement already satisfied: spacy-loggers<2.0.0,>=1.0.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/py
thon3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.0.5)
Requirement already satisfied: murmurhash<1.1.0,>=0.28.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pyth
on3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.0.12)
Requirement already satisfied: cymem<2.1.0,>=2.0.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.0.11)
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.
10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.0.9)
Requirement already satisfied: thinc<8.3.0,>=8.1.8 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (8.2.5)
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.1
0/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.1.3)
Requirement already satisfied: srsly<3.0.0,>=2.4.3 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.5.1)
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python
3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.0.10)
Requirement already satisfied: weasel<0.4.0,>=0.1.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.1
0/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.3.4)
Requirement already satisfied: typer<0.10.0,>=0.3.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.1
0/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.9.4)
Requirement already satisfied: smart-open<7.0.0,>=5.2.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pytho
n3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (6.4.0)
Requirement already satisfied: tqdm<5.0.0,>=4.38.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (4.66.1)
Requirement already satisfied: requests<3.0.0,>=2.13.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python
3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.31.0)
Requirement already satisfied: pydantic!=1.8,!<1.8.1,<3.0.0,>=1.7.4 in /home/ec2-user/anaconda3/envs/pytorch_p3
10/lib/python3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.10.13)
Requirement already satisfied: jinja2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-package
s (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.1.2)
Requirement already satisfied: setuptools in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-pac
kages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (68.2.2)
Requirement already satisfied: packaging>=20.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/sit
e-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (21.3)
Requirement already satisfied: langcodes<4.0.0,>=3.2.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python
3.10/site-packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.5.0)
Requirement already satisfied: numpy>=1.19.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-
packages (from spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.23.5)
Requirement already satisfied: language-data>=1.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/
site-packages (from langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.3.0)
Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pytho
n3.10/site-packages (from packaging>=20.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.0.9)
Requirement already satisfied: typing-extensions>=4.2.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pytho
n3.10/site-packages (from pydantic!=1.8,!<1.8.1,<3.0.0,>=1.7.4->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (4.
8.0)
Requirement already satisfied: charset-normalizer<4,>=2 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pytho
n3.10/site-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.3.1)
Requirement already satisfied: idna<4,>=2.5 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-p
ackages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (3.4)
Requirement already satisfied: urllib3<3,>=1.21.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/
site-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1.26.18)
Requirement already satisfied: certifi>=2017.4.17 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/
site-packages (from requests<3.0.0,>=2.13.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2023.7.22)
Requirement already satisfied: blis<0.8.0,>=0.7.8 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/
site-packages (from thinc<8.3.0,>=8.1.8->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.7.11)
Requirement already satisfied: confection<1.0.0,>=0.0.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/pytho
n3.10/site-packages (from thinc<8.3.0,>=8.1.8->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.1.5)
Requirement already satisfied: click<9.0.0,>=7.1.1 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10
/site-packages (from typer<0.10.0,>=0.3.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (8.1.7)
Requirement already satisfied: cloudpathlib<0.17.0,>=0.7.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/py
thon3.10/site-packages (from weasel<0.4.0,>=0.1.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (0.16.0)
Requirement already satisfied: MarkupSafe>=2.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/sit
e-packages (from jinja2->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (2.1.3)
Requirement already satisfied: marisa-trie>=1.1.0 in /home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/
site-packages (from language-data>=1.2->langcodes<4.0.0,>=3.2.0->spacy<3.8.0,>=3.7.2->en-core-web-sm==3.7.1) (1
.2.1)
Installing collected packages: en-core-web-sm
Successfully installed en-core-web-sm-3.7.1
✓ Download and installation successful
You can now load the package via spacy.load('en_core_web_sm')

```

```

In [3]: # Import the dependencies
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
import re, string
import Stemmer
import spacy
from spacy import displacy

```

```
import pandas as pd
```

```
Matplotlib is building the font cache; this may take a moment.
/home/ec2-user/anaconda3/envs/pytorch_p310/lib/python3.10/site-packages/torch/cuda/__init__.py:551: UserWarning
: Can't initialize NVML
warnings.warn("Can't initialize NVML")
```

Next, you need to create a function to preprocess text so that only real words, not special characters and numbers, are displayed.

```
In [4]: # Preprocess text
def preProcessText(text):
    # Lowercase and strip leading and trailing white space
    text = text.lower().strip()

    # Remove HTML tags
    text = re.compile("<.*?>").sub("", text)

    # Remove punctuation
    text = re.compile("[%s]" % re.escape(string.punctuation)).sub(" ", text)

    # Remove extra white space
    text = re.sub("\s+", " ", text)

    # Remove numbers
    text = re.sub(r"[0-9]", "", text)

    return text
```

Word cloud

Word clouds, which are also known as *text clouds* or *tag clouds*, help you visualize text data by highlighting the important words or phrases. Word clouds convey crucial information at a glance by making commonly occurring words bigger and bolder. These clouds are commonly used to compare and contrast two pieces of text. Word clouds are also used to identify the topic of a document.

To create a word cloud, you will use [WordCloud for Python](#).

The following text is from the [What Is Natural Language Processing \(NLP\)?](#) page on [aws.amazon.com](#).

```
In [5]: text = "Natural language processing (NLP) is a machine learning technology that gives computers the \
ability to interpret, manipulate, and comprehend human language. Organizations today have large volumes \
of voice and text data from various communication channels like emails, text messages, social media \
newsfeeds, video, audio, and more. They use NLP software to automatically process this data, analyze \
the intent or sentiment in the message, and respond in real time to human communication. \
Natural language processing (NLP) is critical to fully and efficiently analyze text and speech data. \
It can work through the differences in dialects, slang, and grammatical irregularities typical in \
day-to-day conversations. \
Companies use it for several automated tasks, such as to: \
<li>Process, analyze, and archive large documents</li> \
<li>Analyze customer feedback or call center recordings</li> \
<li>Run chatbots for automated customer service</li> \
<li>Answer who-what-when-where questions</li> \
<li>Classify and extract text</li> \
You can also integrate NLP in customer-facing applications to communicate more effectively with \
customers. For example, a chatbot analyzes and sorts customer queries, responding automatically to \
common questions and redirecting complex queries to customer support. This automation helps reduce \
costs, saves agents from spending time on redundant queries, and improves customer satisfaction."
```

```
In [6]: # Remove stop words before generating the word cloud
wordcloud = WordCloud(stopwords=STOPWORDS, background_color="black", max_words=300)

# Clean up the text to prevent plotting punctuation and duplicate words (for example, 'Natural' and 'natural')
wordcloud.generate(preProcessText(text))

plt.figure(figsize=(15, 10))
plt.axis("off")
plt.imshow(wordcloud);
```


PRON	Pronoun
PROPN	Proper Noun
PUNCT	Punctuation
SCONJ	Subordinating Conjunction
SYM	Symbol
VERB	Verb
X	Other

Now you can use the tagger to tag each token or word in the following text.

Important: Always remember to preprocess the text before tagging, as we have done before in this notebook.

```
In [8]: # Text sample
text
```

```
Out[8]: 'Natural language processing (NLP) is a machine learning technology that gives computers the ability to interpret, manipulate, and comprehend human language. Organizations today have large volumes of voice and text data from various communication channels like emails, text messages, social media newsfeeds, video, audio, and more. They use NLP software to automatically process this data, analyze the intent or sentiment in the message, and respond in real time to human communication. Natural language processing (NLP) is critical to fully and efficiently analyze text and speech data. It can work through the differences in dialects, slang, and grammatical irregularities typical in day-to-day conversations. Companies use it for several automated tasks, such as to: <li>Process, analyze, and archive large documents</li> <li>Analyze customer feedback or call center recordings</li> <li>Run chatbots for automated customer service</li> <li>Answer who-what-when-where questions</li> <li>Classify and extract text</li> You can also integrate NLP in customer-facing applications to communicate more effectively with customers. For example, a chatbot analyzes and sorts customer queries, responding automatically to common questions and redirecting complex queries to customer support. This automation helps reduce costs, saves agents from spending time on redundant queries, and improves customer satisfaction.'
```

Try it yourself!



To use the spaCy part-of-speech tagger, run the following cell.

Observe the tags that are assigned to each word, and use the table from a previous cell to understand the meaning of each tag.

```
In [9]: # Part-of-speech tagging
nlp = spacy.load("en_core_web_sm")
doc = nlp(text)
token_text = [token.orth_ for token in doc]
token_lemma = [token.lemma_ for token in doc]
token_pos = [token.pos_ for token in doc]
pd.DataFrame(zip(token_text, token_pos),
              columns=['token_text', 'token_lemma'])
```

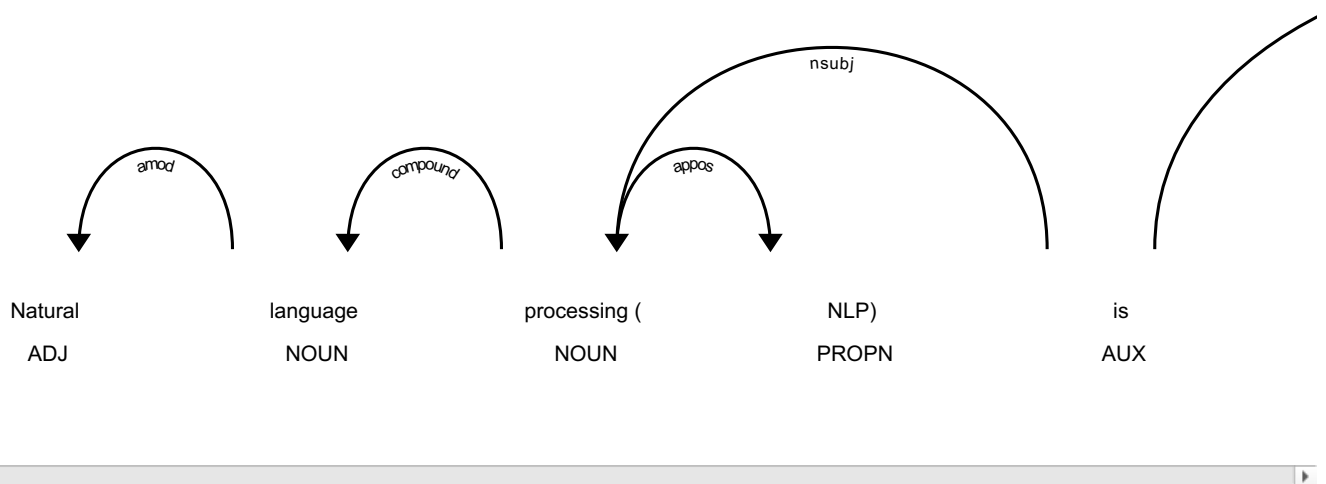
```
Out[9]:
```

	token_text	token_lemma
0	Natural	ADJ
1	language	NOUN
2	processing	NOUN
3	(PUNCT
4	NLP	PROPN
...
248	and	CCONJ
249	improves	VERB
250	customer	NOUN
251	satisfaction	NOUN
252	.	PUNCT

253 rows × 2 columns

Refer to the table in a previous cell to identify the tags that the spaCy tagger produces.

```
In [10]: # uncomment the code below to display the dependency parse or named entities and their tags
displacy.render(doc, style="dep")
```



Stemming and lemmatization

Stemming and lemmatization are two ways to process words so that a model will be more efficient. Both methods remove parts of words so that they can be grouped together.

For example, in the following sentence, "ning" would be removed from "running" so that "running" and "run" would be categorized the same.

The child enjoys **running**, so they **run** every day.

What could make stemming and lemmatization difficult to do properly?

Try it yourself!

Activity



In the next few sections, you will compare stemming and lemmatization.
Consider which text processing method is more suitable for the use case that is provided.

Stemming

Stemming is a rule-based system to convert words into their root forms by removing suffixes. This method helps to enhance similarities (if any) between sentences.

Examples:

"jumping", "jumped" -> "jump"

"cars" -> "car"

```
In [11]: # let's list the language available
print(Stemmer.algorithms())

['arabic', 'armenian', 'basque', 'catalan', 'danish', 'dutch', 'english', 'finnish', 'french', 'german', 'greek', 'hindi', 'hungarian', 'indonesian', 'irish', 'italian', 'lithuanian', 'nepali', 'norwegian', 'porter', 'portuguese', 'romanian', 'russian', 'serbian', 'spanish', 'swedish', 'tamil', 'turkish', 'yiddish']
```

```
In [12]: # we will use english for this example
stemmer = Stemmer.Stemmer('english')
```

```
In [13]: original_text = "  This is a message to be cleaned. It may involve some things like: <br>, ?, :, ' ' adjacent
print(original_text)
```

 This is a message to be cleaned. It may involve some things like:
, ?, :, ' ' adjacent spaces and tabs
.

```
In [14]: # Cleaned text
cleaned_text = preProcessText(original_text)
print(cleaned_text)
```

this is a message to be cleaned it may involve some things like adjacent spaces and tabs

```
In [15]: # Use a tokenizer (nlp) and stemmer from the PyStemmer library
nlp = spacy.load("en_core_web_sm")
doc = nlp(original_text)
stemmed_sentence = []
original_sentence = []
# Tokenize the sentence
for token in doc:
    original_sentence.append(token.text)
    stemmed_sentence.append(stemmer.stemWord(token.text))

pd.DataFrame(zip(original_sentence, stemmed_sentence),
             columns=['token_text', 'token_stemmer'])
```


Out[15]:

	token_text	token_stemmer
0		
1	This	This
2	is	is
3	a	a
4	message	messag
5	to	to
6	be	be
7	cleaned	clean
8	.	.
9	It	It
10	may	may
11	involve	involv
12	some	some
13	things	thing
14	like	like
15	:	:
16	<	<
17	br	br
18	>	>
19	,	,
20	?	?
21	,	,
22	:	:
23	,	,
24	"	"
25		
26	adjacent	adjac
27	spaces	space
28	and	and
29	tabs	tab
30		
31	.	.
32		

```
In [16]: stemmed_text = " ".join(stemmed_sentence)
print(stemmed_text)

This is a messag to be clean . It may involv some thing like : < br > , ? , : , ' ' adjac space and tab .
```

From the output of the previous code cell, you can see that stemming isn't perfect. It makes mistakes, such as "messag", "involv", and "adjac". Stemming is a rule-based method that sometimes mistakenly removes suffixes from words. It does run quickly, which makes it appealing to use for massive datasets.

Lemmatization

If you aren't satisfied with the result of stemming, you can use the lemmatization instead. This method usually requires more work but gives better results.

Lemmatization needs to know the correct word position tags, such as "noun", "verb", or "adjective". You need to use another spaCy function to feed this information to the lemmatizer.

The cell below uses part of the full list of position tags listed in the previous session `Part-of-speech tagging`.

```
In [17]: # get the original and the lemmatized the word/token
token_lemma = [token.lemma_ for token in doc]
token_text = [token.orth_ for token in doc]
pd.DataFrame(zip(token_text, token_lemma),
              columns=['token_text', 'token_lemma'])
```

```
Out[17]:
```

	token_text	token_lemma
0		
1	This	this
2	is	be
3	a	a
4	message	message
5	to	to
6	be	be
7	cleaned	clean
8	.	.
9	It	it
10	may	may
11	involve	involve
12	some	some
13	things	thing
14	like	like
15	:	:
16	<	<
17	br	br
18	>	>
19	,	,
20	?	?
21	,	,
22	:	:
23	,	,
24	"	"
25		
26	adjacent	adjacent
27	spaces	space
28	and	and
29	tabs	tab
30		
31	.	.
32		

```
In [18]: lemmatized_text = " ".join(token_lemma)
print(lemmatized_text)
```

this be a message to be clean . it may involve some thing like : < br > , ? , : , ' ' adjacent space and t
ab .

How do the results compare? Is the lemmatized text better than the stemmed text?

Named entity recognition

Named entity recognition involves identification of key information in text and then classifying that information into predefined categories, such as person, organization, place, or date. This is one of the most popular NLP tasks.

For this section, you will use [spaCy](#). The following table lists the categories and meanings of the category labels that the spaCy module uses.

Category	Meaning
CARDINAL	Numerals that don't fall under another type
DATE	Absolute or relative dates or periods
EVENT	Named hurricanes, battles, wars, sports events, and so on
FAC	Buildings, airports, highways, bridges, and so on
GPE	Countries, cities, states
LANGUAGE	Any named language

LAW	Named documents made into laws
LOC	Non-GPE locations, mountain ranges, bodies of water
MONEY	Monetary values, including unit
NORP	Nationalities, or religious or political groups
ORDINAL	"first", "second", and so on
ORG	Companies, agencies, institutions, and so on
PERCENT	Percentage, including "%"
PERSON	People, including fictional
PRODUCT	Objects, vehicles, foods, and so on (not services)
QUANTITY	Measurements, as of weight or distance
WORK_OF_ART	Titles of books, songs, and so on

The following text was retrieved from the [Amazon - The Climate Pledge](#) page of the About Amazon website.

```
In [19]: # Sample text for named entity recognition
ner_text = "Amazon and Global Optimism co-founded The Climate Pledge, \
a commitment to net-zero \
carbon by 2040."
```

```
In [20]: # Load the spaCy English pipeline for named entity recognition
NER = spacy.load("en_core_web_sm")

# Tag entities in the text
for word in NER(ner_text).ents:
    print(word.text, word.label_)
```

```
Amazon ORG
Global Optimism ORG
The Climate Pledge WORK_OF_ART
2040 DATE
```

Visualizing the Tags

spaCy has a visualizer called displaCy that you can use to visualize tags. Run the following cell to see it working.

```
In [21]: # Visual tag with text
displacy.render(NER(ner_text), style="ent", jupyter=True)
```

Amazon **ORG** and Global Optimism **ORG** co-founded The Climate Pledge **WORK_OF_ART**, a commitment to net-zero carbon by 2040 **DATE**.

As you can see, named entity recognition can help you identify different entities in text. The process isn't always correct, but it can process large sections of text faster than a human can.

Conclusion

In this lab, you practiced using text processing techniques.

Next lab

In the next lab, you will learn about the bag-of-words (BoW) method to convert text data into numerical values.