

Task 1: Spacy Installation

Task 2: spaCy Hello World

```
In [1]: from spacy.tokens.doc import Doc

from spacy.vocab import Vocab

doc = Doc(Vocab(), words = ['u'Hello', 'u'World!'])
print(doc)

Hello World!

In [2]: print(type(doc))
print(doc.vocab)
for token in doc:
    lexeme = doc.vocab[token.text]
    print(lexeme.text)

<class 'spacy.tokens.doc.Doc'>
<spacy.vocab.Vocab object at 0x7f688e4569e0>
Hello
World!
```

Questions:

1. The `Vocab()` object belongs to which class?

Belongs to Class:

- The `Vocab()` object belongs to the `spacy.vocab.Vocab` class.

Purpose:

- Serves as a memory-efficient container.
- Manages the vocabulary of the language model.
- Facilitates quick access to word representations.

2. What is a lexeme object (You need to check the API)?

Definition:

- The Lexeme object allows access to properties and features associated with individual words.

Purpose:

- Provides a foundation for linguistic analysis.
- Essential for various natural language processing tasks.

API Check:

- The Lexeme object is an instance of the `spacy.lexeme.Lexeme` class.

```
In [3]: import spacy

nlp = spacy.load('en_core_web_sm')

doc = nlp('u'I want to learn spacy')

token_text1 = [token.text for token in doc]
token_text2 = [doc[i].text for i in range(len(doc))]

print(token_text1)
print(token_text2)

['I', 'want', 'to', 'learn', 'spacy']
['I', 'want', 'to', 'learn', 'spacy']
```

Questions:

1. What is `en_core_web_sm`?

Answer: `en_core_web_sm` is a pre-trained English language model in Spacy. It belongs to the "web" models and offers various functions for performing natural language tasks.

2. What is the size of `en_core_web_sm`?

Answer: `en_core_web_sm` is an English multi-task CNN (Convolutional Neural Network) model trained on OntoNotes, and its size is approximately 11 MB.

3. What other variations can be used?

- en_core_web_md:**
 - A medium-sized English model with word vectors. It includes more features and is larger than `en_core_web_sm`.
- en_core_web_lg:**
 - A large-sized English model with word vectors. It is the largest among the English models provided by Spacy and includes even more features.
- Specialized Models:**
 - Spacy offers models trained for specific tasks, such as `en_core_web_trf`, which is a transformer-based model. Transformer-based models may provide improved performance in certain scenarios.

```
In [4]: doc = nlp('u'I want to learn spacy.')
for i in range(len(doc)):
    print([t for t in doc[i].lefts])

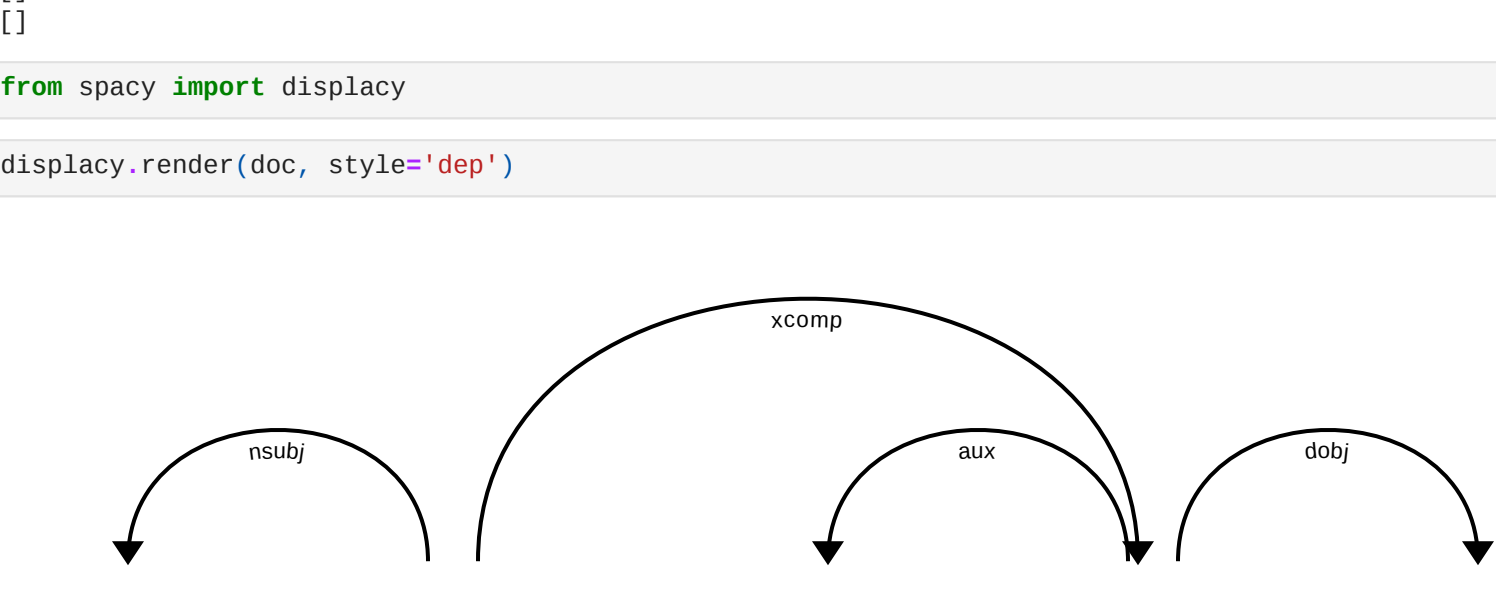
[]
[]
[]
[to]
[]
[]

In [5]: doc = nlp('u'I want to learn spacy.')
for i in range(len(doc)):
    print([t for t in doc[i].rights])
    print([t for t in doc[i].children])

[]
[]
[learn, .]
[to, learn, .]
[]
[]
[spacy]
[to, spacy]
[]
[]
[]
[]
```

```
In [6]: from spacy import displacy

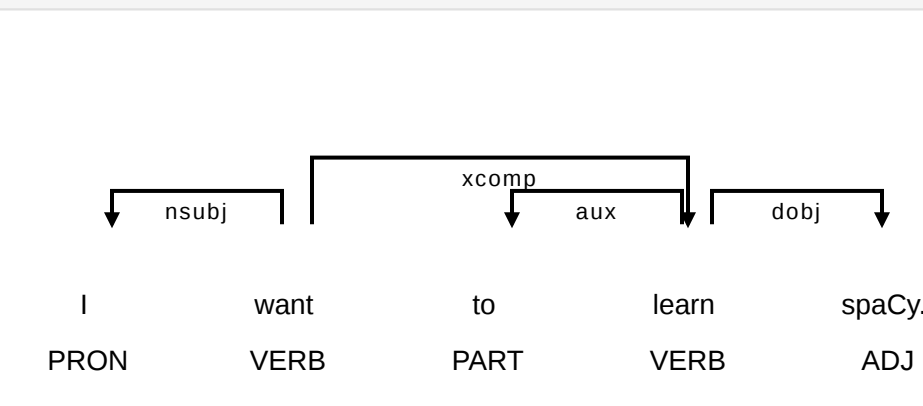
In [7]: displacy.render(doc, style='dep')
```



Questions

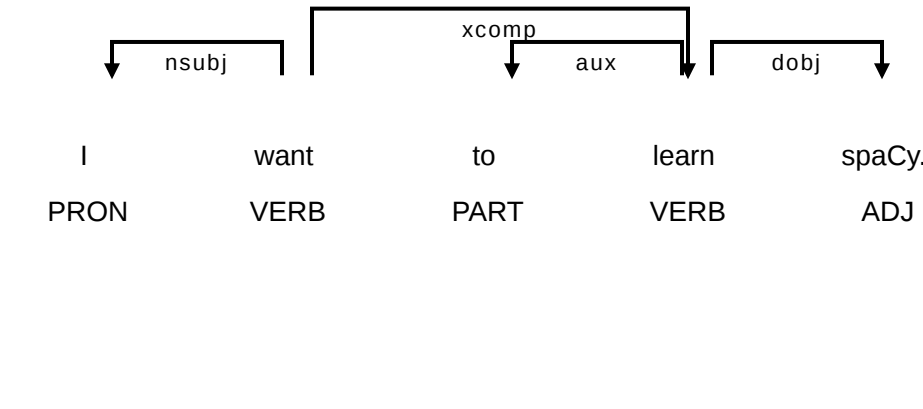
1. Draw the left and right dependencies for the sentence: I want to learn spaCy.

```
In [8]: # Draw left and right dependencies
displacy.render(doc, style='dep', options={'distance': 100, 'compact': True, 'jupyter': True})
```



2. Draw the children for the sentence: I want to learn spaCy.

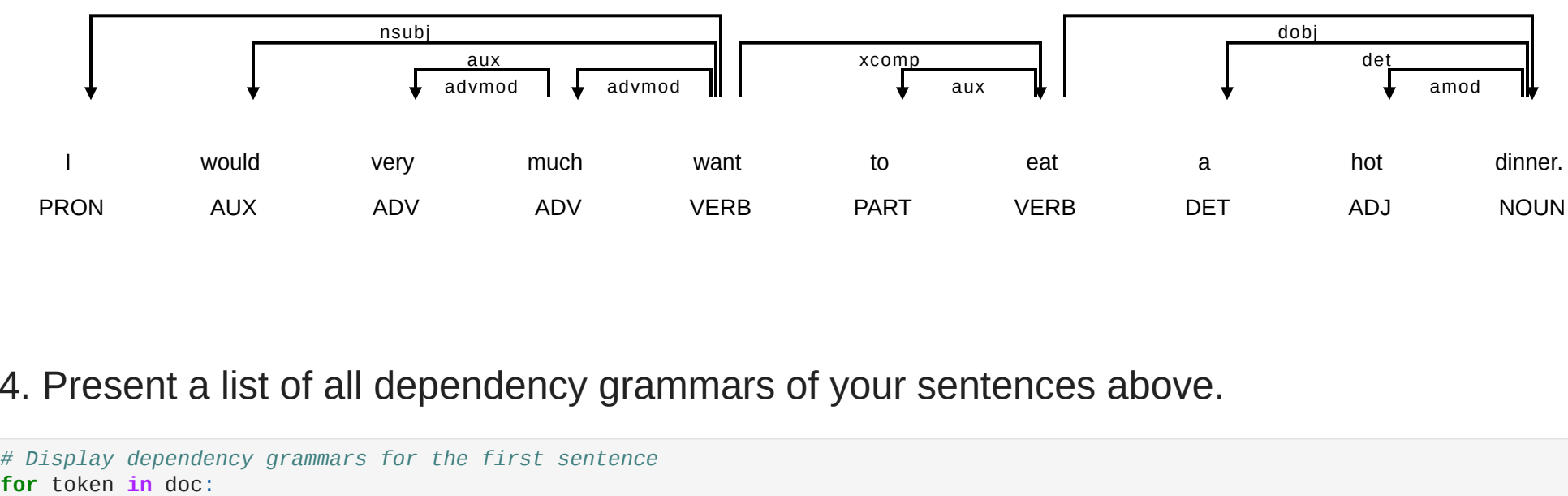
```
In [9]: # Draw children
displacy.render(doc, style='dep', options={'distance': 100, 'compact': True, 'mode': 'dep'}, jupyter=True)
```



3. Draw the left and right dependencies for the sentence: I would very much want to eat a hot dinner.

```
In [10]: # Process the second sentence
doc2 = nlp('I would very much want to eat a hot dinner.')

# Draw left and right dependencies
displacy.render(doc2, style='dep', options={'distance': 100, 'compact': True, 'jupyter': True})
```



4. Present a list of all dependency grammars of your sentences above.

```
In [11]: # Display dependency grammars for the first sentence
for token in doc:
    print("({token.text}): {token.dep}")
print("\n")

# Display dependency grammars for the second sentence
for token in doc2:
    print("({token.text}): {token.dep}")
print("\n")

I: nsubj
want: ROOT
to: aux
learn: xcomp
spacy: dobj
.: punct

I: nsubj
would: aux
very: advmod
much: advmod
want: ROOT
to: aux
eat: xcomp
a: det
hot: amod
dinner: amod
.: punct
```

Task 3: NLTK vs spaCy Pipelines

Step 01: Take Raw Text

```
In [12]: texts = ['uWe are nearing the end of the semester at Peshawar. Final exams of the Fall 2023 semester will start soon.']
```

Step 02: Do Sentence Level Segmentation

```
In [13]: import nltk

In [14]: for text in texts:
    sentences = nltk.sent_tokenize(text)
    print(sentences)

['We are nearing the end of the semester at Peshawar.', 'Final exams of the Fall 2023 semester will start soon.']
```

Step03: Extract Tokens in Sentences

```
In [15]: for sentence in sentences:
    words = nltk.word_tokenize(sentence)
    print(words)

['We', 'are', 'nearing', 'the', 'end', 'of', 'the', 'semester', 'at', 'Peshawar', '.']
['Final', 'exams', 'of', 'the', 'Fall', '2023', 'semester', 'will', 'start', 'soon', '.']
```

Step 04: Perform Part of Speech Tagging

```
In [16]: tagged_words = nltk.pos_tag(words)

print(tagged_words)

[('Final', 'JJ'), ('exams', 'NN'), ('of', 'IN'), ('the', 'DT'), ('Fall', 'NN'), ('2023', 'CD'), ('semester', 'NN'), ('will', 'MD'), ('start', 'VB'), ('soon', 'RB'), ('.', '.')]

In [17]: ne_tagged_words = nltk.ne_chunk(tagged_words)
print(ne_tagged_words)

(S
  Final/JJ
  exams/NN
  of/IN
  the/DT
  Fall/NN
  2023/CD
  semester/NN
  will/MD
  start/VB
  soon/RB
  ./.)
```

Step 05: Detect Named Entity of Tagged Words

```
In [18]: nltk.download('punkt') # Sentence Tokenize

[nltk_data] Downloading package punkt to /home/jadi/nltk.data...
[nltk_data] Package punkt is already up-to-date!

Out[18]: True
```

```
In [19]: nltk.download('averaged_perceptron_tagger') # POS Tagging

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /home/jadi/nltk.data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!

Out[19]: True
```

```
In [20]: nltk.download('maxent_ne_chunker') # Named Entity Chunking

[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data] /home/jadi/nltk.data...
[nltk_data] Package maxent_ne_chunker is already up-to-date!

Out[20]: True
```

```
In [21]: nltk.download('words') # Word Tokenize

[nltk_data] Downloading package words to /home/jadi/nltk.data...
[nltk_data] Package words is already up-to-date!

Out[21]: True
```

```
In [22]: from spacy import displacy
doc = nlp('uWe are nearing the end of the semester at Peshawar. Final exams of the Fall 2023 semester will start soon.')
displacy.render(doc, style='ent')
```

We are nearing **the end of the semester** **DATE** at **Peshawar** **GPE**. Final exams of the Fall 2023 semester will start soon.

```
In [23]: for ent in doc.ents:
    print(ent.text, ent.label_)

the end of the semester DATE
Peshawar GPE
```

Questions

1. How did the Named Entity Output of the NLTK pipeline look like? Present its output.

```
In [24]: import nltk

texts = ['uWe are nearing the end of the semester at Peshawar. Final exams of the Fall 2023 semester will start soon.']

for text in texts:
    sentences = nltk.sent_tokenize(text)
    print(sentences)

    for sentence in sentences:
        words = nltk.word_tokenize(sentence)
        print(words)

        tagged_words = nltk.pos_tag(words)
        print(tagged_words)

        ne_tagged_words = nltk.ne_chunk(tagged_words)
        print(ne_tagged_words)

['We', 'are', 'nearing', 'the', 'end', 'of', 'the', 'semester', 'at', 'Peshawar', '.']
[('We', 'PRP'), ('are', 'VBP'), ('nearing', 'VBG'), ('the', 'DT'), ('end', 'NN'), ('of', 'IN'), ('the', 'DT'), ('semester', 'NN'), ('at', 'IN'), ('Peshawar', 'NNP'), ('.', '.')]
(S
  we/PRP
  are/VBP
  nearing/VBG
  the/DT
  end/NN
  of/IN
  the/DT
  semester/NN
  at/IN
  (ORGANIZATION Peshawar/NNP)
  ./.)

['Final', 'exams', 'of', 'the', 'Fall', '2023', 'semester', 'will', 'start', 'soon', '.']
[('Final', 'JJ'), ('exams', 'NN'), ('of', 'IN'), ('the', 'DT'), ('Fall', 'NN'), ('2023', 'CD'), ('semester', 'NN'), ('will', 'MD'), ('start', 'VB'), ('soon', 'RB'), ('.', '.')]
(S
  Final/JJ
  exams/NN
  of/IN
  the/DT
  Fall/NN
  2023/CD
  semester/NN
  will/MD
  start/VB
  soon/RB
  ./.)
```

2. How did the Named Entity Output of the spaCy pipeline look like? Present its output.

```
In [25]: from spacy import displacy

doc = nlp('uWe are nearing the end of the semester at Peshawar. Final exams of the Fall 2023 semester will start soon.')
```

We are nearing **the end of the semester** **DATE** at **Peshawar** **GPE**. Final exams of the Fall 2023 semester will start soon.

```
In [26]: import spacy
nlp = spacy.load('en_core_web_sm')
nlp.pipe_names

Out[26]: ['tok2vec', 'tagger', 'parser', 'attribute_ruler', 'lemmatizer', 'ner']
```

What is the default pipeline structure of spaCy?

- Tokenizer:** Splits text into individual tokens (words or subwords).
- Part-of-Speech Tagger:** Assigns parts of speech (e.g. noun, verb) to tokens.
- Dependency Parser:** Establishes syntactic relationships between words.
- Named Entity Recognizer (NER):** Identifies entities like persons, locations, dates, etc.
- Sentence Classification:** Assigns labels to the entire text (e.g., sentiment analysis).
- Entity Linker:** Connects named entities to knowledge base entries.
- Text Vectorization:** Generates vectors representing word and document semantics.
- Sentence Boundary Detector:** Identifies sentence boundaries.
- Morphological Analysis:** Provides information about word morphology (e.g., lemma).

```
In [27]: nlp = spacy.load('en_core_web_sm', disable=['parser'])
```

```
In [28]: # Adding a custom component
from spacy.language import Language
def my_component(doc):
    # Do something to the doc here
    return doc
```

```
In [29]: nlp.add_pipe('my_component')
```

```
Out[29]: <function __main__.my_component(doc)>
```

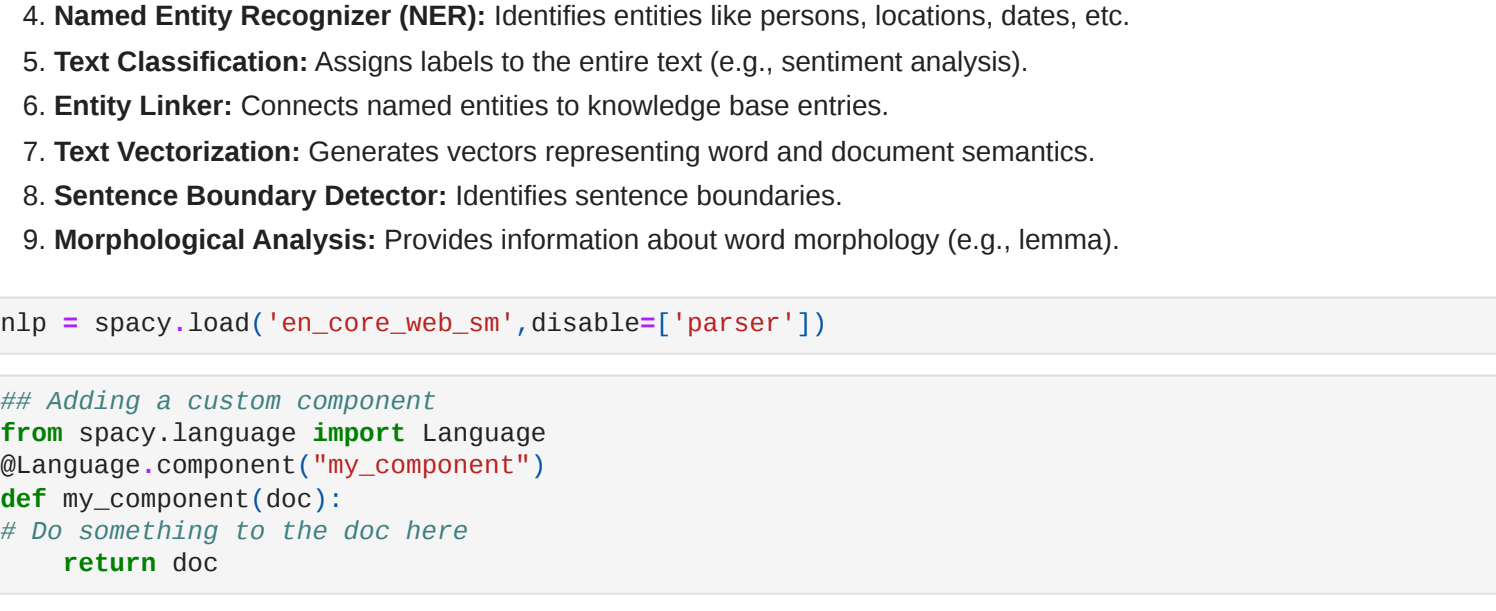
Task 4: Finding Patterns in Sentences

```
In [30]: import spacy
from spacy import displacy

# Load English model with a dependency parser
nlp = spacy.load('en_core_web_sm')

# Process the text
doc = nlp('uHow do I learn spacy.')

# Visualize the dependency tree
displacy.render(doc, style='dep', jupyter=True)
```



```
In [31]: import spacy
from spacy import displacy

# Load English model with a dependency parser
nlp = spacy.load('en_core_web_sm')
```

Process the text

```
doc = nlp('uHow do I learn spacy.')

# Check for the dependency pattern
if dep_pattern(doc):
    print('Found')
else:
    print('Not Found')
```

```
In [32]: nsubj
aux
xcomp
dobj
Not Found
```