

AnalyzingSentenceStructure

November 15, 2021

```
In [84]: import markdown
```

1 My Jupyter Notebook

1.1 Installing necessary packages

```
In [85]: # Installing packages
        #!pip install spacy
        #!pip install markdown
        #!python -m spacy download en_core_web_sm
```

1.2 Importing spacy and pandas

```
In [86]: import spacy
        import pandas as pd
```

1.3 Reading wiki_sentences dataset

Visit source

```
In [87]: candidate_sentences = pd.read_csv("wiki_sentences_v2.csv")
```

```
In [88]: # Getting the size of pandas dataframe and looking at the first few rows for sanity c
        candidate_sentences.shape
        candidate_sentences['sentence'].sample(5)
```

```
Out[88]: 3140    lost is a music festival driven by culture, tr...
        3657    nine rewrites were performed by nine separate ...
        3612    colorization is sometimes used on documentary ...
        3216                                ma was released on may 31, 2019.
        3309    some established graduate programs in the fiel...
        Name: sentence, dtype: object
```

1.3.1 Loading en_core_web_sm, which is a small English pipeline trained on written web text (blogs, news, comments), that includes vocabulary, syntax and entities.

```
In [89]: nlp = spacy.load('en_core_web_sm')
        doc = nlp("the drawdown process is governed by astm standard d823")
```

```

        for tok in doc:
            print(tok.text, "...", tok.dep_)

the ... det
drawdown ... compound
process ... nsubjpass
is ... auxpass
governed ... ROOT
by ... agent
astm ... compound
standard ... compound
d823 ... pobj

```

1.4 Spacy.explain

```
In [90]: spacy.explain("auxpass")
```

```
Out[90]: 'auxiliary (passive)'
```

1.5 Accessing token attributes

```
In [91]: doc = nlp("Airstrikes continued into the early hours of Monday morning in Gaza.")
        # Token texts
        [token.text for token in doc]
```

```
Out[91]: ['Airstrikes',
          'continued',
          'into',
          'the',
          'early',
          'hours',
          'of',
          'Monday',
          'morning',
          'in',
          'Gaza',
          '.']
```

1.6 Accessing spans

```
In [92]: doc = nlp("I argue that states built environments of conflict are material manifestat.
        span = doc[2:4]
        span.text
```

```
Out[92]: 'that states'
```

1.7 Creating a span manually

```
In [93]: from spacy.tokens import Span
         # Create a Doc object
         doc = nlp("Airstrikes continued into the early hours of Monday morning in Gaza.")
         span = Span(doc, 4, 11, label="NORP")
         span.text
```

```
Out[93]: 'early hours of Monday morning in Gaza'
```

1.8 Part-of-speech tags

```
In [94]: doc = nlp("Airstrikes continued into the early hours of Monday morning in Gaza.")
         # Coarse-grained part-of-speech tags
         [token.pos_ for token in doc]
         # Fine-grained part-of-speech tags
         [token.tag_ for token in doc]
```

```
Out[94]: ['NNS', 'VBD', 'IN', 'DT', 'JJ', 'NNS', 'IN', 'NNP', 'NN', 'IN', 'NNP', '.']
```

```
In [95]: spacy.explain("NNP")
```

```
Out[95]: 'noun, proper singular'
```

1.9 Syntactic dependencies

```
In [96]: doc = nlp("Airstrikes continued into the early hours of Monday morning in Gaza.")
         # Dependency labels
         [token.dep_ for token in doc]
         [token.head.text for token in doc]
```

```
Out[96]: ['continued',
          'continued',
          'continued',
          'hours',
          'hours',
          'into',
          'hours',
          'morning',
          'of',
          'continued',
          'in',
          'continued']
```

1.10 Named entities

```
In [97]: doc = nlp("Israel's new plan is to 'shrink,' not solve, the Palestinian conflict.")
         # Text and label of named entity span
         [(ent.text, ent.label_) for ent in doc.ents]
```

```
Out[97]: [('Israel', 'GPE'), ('Palestinian', 'NORP')]
```

```
In [98]: #spacy.explain("NORP") --> 'Nationalities or religious or political groups'
        #spacy.explain("GPE") --> 'Countries, cities, states'
```

1.11 Syntax iterators – Sentences

```
In [99]: doc = nlp("This is more than just eavesdropping, its terrifying. The spyware takes control of the phone. It can make calls to anybody, send messages and it can download content, Aboudi told me.")
        # doc.sents is a generator that yields sentence spans
        [sent.text for sent in doc.sents]
```

```
Out[99]: ['This is more than just eavesdropping, its terrifying.',
          'The spyware takes complete control over the phone.',
          'It can make calls to anybody, send messages and it can download content, Aboudi told me.']
```

1.12 Base noun phrases

```
In [100]: doc = nlp("Ghassan Halaika, a Jerusalem-based researcher with Al Haq, recently noticed some strange things in his phone.")
        # doc.noun_chunks is a generator that yields spans
        [chunk.text for chunk in doc.noun_chunks]
```

```
Out[100]: ['Ghassan Halaika',
           'a Jerusalem-based researcher',
           'Al Haq',
           'strange things',
           'his phone']
```

1.13 Label Explanations

```
In [101]: #Exploring label explanations in spacy
        spacy.explain("NORP")
        # 'Nationalities or religious or political groups'
        # spacy.explain("GPE")
        # 'Countries, cities, states'
```

```
Out[101]: 'Nationalities or religious or political groups'
```

2 Visualizing dependencies from various news sources

2.1 Importing displacy for visualizing sentence structure

```
In [102]: from spacy import displacy
```

2.2 Sentence from Al Jazeera article

Palestinian rights activists defiant over Israeli spyware hacks

```
In [103]: doc = nlp("The rights groups deny any links to the PFLP and Israel has failed to publish details of the spyware hacks.")
        displacy.render(doc, style="dep")
```

<IPython.core.display.HTML object>

2.3 Sentence from NYT article

Conflict Spirals Across Israel and the Palestinian Territories

```
In [104]: doc = nlp("An American envoy landed in Israel for cease-fire talks with Palestinians  
displacy.render(doc, style="dep")
```

<IPython.core.display.HTML object>

2.4 Visualize named entities

```
In [105]: doc = nlp("Internment, Torture and Pro-government Militia in Northern Ireland (with S  
displacy.render(doc, style="ent")
```

<IPython.core.display.HTML object>