



© 2020 Trilogy Education Services, a 2U, Inc. brand. All Rights Reserved.

Class Objectives

By the end of the class, you will be able to:

- Understand spaCy capabilities and where to find documentation.
- Be able to use POS tagged text to extract specific words.
- Use dependency parsed text to extract descriptors.
- Extract specific types of entities from text.
- Correlate text features to real-world series like stock prices.
- Create a dashboard from NLP sentiment features.



spaCy

3

spaCy

spaCy

- Core functions depend on language models learned from tagged text
- Fast and flexible
- Designed specifically for production use

NLTK

- Core functions depend on language models learned from programmed rules
- Accurate
- Intended for educational and prototyping purposes

4

spaCy

We will be using spaCy for:



Part of speech tagging



Named entity recognition



Dependency parsing



These tasks are more suitable for **model-based solutions** because they are complex and depend highly on context.

5

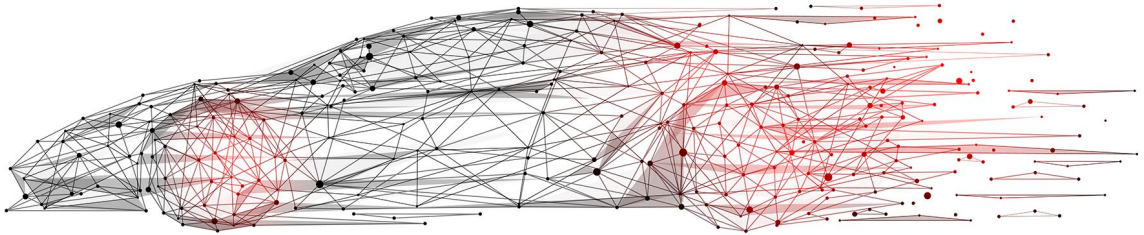


spaCy also provides tools for tasks like **tokenization** and **lemmatization**, which we've already learned with NLTK, and creating word vectors.

6

spaCy

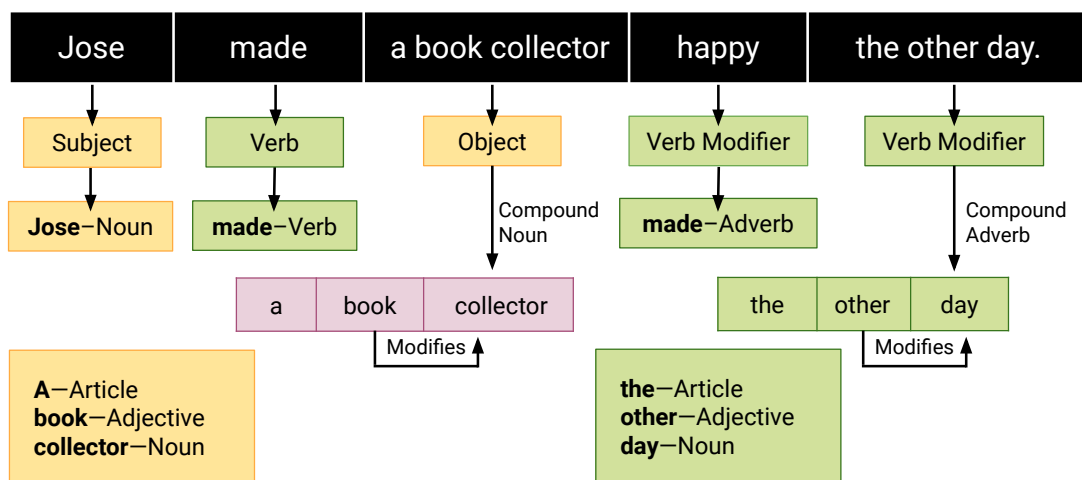
In comparison to NLTK, spaCy's language models trades off accuracy for speed, so if the corpus is large then you may prefer a simpler, rule-based solution.



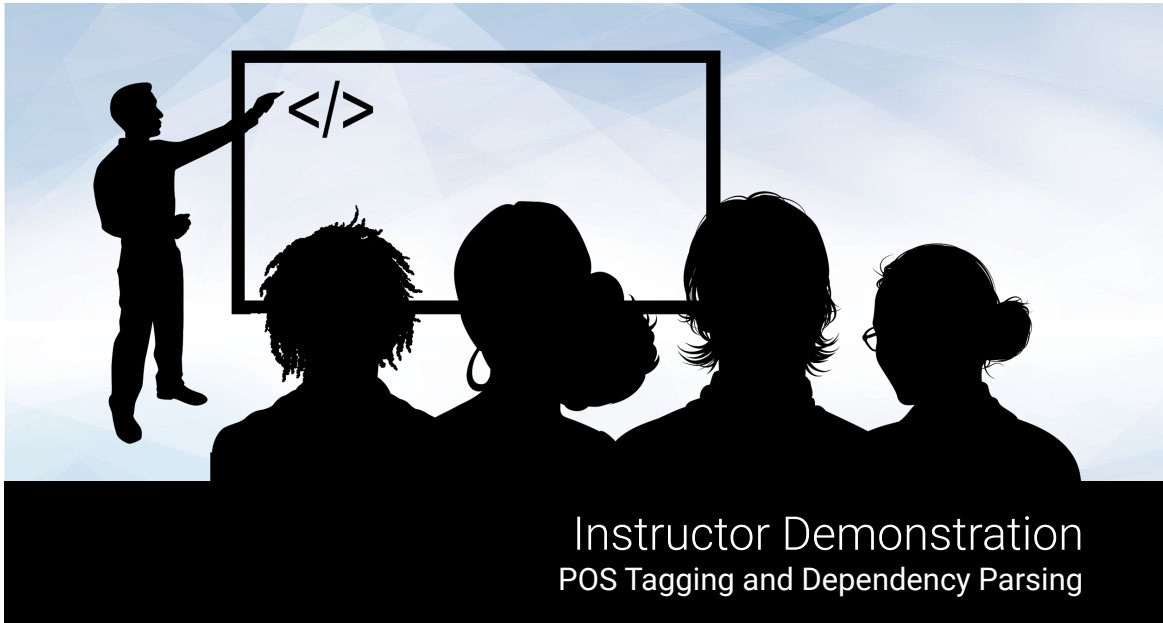
7

Part of speech tagging

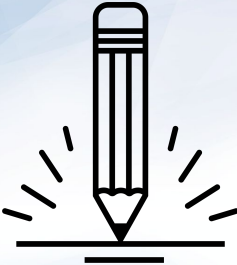
Categorizing each word in a sentence by its grammatical role in the sentence.



8




9

An icon of a pencil with motion lines, indicating writing or drawing.

Activity: Describing America

In this activity you will use the inaugural address corpus from NLTK and spacy's parsing and tagging modules to analyze the text that presidents have used to describe America.

Suggested Time:
15 minutes

An icon of an alarm clock, representing time or duration.

10



Time's Up! Let's Review.

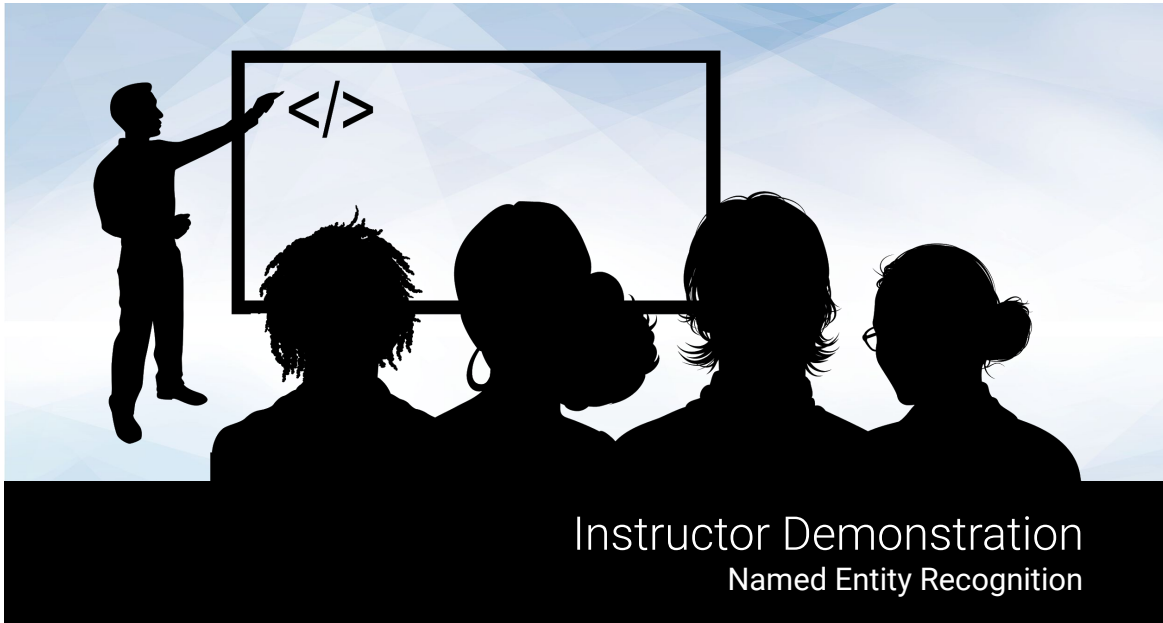
11

Named Entity Recognition

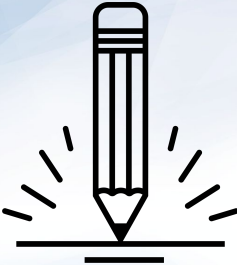
Extracting named entities, which include proper nouns and other specific types of nouns such as currencies, from a text.

US GPE unveils world's most powerful supercomputer, beats **China GPE**. The **US GPE** has unveiled the world's most powerful supercomputer called 'Summit', beating the previous record-holder **China GPE**'s **Sunway TaighuLight ORG**. With a peak performance of **200,000 CARDINAL** trillion calculations per **Second ORDINAL**, it is over twice as fast as **Sunway TaighuLight ORG**, which is capable of **93,000 CARDINAL** Trillion calculations per second. Summit has **4,608 CARDINAL** servers, which reportedly take up the size of **two CARDINAL** tennis courts.

12




13

An icon of a pencil with motion lines around its tip, positioned over a horizontal line.

Activity: NER Clouds

In this activity you will extract named entities of their own choosing from the Reuters corpus and build a wordcloud from these entities.

Suggested Time:
15 minutes

An icon of an alarm clock.

14



Time's Up! Let's Review.

15



Text as Feature

17

Tools and Techniques

Tools and techniques used to create numerical features (structured data) from text (unstructured data):

Tools

- NLTK
- wordcloud
- spaCy



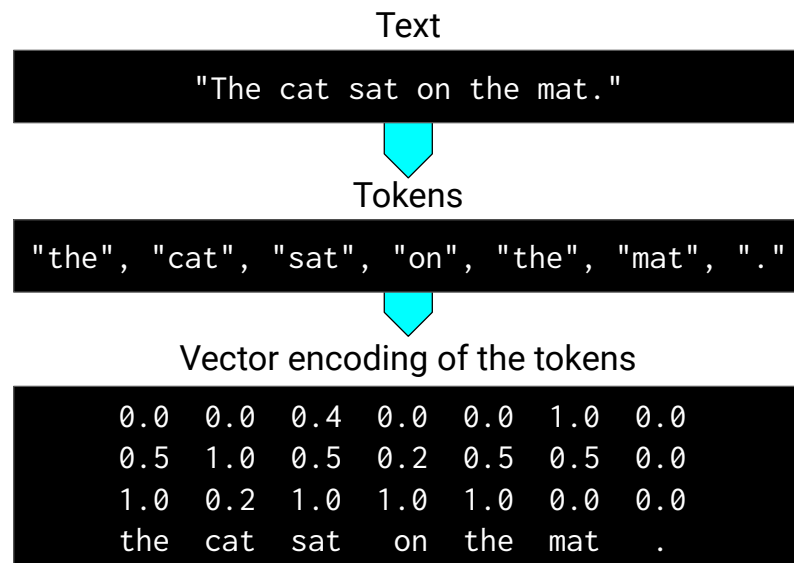
Techniques

- preprocessing
- tokenizing text
- lemmatizing text
- aggregating word counts
- creating n-grams
- normalizing to tf-idf weights
- sentiment analysis
- parsing and pos-tagging text
- named entity recognition

18

Text as Feature

In order to use this data for classification or prediction, we need to make them features—numerical representations of unstructured text.



19

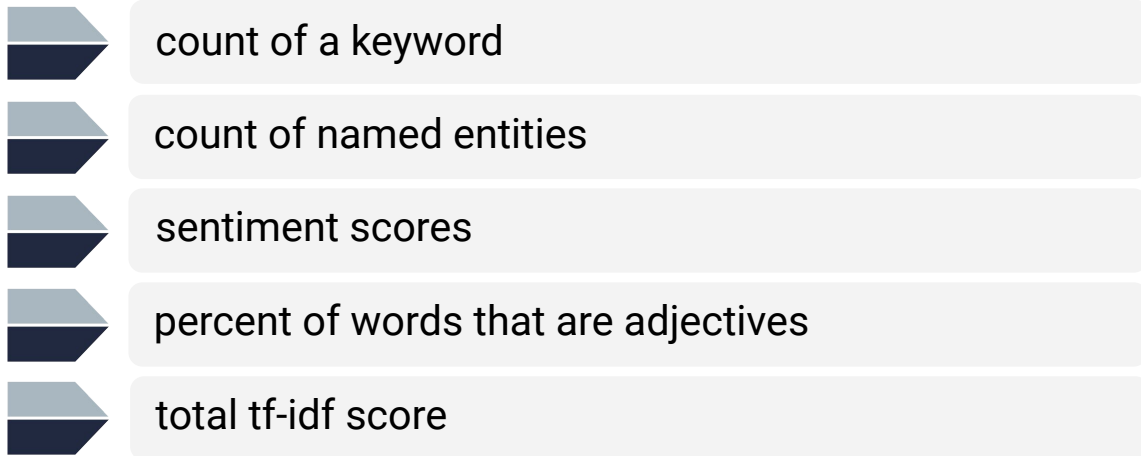


What are some examples of features that can be created from text documents?

20

Text as Feature

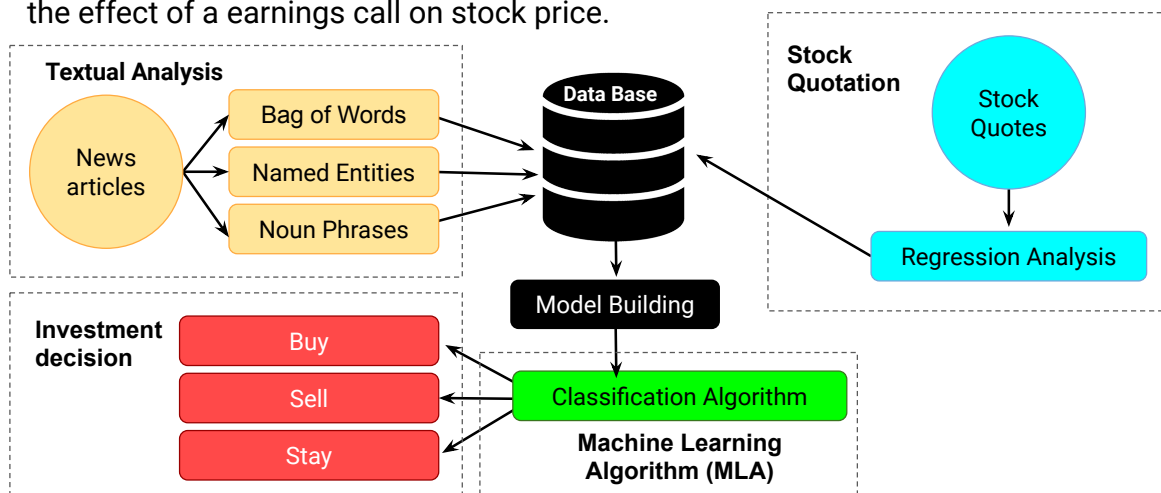
Examples:



21

Text as Feature

These features can then be used to classify a document to a category or predict the effect of a earnings call on stock price.



22



Activity: Correlating Returns

In this activity you will create a sentiment index from newsapi headlines and correlated it to S&P 500 daily returns, looking for text topic that generates the highest correlation.

Suggested Time:
15 minutes



23



Time's Up! Let's Review.

24

Machine Learning Review

Questions?