

# NLTK II

Marina Sedinkina  
- Folien von Desislava Zhekova -

CIS, LMU  
`marina.sedinkina@campus.lmu.de`

December 11, 2018

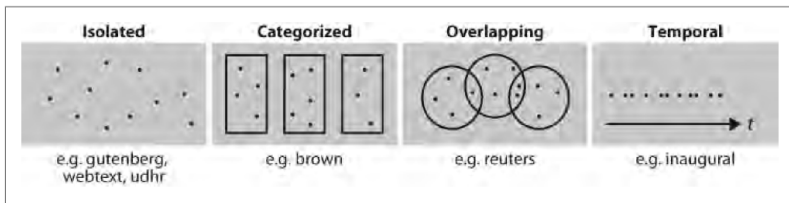
# Outline

- 1 Corpora
- 2 Preprocessing
  - Normalization
- 3 spaCy
  - Tokenization with spaCy
- 4 References

# NLP and Corpora

- **Corpora** are large collections of linguistic data
- designed to achieve specific goal in NLP: data should provide best representation for the task. Such tasks are for example:
  - word sense disambiguation:
  - sentiment analysis
  - text categorization
  - part of speech tagging

# Corpora Structure



# Corpora

- When the `nltk.corpus` module is imported, it automatically creates a set of corpus reader instances that can be used to access the corpora in the NLTK data distribution
- The corpus reader classes may be of several subtypes:  
`CategorizedTaggedCorpusReader`,  
`BracketParseCorpusReader`,  
`WordListCorpusReader`, `PlaintextCorpusReader`  
...

```
1 from nltk.corpus import brown
2
3 print(brown)
4
5 # prints
6 # <CategorizedTaggedCorpusReader in '... /corpora/brown' (not
   # loaded yet)>
```

# Corpus functions

Objects of type `CorpusReader` support the following functions:

Example	Description
<code>fileids()</code>	The files of the corpus
<code>fileids([categories])</code>	The files of the corpus corresponding to these categories
<code>categories()</code>	The categories of the corpus
<code>categories([fileids])</code>	The categories of the corpus corresponding to these files
<code>raw()</code>	The raw content of the corpus
<code>raw(fileids=[f1, f2, f3])</code>	The raw content of the specified files
<code>raw(categories=[c1, c2])</code>	The raw content of the specified categories
<code>words()</code>	The words of the whole corpus
<code>words(fileids=[f1, f2, f3])</code>	The words of the specified fileids
<code>words(categories=[c1, c2])</code>	The words of the specified categories

## Corpus functions

<code>sents()</code>	The sentences of the specified categories
<code>sents(fileids=[f1, f2, f3])</code>	The sentences of the specified fileids
<code>sents(categories=[c1, c2])</code>	The sentences of the specified categories
<code>abspath(fileid)</code>	The location of the given file on disk
<code>encoding(fileid)</code>	The encoding of the file (if known)
<code>open(fileid)</code>	Open a stream for reading the given corpus file
<code>root()</code>	The path to the root of locally installed corpus
<code>readme()</code>	The contents of the README file of the corpus

---

# Brown Corpus

```
1 from nltk.corpus import brown
2
3 print(brown.categories())
4 # ["adventure", "belles_lettres", "editorial", "
   #  fiction", "government", "hobbies", "humor", "
   #  learned", "lore", "mystery", "news", "religion",
   #  "reviews", "romance", "science_fiction"]
```



## Brown Corpus

```
1 from nltk.corpus import brown
2
3 print(brown.categories())
4 print(brown.words(categories="news"))
5 # ["The", "Fulton", "County", "Grand", "Jury", "said",
    #  ..., ...]
```

Access the list of words, but restrict them to a specific category.

# Brown Corpus

```
1 from nltk.corpus import brown
2
3 print(brown.categories())
4 print(brown.words(categories="news"))
5
6 print(brown.words(fileids=["cg22"]))
7 # ["Does", "our", "society", "have", "a", "runaway",
   #  ", ", ... ]
```

Access the list of words, but restrict them to a specific file.

# Brown Corpus

```
1 from nltk.corpus import brown
2
3 print(brown.categories())
4 print(brown.words(categories="news"))
5 print(brown.words(fileids=["cg22"]))
6
7 print(brown.sents(categories=["news", "editorial", "
    reviews"]))
8 # [[ "The", "Fulton", "County" ... ], [ "The", "jury", "
    further" ... ], ... ]
```

Access the list of sentences, but restrict them to a given list of categories.


# Brown Corpus

We can compare genres in their usage of modal verbs:

```
1 import nltk
2 from nltk.corpus import brown
3
4 news_text = brown.words(categories="news")
5 fdist = nltk.FreqDist([w.lower() for w in news_text])
6 modals = ["can", "could", "may", "might", "must", "will"]
7
8 for m in modals:
9     print(m + ":", fdist[m])
10
11 # can: 94
12 # could: 87
13 # may: 93
14 # might: 38
15 # must: 59
```

# Gutenberg Corpus

NLTK includes a small selection of texts from the Project Gutenberg electronic text archive, which contains more than 50 000 free electronic books, hosted at `http://www.gutenberg.org`.



The screenshot shows the Project Gutenberg website interface. On the left is a navigation sidebar with sections for 'search book catalog', 'search website', and 'donate'. The main content area features a header 'Free ebooks - Project Gutenberg', a search bar, and a list of 'Some of Our Latest Books' with book covers. A 'Welcome' message is visible at the bottom of the main content area.

**Project Gutenberg**

search book catalog

search website

donate  
Project Gutenberg needs your donation!

Free ebooks - Project Gutenberg

From Project Gutenberg, the first producer of free ebooks.

Book search · Book categories · Browse catalog · Mobile site · Report errors · Terms of use

New Kindle Fire Review

Before you buy: [Read our Webmaster's review of the new Kindle Fire.](#)

Some of Our Latest Books

Welcome

# Gutenberg Corpus

```
1 >>> import nltk
2 >>> nltk.corpus.gutenberg.fileids()
3 ["austen-emma.txt", "austen-persuasion.txt", "austen-
  sense.txt", "bible-kjv.txt", "blake-poems.txt", "
  bryant-stories.txt", "burgess-busterbrown.txt", "
  carroll-alice.txt", "chesterton-ball.txt", "
  chesterton-brown.txt", "chesterton-thursday.txt",
  "edgeworth-parents.txt", "melville-moby_dick.txt",
  "milton-paradise.txt", "shakespeare-caesar.txt",
  "shakespeare-hamlet.txt", "shakespeare-macbeth
  .txt", "whitman-leaves.txt"]
```

# Inaugural Address Corpus

Time dimension property:

```
1 >>> from nltk.corpus import inaugural
2 >>> inaugural.fileids()
3 ["1789–Washington.txt", "1793–Washington.txt", "1797–
   Adams.txt", ... ]
4 >>> [fileid[:4] for fileid in inaugural.fileids()]
5 ["1789", "1793", "1797", "1801", "1805", "1809", "
   1813", "1817", "1821", ... ]
```

# Gutenberg Corpus

Extract statistics about the corpus:

```
1 for fileid in gutenbergs.fileids():
2
3     num_vocab=len(set([w.lower() for w in gutenbergs.words(fileid)]))
4
5     x1= len(gutenbergs.raw(fileid))/len(gutenbergs.words(fileid))
6     x2= len(gutenbergs.words(fileid))/len(gutenbergs.sents(fileid))
7     x3= len(gutenbergs.words(fileid))/num_vocab
```

???

What is calculated here?



# Annotated Text Corpora

Many text corpora contain linguistic annotations:

- part-of-speech tags
- named entities
- syntactic structures
- semantic roles

# Annotated Text Corpora

```
#begin document <document ID>
<sentence>

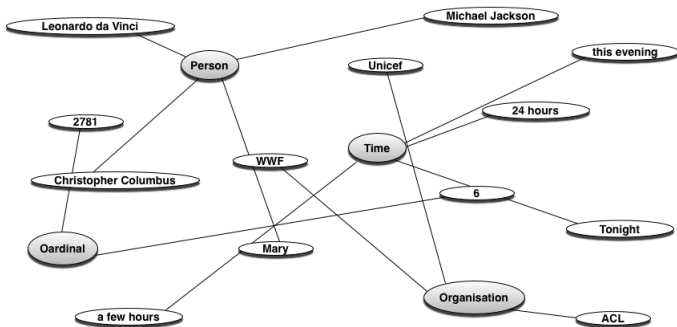
<sentence>
...
<sentence>

#end document <document ID>
...
#begin document <document ID>
<sentence>

<sentence>
...
<sentence>

#end document <document ID>
```

# Annotated Text Corpora



## Annotated Text Corpora

Word#	Word	POS	ParseBit	PredLemma	PFID	WS	SA	NE	PredArgs	PredArgs	Coref
0	It	PRP	(TOP(S(NP*))	-	-	-	Speaker#1 *	*	(ARG1*)	(22)	
1	is	VBZ	(VP*	-	03	-	Speaker#1 *	(V*)	*	-	
2	composed	VBN	(VP*	-	01	2	Speaker#1 *	*	(V*)	-	
3	of	IN	(PP*	-	-	-	Speaker#1 *	*	(ARG2*	-	
4	a	DT	(NP(NP*	-	-	-	Speaker#1 *	*	*	(24)	
5	primary	JJ	*	-	-	-	Speaker#1 *	*	*	-	
6	stele	NN	*)	-	-	-	Speaker#1 *	*	*	(24)	
7	,	,	*	-	-	-	Speaker#1 *	*	*	-	
8	secondary	JJ	(NP*	-	-	-	Speaker#1 *	*	*	(13)	
9	steles	NNS	*)	-	-	-	Speaker#1 *	*	*	(13)	
10	,	,	*	-	-	-	Speaker#1 *	*	*	-	
11	a	DT	(NP*	-	-	-	Speaker#1 *	*	*	-	
12	huge	JJ	*	-	-	-	Speaker#1 *	*	*	-	
13	round	NN	*	-	-	-	Speaker#1 *	*	*	-	
14	sculpture	NN	(NML(NML*)	-	-	-	Speaker#1 *	*	*	-	
15	and	CC	*	-	-	-	Speaker#1 *	*	*	-	
16	beacon	NN	(NML*	-	-	-	Speaker#1 *	*	*	-	
17	tower	NN	*)	-	-	-	Speaker#1 *	*	*	-	
18	,	,	*	-	-	-	Speaker#1 *	*	*	-	
19	and	CC	*	-	-	-	Speaker#1 *	*	*	-	
20	the	DT	(NP*	-	-	-	Speaker#1 (WORK_OF_ART*	*	*	-	
21	Great	NNP	*	-	-	-	Speaker#1 *	*	*	-	
22	Wall	NNP	*)	-	-	-	Speaker#1 *)	*	*	-	
23	,	,	*	-	-	-	Speaker#1 *	*	*	-	
24	among	IN	(PP*	-	-	-	Speaker#1 *	*	*	-	
25	other	JJ	(NP*	-	-	-	Speaker#1 *	*	*	-	
26	things	NNS	*)	-	-	-	Speaker#1 *)	*	*	-	
27	.	.	*)	-	-	-	Speaker#1 *)	*	*	-	

# Annotated Text Corpora

download required corpus via `nltk.download()`

Corpus	Compiler	Contents
Brown Corpus	Francis, Kucera	15 genres, 1.15M words, tagged, categorized
CESS Treebanks	CLiC-UB	1M words, tagged and parsed (Catalan, Spanish)
Chat-80 Data Files	Pereira & Warren	World Geographic Database
CMU Pronouncing Dictionary	CMU	127k entries
CoNLL 2000 Chunking Data	CoNLL	270k words, tagged and chunked
CoNLL 2002 Named Entity	CoNLL	700k words, POS and named entity tagged (Dutch, Spanish)
CoNLL 2007 Dependency Parsed Treebanks (selections)	CoNLL	150k words, dependency parsed (Basque, Catalan)
Dependency Treebank	Narad	Dependency parsed version of Penn Treebank sample
Floresta Treebank	Diana Santos et al.	9k sentences, tagged and parsed (Portuguese)
Gazetteer Lists	Various	Lists of cities and countries

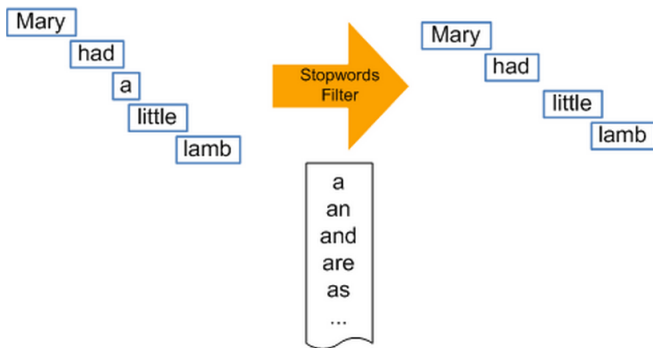
## Wordlist Corpora

Word lists are a type of lexical resources. NLTK includes some examples:

- `nltk.corpus.stopwords`
- `nltk.corpus.names`
- `nltk.corpus.swadesh`
- `nltk.corpus.words`

# Stopwords

Stopwords are high-frequency words with little lexical content such as **the, to, and**.



## Wordlists: Stopwords

```
1 >>> from nltk.corpus import stopwords
2 >>> stopwords.words("english")
3 ["a", "a's", "able", "about", "above", "according", "accordingly",
   "across", "actually", "after", "afterwards", "again", "
   against", "ain't", "all", "allow", "allows", "almost", "alone",
   "along", "already", "also", "although", "always", ... ]
```

Also available for: Danish, Dutch, English, Finnish, French, German, Hungarian, Italian, Norwegian, Portuguese, Russian, Spanish, Swedish and Turkish



## Wordlists: Names

- Names Corpus is a wordlist corpus, containing 8,000 first names categorized by gender.
- The male and female names are stored in separate files.

# Wordlists

```
1 import nltk
2
3 names = nltk.corpus.names
4 print(names.fileids())
5 # ["female.txt", "male.txt"]
6
7 female_names = names.words(names.fileids()[0])
8 male_names = names.words(names.fileids()[1])
9
10 print([w for w in male_names if w in female_names])
11 #["Abbey", "Abbie", "Abby", "Addie", "Adrian", "Adrien", "Ajay", "
    Alex", "Alexis", "Alfie", "Ali", "Alix", "Allie", "Allyn", "
    Andie", "Andrea", "Andy", "Angel", "Angie", "Ariel", "Ashley",
    ", "Aubrey", "Augustine", "Austin", "Averil", ... ]
```

## Wordlists

NLP application for which gender information would be helpful

### Anaphora Resolution:

**Adrian** drank from the cup. **He** liked the tea.

### Note

Both **he** as well as **she** will be possible solutions when Adrian is the antecedent, since this name occurs in both lists: female and male names.

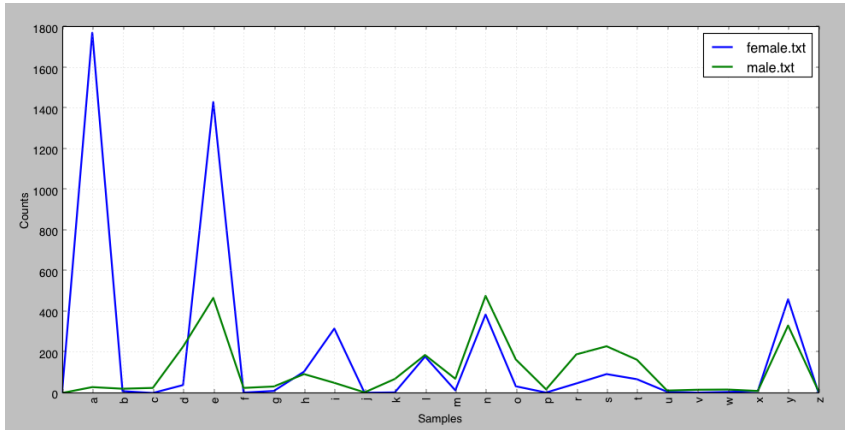
# Wordlists

```
1 import nltk
2 names = nltk.corpus.names
3
4 cfd = nltk.ConditionalFreqDist(
5     (fileid, name[-1])
6     for fileid in names.fileids()
7     for name in names.words(fileid))
```

???

What will be calculated for the conditional frequency distribution stored in `cfd`?

# Wordlists



## Wordlists: Swadesh

- comparative wordlist
- lists about 200 common words in several languages.

# Comparative Wordlists

```
1 >>> from nltk.corpus import swadesh
2 >>> swadesh.fileids()
3 ["be", "bg", "bs", "ca", "cs", "cu", "de", "en", "es", "fr", "hr",
   "it", "la", "mk", "nl", "pl", "pt", "ro", "ru", "sk", "sl",
   "sr", "sw", "uk"]
4 >>> swadesh.words("en")
5 ["I", "you (singular), thou", "he", "we", "you (plural)", "they",
   "this", "that", "here", "there", "who", "what", "where", "
   when", "how", "not", "all", "many", "some", "few", "other", "
   one", "two", "three", "four", "five", "big", "long", "wide",
   ... ]
```

## Comparative Wordlists

```
1 >>> fr2en = swadesh.entries(["fr", "en"])
2 >>> fr2en
3 [ ("je", "I"), ("tu", "vous", "you (singular), thou"), ("il", "he"),
   ... ]
4 >>> translate = dict(fr2en)
5 >>> translate["chien"]
6 "dog"
7 >>> translate["jeter"]
8 "throw"
```



## Comparative Wordlists

```
1 >>> de2en = swadesh.entries(["de", "en"]) # German-English
2 >>> es2en = swadesh.entries(["es", "en"]) # Spanish-English
3 >>> translate.update(dict(de2en))
4 >>> translate.update(dict(es2en))
5 >>> translate["Hund"] "dog"
6 >>> translate["perro"] "dog"
```

## Words Corpus

- NLTK includes some corpora that are nothing more than wordlists.
- The Words Corpus `/usr/share/dict/words` from Unix is used by some spell checkers.

```
1
2 def method_x(text):
3     text_vocab=set(w.lower() for w in text if w.isalpha())
4     english_vocab=set(w.lower() for w in nltk.corpus.words.words())
5     x=text_vocab - english_vocab
6     return sorted(x)
7
8 >>> method_x(nltk.corpus.gutenberg.words('austen-sense.txt'))
```

???

What is calculated in this function?

# Words Corpus

```
1
2 def method_x(text):
3     text_vocab=set(w.lower() for w in text if w.isalpha())
4     english_vocab=set(w.lower() for w in nltk.corpus.words.words())
5     x= text_vocab - english_vocab
6     return sorted(x)
7
8 >>> method_x(nltk.corpus.gutenberg.words('austen-sense.txt'))
9 ['abbeyland', 'abhorred', 'abilities', 'abounded', ... ]
```

# Preprocessing

Original	The boy's cars are different colors.
Tokenized	["The", "boy's", "cars", "are", "different", "colors."]
Punctuation removal	["The", "boy's", "cars", "are", "different", "colors"]
Lowecased	["the", "boy's", "cars", "are", "different", "colors"]
Stemmed	["the", "boy's", "car", "are", "differ", "color"]
Lemmatized	["the", "boy's", "car", "are", "different", "color"]
Stopword removal	["boy's", "car", "different", "color"]

# Tokenization

- Tokenization is the process of breaking raw text into its building parts: words, phrases, symbols, or other meaningful elements called tokens.
- A list of tokens is almost always the first step to any other NLP task, such as part-of-speech tagging and named entity recognition.

# Tokenization

- **token** – is an instance of a sequence of characters in some particular document that are grouped together as a useful semantic unit for processing
- **type** – is the class of all tokens containing the same character sequence

# Tokenization

- What is Token?
- Fairly trivial: chop on whitespace and throw away punctuation characters.
- Tricky cases: various uses of the apostrophe for possession and contractions?

# Tokenization

Mrs. O'Reilly said that the girls' bags from  
H&M's shop in New York aren't cheap.

Mrs.	"Mrs."; "Mrs", "."
O'Reilly	"O'Reilly"; "OReilly"; "O'", "Reilly"; "O", "'", "Reilly";
aren't	"aren't"; "arent"; "are", "n't"; "aren", "t"
...	...



# Tokenization

???

Tokenize manually the following sentence. How many tokens do you get?

Mrs. O'Reilly said that the girls' bags from H&M's shop in New York aren't cheap.

# Tokenization

???

Tokenize manually the following sentence:

Mrs. O'Reilly said that the girls' bags from H&M's shop in New York aren't cheap.

Answer

NLTK returns the following 20 tokens:

```
["Mrs.", "O'Reilly", "said", "that", "the",  
"girls", "'", "bags", "from", "H", "&", "M",  
"'s", "shop", "in", "New", "York", "are",  
"n't", "cheap", "."]
```

# Tokenization

Most decisions need to be met depending on the language at hand.  
Some problematic cases for English include:

- hyphenation – *ex-wife, Cooper-Hofstadter, the bazinga-him-again maneuver*
- internal white spaces – *New York, +49 89 21809719, January 1, 1995, San Francisco-Los Angeles*
- apostrophe – *O'Reilly, aren't*
- other cases – *H&M's*

# Sentence Segmentation

Tokenization can be approached at any level:

- word segmentation
- sentence segmentation
- paragraph segmentation
- other elements of the text

# Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1 >>> from nltk import word_tokenize ,  
      wordpunct_tokenize  
2 >>> s = "Good muffins cost $3.88\nin New York.  
      Please buy me\n      two of them.\n\nThanks."  
3 >>> word_tokenize(s)  
4 ['Good', 'muffins', 'cost', '$', '3.88', 'in', 'New',  
   'York', '.', 'Please', 'buy', 'me', 'two', 'of',  
   'them', '.', 'Thanks', '.']  
5 >>> wordpunct_tokenize(s)  
6 ['Good', 'muffins', 'cost', '$', '3', '.', '88', 'in',  
   'New', 'York', '.', 'Please', 'buy', 'me', 'two',  
   'of', 'them', '.', 'Thanks', '.']
```

# Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1 >>> from nltk.tokenize import *
2 >>> # same as s.split():
3 >>> WhitespaceTokenizer().tokenize(s)
4 ['Good', 'muffins', 'cost', '$3.88', 'in', 'New', 'York.', 'Please', 'buy', 'me', 'two', 'of', 'them.', 'Thanks.']
5 >>> # same as s.split(" "):
6 >>> SpaceTokenizer().tokenize(s)
7 ['Good', 'muffins', 'cost', '$3.88\nin', 'New', 'York.', '', 'Please', 'buy', 'me\ntwo', 'of', 'them.\n\nThanks.']
```

# Segmentation

NLTK comes with a whole bunch of tokenization possibilities:

```
1 >>> # same as s.split('\n'):
2 >>> LineTokenizer(blanklines='keep').tokenize(s)
3 ['Good muffins cost $3.88', 'in New York. Please buy
   me', 'two of them.', '', 'Thanks. ']
4 >>> # same as [l for l in s.split('\n') if l.strip()
   ]:
5 >>> LineTokenizer(blanklines='discard').tokenize(s)
6 ['Good muffins cost $3.88', 'in New York. Please buy
   me', 'two of them.', 'Thanks. ']
7 >>> # same as s.split('\t'):
8 >>> TabTokenizer().tokenize('a\tb c\n\t d')
9 ['a', 'b c\n', ' d']
```

# Segmentation

NLTK PunktSentenceTokenizer: divides a text into a list of sentences

```
1 >>> import nltk.data
2 >>> text = "Punkt knows that the periods in Mr. Smith and
  Punkt knows that the periods in Mr. Smith and
  Johann S. Bach do not mark sentence boundaries. And
  sometimes sentences can start with non-capitalized
  words. i is a good variable name."
3 >>> sent_detector = nltk.data.load('tokenizers/punkt/
  english.pickle')
4 >>> print '\n———\n'.join(sent_detector.tokenize(text.
  strip()))
5 # Punkt knows that the periods in Mr. Smith and Johann S.
  Bach do not mark sentence boundaries.
6 # ——
7 # And sometimes sentences can start with non-capitalized
  words.
8 # ——
9 # i is a good variable name.
```



# Normalization

Once the text has been segmented into its tokens (paragraphs, sentences, words), most NLP pipelines do a number of other basic procedures for text normalization, e.g.:

- lowercasing
- stemming
- lemmatization
- stopword removal

# Lowercasing

Lowercasing:

```
1 import nltk
2
3 string = "The boy's cars are different colors."
4 tokens = nltk.word_tokenize(string)
5 lower = [x.lower() for x in tokens]
6 print(" ".join(lower))
7
8 # prints
9 # the boy 's cars are different colors .
```

- Often, however, instead of working with all word forms, we would like to extract and work with their base forms (e.g. lemmas or stems)
- Thus with **stemming** and **lemmatization** we aim to reduce inflectional (and sometimes derivational) forms to their base forms.

# Stemming

**Stemming:** removing morphological affixes from words, leaving only the word stem.

```
1 import nltk
2
3 string = "The boy's cars are different colors."
4 tokens = nltk.word_tokenize(string)
5 lower = [x.lower() for x in tokens]
6 stemmed = [stem(x) for x in lower]
7 print(" ".join(stemmed))
8
9 def stem(word):
10     for suffix in ["ing", "ly", "ed", "ious", "ies", "ive", "es",
11                  "s", "ment"]:
12         if word.endswith(suffix):
13             return word[:-len(suffix)]
14     return word
15 # prints
16 # the boy 's car are different color .
```

# Stemming

Stemming:

```
1 import nltk
2 import re
3
4 string = "The boy's cars are different colors."
5 tokens = nltk.word_tokenize(string)
6 lower = [x.lower() for x in tokens]
7 stemmed = [stem(x) for x in lower]
8 print(" ".join(stemmed))
9
10 def stem(word):
11     regexp = r"^(.*?)(ing|ly|ed|ious|ies|ive|es|s|ment)?$"
12     stem, suffix = re.findall(regexp, word)[0]
13     return stem
14
15 # prints
16 # the boy 's car are different color .
```

# Stemming

NLTK's stemmers:

- **Porter Stemmer** is the oldest stemming algorithm supported in NLTK, originally published in 1979.

`http:`

`//www.tartarus.org/~martin/PorterStemmer/`

- **Lancaster Stemmer** is much newer, published in 1990, and is more aggressive than the Porter stemming algorithm.
- **Snowball stemmer** currently supports several languages: Danish, Dutch, English, Finnish, French, German, Hungarian, Italian, Norwegian, Porter, Portuguese, Romanian, Russian, Spanish, Swedish.
- **Snowball stemmer**: slightly faster computation time than porter.

# Stemming

NLTK's stemmers:

```
1 import nltk
2
3 string = "The boy's cars are different colors."
4 tokens = nltk.word_tokenize(string)
5 lower = [x.lower() for x in tokens]
6
7 porter = nltk.PorterStemmer()
8 stemmed = [porter.stem(t) for t in lower]
9 print(" ".join(stemmed))
10 # prints
11 # the boy 's car are differ color .
12
13 lancaster = nltk.LancasterStemmer()
14 stemmed = [lancaster.stem(t) for t in lower]
15 print(" ".join(stemmed))
16 # prints
17 # the boy 's car ar diff col .
```

# Stemming

NLTK's stemmers:

```
1 import nltk
2
3 string = "The boy's cars are different colors."
4 tokens = nltk.word_tokenize(string)
5 lower = [x.lower() for x in tokens]
6
7 snowball = nltk.SnowballStemmer("english")
8 stemmed = [snowball.stem(t) for t in lower]
9 print(" ".join(stemmed))
10 # prints
11 # the boy 's car are differ color .
```



# Lemmatization

- stemming can often create non-existent words, whereas lemmas are actual words
- **NLTK WordNet Lemmatizer** uses the WordNet Database to lookup lemmas

```
1 import nltk
2 string = "The boy's cars are different colors."
3 tokens = nltk.word_tokenize(string)
4 lower = [x.lower() for x in tokens]
5 porter = nltk.PorterStemmer()
6 stemmed = [porter.stem(t) for t in lower]
7 print(" ".join(stemmed))
8 # prints the boy 's car are differ color .
9 wnl = nltk.WordNetLemmatizer()
10 lemmatized = [wnl.lemmatize(t) for t in lower]
11 print(" ".join(lemmatized))
12 # prints the boy 's car are different color .
```

# Stopword removal:

Stopword removal:

```
1 import nltk
2
3 string = "The boy's cars are different colors."
4 tokens = nltk.word_tokenize(string)
5 lower = [x.lower() for x in tokens]
6 wnl = nltk.WordNetLemmatizer()
7 lemmatized = [wnl.lemmatize(t) for t in lower]
8
9 content = [x for x in lemmatized if x not in nltk.
            corpus.stopwords.words("english")]
10 print(" ".join(content))
11 # prints
12 # boy 's car different color .
```

# Industrial-Strength Natural Language Processing

IN PYTHON

## Fastest in the world

spaCy excels at large-scale information extraction tasks. It's written from the ground up in carefully memory-managed Cython. Independent research has confirmed that spaCy is the fastest in the world. If your application needs to process entire web dumps, spaCy is the

## Get things done

spaCy is designed to help you do real work — to build real products, or gather real insights. The library respects your time, and tries to avoid wasting it. It's easy to install, and its API is simple and productive. We like to think of spaCy as the Ruby on Rails of Natural Language

## Deep learning

spaCy is the best way to prepare text for deep learning. It interoperates seamlessly with TensorFlow, PyTorch, scikit-learn, Gensim and the rest of Python's awesome AI ecosystem. With spaCy, you can easily construct linguistically sophisticated statistical models for a

- a free, open-source library for advanced **Natural Language Processing** (NLP) in Python
- you can build applications that process and "understand" large volumes of text:
  - information extraction systems
  - natural language understanding systems
  - **pre-process text**
  - deep learning
- **NLTK** - a platform for teaching and research
- **spaCY** - designed specifically for **production use**

## ADAM: Question Answering System



A question answering system that extracts answers from Wikipedia to questions posed in natural language.

## CleanNLP

A tidy data model for NLP in R

## displaCy ENT



A modern named entity visualizer

## EpiTator



Extracts case counts, resolved location/species/disease names, date ranges and more

## AllenNLP



An open-source NLP research library, built on PyTorch and spaCy

## displaCy



A modern syntactic dependency visualizer

## Dragonfire



An open-source virtual assistant for Ubuntu based Linux distributions

## ExcelCy

Excel Integration with spaCy. Training NER using XLSX from PDF, DOCX, PPT, PNG or JPG.

# spaCy Features

NAME	DESCRIPTION
<b>Tokenization</b>	Segmenting text into words, punctuations marks etc.
<b>Part-of-speech (POS) Tagging</b>	Assigning word types to tokens, like verb or noun.
<b>Dependency Parsing</b>	Assigning syntactic dependency labels, describing the relations between individual tokens, like subject or object.
<b>Lemmatization</b>	Assigning the base forms of words. For example, the lemma of "was" is "be", and the lemma of "rats" is "rat".
<b>Sentence Boundary Detection (SBD)</b>	Finding and segmenting individual sentences.
<b>Named Entity Recognition (NER)</b>	Labelling named "real-world" objects, like persons, companies or locations.
<b>Similarity</b>	Comparing words, text spans and documents and how similar they are to each other.
<b>Text Classification</b>	Assigning categories or labels to a whole document, or parts of a document.
<b>Rule-based Matching</b>	Finding sequences of tokens based on their texts and linguistic annotations, similar to regular expressions.

## Statistical models

- some of spaCy's features require **statistical models** to be loaded (e.g. to predict whether a word is a verb or a noun)
- spaCy currently offers statistical models for **8 languages**
- Components:
  - **Binary weights** for the part-of-speech tagger, dependency parser and named entity recognizer
  - **Lexical entries** in the vocabulary, i.e. words and e.g. their shape or spelling.
  - **Word vectors**, i.e. multi-dimensional meaning representations of words
  - **Configuration** options, like the language and processing pipeline settings

# Linguistic annotations

spaCy provides a variety of linguistic annotations to give you **insights into a text's grammatical structure.**

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'Apple is looking at buying U.K. startup for $1
5 billion')
6 for token in doc:
7     print(token.text, token.pos_, token.dep_)
8
9 >>> Apple PROPN nsubj
10    is VERB aux
11    looking VERB ROOT
12    at ADP prep
13    buying VERB pcomp
14    ...
```



# Part-of-speech tags and dependencies

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'Apple is looking at buying U.K. startup for $1 billion
5         ')
6 for token in doc:
7     print(token.text, token.lemma_, token.pos_, token.tag_,
8           token.dep_, token.shape_, token.is_alpha, token.is_stop)
9 >>>
10 Apple apple PROPN NNP nsubj Xxxxx True False
11 is be VERB VBZ aux xx True True
12 looking look VERB VBG ROOT xxxx True False
13 at at ADP IN prep xx True True
14 buying buy VERB VBG pcomp xxxx True False
15 U.K. u.k. PROPN NNP compound X.X. False False
16 startup startup NOUN NN dobj xxxx True False
17 for for ADP IN prep xxx True True
18 $ $ SYM $ quantmod $ False False
19 1 1 NUM CD compound d False False
20 billion billion NUM CD pobj xxxx True False
```

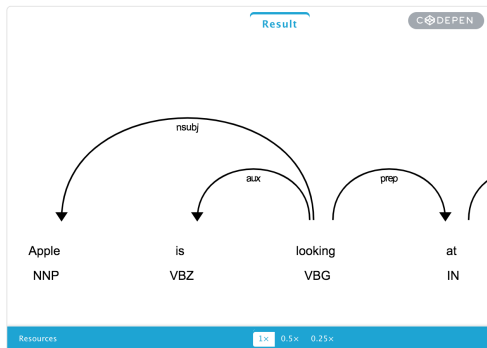
# Part-of-speech tags and dependencies

TEXT	LEMMA	POS	TAG	DEP	SHAPE	ALPHA	STOP
Apple	apple	PROPN	NNP	nsubj	Xxxxx	True	False
is	be	VERB	VBZ	aux	xx	True	True
looking	look	VERB	VBG	ROOT	xxxx	True	False
at	at	ADP	IN	prep	xx	True	True
buying	buy	VERB	VBG	pcomp	xxxx	True	False
U.K.	u.k.	PROPN	NNP	compound	X.X.	False	False
startup	startup	NOUN	NN	doobj	xxxx	True	False
for	for	ADP	IN	prep	xxx	True	True
\$	\$	SYM	\$	quantmod	\$	False	False
1	1	NUM	CD	compound	d	False	False
billion	billion	NUM	CD	pobj	xxxx	True	False

# Part-of-speech tags and dependencies

## Use built-in **displaCy** visualizer

- 1 `from spacy import displacy`
- 2 `displacy.serve(doc, style='dep')`



# Named Entities

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'Apple is looking at buying U.K. startup for $1 billion
5         ')
6 for ent in doc.ents:
7     print(ent.text, ent.start_char, ent.end_char, ent.label_)
8 >>>
9 Apple 0 5 ORG
10 U.K. 27 31 GPE
11 $1 billion 44 54 MONEY
```

# Named Entities

TEXT	START	END	LABEL	DESCRIPTION
Apple	0	5	ORG	Companies, agencies, institutions.
U.K.	27	31	GPE	Geopolitical entity, i.e. countries, cities, states.
\$1 billion	44	54	MONEY	Monetary values, including unit.

# Named Entities

## Use built-in **displaCy** visualizer

```
1 from spacy import displacy
2 displacy.serve(doc, style='ent')
```



# Word vectors and similarity

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_md') # make sure to use larger
   model! python -m spacy download en_core_web_md
4 tokens = nlp(u'dog cat banana')
5
6 for token1 in tokens:
7     for token2 in tokens:
8         print(token1.text, token2.text, token1.similarity(token2))
9 >>>
10 dog dog 1.0
11 dog cat 0.80168545
12 dog banana 0.24327643
13 cat dog 0.80168545
14 cat cat 1.0
15 cat banana 0.28154364
16 banana dog 0.24327643
17 banana cat 0.28154364
18 banana banana 1.0
```

# Word vectors and similarity

	DOG	CAT	BANANA
DOG	1.00	0.80	0.24
CAT	0.80	1.00	0.28
BANANA	0.24	0.28	1.00

## BANANA.VECTOR

```
array([ 2.02280000e-01, -7.66180009e-02,  3.70319992e-01,  
       3.28450017e-02, -4.19569999e-01,  7.20689967e-02,  
       -3.74760002e-01,  5.74599989e-02, -1.24009997e-02,  
       5.29489994e-01, -5.23800015e-01, -1.97710007e-01,  
       -3.41470003e-01,  5.33169985e-01, -2.53309999e-02,  
       1.73800007e-01,  1.67720005e-01,  8.39839995e-01,  
       5.51070012e-02,  1.05470002e-01,  3.78719985e-01,  
       2.42750004e-01,  1.47449998e-02,  5.59509993e-01,  
       1.25210002e-01,  6.75960004e-01,  3.58420014e-01])
```

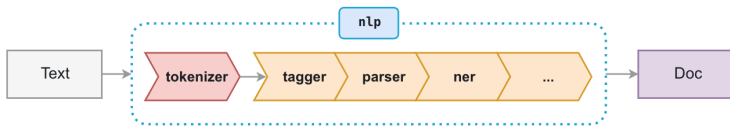


# Word vectors and similarity

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_md')
4 tokens = nlp(u'dog cat banana afskfsd')
5
6 for token in tokens:
7     print(token.text, token.has_vector, token.vector_norm,
8           token.is_oov)
9 >>>
10 dog True 7.0336733 False
11 cat True 6.6808186 False
12 banana True 6.700014 False
13 afskfsd False 0.0 True
```

`en_vectors_web_lg` includes over 1 million unique vectors

# Pipelines



"pipeline": ["tagger", "parser", "ner"]

# Vocab, hashes and lexemes

- spaCy stores data in a vocabulary **Vocab**
- encodes all string to hash values
- uses **lookup table that works in both directions** - you can look up a string to get its hash, or a hash to get its string

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'I love coffee')
5 print(doc.vocab.strings[u'coffee']) # 3197928453018144401
6 print(doc.vocab.strings[3197928453018144401]) # 'coffee'
7 >>>
8 3197928453018144401
9 coffee
```

# Vocab, hashes and lexemes

- Each entry in the vocabulary is called **Lexeme**

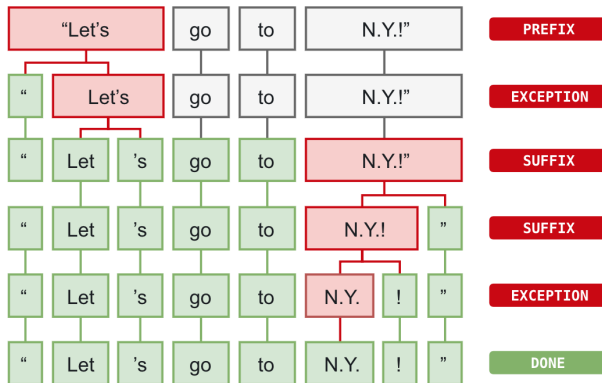
```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'I love coffee')
5 for word in doc:
6     lexeme = doc.vocab[word.text]
7     print(lexeme.text, lexeme.orth, lexeme.shape_, lexeme.prefix_,
8           lexeme.suffix_, lexeme.is_alpha, lexeme.is_digit,
9           lexeme.is_title, lexeme.lang_)
10 >>> I 4690420944186131903 X | | True False True en
11 love 3702023516439754181 xxxx | ove True False False en
12 coffee 3197928453018144401 xxxx c fee True False False en
```

# Tokenization with spaCy

```
1 import spacy
2
3 nlp = spacy.load('en_core_web_sm')
4 doc = nlp(u'Apple is looking at buying U.K. startup for $1 billion
5         ')
6 for token in doc:
7     print(token.text)
8 >>>
9     Apple
10    is
11    looking
12    at
13    buying
14    U.K.
15    ...
```

# Tokenization with spaCy

- Does the substring match a tokenizer exception rule? (U.K.)
- Can a prefix, suffix or infix be split off? (e.g. punctuation)



# Tokenization with spaCy

- Tokenizer exceptions strongly depend on the specifics of the individual language
- Global and language-specific tokenizer data is supplied via the language data in **spacy/lang**

LANGUAGE	CODE	LANGUAGE DATA	MODELS
English	en	<code>lang/en</code> </>	<a href="#">4 models</a>
German	de	<code>lang/de</code> </>	<a href="#">1 model</a>
Spanish	es	<code>lang/es</code> </>	<a href="#">2 models</a>
Portuguese	pt	<code>lang/pt</code> </>	<a href="#">1 model</a>
French	fr	<code>lang/fr</code> </>	<a href="#">2 models</a>
Italian	it	<code>lang/it</code> </>	<a href="#">1 model</a>

## Adding special case tokenization rules

- The tokenizer exceptions define special cases like "don't" in English, which needs to be split into two tokens: {ORTH: do} and {ORTH: n't, LEMMA: not}
- Later: also possibility to customize spaCy Tokenizer class

```
1 import spacy
2 from spacy.symbols import ORTH, LEMMA, POS, TAG
3
4 # add special case rule
5 special_case = [{ORTH: u'gim', LEMMA: u'give', POS: u'VERB'},
6                 {ORTH: u'me'}]
7 nlp.tokenizer.add_special_case(u'gimme', special_case)
8
9 # check new tokenization
10 print([w.text for w in nlp(u'gimme that')]) #['gim', 'me', 'that']
```



# Other Tokenizers

Tweet NLP



Carnegie Mellon

We provide a tokenizer, a part-of-speech tagger, hierarchical word clusters, and a dependency parser for tweets, along with annotated corpora and web-based annotation tools.

Contributors: [Archna Bhatia](#), [Dipanjan Das](#), [Chris Dyer](#), [Jacob Eisenstein](#), [Jeffrey Flanigan](#), [Kevin Gimpel](#), [Michael Heilman](#), [Lingpeng Kong](#), [Daniel Mills](#), [Brendan O'Connor](#), [Olutobi Owoputi](#), [Nathan Schneider](#), [Noah Smith](#), [Swabha Swayamdipta](#) and [Dani Yogatama](#).

## Quick Links

- [Part-of-speech Tagger and POS annotated data](#) - also [Ttokenizer](#): tokenizer software (part of tagger package) and [Tagging Models](#) -- [Download Link](#)
- [Tweeboparser](#) and [Tweebank](#): Dependency parser software and dependency annotated data -- [Download Link](#)
- [Documentation, annotation guidelines, and papers](#) describing this work
- [Hierarchical Twitter Word Clusters](#)

# References

- <http://www.nltk.org/book/>
- <https://github.com/nltk/nltk>
- <https://spacy.io/>