

03

Python for (Corpus) Linguists

Python Programming for Linguists

Ingo Kleiber, 2020/2021

List Comprehensions

```
numbers = [10, 20, 30]
```

```
times_ten = [n * 10 for n in numbers]
```



List Comprehension

This is equal to:

```
times_ten = []  
for n in numbers:  
    times_ten.append(n * 10)
```

List Comprehensions

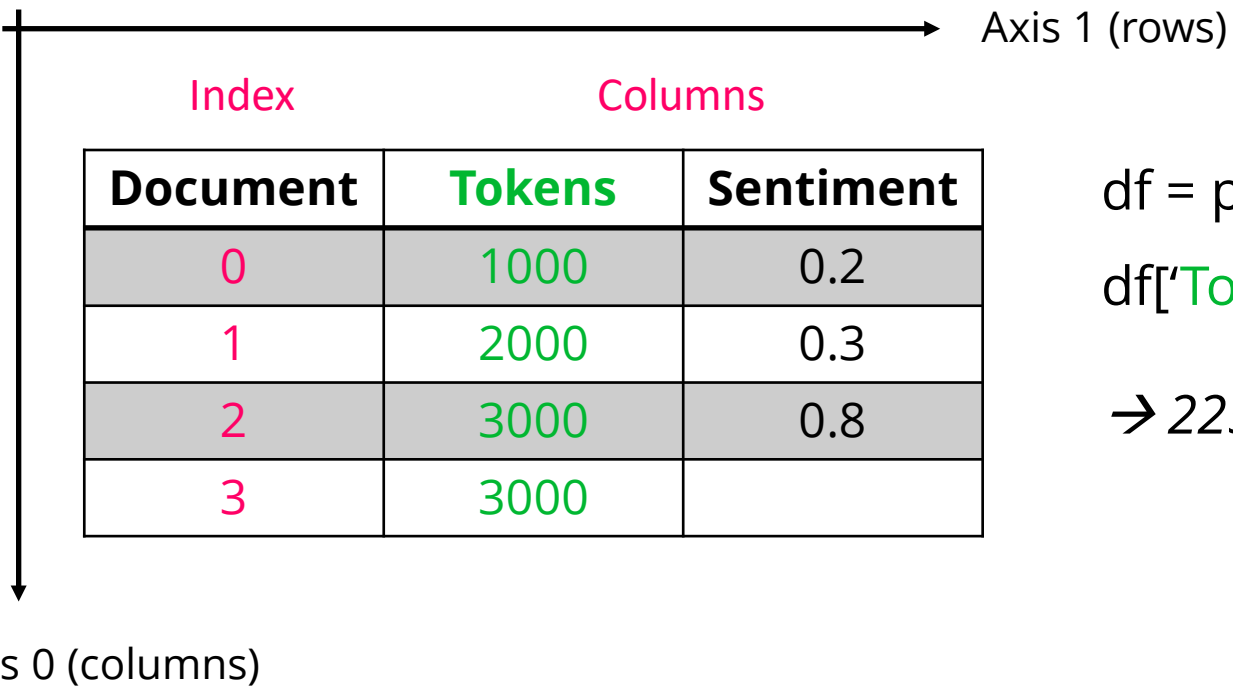
```
lol = [  
    ['A', 1],  
    ['B', 2],  
    ['C', 3]  
]
```

```
only_first_element= [n[1] for n in lol]
```

```
[1, 2, 3]
```

Pandas and DataFrames

Pandas is a very powerful data analysis and manipulation tool/library. The key component are DataFrame objects which are essentially very powerful tables.



The diagram shows a table representing a DataFrame. A horizontal arrow at the top points to the right and is labeled 'Axis 1 (rows)'. A vertical arrow on the left points downwards and is labeled 'Axis 0 (columns)'. The table has three columns: 'Document' (labeled 'Index' in pink), 'Tokens' (labeled 'Columns' in pink), and 'Sentiment'. The 'Document' column contains values 0, 1, 2, and 3, with 0, 2, and 3 highlighted in pink. The 'Tokens' column contains values 1000, 2000, 3000, and 3000, with 1000, 2000, and 3000 highlighted in green. The 'Sentiment' column contains values 0.2, 0.3, 0.8, and an empty cell. The first row (Document 0) has a gray background.

Document	Tokens	Sentiment
0	1000	0.2
1	2000	0.3
2	3000	0.8
3	3000	

```
df = pd.DataFrame(...)
```

```
df['Tokens'].mean()
```

→ 2250.0

TextDirectory Refresher

TextDirectory is a library that is useful when working with multiple text files in one directory. We can filter files based on various criteria and also run transformations (e.g., transforming the corpus to lowercase) on the texts.

```
wikipedia = textdirectory.TextDirectory(directory='data/wikipedia', autoload=True)
```

Load all files in the directory data/wikipedia.

```
wikipedia.filter_by_random_sampling(10)
```

Reduce the selection to 10 randomly sampled files

```
wikipedia.stage_transformation(['transformation_lowercase'])
```

Schedule/stage that all files (texts) are being transformed to lowercase

```
text = wikipedia.aggregate_to_memory()
```

Run the transformation and aggregate all documents into one string

Exercise 8 – Concordancer

Write a basic concordancer that can generate concordances based on a given file and a given search term. If you want to challenge yourself, try to format the concordances in KWIC format.

RegEx-Based Approach

We will use a **regular expression** to find all instances of the search term as well as 25 characters before and after (left and right).

Token-Based Approach

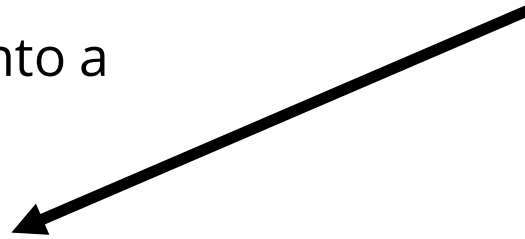
We will **tokenize** the text so that we can define a window/span in terms of tokens (words) instead of characters. We will then generate a left and right window to print KWIC concordances.

join

We can use `.join()` to turn an *iterable* into a string.

Iterables are sequences that can be iterated over using, for example, a for-loop.

These include, for example, *lists*, *sets*, and *strings*.



```
tokens = ['The', 'cat', 'is', 'grey']
```

```
s1 = ''.join(tokens)   The cat is grey
```

```
s2 = '-'.join(tokens)  The-cat-is-grey
```

Slicing Tokens

`text_tokenized` = ['the', 'cat', 'is', 'grey', 'and', 'likes', 'mice']

`search_word` = 'grey'

`lr` = 2

Let's call this index (for the search term) *id*



0	1	2	3	4	5	6
the	cat	is	grey	and	likes	mice

`text_tokenized[id - lr: id]`
1 3

`text_tokenized[id + 1: id + lr]`
4 5

Exercise 9 – N-Grams

Write a function that produces all n-grams based on a given text file and an n . *Hint:* The NLTK provides a fairly easy solution to generating n-grams.

NLTK Approach

NLTK has an `ngram` method that allows us to generate n-grams very easily.

Plain Old Python

In order to generate n-grams ourselves, we need to know that the number of n-grams will be *the number of tokens + 1 - n*. Once we know how many n-grams there are, we can create a loop that appends the n-grams, which we get by slicing the tokenized text, to a list of n-grams.

Plain Old Python

`text` = 'I really like Python, it is pretty awesome.' There are **six trigrams** here. $n = 3$

6

```
for i in range(no_of_ngrams):  
    print(tokenized_text[i:i+n])
```

i	<code>tokenized_text[i:i+n]</code>
0	['I', 'really', 'like']
1	['really', 'like', 'Python']
2	['like', 'Python', 'it']
3	['Python', 'it', 'is']
...	...

1 2 3

→ ['I', 'really', 'like', 'Python', 'it', 'is', 'pretty', 'awesome']

└──────────────────┘

1: 1+n

1: 4

Exercise 10 – Frequency Analysis

Write a script that generates a frequency table for a given text. The list should contain all types and their frequencies. *Hint:* Have a look at Python's Counter capabilities.

NLTK Approach

We can use *NLTK*'s FreqDist to generate frequency distributions of tokenized texts.

Counter Approach

We can also use Python's Counter to count all elements in an iterable (i.e., a list of tokens).

spaCy Approach

After creating a *spaCy* document (see Exercise 13), we can use the `.count_by()` method to get frequency distributions.

Counter

Counter can be used to count hashable objects (e.g., a list). The resulting **counter object** behaves a lot like a dictionary and contains the individual elements as well as their counts.

```
numbers = [1, 1, 2, 3, 3, 4]
```

```
counts = Counter(numbers)
```

```
counts[1] → 2
```

```
counts.most_common(2) → [(1, 2), (3, 2)]
```

Exercise 11 – Computing Basic Statistics

Write a script that generates the following statistics for a given search term and a set of text files (a corpus): The absolute and relative frequencies; the mean frequency; the standard deviation. Also try to plot the frequency distribution across files.

Basic Approach

We define two functions for getting the absolute and relative frequencies of a given text. Then we are using a third function to generate frequencies for a number of texts which we will store in a list. Finally, we can use Python's statistics functions to get the required statistics.

Pandas DataFrame Approach

After getting the **vocabulary** of the corpus, we use one of the functions from above to populate two frequency tables. Then we create *Pandas* DataFrames from these tables.

Vocabulary

In NLP it is very common to store the **vocabulary** (essentially a list of types) in a data structure separate from everything else. Aside from some other benefits, this avoids duplication and reduces memory cost. *Here's a very simplistic example:*

Index	0	1	2	3	4
vocabulary =	['the', 'grey', 'cat', 'is', 'black']				

v = vocabulary

sentence = [v[0], v[2], v[3], v[4]]

Lists and Sets

Sets, in the mathematical sense, are well-defined collections of distinct elements.

```
list_with_duplicate = ['A', 'B', 'B', 'C']  
s = set(list_with_duplicate)    → {'A', 'B', 'C'}
```

In Python, sets are **unordered** and only **contain unique elements**.

While sets can be used for many things (especially when leveraging set theory), we will simply use them to turn a *list of tokens* (containing duplicates) into a *set of types*.

Vocabulary

In NLP it is very common to store the **vocabulary** (essentially a list of types) in a data structure separate from everything else. Aside from some other benefits, this avoids duplication and reduces memory cost. *Here's a very simplistic example:*

Index	0	1	2	3	4
vocabulary =	['the', 'grey', 'cat', 'is', 'black']				

We can use the same vocabulary to index frequency tables.

document = 'The cat is grey. The cat is **black**'

frequencies = {0: 2, 1: 1, 2: 2, 3: 2, 4: 1}

Type <i>vocabulary[i]</i>	Frequency	
0	2	<i>the</i>
1	1	<i>grey</i>
2	2	<i>cat</i>
3	2	<i>is</i>
4	1	<i>black</i>

Exercise 12 – NLTK Stemming, Lemmatization, and WordNet

Use NLTK to stem and lemmatize the following words. Use the PorterStemmer, the LancasterStemmer, and the WordNetLemmatizer and compare your results. What are the pros and cons of these approaches?

```
words = ['connection', 'become', 'caring', 'are', 'women', 'driving']
```

Of course, feel free to add more examples! Since you already have WordNet, try to find the synonyms for *fantastic* using WordNet.

Stemming and Lemmatizing

We are using *NLTK* to compare three stemmers and/or lemmatizers. After looking at them qualitatively, we are testing how fast they can lemmatize a large number of words.

WordNet Synsets

We are using *NLTK* to access *WordNet* data. More precisely, we are accessing the synsets for *fantastic* in order to find possible synonyms.

Exercise 13 – spaCy Tagging

Use spaCy to automatically tag/annotate a text file of your choice for PoS, NERs, and Universal Dependencies.

Here we are using *spaCy* and a small **language model** (*en_core_web_sm*) to tag a given text. After creating a *spaCy* document – using the model – we can loop over the tokens (and entities) to access their tags.

Using *displaCy*, *spaCy*'s visualizer library, we can also generate graphs for the dependencies.

Exercise 14 – Parsing XML

Write a function that allows you to extract all elements with a given attribute from an XML file.
For example, the function should be able to produce the following output for the file `data/xml/bnc_style.xml` and the attribute `pos= "VERB"`: have, bought

RegEx-Based Approach

We are using a rather simple regular expression to find XML elements that contain the desired attribute and value. This solution, while being very straightforward, is not very robust if, for example, the underlying XML changes slightly.

Parsing Approach (LXML)

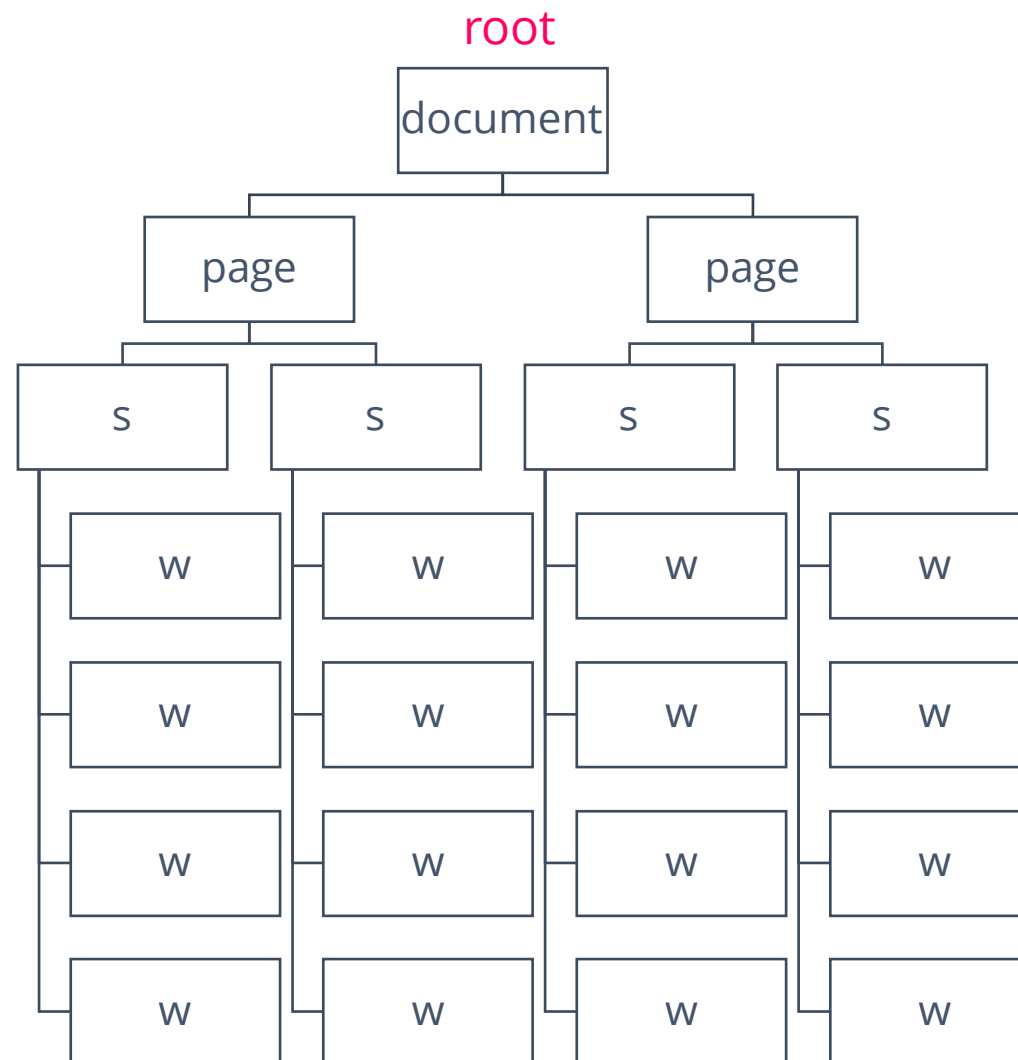
Here we are using an XML library (*LXML*) to parse and then navigate the XML structure/tree. We can also use **XPath** to navigate the document comfortably.

```

1 <document>
2   <page pg_nr="1">
3     <s>
4       <w pos="determiner">The</w>
5       <w pos="noun">flower</w>
6       <w pos="verb">was</w>
7       <w pos="adjective">red.</w>
8     </s>
9     <s>
10      <w pos="pronoun">It</w>
11      <w pos="verb">smelled</w>
12      <w pos="preposition">of</w>
13      <w pos="noun">summer.</w>
14    </s>
15  </page>
16  <page pg_nr="2">
17    <s>
18      <w pos="pronoun">She</w>
19      <w pos="verb">enjoyed</w>
20      <w pos="det">the</w>
21      <w pos="noun">trip.</w>
22    </s>
23    <s>
24      <w pos="pronoun">They</w>
25      <w pos="verb">took</w>
26      <w pos="det">a</w>
27      <w pos="noun">bus.</w>
28    </s>
29  </page>
30 </document>

```

XML as a Tree



XML XPath

XPath is a query language used for selecting nodes in XML documents.

`/page[@pg_nr='2']/s[2]/w[1]`

/ Select from the root node
@ Select attribute

```
1 <document>
2   <page pg_nr="1">
3     <s>
4       <w pos="determiner">The</w>
5       <w pos="noun">flower</w>
6       <w pos="verb">was</w>
7       <w pos="adjective">red.</w>
8     </s>
9     <s>
10      <w pos="pronoun">It</w>
11      <w pos="verb">smelled</w>
12      <w pos="preposition">of</w>
13      <w pos="noun">summer.</w>
14    </s>
15  </page>
16  <page pg_nr="2">
17    <s>
18      <w pos="pronoun">She</w>
19      <w pos="verb">enjoyed</w>
20      <w pos="det">the</w>
21      <w pos="noun">trip.</w>
22    </s>
23    <s>
24      <w pos="pronoun">They</w>
25      <w pos="verb">took</w>
26      <w pos="det">a</w>
27      <w pos="noun">bus.</w>
28    </s>
29  </page>
30 </document>
```

s[2]
w[1]

`/page[@pg_nr='2']`

`/page[@pg_nr='2']/s[2]`

`/page[@pg_nr='2']/s[2]/w[1]`

Exercise 15 – Web Scraping

Write a function that scrapes the text from a given website. The function should take a URL as its input and return the text present on the given website (e.g., Wikipedia). If you want to challenge yourself even further, try to remove boilerplate (everything that is not the main text) from the result.

HTML and BeautifulSoup Parsing

The first function will use *requests* to get the HTML for the article. We are then using *BeautifulSoup* to parse the HTML and only return the content of the *bodyContent* div of the Wikipedia article.

HTML and BeautifulSoup Parsing

The second function also retrieves the HTML using *requests*. Instead of parsing the site ourselves, we are using *jusText* to identify non-boilerplate paragraphs which we then combine into one string.

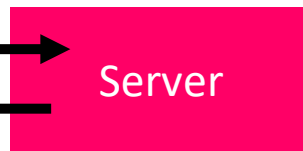
HTTP GET and *Requests*

For example a **browser**
or a **Python script**

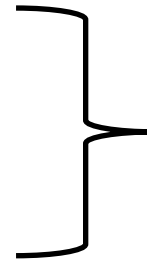


(1) HTTP GET

For example Wikipedia's
webserver



(2)
Code 200
Content / Files



In order to retrieve a website, a **client** (e.g., a browser) sends an HTTP GET request to a **webserver**. The server then responds with the website by sending **HTML** and possibly other content/files such as images.

The server also sends a status code indicating whether the requests worked.

We can use the Requests library to send HTTP requests:

```
r = requests.get('https://en.wikipedia.org/wiki/Linguistics')
```

`r.status_code` → 200

`r.text` → *text/HTML*

`r.content` → *Binary/non-text content*

Code	Meaning
200	OK
403	Forbidden
404	Not Found
5XX	Server Error
...	...

Exercise 16 – Putting Everything Together

The ultimate goal of this exercise is to write a system which can perform basic (comparative) keyword analysis on two corpora.

1. Use your web scraper to build a small Wikipedia corpus of about three to five articles. Ideally, they will belong to a similar topic, e.g., politics.
2. Find a suitable reference corpus to compare your Wikipedia corpus with.
3. Use your new skills to generate frequency lists for both corpora.
4. Implement any keyness statistic (e.g., simple maths or log-likelihood) and determine the keywords

Hint: To download the COCA sampler, run the following command in a Google Colab cell:

```
!cd python-programming-for-linguists/2020/data && sh download_coca.sh
```

This will download and extract the *COCA sampler* to your `/data/corpora/coca` folder.

Enumerate

```
l = ['A', 'B', 'C']
```

```
for index, value in enumerate(l):  
    print(index, value)
```

0 A

1 B

2 C

Lambda Functions

aka. Anonymous Functions

Lambda functions are very powerful but quite hard to comprehend. On the surface level, and we will not go any deeper, these are functions without a name. They are used when we only require a function for a short period of time.

```
x = lambda a: a + 10
```

```
x(5) → 15
```

We're only going to use them once!

They are, for example, useful when `.apply`-ing functions to a DataFrame.

Simple Maths Parameter

The k parameter works almost as a filter. The lower we set the parameter, the more low-frequency items we will identify as keywords.

$$SMP = \frac{RF_T + k}{RF_R + k} \quad k = 100$$

See Kilgarriff, Adam. (2009). Simple Maths for Keywords. In Proceedings of the Corpus Linguistics Conference, Liverpool, July.