

In [2]:

```
1 %%html
2 <style>
3 table {display: block;}
4 td {
5     font-size: 18px
6 }
7 .rendered_html { font-size: 18px; }
8 *{ line-height: 200%; }
9 </style>
```

# Natural Language Processing and the Web

## WS 2022/23 - Practice Class - **Tutorial 2**

### Content

In this practice class, we will discuss the following main points

- How to [read/write](#) text data from/to a file system and web page
- Text [Segmentation](#), [lemmatization](#), [stemming](#)
- [Part of speech tagging](#) using NLTK

To run this notebook, login to Moodle, click [nlp4web tutorial2](#). This will download ch2.zip that you have to unzip to your preferred location. Start [Anaconda navigation](#), open [Jupyter notebook](#) and navigate to the ch2 folder you have unzipped and open [ch2.ipynb](#) file

### How to read/read data from a file system

To open an existing file, you can do it as follows

```
file = open("filename", "r")
```

You can then get the content as

```
file.read()
```

If you want to handle some special characters encodings, you need to open the file with the given encodings. For example, to open the file with [UTF-8](#) encoding

```
import io  
file = io.open("filename", mode="r", encoding=  
"utf-8")
```

If you want to read the file line by line, you can do it as follows

```
file = open("filename", "r")  
for line in file:  
    print (line)
```

To read multiple files from a folder or sub folders, you can do it as follows

```
import os  
files = os.listdir("foldername")
```

Or using `os.scandir`

```
import os

with os.scandir("foldername") as files:
    for file in files:
        print(file.name)
```

In [1]:

```
1 # Install chardet library to detect the encoding of a text
2 !pip install chardet
```

Requirement already satisfied: chardet in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (5.0.0)

[notice] A new release of pip available: 22.2.2 -> 22.3

[notice] To update, run: `pip install --upgrade pip`

In [2]:

```
1 #Checking the encoding of a file
2 import chardet
3 def detect_encoding(file):
4     with open(file, 'rb') as raw_data:
5         result = chardet.detect(raw_data.read())
6         return(result)
```

In [3]:

```
1 print(detect_encoding("data/news/news.txt"))
2 print(detect_encoding("data/news/news4.txt"))
```

```
{'encoding': 'ISO-8859-1', 'confidence': 0.73, 'language': ''}
{'encoding': 'utf-8', 'confidence': 0.99, 'language': ''}
```

In [4]:

```
1
2 # Open file for reading ("r")
3 file = open("data/test.txt", "r")
4 # Print all the content
5 print("All content:", file.read())
6 file.close()
7 # Iterate line by line
8 file = open("data/test.txt", "r")
9 print("Printng line by line ...")
10 for line in file:
11     print(line)
12 file.close()
```

All content: This is line one.

And this is line two.

Printng line by line ...

This is line one.

And this is line two.

In [5]:

```
1 import io
2 file = io.open("data/news/news4.txt", mode="r", encoding="utf-8")
3 file.read()
```

Out[5]:

'የፋይናንስ ኢንዱስትሪውን በ500 ሚሊዮን ብር የተከፈለ ካፒታል የሚቀላቀለው አማራ ባንክ ከ21 August 2019 በኋላ ታዲያ በፋይናንስ ኢንዱስትሪው ውስጥ 2011 ዓ.ም. በተለየ የሚታይበትን ክስተት አስተናግዷል። ይህም የኢትዮጵያ ብሔራዊ ባንክ ለማቋቋም የሚጠይቀውን የተከፈለ ካፒታል መጠን ከ100 ሚሊዮን ብር ወደ 500 ሚሊዮን ብር ካሳደገ ወዲህ፣ ላለፉት ሰባት ዓመታት አንድም ባንክ አሳይቋቋም ቆይቶ ዘንድሮ ግን ሰባት ያህል ባንኮች ለምሥረታ መዘጋጀታቸው ነው። እነዚህ ባንኮች በአዲሱ ዓመት ወደ ሥራ እንዲመገቡ በማስታወቅ የአክሲዮን ሽያጭ ውስጥ ገብተዋል። በዕቅዳቸው መሠረት ከተጓዙ በቀጣዩ ዓመት የአገሪቱ ባንኮች ቁጥር ወደ 25 ያድጋል። እስከ ቀጣዩ ዓመት አጋማሽ ድረስ ወደ ሥራ ለመግባት ካቀዱትና በወራት ውስጥ የአክሲዮን ሽያጫቸውን ጨርሰው የባንክ ኢንዱስትሪውን እንዲሟሟ ቀላቀሉ ካስታወቁት ውስጥ አንዱ፣ አማራ ባንክ አክሲዮን ማሳበር ነው። በሁለት ቢሊዮን ብር የተፈቀደ ካፒታልና በ500 ሚሊዮን ብር የተከፈለ ካፒታል ወደ ሥራ እንዲገቡ ያለውን አማራ ባንክ አደራጅ ኮሚቴውን በሰብሳቢነት የሚመሩት፣ የገቢዎችና ጉምሩክ ባለሥልጣን የቀድሞ ዋና ዳይሬክተር አቶ መለኮ ፈንታ ናቸው። የባንኩ የፕሮጀክት አስተባባሪ በመሆን የተሰየሙት ደግሞ የቀድሞ የልማት ባንክ ምክትል ፕሬዚዳንት፣ ከዚያም የዓባይ ባንክ የመጀመርያ ዋ ፕሬዚዳንት ወ/ሮ መሰንበት ሸንቁጥ ናቸው። ወ/ሮ መሰንበት የአዲስ አበባ ንግድና ዘርፍ ማሳበራት ምክር ቤት ፕሬዚዳንት በመሆንም እያገለገሉ ነው። አማራ ባንክ ቅዳሜ ነሐሴ 11 ቀን 2011 ዓ.ም. የአክሲዮን ሽያጩን በይፋ የጀመረበት ሥነ ሥርዓት ላይ በክብር እንግ

In [6]:

```
1 import os
2 folder = "./data/"
3 with os.scandir(folder) as files:
4     for file in files:
5         print(file.name)
```

```
.DS_Store
names2.csv
simple-documents
names.csv
news
testout.txt
test.txt
```

## Writing to a file

You can write text data to a file as follows:

```
file = open("filename", "w")
file.write("content")
file.close()
```

If you want to append to an existing file,

```
file = open("filename", "a")
file.write("content to append")
```

In [8]:

```
1 # Write to a file "testout.txt", if the file does not exist, create it in the cu
2 file = open("testout.txt", "w")
3 file.write("This is the first line\n")
4 file.close()
5
6 # Append content to the "testout.txt" file
7 file = open("testout.txt", "a")
8 file.write("This is the second appended line\n")
9 file.close()
10
```

## How to read data from web page content

The standard [requests](#) library can be used to retrieve content from a web page. Lets see some of the methods for the [GET](#) HTTP action. A better way of retrieving web content will be discussed later, mainly using different packages.

In [10]:

```
1 # requests is the standard html library to make HTTP requests
2 import requests
3 link1 = "https://en.wikipedia.org/wiki/Natural_language_processing"
4
5 # get is an HTTP method to retrieve data from a specified resource
6 # response is an HTTP "Response" object to inspect the request
7 response = requests.get(link1)
8 # HTTP status code, example 200 is an ok status, 404 is not found status code
9 status = response.status_code
10 print(status)
```

200



In [12]:

```
1 # the information obtained from the resource, a payload in the form of string
2 content = response.text
3
4 print (type(content))
5 print("====")
6 print (content)
```

```
<class 'str'>
====
<!DOCTYPE html>
<html class="client-nojs" lang="en" dir="ltr">
<head>
<meta charset="UTF-8"/>
<title>Natural language processing - Wikipedia</title>
<script>document.documentElement.className="client-js";RLCONF={"wgBreakFrames":false,"wgSeparatorTransformTable":["",""],"wgDigitTransformTable":["",""],"wgDefaultDateFormat":"dmy","wgMonthNames":["","January","February","March","April","May","June","July","August","September","October","November","December"],"wgRequestId":"612780b3-202e-491b-bd00-172fa9fb7437","wgCSPNonce":false,"wgCanonicalNamespace":"","wgCanonicalSpecialPageName":false,"wgNamespaceNumber":0,"wgPageName":"Natural_language_processing","wgTitle":"Natural language processing","wgCurRevisionId":1112038569,"wgRevisionId":1112038569,"wgArticleId":21652,"wgIsArticle":true,"wgIsRedirect":false,"wgAction":"view","wgUserName":null,"wgUserGroups":["*"],"wgCategories":["All accuracy disputes","Accuracy disputes from December 2013","CS1 maint: location","Arti
```

In [13]:

```
1 # Examine which encoding the request determines for the response
2 encoding = response.encoding
3 print(encoding)
```

UTF-8





In [16]:

```
1 link2 = "https://api.github.com"
2
3 # get is an HTTP method to retrieve data from a specified resource
4 # response is an HTTP "Response" object to inspect the request
5 response = requests.get(link2)
6 # HTTP status code, example 200 is an ok status, 404 is not found status code
7 status = response.status_code
8 print(status)
```

200

In [17]:

```
1 # the information obtained from the resource, a payload in the form of bytes
2 content = response.content
3 print (content)
```

```
b'{"current_user_url": "https://api.github.com/user",\n "current_user_authorized_applications_html_url": "https://github.com/settings/connections/applications{/client_id}",\n "authorizations_url": "https://api.github.com/authorizations",\n "code_search_url": "https://api.github.com/search/code?q={query}{&page,per_page,sort,order}",\n "commit_search_url": "https://api.github.com/search/commits?q={query}{&page,per_page,sort,order}",\n "emails_url": "https://api.github.com/user/emails",\n "emojis_url": "https://api.github.com/emojis",\n "events_url": "https://api.github.com/events",\n "feeds_url": "https://api.github.com/feeds",\n "followers_url": "https://api.github.com/user/followers",\n "following_url": "https://api.github.com/user/following{/target}",\n "gists_url": "https://api.github.com/gists{/gist_id}",\n "hub_url": "https://api.github.com/hub",\n "issue_search_url": "https://api.github.com/search/issues?q={query}{&page,per_page,sort,order}",\n "issues_url": "https://api.github.com/issues",\n "keys_url": "https://api.github.com/user/keys",\n "label_search_url": "https://api.github.com/search/labels?q={query}&repository_id={repository_id}{&page,per_page}",\n "notifications_url": "https://api.github.com/notifications",\n "organization_url": "https://api.github.com/orgs/{org_id}",\n "organization_repositories_url": "https://api.github.com/orgs/{org_id}/repositories"}
```

In [18]:

```
1 headers = response.headers
2 print(headers)
3 print("====")
4 contentType = response.headers.get('content-type')
5 print(contentType)
```

```
{'Server': 'GitHub.com', 'Date': 'Tue, 25 Oct 2022 08:12:04 GMT', 'Cache-Control': 'public, max-age=60, s-maxage=60', 'Vary': 'Accept, Accept-Encoding, Accept, X-Requested-With', 'ETag': '"4f825cc84e1c733059d46e76e6df9db557ae5254f9625dfe8e1b09499c449438"', 'Access-Control-Expose-Headers': 'ETag, Link, Location, Retry-After, X-GitHub-OTP, X-RateLimit-Limit, X-RateLimit-Remaining, X-RateLimit-Used, X-RateLimit-Resource, X-RateLimit-Reset, X-OAuth-Scopes, X-Accepted-OAuth-Scopes, X-Poll-Interval, X-GitHub-Media-Type, X-GitHub-SSO, X-GitHub-Request-Id, Deprecation, Sunset', 'Access-Control-Allow-Origin': '*', 'Strict-Transport-Security': 'max-age=31536000; includeSubdomains; preload', 'X-Frame-Options': 'deny', 'X-Content-Type-Options': 'nosniff', 'X-XSS-Protection': '0', 'Referrer-Policy': 'origin-when-cross-origin, strict-origin-when-cross-origin', 'Content-Security-Policy': "default-src 'none'", 'Content-Type': 'application/json; charset=utf-8', 'X-GitHub-Media-Type': 'github.v3; format=json', 'Content-Encoding': 'gzip', 'X-RateLimit-Limit': '60', 'X-RateLimit-Remaining': '59', 'X-RateLimit-Reset': '1666689126', 'X-RateLimit-Resource': 'core', 'X-RateLimit-Used': '1', 'Accept-Ranges': 'bytes', 'Content-Length': '530', 'X-GitHub-Request-Id': 'E7A2:F493:9B67D5:9DD7F4:63579A56'}
====
application/json; charset=utf-8
```

## print a formatted content using pprint





In [23]:

```
1 # Get only the text content of the page
2 print(soup.text)
```

Beautiful Soup: We called him Tortoise because he taught us.

[ Download | Documentation | Hall of Fame | For enterprise | Source |  
Changelog | Discussion group | Zine ]

Beautiful Soup

In [24]:

```
1 # list all the available elements in the page
2 [type(item) for item in list(soup.children)]
```

Out[24]:

```
[bs4.element.Doctype,
bs4.element.NavigableString,
bs4.element.Tag,
bs4.element.NavigableString,
bs4.element.Tag,
bs4.element.NavigableString,
bs4.element.Tag,
bs4.element.Comment,
bs4.element.Tag,
bs4.element.NavigableString,
bs4.element.NavigableString,
bs4.element.NavigableString,
bs4.element.NavigableString,
bs4.element.NavigableString]
```







```
rote about what I learned about software development from working on
Beautiful Soup. Thanks!</p>
</p></p></p></p></p></body>]
```

In [26]:

```
1 body = list(html.children)[3]
```

In [27]:

```
1 # Get all the sub-tags under body
2 for tag in list(body.children):
3     if tag.name:
4         print (tag.name)
```

```
style
img
br
p
div
p
p
```

In [28]:

```
1 # Get all the text under the p tags
2 for tag in list(body.children):
3     if tag.name and tag.name == 'p':
4         #pass # uncomment the following line
5         print (tag.text)
6 # OR use findall by the tag name
7 print("====Using find all====")
8 #for p in soup.find_all('p'):
9 #     print(p.text)
```

[ [Download](#) | [Documentation](#) | [Hall of Fame](#) | [For enterprise](#) | [Source](#) | [Changelog](#) | [Discussion group](#) | [Zine](#) ]  
You didn't write that awful page. You're just trying to get some data out of it. Beautiful Soup is here to help. Since 2004, it's been saving programmers hours or days of work on quick-turnaround screen scraping projects.  
Beautiful Soup is a Python library designed for quick turnaround projects like screen-scraping. Three features make it powerful:

Beautiful Soup provides a few simple methods and Pythonic idioms for navigating, searching, and modifying a parse tree: a toolkit for dissecting a document and extracting what you need. It doesn't take much code to write an application

Beautiful Soup automatically converts incoming documents to Unicode and outgoing documents to UTF-8. You don't have to think about encodings, unless the document doesn't specify an encoding and Beautiful Soup can't detect one. Then you just have to specify the

## Reading and writing **CSV** files.

CSV files are comma delimited files that can be used to store structured files in the form of lists. It might contain headers, and the header as well as individual records are separated by new lines. In general delimiters can be also user defined.

You can use the [csv library](#) to process CSV files.

```
import csv
```

In [29]:

```
1 # Read personal information: name.csv is obtained from https://github.com/Coreym  
2 import csv  
3 with open('data/names.csv','r') as csv_names:  
4     name_reader = csv.reader(csv_names,delimiter=',')  
5     # skip the header  
6     next(name_reader)  
7     for name in name_reader:  
8         print(name)
```

```
['John', 'Doe', 'john-doe@bogusemail.com']  
['Mary', 'Smith-Robinson', 'maryjacobs@bogusemail.com']  
['Dave', 'Smith', 'davesmith@bogusemail.com']  
['Jane', 'Stuart', 'janestuart@bogusemail.com']  
['Tom', 'Wright', 'tomwright@bogusemail.com']  
['Steve', 'Robinson', 'steverobinson@bogusemail.com']  
['Nicole', 'Jacobs', 'nicolejacobs@bogusemail.com']  
['Jane', 'Wright', 'janewright@bogusemail.com']  
['Jane', 'Doe', 'janedoe@bogusemail.com']  
['Kurt', 'Wright', 'kurtwright@bogusemail.com']  
['Kurt', 'Robinson', 'kurtrobinson@bogusemail.com']  
['Jane', 'Jenkins', 'janejenkins@bogusemail.com']  
['Neil', 'Robinson', 'neilrobinson@bogusemail.com']  
['Tom', 'Patterson', 'tompatterson@bogusemail.com']  
['Sam', 'Jenkins', 'samjenkins@bogusemail.com']  
['Steve', 'Stuart', 'stevestuart@bogusemail.com']  
['Maggie', 'Patterson', 'maggiepatterson@bogusemail.com']  
['Maggie', 'Stuart', 'maggiestuart@bogusemail.com']  
['Jane', 'Doe', 'janedoe@bogusemail.com']  
['Steve', 'Patterson', 'stevepatterson@bogusemail.com']  
['Dave', 'Smith', 'davesmith@bogusemail.com']  
['Sam', 'Wilks', 'samwilks@bogusemail.com']  
['Kurt', 'Jefferson', 'kurtjefferson@bogusemail.com']  
['Sam', 'Stuart', 'samstuart@bogusemail.com']  
['Jane', 'Stuart', 'janestuart@bogusemail.com']  
['Dave', 'Davis', 'davedavis@bogusemail.com']  
['Sam', 'Patterson', 'sampatterson@bogusemail.com']  
['Tom', 'Jefferson', 'tomjefferson@bogusemail.com']  
['Jane', 'Stuart', 'janestuart@bogusemail.com']  
['Maggie', 'Jefferson', 'maggiejefferson@bogusemail.com']  
['Mary', 'Wilks', 'marywilks@bogusemail.com']  
['Neil', 'Patterson', 'neilpatterson@bogusemail.com']  
['Corey', 'Davis', 'coreydavis@bogusemail.com']  
['Steve', 'Jacobs', 'stevejacobs@bogusemail.com']  
['Jane', 'Jenkins', 'janejenkins@bogusemail.com']  
['John', 'Jacobs', 'johnjacobs@bogusemail.com']  
['Neil', 'Smith', 'neilsmith@bogusemail.com']  
['Corey', 'Wilks', 'coreywilks@bogusemail.com']  
['Corey', 'Smith', 'coreysmith@bogusemail.com']  
['Mary', 'Patterson', 'marypatterson@bogusemail.com']  
['Jane', 'Stuart', 'janestuart@bogusemail.com']  
['Travis', 'Arnold', 'travisarnold@bogusemail.com']
```

```
['John', 'Robinson', 'johnrobinson@bogusemail.com']  
['Travis', 'Arnold', 'travisarnold@bogusemail.com']
```

In [30]:

```
1 # Re-write the csv file with a different delimiter  
2 with open('data/names.csv', 'r') as csv_names:  
3     name_reader = csv.reader(csv_names, delimiter=',')  
4     with open('data/names2.csv', 'w') as csv_names:  
5         name_writer = csv.writer(csv_names, delimiter="\t")  
6         for name in name_reader:  
7             name_writer.writerow(name)  
8
```

## Excercise\_1 ( 3pt)

1. List all the files in the subdirectories of the [data](#) folder that you have downloaded from Moodle. (0.5 pt)
2. Merge all the news content into a single file and save into a different file named [allnews.txt](#) using the ISO-8859-1 encoding. (0.5 pt)
3. Read the country name and capital city from the [this](#) (<https://geographyfieldwork.com/WorldCapitalCities.htm>) page, which lists the world capital cities with their country. Save the result as a [comma separated value \(csv\)](#) file format. (2 pts)

## Tokenization

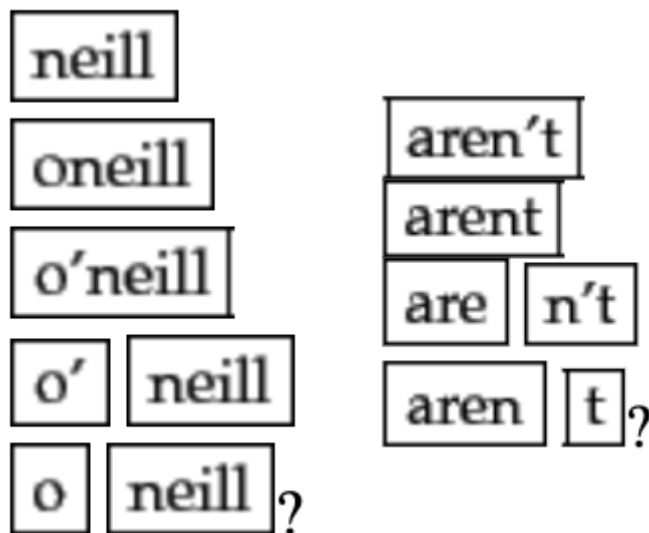
Tokenization is the first step in almost all NLP applications. It is an important prerequisite for a number of downstream NLP tasks. Consider the following example:

Mr. O'Neill thinks that the boys' stories about Chile's capital aren't amusing.``

For O'Neill and aren't, which of the following are the desired tokenization?

O'Neill

aren't



## Tokenizing with regex

### Sentence segmentation

Lets assume that the text is properly written English documents. We can build a simple segmenter, that will split the sentence by [a full stop](#).


```
m = re.split('(\.)', text)
```

However, if the text contains . which is not the end of the sentence by itself, it will wrongly split them.

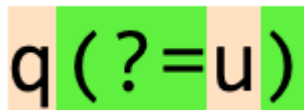
Better solution is to use [lookahead](#) and [lookbehind](#) regex patterns.

[Lookahead](#) mechanism asserts the presence or absence of certain expressions without consuming the pattern itself. Positive [lookahead](#) is defined with [\(?=REGEX\)](#) while negative [lookahead](#) is defined as [\(?!REGEX\)](#)

Example: Match a `q` not followed by a `u` .

The image shows the regex pattern `q(?:!u)`. The character `q` is in a grey box. The opening parenthesis `(` is in a green box. The characters `?!u` are in a green box. The closing parenthesis `)` is in a grey box.

Match a `q` that is followed by a `u`

The image shows the regex pattern `q(?=u)`. The character `q` is in an orange box. The opening parenthesis `(` is in a green box. The characters `=u` are in a green box. The closing parenthesis `)` is in an orange box.

Positive [lookbehind](#) is defined with [\(?<=REGEX\)](#) while negative [lookbehind](#) is defined as [\(?<!REGEX\)](#)

In [31]:

```
1 import re
2 #Split sentence with .
3 text = """Fruits like apple, orange, and mango are healthy. But they are expensive
4 ($1.5/kililo), i.e Mrs. Bean can't afford them! However, One can order some online from
5 www.rewe.de
6
7 sents = re.split(r"\.", text) # Add space here and see the effect
8
9 for sent in sents:
10     print (sent) # What is the problem here?
11
12 print("===")
13 # Find all words that are followed by , or .
14 # \w+ --> one or more words
15 # (?=,|\.) --> followed by , or . or !
16 p = re.compile("\w+(?=,|\.|\!)")
17 print("Lookahead matches", p.findall(text))
```

Fruits like apple, orange, and mango are healthy

But they are expensive

(\$1.5/kililo), i

e Mrs

Bean can't afford them! However, One can order some online from www

rewe

de

===

Lookahead matches ['apple', 'orange', 'healthy', '1', 'i', 'Mrs', 'the', 'm', 'However', 'www', 'rewe', 'de']

In [32]:

```
1 # Sentence splitting using lookbehinde and lookahead
2 print("unsegmented text:", text)
3 print("===")
4 # lookahead/lookbehind matchs
5 # (?<![A-Z][a-z]\.) --> If . is Not behind two characters (first cap, second sma
6 # (?<=[\./?|\!])\s+ --> If . is before one or more spaces
7 # After the space, the first character should be Capitalized
8 sents = re.split("(?<![A-Z][a-z]\.)(?<=[\./?|\!])\s+(?=[A-Z])", text)
9 for sent in sents:
10     print (sent)
```

unsegmented text: Fruits like apple, orange, and mango are healthy. But  
t they are expensive  
(\$1.5/killo), i.e Mrs. Bean can't afford them! However, One can order  
some online from www.rewe.de.

===

Fruits like apple, orange, and mango are healthy.  
But they are expensive  
(\$1.5/killo), i.e Mrs.  
Bean can't afford them!  
However, One can order some online from www.rewe.de.

## Exercise\_2 ( 3pt)

Modify the [regex](#) above for sentence segmentation so that the following text are split  
into correct sentences.

Fruits like apple, orange, and mango are health. But  
they are expensive, i.e Mr. Bean can't afford them!  
One can order some online from www.rewe.de. Prof.  
Karl, Dep. of Plant Science. Email:  
karl@plant.science.de. Regards!



# Word tokenization with regex

Word tokenization for English languages can be attained using white space characters. However, there are still special cases that need to be addressed.

In [33]:

```
1 # English tokenization,  
2 tokens = re.split("(\w+|\$[\d\.]+\S+)", text)  
3 for token in tokens:  
4     if token.strip():  
5         print (token) # Can you fix some of the issues?
```

```
Fruits  
like  
apple  
,  
orange  
,  
and  
mango  
are  
healthy  
.   
But  
they  
are  
expensive  
($1.5/killo),  
i  
.e  
Mrs  
.   
Bean  
can  
't  
afford  
them  
!  
However  
,  
One  
can  
order  
some  
online  
from  
www  
.rewe.de.
```

## Excercise\_2( 3pt)

Modify/re-write the word tokenization pattern given above so that you can achieve near ideal tokenization for the following text

```
"I said, 'what're you? Crazy?'" said Sandowsky. "I  
can't afford to do that."
```

See the ideal tokenization result from the Exercise\_2 - Ideal tokenization - file in Moodle.

## Tokenization with NLTK

NLTK includes tokenizers for different languages

In [34]:

```
1 import nltk
2 sentence = "These are simply sentences obtained from the treminals."
3 tokens = nltk.word_tokenize(sentence, language='english')
4 print(tokens)
5 sentence2 = ""
6 Die Brände in Brasilien setzen erhebliche Mengen an klimaschädlichen Treibhausga
7 Die Nasa hat nun simuliert, wie sich Kohlenmonoxid über Südamerika ausbreitet. A
8 das Gas der Gesundheit erheblich.""
9 tokens2 = nltk.word_tokenize(sentence2, language='german')
10 print (tokens2)
```

```
['These', 'are', 'simply', 'sentences', 'obtained', 'from', 'the', 'tr  
eminals', '.']  
['Die', 'Brände', 'in', 'Brasilien', 'setzen', 'erhebliche', 'Mengen',  
'an', 'klimaschädlichen', 'Treibhausgasen', 'frei', '.', 'Die', 'Nas  
a', 'hat', 'nun', 'simuliert', ',', 'wie', 'sich', 'Kohlenmonoxid', 'ü  
ber', 'Südamerika', 'ausbreitet', '.', 'Am', 'Boden', 'schadet', 'da  
s', 'Gas', 'der', 'Gesundheit', 'erheblich', '.']
```

# Stemmer in NLTK

**Stemming** tries to remove or chops off the end of a word to reduce the word to its base form, somehow a crude approach.

Different stemmers use different rules.

In [35]:

```
1 porter = nltk.PorterStemmer()
2 lancaster = nltk.LancasterStemmer()
3 [porter.stem(token) for token in tokens]
4
```

Out[35]:

```
['these', 'are', 'simpli', 'sentenc', 'obtain', 'from', 'the', 'tremi  
n', '.']
```

In [36]:

```
1 [lancaster.stem(token) for token in tokens]
```

Out[36]:

```
['thes', 'ar', 'simply', 'sent', 'obtain', 'from', 'the', 'tremi  
'n', '.']
```

In [37]:

```
1 #Example from: https://www.datacamp.com/community/tutorials/stemming-lemmatizati
2 #A list of words to be stemmed
3 word_list = ["friend", "friendship", "friends", "friendships", "stabil", "destabil
4             "railroad", "moonlight", "football", "ill", "illness", "sick", "sickness"
5 print("{0:20}{1:20}{2:20}".format("Word", "Porter Stemmer", "lancaster Stemmer"))
6 for word in word_list:
7     print("{0:20}{1:20}{2:20}".format(word, porter.stem(word), lancaster.stem(word))
```

Word	Porter Stemmer	lancaster Stemmer
friend	friend	friend
friendship	friendship	friend
friends	friend	friend
friendships	friendship	friend
stabil	stabil	stabl
destabilize	destabil	dest
misunderstanding	misunderstand	misunderstand
railroad	railroad	railroad
moonlight	moonlight	moonlight
football	footbal	footbal
ill	ill	il
illness	ill	il
sick	sick	sick
sickness	sick	sick
hope	hope	hop

In [38]:

```
1 #CISTEM Stemmer for German
2 de_stem = nltk.stem.cistem.Cistem()
3 #Show the top 6 stems from the list
4 [de_stem.stem(t) for t in tokens2][:6]
```

Out[38]:

```
['die', 'brand', 'in', 'brasilie', 'setz', 'erheblich']
```

## Tokenization spaCy

In [39]:

```
1 import spacy
2 spacy = spacy.load('en_core_web_sm')
3 doc = spacy(text)
4 print("{0:15}{1:15}".format("Token", "Lemma"))
5 for token in doc:
6     # print the token and its lemma
7     print("{0:15}{1:15}".format(token.text, token.lemma_)) #How does it compare
```

Token	Lemma
Fruits	fruit
like	like
apple	apple
'	'
orange	orange
'	'
and	and
mango	mango
are	be
healthy	healthy
.	.
But	but
they	they
are	be
expensive	expensive

(	(
\$	\$
1.5	1.5
/	/
killo	killo
)	)
'	'
i.e	i.e
Mrs.	Mrs.
Bean	Bean
ca	ca
n't	n't
afford	afford
them	they
!	!
However	however
'	'
One	one
can	can
order	order
some	some
online	online
from	from
www.rewe.de	www.rewe.de
.	.

In [40]:

```
1 for sent in doc.sents:
2     print (sent) #How does it compare with the sentences in the regex and NLTK?
```

```
Fruits like apple, orange, and mango are healthy.
But they are expensive
($1.5/killo), i.e Mrs. Bean can't afford them!
However, One can order some online from www.rewe.de.
```

## Lemmatization

A lemma is the canonical, uninflected or dictionary form of a word. For example, the lemma of `smallest` is `small`, and the lemma of `eating` is `eat`. In many languages, the lemma for nouns is the `nominative singular` form, the lemma for adjectives is the `nominative singular positive` form, and the lemma for verbs is the `infinitive`. But given an inflected form, finding the lemma (a process called `lemmatization`) is not always as easy. Words often undergo regular spelling changes when inflected for example, in English, verbs and adjectives ending in `-e` often drop this letter when inflecting: `bake` → `baking`. Sometimes final consonants are doubled, as in (British) English `travel` → `travelling`.

An accurate algorithm for lemmatization must be aware of these sorts of inflectional rules and know how to undo them to arrive at the `base form` of the word. It must also know about completely irregular cases, such as `go` → `went`, `mouse` → `mice`, and `good` → `better`. Lemmatization is a difficult task for computers, and requires some basic understanding of the grammatical context and properties of the word. For example, the lemma of `dove` depends on whether the word is being used as a noun

(as in the `bird`) or a verb (as in the past tense of `dive`). However, lemmatization is an important task because, as with part-of-speech tagging, many NLP applications rely on lemmatized text.

Examples of lemmatization:

```
corpora : corpus
```

```
better : good
```

## NLTK Lemmatizer

In [41]:

```
1 import nltk
2 nltk.download('wordnet')
3 # Lemmatize using WordNet's built-in morphy function
4 # Returns the input unchanged if it cannot be found in WordNet
5 from nltk.stem import WordNetLemmatizer
6 lemmatizer = WordNetLemmatizer()
7 print("rocks :", lemmatizer.lemmatize("rocks"))
8 print("corpora :", lemmatizer.lemmatize("corpora"))
9 #Give the POS tag as a context to the tager, a denotes adjective in "pos"
10 print("better :", lemmatizer.lemmatize("better", pos ="a"))
11 #Lemmatizing sentence
12 sentence = "The striped bats are hanging on their feet for best"
13 word_list = nltk.word_tokenize(sentence)
14 print("words:",word_list)
15 # Lemmatize list of words and join
16 lemmatized_output = ' '.join([lemmatizer.lemmatize(w) for w in word_list])
17 print("lemma:",lemmatized_output)
```

```
[nltk_data] Downloading package wordnet to
[nltk_data]    /Users/abhikjana/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

```
rocks : rock
corpora : corpus
better : good
words: ['The', 'striped', 'bats', 'are', 'hanging', 'on', 'their', 'feet', 'for', 'best']
lemma: The striped bat are hanging on their foot for best
```

## spaCy Lemmatizzer



In [42]:

```
1 import spacy
2 # Initialize spacy 'en_core_web_sm' model, keeping only tagger component needed
3 nlp = spacy.load('en_core_web_sm', disable=['parser', 'ner'])
4 sentence = "The striped bats are hanging on their feet for best"
5 # Parse the sentence using the loaded 'English' model object `nlp`
6 doc = nlp(sentence)
7 # Extract the lemma for each token and join
8 " ".join(["[" + token.lemma_ + "]" for token in doc])
```

Out[42]:

```
'[the] [striped] [bat] [be] [hang] [on] [their] [foot] [for] [good]'
```

## TextBlob Lemmatizer

In [43]:

```
1 !pip install textblob
```

```
Requirement already satisfied: textblob in /Library/Frameworks/Python.
framework/Versions/3.8/lib/python3.8/site-packages (0.17.1)
Requirement already satisfied: nltk>=3.1 in /Library/Frameworks/Pytho
n.framework/Versions/3.8/lib/python3.8/site-packages (from textblob)
(3.5)
Requirement already satisfied: click in /Library/Frameworks/Python.fra
mework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textb
lob) (8.0.3)
Requirement already satisfied: regex in /Library/Frameworks/Python.fra
mework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textb
lob) (2020.11.13)
Requirement already satisfied: tqdm in /Library/Frameworks/Python.fram
ework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textbl
ob) (4.54.1)
Requirement already satisfied: joblib in /Library/Frameworks/Python.fr
amework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->text
blob) (1.0.0)
```

[notice] A new release of pip available: 22.2.2 -> 22.3

[notice] To update, run: pip install --upgrade pip

In [44]:

```
1 from textblob import TextBlob, Word
2 # Lemmatize a word, use the WordNet's morphy function
3 word = 'stripes'
4 w = Word(word)
5 w.lemmatize()
```

Out[44]:

'stripe'

In [45]:

```
1 # Lemmatize a sentence
2 sentence = "The striped bats are hanging on their feet for best"
3 sent = TextBlob(sentence)
4 " ".join(["["+ w.lemmatize()+"]" for w in sent.words])
```

Out[45]:

'[The] [striped] [bat] [are] [hanging] [on] [their] [foot] [for] [best]'

## Parts of speech tagging with NLTK

Part-of-speech tagging (POS tagging) is the process of marking up the words in a text with their corresponding part of speech (e.g., noun, verb, adjective). For example, take the following sentence:

A dog had seen the cutest ferrets.

A tokenizer would split it into the following tokens:

A	dog	had	seen	the	cutest	ferrets
---	-----	-----	------	-----	--------	---------

A part-of-speech tagger could then assign labels, or tags, to the tokens according to their respective parts of speech:

A<sub>DT</sub> dog<sub>NN</sub> had<sub>VBD</sub> seen<sub>VBN</sub> the<sub>DT</sub> cutest<sub>JJS</sub> ferrets<sub>NNS</sub>

The Penn Treebank tags used here are as follows: DT determiner NN noun, singular or mass VBD verb, past tense JJS adjective, superlative NNS noun, plural VBN verb, past participle

The inventory from which these POS tags are drawn varies from language to language , and from application to application .

NLTK includes a Part-of-speech tagger, which assign a [tag](#), or [word class](#), or [lexical category](#) for a given token in a text. The default POS tagset for English is based on [PennTreebank tagset](#)

([https://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn\\_treebank\\_pos.html](https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html)).

NLTK also include the [Universal POS tagset](#)

(<https://universaldependencies.org/u/pos/>).

In [46]:

```
1 from nltk.tokenize import sent_tokenize, word_tokenize
2 from nltk.tag import pos_tag
3 import nltk
4 nltk.download('tagsets')
5 text = "I saw a man sawing the tree with a saw. He can't finish it ontime."
6 sentences = sent_tokenize(text)
7 print("tagged", pos_tag(word_tokenize(sentences[0])))
8 print("=====")
9 for sentence in sentences:
10     for token, pos in pos_tag(word_tokenize(sentence)):
11         print(token, pos)
12 # to get information about a given tag
13 print("=====")
14 nltk.help.upenn_tagset("VB")
```

```
tagged [('I', 'PRP'), ('saw', 'VBD'), ('a', 'DT'), ('man', 'NN'), ('sa
wing', 'VBG'), ('the', 'DT'), ('tree', 'NN'), ('with', 'IN'), ('a', 'D
T'), ('saw', 'NN'), ('.', '.')]
=====
```

```
I PRP
saw VBD
a DT
man NN
sawing VBG
the DT
tree NN
with IN
a DT
saw NN
. .
He PRP
ca MD
n't RB
finish VB
it PRP
ontime RB
. .
=====
```

```
VB: verb, base form
ask assemble assess assign assume atone attention avoid bake balka
nize
bank begin behold believe bend benefit bevel beware bless boil bom
b
boost brace break bring broil brush build ...
```

```
[nltk_data] Downloading package tagsets to
[nltk_data] /Users/abhikjana/nltk_data...
[nltk_data] Package tagsets is already up-to-date!
```

In [47]:

```
1 # you can also decide to use the Universal POS tagset
2 nltk.download('universal_tagset')
3 print("{0:15}{1:15}".format("Token", "POS tag"))
4 for sentence in sentences:
5     for token, pos in pos_tag(word_tokenize(sentence), tagset='universal'):
6         print("{0:15}{1:15}".format(token, pos) )
```

Token	POS tag
I	PRON
saw	VERB
a	DET
man	NOUN
sawing	VERB
the	DET
tree	NOUN
with	ADP
a	DET
saw	NOUN
.	.
He	PRON
ca	VERB
n't	ADV
finish	VERB
it	PRON
ontime	ADV
.	.

```
[nltk_data] Downloading package universal_tagset to
[nltk_data]    /Users/abhikjana/nltk_data...
[nltk_data]    Package universal_tagset is already up-to-date!
```

## WordNet Lemmatizer with appropriate POS tag

wordnet pos tags are: noun (n), verb(v) , adj(a) and adv (r)

In [48]:

```
1 # Lemmatize with POS Tag
2 from nltk.corpus import wordnet
3 # wordnet pos tags are: noun (n), verb(v) , adj(a) and adv (r)
4 def get_wordnet_pos(word):
5     """Map POS tag to first character lemmatize() accepts"""
6     tag = nltk.pos_tag([word])[0][1][0].upper()
7     tag_dict = {"J": wordnet.ADJ,
8                 "N": wordnet.NOUN,
9                 "V": wordnet.VERB,
10                "R": wordnet.ADV}
11     return tag_dict.get(tag, wordnet.NOUN)
12
13 # 1. Init Lemmatizer
14 lemmatizer = WordNetLemmatizer()
15 # 2. Lemmatize Single Word with the appropriate POS tag
16 word = 'feet'
17 print(lemmatizer.lemmatize(word, get_wordnet_pos(word)))
18 # 3. Lemmatize a Sentence with the appropriate POS tag
19 sentence = "The striped bats are hanging on their feet for best"
20 print([lemmatizer.lemmatize(w, get_wordnet_pos(w)) for w in nltk.word_tokenize(s
21
```

```
foot
['The', 'strip', 'bat', 'be', 'hang', 'on', 'their', 'foot', 'for', 'best']
```

## Part of speech tagging with spaCy

In [47]:

```
1 import spacy
2 import pprint
3 # Load English tokenizer, tagger,
4 # parser, NER and word vectors
5 nlp = spacy.load("en_core_web_sm")
6 text = ("I saw a man sawing the tree with a saw. He can't finish it ontime!")
7 doc = nlp(text)
8 # Print token and Tag
9 for token in doc:
10     print(token, token.pos_)
11 # Example list of Verb tokens
12 print("Verbs:", [token.text for token in doc if token.pos_ == "VERB"])
```

```
I PRON
saw VERB
a DET
man NOUN
sawing VERB
the DET
tree NOUN
with ADP
a DET
saw NOUN
. PUNCT
He PRON
ca AUX
n't PART
finish VERB
it PRON
ontime ADV
! PUNCT
Verbs: ['saw', 'sawing', 'finish']
```

## Part of speech tagging with TextBlob

In [49]:

```
1 !pip install -U textblob
```

```
Requirement already satisfied: textblob in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (0.17.1)
Requirement already satisfied: nltk>=3.1 in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from textblob) (3.5)
Requirement already satisfied: tqdm in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textblob) (4.54.1)
Requirement already satisfied: click in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textblob) (8.0.3)
Requirement already satisfied: regex in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textblob) (2020.11.13)
Requirement already satisfied: joblib in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk>=3.1->textblob) (1.0.0)
```

[notice] A new release of pip available: 22.2.2 -> 22.3

[notice] To update, run: `pip install --upgrade pip`

In [50]:

```
1 !python -m textblob.download_corpora
```

```
[nltk_data] Downloading package brown to /Users/abhikjana/nltk_data...
[nltk_data]   Package brown is already up-to-date!
[nltk_data] Downloading package punkt to /Users/abhikjana/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   /Users/abhikjana/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /Users/abhikjana/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-date!
[nltk_data] Downloading package conll2000 to
[nltk_data]   /Users/abhikjana/nltk_data...
[nltk_data]   Package conll2000 is already up-to-date!
[nltk_data] Downloading package movie_reviews to
[nltk_data]   /Users/abhikjana/nltk_data...
[nltk_data]   Package movie_reviews is already up-to-date!
Finished.
```



In [51]:

```
1 from textblob import TextBlob
2 text = ("I saw a man sawing the tree with a saw. He can't finish it ontime!")
3 # create a textblob object
4 blob_object = TextBlob(text)
5 # print word with pos tag.
6 print(blob_object.tags)
```

```
[('I', 'PRP'), ('saw', 'VBD'), ('a', 'DT'), ('man', 'NN'), ('sawing', 'VBG'), ('the', 'DT'), ('tree', 'NN'), ('with', 'IN'), ('a', 'DT'), ('saw', 'NN'), ('He', 'PRP'), ('ca', 'MD'), ('n't', 'RB'), ('finish', 'VB'), ('it', 'PRP'), ('ontime', 'RB')]
```

## TextBlob Lemmatizer with appropriate POS tag

In [52]:

```
1 #Define function to lemmatize each word with its POS tag
2 def lemmatize_with_postag(sentence):
3     sent = TextBlob(sentence)
4     tag_dict = {"J": 'a',
5                 "N": 'n',
6                 "V": 'v',
7                 "R": 'r'}
8     words_and_tags = [(w, tag_dict.get(pos[0], 'n')) for w, pos in sent.tags]
9     lemmatized_list = ["["+wd.lemmatize(tag) +"]" for wd, tag in words_and_tags]
10    return " ".join(lemmatized_list)
11 # Lemmatize
12 sentence = "The striped bats are hanging on their feet for best"
13 lemmatize_with_postag(sentence)
```

Out[52]:

```
'[The] [striped] [bat] [be] [hang] [on] [their] [foot] [for] [best]'
```

## Excercise\_3 ( 3pt)

# Lemmatization for German

There is no lemmatization library in NLTK for German. However, the [GermaLemma](https://github.com/WZBSocialScienceCenter/germalemma) (<https://github.com/WZBSocialScienceCenter/germalemma>) (<https://github.com/WZBSocialScienceCenter/germalemma>) (<https://github.com/WZBSocialScienceCenter/germalemma>) library is an open source lemmatizer for German. To lemmatize a word, you need to pass the POS tag as a secondary argument. In this exercise, you can use the POS tagger for German from [pattern.de](https://github.com/WZBSocialScienceCenter/germalemma) but then you have to convert tags into N , V , ADJ , or ADV . So your task is, when the word category is in one of the four tags, map them and pass to the lemmatizer. If the POS tag is not in the four categories, return the word itself as the lemma. See the cells below on how to execute the lemmatizer and pos tager for German.

You can install [GermaLemma](https://github.com/WZBSocialScienceCenter/germalemma) as

```
pip install -U germalemma
```

Also make sure mysql and related packages are installed

In [52]:

```
1 #uncomment the following for Mysql cleint dev in Linux
2 # !sudo apt install default-libmysqlclient-dev
3 # Installing GermaLemma
4 !pip install -U germalemma
```

Collecting germalemma

Using cached germalemma-0.1.3-py3-none-any.whl (2.3 MB)

Collecting PatternLite>=3.6

Using cached PatternLite-3.6-py3-none-any.whl (22.1 MB)

Collecting Pyphen>=0.9.5

Downloading pyphen-0.13.0-py3-none-any.whl (2.0 MB)

---

2.0/2.0 MB 20.5 MB/s eta

0:00:0000:0100:01

Requirement already satisfied: numpy in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from PatternLite>=3.6->germalemma) (1.17.4)

Requirement already satisfied: scipy in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from PatternLite>=3.6->germalemma) (1.4.1)

Requirement already satisfied: nltk in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from PatternLite>=3.6->germalemma) (3.5)

Requirement already satisfied: regex in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk->PatternLite>=3.6->germalemma) (2020.11.13)

Requirement already satisfied: click in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk->PatternLite>=3.6->germalemma) (8.0.3)

Requirement already satisfied: tqdm in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk->PatternLite>=3.6->germalemma) (4.54.1)

Requirement already satisfied: joblib in /Library/Frameworks/Python.framework/Versions/3.8/lib/python3.8/site-packages (from nltk->PatternLite>=3.6->germalemma) (1.0.0)

Installing collected packages: Pyphen, PatternLite, germalemma

Successfully installed PatternLite-3.6 Pyphen-0.13.0 germalemma-0.1.3

[notice] A new release of pip available: 22.2.2 -> 22.3

[notice] To update, run: pip install --upgrade pip

In [53]:

```
1 # POS taggin for German (ISSUE - run this cell multiple times if it raises excep
2 from germalemma import GermaLemma
3 from pattern.de import parse, split
4 s = parse('Die Katze liegt auf der Matte.')
5 for sentence in split(s):
6     for token in sentence:
7         print (token, token.pos)
```

```
Die DT
Katze NN
liegt VB
auf IN
der DT
Matte NN
. .
```

**Address the TODO part below to complete your exercise**

In [54]:

```
1 # Lemmatizer for German using GermaLemma lematizer -
2 #(ISSUE - run this cell multiple times, it may happen that it raises exceptions,
3
4 sentence2 = ""Die Brände in Brasilien setzen erhebliche Mengen an klimaschädlic
5 Die Nasa hat nun simuliert, wie sich Kohlenmonoxid über Südamerika ausbreitet.
6 Am Boden schadet das Gas der Gesundheit erheblich.""
7 de_lemma = GermaLemma()
8 # POS tagger
9 poses = parse(sentence2)
10 for sentence in split(poses):
11     print ("===")
12     for token in sentence:
13         print (token, token.pos)
14         #TODO: Here map the POS tag to V, N, ADJ, or ADV as MAPD_POS. MAPD_POS =
15         # If the POS tag is not in V, N, ADJ, or ADV, no need to lemmatize
16         # Print the lemma here
17         # Print([de_lemma.find_lemma(token, MAPED_POS) for token in split(poses),
18     print ("===")
19 #
```

===

Die DT  
Brände NNS  
in IN  
Brasilien NNP  
setzen VB  
erhebliche JJ  
Mengen NN  
an IN  
klimaschädlichen JJ  
Treibhausgasen NN  
frei JJ

. .

===

===

Die DT  
Nasa NN  
hat VB  
nun IN  
simuliert NN  
, ,  
wie IN  
sich PRP  
Kohlenmonoxid NN  
über IN

Südamerika NNP  
ausbreitet NNP  
.  
.  
===  
===  
Am IN  
Boden NN  
schadet VB  
das DT  
Gas NN  
der DT  
Gesundheit NN  
erheblich JJ  
.  
.  
===

## Resources

- [HTTP request \(https://realpython.com/python-requests/\)](https://realpython.com/python-requests/)
- [Regex \(https://www.programiz.com/python-programming/regex\)](https://www.programiz.com/python-programming/regex)
- [Regex lookahead and lookbehind \(https://www.regular-expressions.info/lookaround.html\)](https://www.regular-expressions.info/lookaround.html)
- [BeautifulSoup tutorial \(https://www.dataquest.io/blog/web-scraping-tutorial-python/\)](https://www.dataquest.io/blog/web-scraping-tutorial-python/)
- [Online Regex editor \(https://regex101.com/r/nG1gU7/\)](https://regex101.com/r/nG1gU7/)
- [Processing Raw Text NLTK \(https://www.nltk.org/book/ch03.html\)](https://www.nltk.org/book/ch03.html)
- [Different Stemmers in NLTK \(https://www.nltk.org/api/nltk.stem.html\)](https://www.nltk.org/api/nltk.stem.html)
- [German lemmatizer \(https://github.com/WZBSocialScienceCenter/germalemma\)](https://github.com/WZBSocialScienceCenter/germalemma)

In [ ]:

1	
---	--