

PDF & Twitter Scraping using Python

Sunbin Cho | Tom Bugdalle | Sergio López | Loukas Sfountouris

Université Paris I Panthéon-Sorbonne
Methods in Empirical Development Economics
Python Project

January 16, 2022



Overview

1 Topic

2 Methodology

- PDF Scraping + NLP
- Twitter Scraping + NLP

3 Data visualization

- PDF Graphs
- Twitter Graphs

4 Discussion

- Word-Cloud

5 References

The topic

- ↔ Public awareness due to climate change: The case of COP26 agenda
- ↔ Urgency of the subjects discussed on the annual COP, and the implications of the decisions that are taken there (Official statements)
- ↔ Topic is now subject of public debate in different spheres (e.g. Twitter)

PDF Scrapping: Web Crawling

```
9 # In[6]:
10 #pip install selenium
11 import requests
12 from bs4 import BeautifulSoup
13
14
15 # In[7]:
16
17 from selenium import webdriver
18 import time
19
20
21
22 # In[68]:
23
24 #Call the web page and maximise
25 driver = webdriver.Chrome()
26 web_url = 'https://www.unfccc.int/sites/submissionsstaging/Pages/Home.aspx'
27 time.sleep(1)
28 driver.get(web_url)
29 driver.maximize_window()
30 time.sleep(5)
```

Web Crawling (UNFCCC)
for downloading all the
opening statements

```
35 # In[69]:
36
37
38 #Click and select years
39 x_path = '//*[@id="searchyear"]/button'
40 year_button = driver.find_element_by_xpath(x_path).click()
41 driver.find_element_by_xpath('//*[@id="open"]').click()
42 #driver.find_element_by_xpath('//*[@id="2022"]').click()
43 #driver.find_element_by_xpath('//*[@id="2021"]').click()
44 driver.find_element_by_xpath('//*[@id="2020"]').click()
45 driver.find_element_by_xpath('//*[@id="2019"]').click()
46 driver.find_element_by_xpath('//*[@id="2018"]').click()
47 driver.find_element_by_xpath('//*[@id="2017"]').click()
48 driver.find_element_by_xpath('//*[@id="2016"]').click()
49 driver.find_element_by_xpath('//*[@id="2015"]').click()
50
51
52
53 # ### Search "opening"
54
55 # In[70]:
56
57
58 search_box = driver.find_element_by_xpath('//*[@id="accordion"]/div[2]/div/div/div[2]/div/div/input')
59 search_box.click()
60 search_box.send_keys("opening")
61 driver.find_element_by_xpath('//*[@id="accordion"]/div[2]/div/div/div[2]/div/div/button[1]').click()
```

Select all years and search
“Opening”

Web Crawling (UNFCCC) for downloading all the opening statements

The screenshot shows a web browser window with the URL www4.unfccc.int/sites/submissionsstaging/Pages/Home.aspx. The page features a search bar with the term 'opening' entered. Below the search bar, there are filters for 'SELECTED TAGS' including years from 2022 to 2015. The main content area displays two sections: 'CALLS FOR SUBMISSIONS, ELECTIONS AND STATEMENTS FOR CONSIDERATION AT UPCOMING SESSIONS (8)' and 'CALLS FOR SUBMISSIONS, ELECTIONS AND STATEMENTS CONSIDERED AT PAST SESSIONS (8)'. The first section shows details for 'General statements' with a deadline of 25/05/2016 and a session name of 'SBI.44'. A 'START SUBMISSION' button is visible. The second section shows details for 'General statements' with a deadline of 25/05/2016 and a session name of 'SBI.44'. A table below lists parties and their submission dates.

Parties	Submission date
on behalf of Climate Action Network on behalf of Environmental NGOs Constituency	16/05/2016

Download and convert into pdf

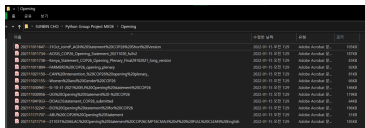
Download all the COP26 opening statement pdf

```
89 # In[ ]:
90
91
92 pdf_link = []
93
94 for a in upcoming_soup.find_all('a', href=True):
95     a_string = a['href']
96     if a_string.find('.pdf') >= 0:
97         if a_string.find('COP26') >= 0:
98             pdf_link.append(a_string)
99
```

Line 136: Read web file and save opening statements into pdf files

```
132 # In[ ]:
133 ## Save opening statements pdf files
134 ## Read web file
135 #pip install urllib3
136 import urllib.request
137
138 for text in sites:
139     report = text.replace(' ', '%20')
140     webFile = urllib.request.urlopen(report)
141     pdfFile = open(report.split('/')[1], 'wb')
142     pdfFile.write(webFile.read())
143     webFile.close()
144     pdfFile.close()
145
```

From 14 web files, we excluded 2 of spanish version. We'll do the analysis with 12 Opening Statements



File Name	Size	Modified	Owner	Group
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin
202111011847 - COP26 Opening Statement.pdf	10000	2021-11-01 18:47	Admin	Admin

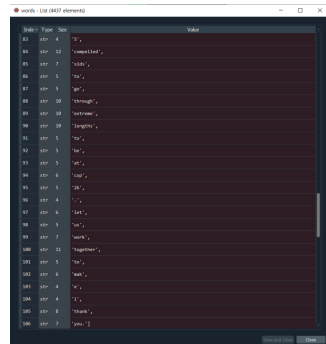
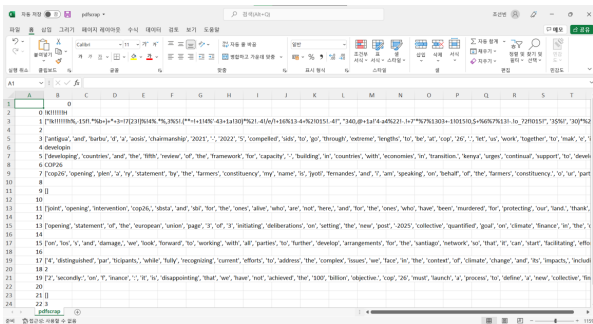
Then we use those ones, to scrap the text inside

```
147 ## Scraping PDF
148
149 path = r"YOUR DIRECTORY"
150 dirs = os.listdir(path)
151 mytext=[]
152 for item in dirs:
153     item_path = os.path.join(path, item)
154     with open(item_path, mode='rb') as f:
155         reader = PdfFileReader(f)
156         for page in reader.pages:
157             pass
158             text = page.extractText().encode('utf-8')
159             mytext.append(text)
160             print(text)
161             text = text.split()
162             print(text)
163
```

Extract/Split text + Word list

Line 144: Extract and Split the text from the Opening statements

Line 175: import CSV file
and create word list



Download required packages & Data cleaning (NLP)

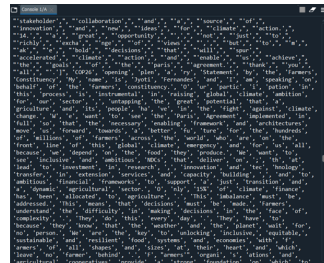
Line 187: Download all packages needed

```
187 # In[ ]: From Loukas Download packages for analysis
188 #ANALYSIS
189 import nltk
190 nltk.download('punkt')
191 nltk.download('wordnet')
192 from nltk import sent_tokenize, word_tokenize
193 from nltk.stem.snowball import SnowballStemmer
194 from nltk.stem.wordnet import WordNetLemmatizer
195 from nltk.corpus import stopwords
196 from nltk.tag import pos_tag
197
198 import pandas as pd
199 import numpy as np
200 import re
201 #pip install spacy
202 import spacy
203
204 nlp = spacy.load('en_core_web_lg')
205
206 import spacy
207 from spacy.lang.en.examples import sentences
208
209 nlp = spacy.load("en_core_web_sm")
```

Line 211: Data cleaning

```
211 # Data Cleansing
212 df2=pd.DataFrame(words)
213 df3=df2[0] #To extract the words from the Column 0 (The name of the column that includes words is 0)
214
215 all_sentences = []
216
217 for word in df3:
218     all_sentences.append(word)
219
220 all_sentences
221 #df1 = df.to_string()
222 #df_split = df1.split()
223 #df_split
224
225 lines = list()
226 for line in all_sentences:
227     words = line.split()
228     for w in words:
229         lines.append(w)
230
231 print(lines)
```

To gather all the words from the word list



A screenshot of a document showing a list of words extracted from a text. The words are displayed in a grid-like format, with some words appearing in quotes. The words include: 'stakeholder', 'collaboration', 'and', 'a', 'source', 'of', 'innovation', 'and', 'new', 'ideas', 'for', 'climate', 'action', 'in', 'the', 'great', 'opportunity', 'not', 'just', 'to', 'a', 'richly', 'each', 'age', 'of', 'views', 'but', 'to', 'a', 'new', 'decisions', 'will', 'shape', 'the', 'future', 'accelerated', 'climate', 'action', 'and', 'enable', 'us', 'achieve', 'the', 'goal', 'of', 'the', 'Paris', 'Agreement', 'Thank', 'you', 'all', 'for', 'your', 'commitment', 'to', 'this', 'Statement', 'by', 'the', 'farmers', 'Constituency', 'My', 'name', 'is', 'Jyoti', 'Fernandes', 'and', 'I', 'am', 'speaking', 'on', 'behalf', 'of', 'the', 'farmers', 'Constituency', 'of', 'the', 'Paris', 'Agreement', 'This', 'process', 'is', 'instrumental', 'in', 'raising', 'global', 'climate', 'ambition', 'for', 'our', 'sector', 'untapping', 'the', 'great', 'potential', 'that', 'a', 'agriculture', 'and', 'its', 'people', 'have', 'in', 'the', 'fight', 'against', 'climate', 'change', 'We', 'want', 'to', 'see', 'the', 'Paris', 'Agreement', 'implemented', 'in', 'full', 'on', 'that', 'the', 'necessary', 'enabling', 'framework', 'and', 'architectures', 'have', 'us', 'forward', 'towards', 'a', 'better', 'for', 'future', 'for', 'the', 'hundreds', 'of', 'millions', 'of', 'farmers', 'across', 'the', 'world', 'who', 'are', 'on', 'the', 'front', 'line', 'of', 'this', 'global', 'climate', 'emergency', 'and', 'for', 'us', 'all', 'because', 'we', 'depend', 'on', 'the', 'food', 'they', 'produce', 'We', 'want', 'to', 'see', 'inclusive', 'and', 'ambitious', 'NDCs', 'that', 'deliver', 'on', 'the', 'at', 'lead', 'to', 'investment', 'in', 'research', 'innovation', 'and', 'knowledge', 'transfer', 'in', 'extension', 'services', 'and', 'capacity', 'building', 'and', 'to', 'ambitious', 'financial', 'frameworks', 'to', 'support', 'a', 'just', 'transition', 'and', 'a', 'dynamic', 'agricultural', 'sector', 'to', 'only', 'the', 'of', 'climate', 'finance', 'has', 'been', 'allocated', 'to', 'agriculture', 'This', 'imbalance', 'must', 'be', 'addressed', 'This', 'means', 'that', 'decisions', 'must', 'be', 'made', 'Farmers', 'understand', 'the', 'difficulty', 'in', 'making', 'decisions', 'in', 'the', 'face', 'of', 'complexity', 'They', 'do', 'this', 'every', 'day', 'They', 'have', 'to', 'because', 'they', 'know', 'that', 'the', 'weather', 'and', 'the', 'planet', 'will', 'for', 'us', 'person', 'We', 'are', 'the', 'key', 'to', 'unlocking', 'inclusive', 'equitable', 'sustainable', 'and', 'resilient', 'food', 'systems', 'and', 'economies', 'with', 'farmers', 'of', 'all', 'shapes', 'and', 'sizes', 'at', 'their', 'heart', 'and', 'which', 'have', 'on', 'farm', 'business', 'regain', 'a', 'strong', 'and', 'agricultural', 'cooperatives', 'provide', 'a', 'strong', 'foundation', 'on', 'which', 'to',

N.B. We installed spacy sm and lg model before (source: https://anaconda.org/conda-forge/spacy-model-en_core_web_sm)

From line 288: Data Visualization (NLP)

1. Counting words

```
288 # In[ ]: VISUALIAZION
289
290 #count words again with nltk
291
292 df = df2[0].value_counts()
293
294 #df
295 #df['freq'] = df.groupby(0)[0].transform('count')
296 #df['freq'] = df.groupby(0)[0].transform('count')
297 #df.sort_values(by = ('freq'), ascending=False)
298 #This will give frequencies of our words
299
300 from nltk.probability import FreqDist
301
302 freqdoctor = FreqDist()
303
304 for words in df:
305     freqdoctor[words] += 1
306
307 freqdoctor
```

2. Creating plot with top overall words

```
309 ##plot
310
311 import matplotlib.pyplot as plt; plt.rcdefaults()
312 import numpy as np
313 import matplotlib.pyplot as plt
314 import seaborn as sns
315
316 #This is a simple plot that shows the top 20 words being used
317 #df.plot(20)
318
319 df = df[:20,]
320 plt.figure(figsize=(10,5))
321 sns.barplot(df.values, df.index, alpha=0.8)
322 plt.title('Top Words Overall')
323 plt.ylabel('Word from Opening Statements', fontsize=12)
324 plt.xlabel('Count of Words', fontsize=12)
325 plt.show()
```

Twitter Scrapping: Pre-requisites

To extract tweets you need to apply at <https://developer.twitter.com/en/products/twitter-api/academic-research> and present your project

Which will give you a bearer token, that will allow you to program Python queries to obtain information from Twitter



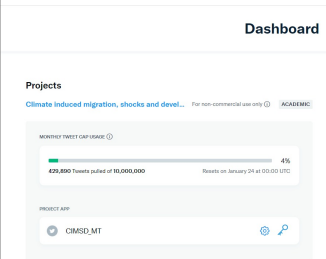
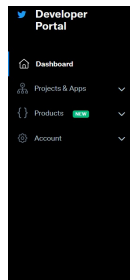
Here are your keys & tokens

For security, this will be the last time we'll display these. If something happens, you can always regenerate them.

API key ⓘ

API secret key ⓘ

Bearer token ⓘ



Packages and Environment variables

Required Packages:

```
# -- coding: utf-8 --  
import pandas as pd  
import requests  
import time
```

Line 24: Put the previously obtained bearer token The 'headers' object uses the token and returns headers we will use for scraping.

```
23 #BEARER TOKEN  
24 bearer_token = "."  
25 #CREATE URL  
26 search_url = "https://api.twitter.com/2/tweets/search/all"  
27 #CREATE HEADERS  
28 headers = {"Authorization": "Bearer {}".format(bearer_token)}  
29
```

Scraping

By using the requests package, execute queries on the desired keywords, and set specific days to concentrate on days of the COP meeting, among other settings

```
query_params = {'query': 'lang:en (cop26)',
                'start_time': {},
                "end_time": yesterday.isoformat("T") + "Z",
                'max_results': 500,
                'expansions': 'author_id,in_reply_to_user_id,geo.place_id',
                'tweet.fields': 'id,text,author_id,in_reply_to_user_id,geo,conversation_id,created_at',
                'user.fields': 'id,name,username,created_at,description,public_metrics,verified',
                'place.fields': 'full_name,id,country,country_code,geo,name,place_type'}
```

Store the results for the analysis

```
58     #Creating a list of responses from the scraping
59     response = requests.request("GET", search_url, headers=headers, params=query_params)
60     #Technical threshold to avoid potential errors in the process
61     if response.status_code == 429:
62         time.sleep(60)
63     response = requests.request("GET", search_url, headers=headers, params=query_params)
64     #Falsification like text (Automatically python sends back an error if response code!=200)
65     if response.status_code != 200:
66         raise Exception(response.status_code, response.text)
67     #The "clean" responses are saved in a different list (in JSON format)
68     r = response.json()
69     #Saving the scrapped data as text
70     for e in r["data"]:
71         #print(e["text"])
72         mytext.append(e["text"])
73     #Dataframing the text, with the help of Panda, in order to be saved in an editing format whether in csv or in xlsx
74     df = pd.DataFrame(r["data"]) #correct one
75     #Exporting data in an Excel and/or csv
76     df.to_excel('twitterscraping.xlsx')
77     df.to_csv('twitterscraping.csv')
```

NLP-Step 1: Importing packages and loading spacy model

Required packages:

```
#ANALYSIS
import nltk
nltk.download('punkt')
nltk.download('wordnet')
from nltk.stem.snowball import SnowballStemmer
import pandas as pd
import re
import spacy
nlp = spacy.load('en_core_web_lg')
```

Split the tweets into lines...

Index	Type	Size	Value
0	str	127	The latest Architects News from piqui! https://t.co/W1ZTwuloQ0 Thanks t ...
1	str	195	Caroline May and Charles Winch contributed an article published in Law ...
2	str	223	If @ScotGovFM was THAT CONCERNED as she claims. The Louisa Jordan woul ...
3	str	278	Tracking today: US @CommerceGov holds civil nuclear trade advisory com ...
4	str	143	The latest Journal of #Sustainability #CSR & #Marketing! https://t.co/...
5	str	101	Syria: How an ecological disaster can unfold in mere decades.

NLP-Step 2: Stemming + Removing Stop words

Remove RT (Retweet acronyms)

```
49 lines = [re.sub(r'^A-Za-z0-9+', '', x) for x in lines]
50
51 lines
52
53 lines2 = []
54
55 for word in lines:
56     if word != '':
57         lines2.append(word)
58
59 #Removing RT (Retweet condition within tweet)
60 for word in lines:
61     if word == 'RT':
62         lines2.remove(word)
63
```

Stem words to their root

```
66 #This is stemming the words to their root
67 from nltk.stem.snowball import SnowballStemmer
68
69 # The Snowball Stemmer requires that you pass a language parameter
70 s_stemmer = SnowballStemmer(language='english')
71
72 stem = []
73 for word in lines2:
74     stem.append(s_stemmer.stem(word))
75
76 stem
77
78 #Part 3
79 #Removing all Stop Words
80 import spacy
81 nlp = spacy.load('en_core_web_lg')
82 stem2 = []
83
84 for word in stem:
85     if word not in nlp.Defaults.stop_words:
86         stem2.append(word)
87
```

Remove all Stop words

```
80 import spacy
81 nlp = spacy.load('en_core_web_lg')
82 stem2 = []
83
84 for word in stem:
85     if word not in nlp.Defaults.stop_words:
86         stem2.append(word)
87
88 stem2
```

NLP-Step 3: Count words and Visualization

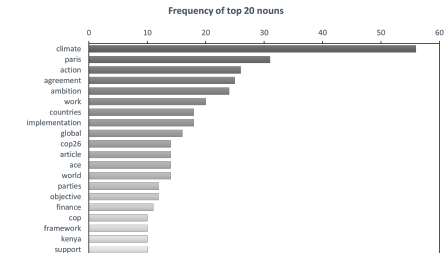
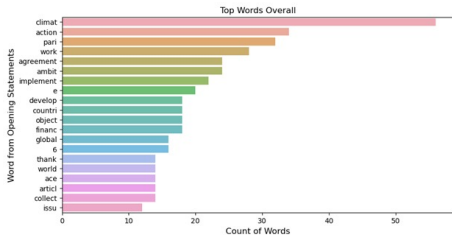
Count word frequency

```
105 from nltk.probability import FreqDist
106
107 freqdoctor = FreqDist()
108
109 for words in df:
110     freqdoctor[words] += 1
111
112 freqdoctor
```

Create and label graphs

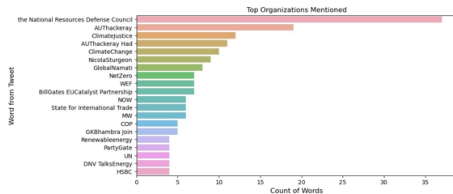
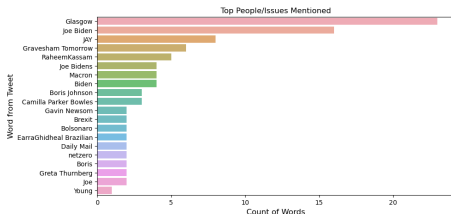
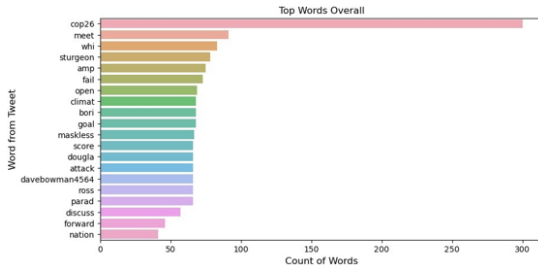
```
152 df7 = df7['Word'].value_counts()
153 df = df7[:20,]
154 plt.figure(figsize=(10,5))
155 sns.barplot(df.values, df.index, alpha=0.8)
156 plt.title('Top Organizations Mentioned')
157 plt.ylabel('Word from Tweet', fontsize=12)
158 plt.xlabel('Count of Words', fontsize=12)
159 plt.show()
```

PDF Graphs



Note: Excel post-editing

Twitter Graphs



PDF Scraping

- Climate (56) is the most used noun with a lead of 25 words, which is in line with one's expectations for speeches at the COP26
- Top words like action (26), agreement (25), implementation (18) and article (14) characterize, at least formally, a strongly solution-oriented approach
- Negative words like devastating (4), damage (4), emergency (2), disaster (2) or instability (2) are used very rarely indicating a reluctance of producing unpleasant emotions
- The relatively frequent use of the country name Kenya (10) seems to be idiosyncratic, as the next country name in the ranking is India (2) with a much lower frequency

Twitter Scraping

- At first glance, we observe a diverse composition of apparently unrelated and incomplete Twitter keywords
- Presumably, this is a result of the high variety of differing motives behind each and every individual tweet
- While top words like “fail” and “attack” can typically characterize a complaining Twitter user, they also create the impression of a rather high degree of pessimism most likely deriving from low expectations for the outcome of cop26
- The relatively frequent use of the word “maskless” in relation to cop26, suggests that twitter users might observe and even criticize a loose enforcement of mask rules among cop26 participants

Word-Cloud: A duality ? (using wordart.com)



PDF



Twitter

Merci!

References

- ❶ Fenniak, M., (2006) PyPDF, All rights reserved
- ❷ Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585, 357–362.
<https://doi.org/10.1038/s41586-020-2649-2>
- ❸ Honnibal, M., Montani, I. (2017). spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing.
- ❹ McKinney, W., et al. (2010). Data structures for statistical computing in python. In *Proceedings of the 9th Python in Science Conference* (Vol. 445, pp. 51–56)
- ❺ Porter, M. F. (2001). Snowball: A language for stemming algorithms
Published online (Accessed 11.03.2008, 15.00h)
- ❻ Richardson, L. (2007). Beautiful soup documentation. April.
- ❼ Van Rossum, G. (2020). The Python Library Reference, release 3.8.2. Python Software Foundation.