Data-X Spring 2019: Homework 7

Webscraping

In this homework, you will do some exercises with web-scraping.

Name: Loi Chii Lek

SID: 3034453674

Fun with Webscraping & Text manipulation

1. Statistics in Presidential Debates

Your first task is to scrape Presidential Debates from the Commission of Presidential Debates website: https://www.debates.org/voter-education/debate-transcripts/ (https://www.debates.org/voter-education/debate-transcripts/)

To do this, you are not allowed to manually look up the URLs that you need, instead you have to scrape them. The root url to be scraped is the one listed above, namely: https://www.debates.org/voter-education/debate-transcripts/ (https://www.debates.org/voter-education/debate-transcripts/)

- 1. By using requests and BeautifulSoup find all the links / URLs on the website that links to transcriptions of **First Presidential Debates** from the years [1988, 1984, 1976, 1960]. In total you should find 4 links / URLs that fulfill this criteria. **Print the urls**.
- 2. When you have a list of the URLs your task is to create a Data Frame with some statistics (see example of output below):
 - A. Scrape the title of each link and use that as the column name in your Data Frame.
 - B. Count how long the transcript of the debate is (as in the number of characters in transcription string). Feel free to include \ characters in your count, but remove any breakline characters, i.e. \n . You will get credit if your count is +/- 10% from our result.
 - C. Count how many times the word war was used in the different debates. Note that you have to convert the text in a smart way (to not count the word warranty for example, but counting war., war!, war, or War etc.
 - D. Also scrape the most common used word in the debate, and write how many times it was used. Note that you have to use the same strategy as in C in order to do this.

Print your final output result.

Tips:

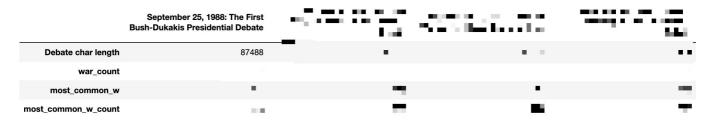
In order to solve the questions above, it can be useful to work with Regular Expressions and explore methods on strings like .strip(), .replace(), .find(), .count(), .lower() etc. Both are very powerful tools to do string processing in Python. To count common words for example I used a Counter object and a Regular expression pattern for only words, see example:

```
from collections import Counter
  import re

counts = Counter(re.findall(r"[\w']+", text.lower()))
```

Read more about Regular Expressions here: https://docs.python.org/3/howto/regex.html)

Example output of all of the answers to Question 1.2:



In [23]:

```
import bs4 as bs
import requests
import re

source = requests.get("https://www.debates.org/voter-education/debate-transcripts/")
soup = bs.BeautifulSoup(source.content, features='html.parser', parse_only=bs.SoupStrainer('a'))
```

In [74]:

September 25, 1988: The First Bush-Dukakis Presidential Debate /voter-education/debate-transcripts/september-25-1988-debate-transcript/

October 7, 1984: The First Reagan-Mondale Presidential Debate /voter-education/debate-transcripts/october-7-1984-debate-transcript/

September 23, 1976: The First Carter-Ford Presidential Debate /voter-education/debate-transcripts/september-23-1976-debate-transcript/

September 26, 1960: The First Kennedy-Nixon Presidential Debate /voter-education/debate-transcripts/september-26-1960-debate-transcript/

In [154]:

```
link_text = {}

for key in links:
    link_source = requests.get("https://www.debates.org" + links[key])
    link_soup = bs.BeautifulSoup(link_source.content, features='html.parser')
    link_text[key] = link_soup.find(id='content-sm')
```

In [155]:

```
link_words = {}

for key in link_text:
    str_append = ""
    for p in link_text[key].find_all('p'):
        str_append += " " + p.get_text().rstrip()
        link_words[key] = str_append
```

In [193]:

```
import pandas as pd
import string
from collections import Counter
df = pd.DataFrame()
char_length = []
war_count = []
most_common_word = []
most common word count = []
for key in link words:
    df[key] = key
    char_length.append(len(link_words[key]))
    no punctuation string = link words[key].translate(str.maketrans('', '', string.punc
tuation)).lower()
    count = sum(1 for _ in re.finditer(r'\b%s\b' % re.escape("war"), no_punctuation_str
ing))
   war_count.append(count)
   wc = Counter(no_punctuation_string.split())
    most common word.append(wc.most common(1)[0][0])
    most_common_word_count.append(wc.most_common(1)[0][1])
df.loc['Debate char length'] = char_length
df.loc['war_count'] = war_count
df.loc['most_common_w'] = most_common_word
df.loc['most common w count'] = most common word count
df.head()
```

Out[193]:

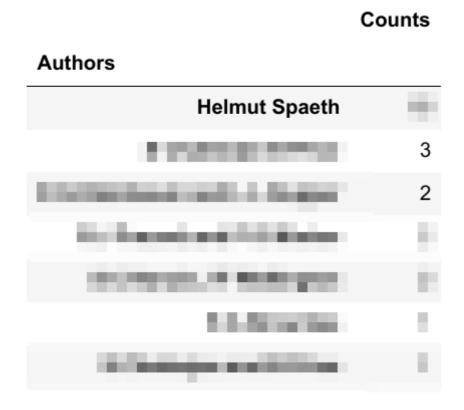
	September 25, 1988: The First Bush-Dukakis Presidential Debate	October 7, 1984: The First Reagan- Mondale Presidential Debate	September 23, 1976: The First Carter-Ford Presidential Debate	September 26, 1960: The First Kennedy-Nixon Presidential Debate
Debate char length	87691	86782	80796	60974
war_count	7	2	7	3
most_common_w	the	the	the	the
most_common_w_count	799	867	856	779

2. Download and read in specific line from many data sets

Scrape the first 27 data sets from this URL http://people.sc.fsu.edu/~jburkardt/datasets/regression/ (i.e. x01.txt - x27.txt). Then, save the 5th line in each data set, this should be the name of the data set author (get rid of the # symbol, the white spaces and the comma at the end).

Count how many times (with a Python function) each author is the reference for one of the 27 data sets. Showcase your results, sorted, with the most common author name first and how many times he appeared in data sets. Use a Pandas DataFrame to show your results, see example. **Print your final output result.**

Example output of the answer for Question 2:



In [249]:

```
import bs4 as bs
import requests
import re
source = requests.get("http://people.sc.fsu.edu/~jburkardt/datasets/regression/")
soup = bs.BeautifulSoup(source.content, features='html.parser', parse_only=bs.SoupStrai
ner('a'))
links = []
for link in soup:
    links.append(link['href'])
links = links[6:33]
authors = []
df2 = pd.DataFrame()
for link in links:
    link_source = requests.get("http://people.sc.fsu.edu/~jburkardt/datasets/regressio
n/" + link)
    link_soup = bs.BeautifulSoup(link_source.content)
    author name = ''
    i = 0
    for char in link_soup.get_text():
        if char == "#":
            i += 1
        if i == 5:
            author_name += char
    author_name = author_name.replace('#', "").replace(',', "").lstrip().rstrip()
    authors.append(author name)
```

In [258]:

```
author_count = Counter(authors)

author_list = []
author_list_count = []
for key, value in author_count.items():
    author_list.append(key)
    author_list_count.append(value)

df2['Author'] = author_list
df2['Count'] = author_list_count

df2.sort_values(by=['Count'],ascending=False)
```

['Helmut Spaeth', 'R J Freund and P D Minton', 'D G Kleinbaum and L L Kupp
er', 'K A Brownlee', 'S Chatterjee and B Price', 'S Chatterjee B Price',
'S C Narula J F Wellington']

Out[258]:

	Author	Count
0	Helmut Spaeth	16
5	S Chatterjee B Price	3
1	R J Freund and P D Minton	2
2	D G Kleinbaum and L L Kupper	2
6	S C Narula J F Wellington	2
3	K A Brownlee	1
4	S Chatterjee and B Price	1