Data-X Spring 2019: Homework 7

Webscraping

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

Name: Jack Xie

https://github.com/JackXie24/xiexiangfie_data_x_s19/tree/m (https://github.com/JackXie24/xiexiangfie_data_x_s19/tree/m

SID: 3032163590

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/

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**, 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.

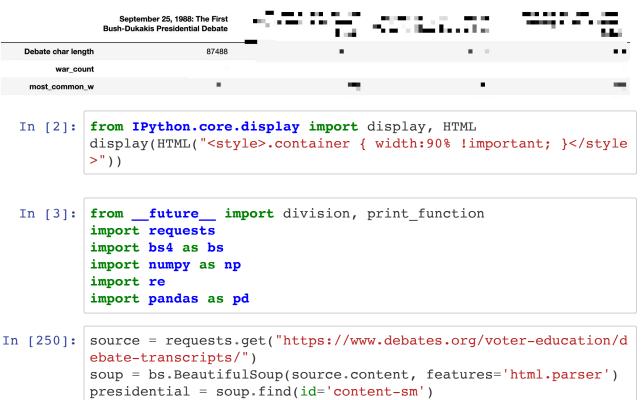
Both are very powerful tools to do string processing in Python. To count common words for example I used a Counter object and a Regular expression pattern for only words, see example:

```
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 (https://docs.python.org/3/howto/regex.html)

Example output of all of the answers to Question 1.2:



```
In [6]: years = [0, 3, 8, 11]
    links = presidential.find_all('a')
    db, dbN = [], []
    for i in links:
        db.append(i.get('href'))
        dbN.append(i.text)
    db = db[-15::]
    dbN = dbN[-15::]
    dbN + db
```

```
Out[6]: ['September 25, 1988: The First Bush-Dukakis Presidential Deba
        te',
         'October 5, 1988: The Bentsen-Quayle Vice Presidential Debat
         'October 13, 1988: The Second Bush-Dukakis Presidential Debat
         'October 7, 1984: The First Reagan-Mondale Presidential Debat
         'October 11, 1984: The Bush-Ferraro Vice Presidential Debat
         'October 21, 1984: The Second Reagan-Mondale Presidential Deb
         'September 21, 1980: The Anderson-Reagan Presidential Debat
        e',
         'October 28, 1980: The Carter-Reagan Presidential Debate',
         'September 23, 1976: The First Carter-Ford Presidential Debat
         'October 6, 1976: The Second Carter-Ford Presidential Debat
        e',
         'October 22, 1976: The Third Carter-Ford Presidential Debat
         'September 26, 1960: The First Kennedy-Nixon Presidential Deb
        ate',
         'October 7, 1960: The Second Kennedy-Nixon Presidential Debat
         'October 13, 1960: The Third Kennedy-Nixon Presidential Debat
         'October 21, 1960: The Fourth Kennedy-Nixon Presidential Deba
         '/voter-education/debate-transcripts/september-25-1988-debate
        -transcript/',
         '/voter-education/debate-transcripts/october-5-1988-debate-tr
        anscripts/',
         '/voter-education/debate-transcripts/october-13-1988-debate-t
        ranscript/',
         '/voter-education/debate-transcripts/october-7-1984-debate-tr
        anscript/',
          '/voter-education/debate-transcripts/october-11-1984-debate-t
        ranscript/',
         '/voter-education/debate-transcripts/october-21-1984-debate-t
        ranscript/',
          '/voter-education/debate-transcripts/september-21-1980-debate
        -transcript/',
         '/voter-education/debate-transcripts/october-28-1980-debate-t
        ranscript/',
         '/voter-education/debate-transcripts/september-23-1976-debate
        -transcript/',
          '/voter-education/debate-transcripts/october-6-1976-debate-tr
        anscript/',
         '/voter-education/debate-transcripts/october-22-1976-debate-t
        ranscript/',
         '/voter-education/debate-transcripts/september-26-1960-debate
        -transcript/',
         '/voter-education/debate-transcripts/october-7-1960-debate-tr
        anscript/',
         '/voter-education/debate-transcripts/october-13-1960-debate-t
        ranscript/',
```

'/voter-education/debate-transcripts/october-21-1960-debate-transcript/']

```
In [7]: | urls = []
          names = []
          for i in np.arange(len(db)):
              if i in years:
                   urls.append(db[i])
                   names.append(dbN[i])
          urls
 Out[7]: ['/voter-education/debate-transcripts/september-25-1988-debate
          -transcript/',
           '/voter-education/debate-transcripts/october-7-1984-debate-tr
          anscript/',
           '/voter-education/debate-transcripts/september-23-1976-debate
          -transcript/',
           '/voter-education/debate-transcripts/september-26-1960-debate
          -transcript/']
 In [8]: names
 Out[8]: ['September 25, 1988: The First Bush-Dukakis Presidential Deba
           'October 7, 1984: The First Reagan-Mondale Presidential Debat
           'September 23, 1976: The First Carter-Ford Presidential Debat
          e',
           'September 26, 1960: The First Kennedy-Nixon Presidential Deb
          ate']
          df = pd.DataFrame(columns=names)
In [31]:
Out[31]:
            September 25, 1988:
                                October 7, 1984:
                                                  September 23,
                                                               September 26, 1960:
                The First Bush-
                               The First Reagan-
                                                  1976: The First
                                                                The First Kennedy-
                     Dukakis
                                      Mondale
                                                    Carter-Ford
                                                                Nixon Presidential
             Presidential Debate Presidential Debate Presidential Debate
                                                                         Debate
```

```
In [75]:
         char_counts = []
         war counts = []
         most_common_words = []
         most_common_words_count = []
         for i in urls:
             db = requests.get("https://www.debates.org"+ i)
             dbsoup = bs.BeautifulSoup(db.content, features='html.parser'
         )
             script = dbsoup.find(id='content-sm')
             char = len(script.text) - script.text.count('\n')
             war = len(re.findall(r"war\b", script.text.lower()))
             counts = Counter(re.findall(r"[\w']+", script.text.lower()))
             most common word = counts.most common()[0][0]
             most common word count = counts.most common()[0][1]
             char_counts.append(char)
             war_counts.append(war)
             most_common_words.append(most_common_word)
             most_common_words_count.append(most_common_word_count)
```

```
In [82]: df.loc["char count"] = char_counts
    df.loc["war count"] = war_counts
    df.loc["most common word"] = most_common_words
    df.loc["most common word count"] = most_common_words_count
```

In [83]: df

Out[83]:

	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
char count	87488	86505	80735	60937
war count	8	2	7	3
most common word	the	the	the	the
most common word count	804	867	857	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:

Authors Helmut Spaeth 3 2

```
In [200]: | source2 = requests.get("http://people.sc.fsu.edu/~jburkardt/data
          sets/regression/")
          soup2 = bs.BeautifulSoup(source2.content, features='html.parser'
          href = soup2.find_all("a")[6:33]
          authors = []
          for i in href:
               txt = requests.get("http://people.sc.fsu.edu/~jburkardt/data
          sets/regression/" + i.get("href"))
               soupT = bs.BeautifulSoup(txt.content, features='html.parser'
           )
               for a in soupT.strings:
                   authors.append(re.split(r'\s{2,}', repr(a).replace(",",
           "").replace("\\n#", ""))[3])
          authors
Out[200]: ['Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'R J Freund and P D Minton',
            'D G Kleinbaum and L L Kupper',
            'Helmut Spaeth',
            'D G Kleinbaum and L L Kupper',
            'K A Brownlee',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'S Chatterjee and B Price',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'R J Freund and P D Minton',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'Helmut Spaeth',
            'S Chatterjee B Price',
            'S Chatterjee B Price',
            'S Chatterjee B Price',
            'S C Narula J F Wellington',
            'S C Narula J F Wellington'
In [237]: data = Counter(authors)
```

```
In [249]: df2 = pd.DataFrame.from_dict(data, orient='index').sort_values(b
    y= 0, ascending=False)
    df2.columns = ['Counts']
    df2.index.names = ['Authors']
    df2
```

Out[249]:

Counts

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

In []: