

Andrew Kim

ADS 509-01: Applied Text Mining

Assignment 1.1: Data Acquisition with APIs and Scrapping

5/15/2023

ADS 509 Module 1: APIs and Web Scrapping

This notebook has two parts. In the first part, you will scrape lyrics from AZLyrics.com. In the second part, you'll run code that verifies the completeness of your data pull.

For this assignment you have chosen two musical artists who have at least 20 songs with lyrics on AZLyrics.com. We start with pulling some information and analyzing them.

General Assignment Instructions

These instructions are included in every assignment, to remind you of the coding standards for the class. Feel free to delete this cell after reading it.

One sign of mature code is conforming to a style guide. We recommend the [Google Python Style Guide](#). If you use a different style guide, please include a cell with a link.

Your code should be relatively easy-to-read, sensibly commented, and clean. Writing code is a messy process, so please be sure to edit your final submission. Remove any cells that are not needed or parts of cells that contain unnecessary code. Remove inessential `import` statements and make sure that all such statements are moved into the designated cell.

Make use of non-code cells for written commentary. These cells should be grammatical and clearly written. In some of these cells you will have questions to answer. The questions will be marked by a "Q:" and will have a corresponding "A:" spot for you. *Make sure to answer every question marked with a Q: for full credit.*

Importing Libraries

```
In [1]: import os
import datetime
import re
```

```
# for the lyrics scrape section
import requests
import time
from bs4 import BeautifulSoup
from collections import defaultdict, Counter
import random
```

```
In [2]: # Use this cell for any import statements you add
import shutil
```

Lyrics Scrape

This section asks you to pull data by scraping www.AZLyrics.com. In the notebooks where you do that work you are asked to store the data in specific ways.

```
In [3]: artists = {'eltonjohn':"https://www.azlyrics.com/j/john.html",
                  'nirvana':"https://www.azlyrics.com/n/nirvana.html"}
# we'll use this dictionary to hold both the artist name and the link on AZlyrics
```

A Note on Rate Limiting

The lyrics site, www.azlyrics.com, does not have an explicit maximum on number of requests in any one time, but in our testing it appears that too many requests in too short a time will cause the site to stop returning lyrics pages. (Entertainingly, the page that gets returned seems to only have the song title to a [Tom Jones song](#).)

Whenever you call `requests.get` to retrieve a page, put a `time.sleep(5 + 10*random.random())` on the next line. This will help you not to get blocked. If you *do* get blocked, which you can identify if the returned pages are not correct, just request a lyrics page through your browser. You'll be asked to perform a CAPTCHA and then your requests should start working again.

Part 1: Finding Links to Songs Lyrics

That general artist page has a list of all songs for that artist with links to the individual song pages.

Q: Take a look at the `robots.txt` page on www.azlyrics.com. (You can read more about these pages [here](#).) Is the scraping we are about to do allowed or disallowed by this page? How do you know?

A:

```
In [4]: # Let's set up a dictionary of lists to hold our links
lyrics_pages = defaultdict(list)
base_url = 'https://www.azlyrics.com'
url_list = []

for artist, artist_page in artists.items() :
    # request the page and sleep
    r = requests.get(artist_page)
    time.sleep(5 + 10*random.random())
    soup = BeautifulSoup(r.text, 'html.parser')

    # now extract the links to lyrics pages from this page
    # store the links `lyrics_pages` where the key is the artist and the
    # value is a list of links.

    for link in soup.find_all('a'):
        # acquire the links
        temp_href = str(link.get('href'))
        temp_url = base_url + temp_href
        url_list.append(temp_url)

    # retain lyric links
    substring = artist + '/'
    url_filtered = [i for i in url_list if substring in i]
    url_filtered

    # save
    lyrics_pages[artist] = url_filtered
```

Let's make sure we have enough lyrics pages to scrape.

```
In [5]: for artist, lp in lyrics_pages.items() :
        assert(len(set(lp)) > 20)
```

```
In [6]: # Let's see how long it's going to take to pull these lyrics
# if we're waiting `5 + 10*random.random()` seconds
for artist, links in lyrics_pages.items() :
    print(f"For {artist} we have {len(links)}.")
    print(f"The full pull will take for this artist will take {round(len(links)*10/3600,2)} hours.")
```

For eltonjohn we have 546.
The full pull will take for this artist will take 1.52 hours.
For nirvana we have 165.
The full pull will take for this artist will take 0.46 hours.

```
In [ ]:
```

Part 2: Pulling Lyrics

Now that we have the links to our lyrics pages, let's go scrape them! Here are the steps for this part.

1. Create an empty folder in our repo called "lyrics".
2. Iterate over the artists in `lyrics_pages`.
3. Create a subfolder in lyrics with the artist's name. For instance, if the artist was Cher you'd have `lyrics/cher/` in your repo.
4. Iterate over the pages.
5. Request the page and extract the lyrics from the returned HTML file using BeautifulSoup.
6. Use the function below, `generate_filename_from_url`, to create a filename based on the lyrics page, then write the lyrics to a text file with that name.

```
In [7]: def generate_filename_from_link(link) :

    if not link :
        return None

    # drop the http or https and the html
    name = link.replace("https","").replace("http","")
    name = link.replace(".html","")

    name = name.replace("/lyrics/", "")

    # Replace useless chareacters with UNDERSCORE
    name = name.replace(":/","").replace(".", "_").replace("/", "_")

    # tack on .txt
    name = name + ".txt"

    return (name)
```

```
In [8]: # Make the lyrics folder here. If you'd like to practice your programming, add functionality
# that checks to see if the folder exists. If it does, then use shutil.rmtree to remove it and create a new one

if os.path.isdir("lyrics") :
    shutil.rmtree("lyrics")

os.mkdir("lyrics")
```

```
In [9]: url_stub = "https://www.azlyrics.com"
start = time.time()

total_pages = 20
for artist in lyrics_pages :
    # Use this space to carry out the following steps:
    # 1. Build a subfolder for the artist
    subfold = "lyrics/" + artist
    if os.path.isdir(subfold):
        shutil.rmtree(subfold + '/')
    os.mkdir(subfold)

    # 2. Iterate over the lyrics pages
    i = 0
    urls = lyrics_pages[artist]

    while i < total_pages:
        temp_url = urls[i]

        # 3. Request the lyrics page.
        # Don't forget to add a line like `time.sleep(5 + 10*random.random())`
        # to sleep after making the request
        r = requests.get(temp_url)
        time.sleep(5 + 10*random.random())

        soup = BeautifulSoup(r.text, 'html.parser')

        # 4. Extract the title and lyrics from the page.
        title = soup.find_all('b')[1].get_text()
        lyrics = soup.find_all('div')[22].get_text()

        song = title + '\n\n' + lyrics

        # 5. Write out the title, two returns ('\n'), and the lyrics. Use `generate_filename_from_url`
        # to generate the filename.
        wd = os.getcwd()
        folders = "\lyrics\\" + artist
        filepath = wd + folders
        filename = generate_filename_from_link(temp_url)
        pathname = filepath + '\\' + filename

        text_file = open(pathname, "w")
        text_file.write(song)
        text_file.close()

        i +=1

    # Remember to pull at least 20 songs per artist. It may be fun to pull all the songs for the artist
```

```
In [10]: print(f"Total run time was {round((time.time() - start)/3600,2)} hours.")
```

Total run time was 0.12 hours.

Evaluation

This assignment asks you to pull data by scraping www.AZLyrics.com. After you have finished the above sections, run all the cells in this notebook. Print this to PDF and submit it, per the instructions.

```
In [11]: # Simple word extractor from Peter Norvig: https://norvig.com/spell-correct.html
def words(text):
    return re.findall(r'\w+', text.lower())
```

Checking Lyrics

The output from your lyrics scrape should be stored in files located in this path from the directory: `/lyrics/[Artist Name]/[filename from URL]`. This code summarizes the information at a high level to help the instructor evaluate your work.

```
In [13]: artist_folders = os.listdir("lyrics/")
artist_folders = [f for f in artist_folders if os.path.isdir("lyrics/" + f)]

for artist in artist_folders :
    artist_files = os.listdir("lyrics/" + artist)
    artist_files = [f for f in artist_files if 'txt' in f or 'csv' in f or 'tsv' in f]

    print(f"For {artist} we have {len(artist_files)} files.")

    artist_words = []

    for f_name in artist_files :
        with open("lyrics/" + artist + "/" + f_name) as infile :
            artist_words.extend(words(infile.read()))

    print(f"For {artist} we have roughly {len(artist_words)} words, {len(set(artist_words))} are unique.")
```

For eltonjohn we have 20 files.
For eltonjohn we have roughly 4889 words, 911 are unique.
For nirvana we have 20 files.
For nirvana we have roughly 4290 words, 574 are unique.

```
In [ ]:
```