

Homework 03

⚠️ Before you start ⚠️

Duplicate this Jupyter Notebook in your week-03 folder (right-click -> Duplicate) and then add your last name to the beginning of it (ie. hw-03-blevins.ipynb - otherwise you risk having all your work overwritten when you try to sync your GitHub repository with your instructor's repository.

⚠️ No, seriously: check the name of this file. Is it the copy you made? (ie. hw-03-blevins.ipynb). If so, you can proceed ⚠️

Student Name: Jiawen Huang

The Data

You're going to analyze several historical documents in this homework. In keeping with the theme of our first unit for the semester, **Slavery and Data**, I've chosen two 19th-century narratives written by formerly enslaved people: [Sojourner Truth](#) and [Henry "Box" Brown](#).

You should have the following files:

- `hw-03-yourlastname.ipynb` (your working version Jupyter Notebook)
- `truth.txt` (Sojourner Truth's narrative)
- `brown.txt` (Henry Brown's narrative)

Load and Process the Data

Use the `open()` and `read()` functions to get the content of each of these files into Python, assigning them the corresponding variable names of `truth_fulltext` and `brown_fulltext`.

```
In [57]: truth_fulltext = open('truth.txt', mode='r', encoding='utf-8').read()
brown_fulltext = open('brown.txt', mode='r', encoding='utf-8').read()
```

In the next two code cells, write `print()` statements that:

- Print the **first 500 characters** of Truth's narrative.
- Print characters **5000 to 6000** of Brown's narrative.

Hint: use the index and slice approaches for strings:

<https://melaniewalsh.github.io/Intro-Cultural-Analytics/02-Python/06-String-Methods.html>.

```
In [59]: print(truth_fulltext[0:500])
```

NARRATIVE OF SOJOURNER TRUTH

HER BIRTH AND PARENTAGE.

THE subject of this biography, SOJOURNER TRUTH, as she now calls herself-but whose name, originally, was Isabella-was born, as near as she can now calculate, between the years 1797 and 1800. She was the daughter of James and Betsey, slaves of one Colonel Ardinburgh, Hurley, Ulster County, New York.

Colonel Ardinburgh belonged to that class of people called Low Dutch.

Of her first master, she can give no account, as she must have been

```
In [60]: print(brown_fulltext[5000:6000])
```

of pity, indignation and
horror.

I first drew the breath of life in Louisa County, Va., forty-five miles from the city of Richmond, in the year 1816. I was born a slave. Not because at the moment of my birth an angel stood by, and declared that such was the will of God concerning me; although in a country whose most honored writings declare that all men have a right to liberty, given them by their Creator, it seems strange that I, or any of my brethren, could have been born without this inalienable right, unless God had thus signified his departure from his usual rule, as described by our fathers. Not, I say, on account of God's willing it to be so, was I born a slave, but for the reason that nearly all the people of this country are united in legislating against heaven, and have contrived to vote down our heavenly father's rules, and to substitute for them, that cruel law which binds the chains of slavery upon one sixth part of the inhabitants of this land. I was born a slave! and wh

In the next code cell complete the following:

- Look at the printed out "slice" of Brown's narrative. Make a new variable and assign it a value of **Brown's birth year**.
- Make a new variable and assign it a value of: **how old Henry Brown would have been in the year 1860**.
- Write a **print statement** using your new variable that says how old Henry Brown would have been in 1860.

```
In [62]: brown_birth=1816
brown_age1860= 1860 - brown_birth
print(f"Henry Brown was {brown_age1860} in 1860")
```

Henry Brown was 44 in 1860

Suppose we want to compare how long each narrative is measured by the number of lines in each text. First, use the `split()` function for each narrative to

break it apart by each new line. The new line character is `\n`. Make two new variables storing a list of the broken apart text: `truth_lines` and `brown_lines`.

```
In [64]: truth_lines= truth_fulltext.split('\n')
brown_lines= brown_fulltext.split('\n')
```

Which narrative has more lines? You can calculate how many lines are in each narrative through the `len()` function which will calculate the **length** of each list of lines you made in the previous section.

- Write two `print()` statements to show **how many lines are in each narrative**.
- Add a third `print()` statement that calculates **the difference between these two narratives measured by their number of lines**.

```
In [66]: print(f"There are {len(truth_lines)} lines in truth.")
print(f"There are {len(brown_lines)} lines in brown.")
print(f"truth have {len(truth_lines)-len(brown_lines)} more lines than brown")
```

There are 3627 lines in truth.
There are 2223 lines in brown.
truth have 1404 more lines than brown

Combine the `len()` and `comparison` functions with an `if statement` to print either `Sojourner Truth's narrative has more lines` or `Henry Brown's narrative has more lines` based on which has more lines.]

```
In [68]: if len(truth_lines)> len(brown_lines):
    print("Sojourner Truth's narrative has more lines")
elif len(truth_lines)== len(brown_lines):
    print("equal")
else :
    print("Henry Brown's narrative has more lines")
```

Sojourner Truth's narrative has more lines

Counting Word Frequency

Look at the code below from Melanie Walsh's "Anatomy of a Python Script" that she used to calculate the most frequently occurring words in a novel "The Yellow Wallpaper." You are going to use this code as a starting point but change it to apply this same approach to the two texts we've been working with. Your goal: **compare the most frequently occurring words in both Truth's narrative and Brown's narrative**.

Note: don't edit Walsh's code cell directly. Instead, copy and paste the code into **the two empty code cells below it** that you can then edit. If you accidentally overwrite it and need to find the original, you can [copy it from the original tutorial](#).

Adjustments you'll need to make to Walsh's code:

- Open the right .txt file.

- Find the most frequent **20 words** instead of 40 words.

```
In [71]: #Walsh's Code - copy this into a new code cell
import re
from collections import Counter

def split_into_words(any_chunk_of_text):
    lowercase_text = any_chunk_of_text.lower()
    split_words = re.split(r"\W+", lowercase_text)
    return split_words

filepath_of_text = "The-Yellow-Wallpaper_Charlotte-Perkins-Gilman.txt"
number_of_desired_words = 40

stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you',
'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her', 'hers',
'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves',
'what', 'which', 'who', 'whom', 'this', 'that', 'these', 'those', 'am', 'is',
'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do',
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
'u', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between',
'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from',
'up', 'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then',
'once', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each',
'few', 'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own',
'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don',
'should', 'now', '']

full_text = open(filepath_of_text, encoding="utf-8").read()

all_the_words = split_into_words(full_text)
meaningful_words = [word for word in all_the_words if word not in stopwords]
meaningful_words_tally = Counter(meaningful_words)
most_frequent_meaningful_words = meaningful_words_tally.most_common(number_of_desired_words)

most_frequent_meaningful_words
```

```
-----  
FileNotFoundError                         Traceback (most recent call last)  
Cell In[71], line 26  
      11     number_of_desired_words = 40  
      12     stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves',  
      13     'you', 'your', 'yours',  
      14     'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her',  
      15     'hers',  
      16     'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves',  
      17     (...)  
      18     'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same',  
      19     'so',  
      20     'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', 'should',  
      21     'now', 've', 'll', 'amp', 'would', 'one']  
---> 26 full_text = open(filepath_of_text, encoding="utf-8").read()  
      27 all_the_words = split_into_words(full_text)  
      28 meaningful_words = [word for word in all_the_words if word not in stopwords]  
  
File ~\anaconda3\Lib\site-packages\IPython\core\interactiveshell.py:324, in _modified_open(file, *args, **kwargs)  
    317 if file in {0, 1, 2}:  
    318     raise ValueError(  
    319         f"IPython won't let you open fd={file} by default "  
    320         "as it is likely to crash IPython. If you know what you are doing,"  
    321         "you can use builtins' open."  
    322     )  
--> 324 return io_open(file, *args, **kwargs)  
  
FileNotFoundError: [Errno 2] No such file or directory: 'The-Yellow-Wallpaper_Ch  
arlotte-Perkins-Gilman.txt'
```

```
In [72]: import re  
from collections import Counter  
  
def split_into_words(any_chunk_of_text):  
    lowercase_text = any_chunk_of_text.lower()  
    split_words = re.split(r"\W+", lowercase_text)  
    return split_words  
  
filepath_of_text = "truth.txt"  
number_of_desired_words = 20  
  
stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you',  
    'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her', 'hers',  
    'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves',  
    'what', 'which', 'who', 'whom', 'this', 'that', 'these', 'those', 'am', 'is',  
    'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do', 'do',  
    'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'u',  
    'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',  
    'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up',  
    'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once',  
    'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'm',  
    'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same', 'so',  
    'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', 'should', 'now',  
    'full_text = open(filepath_of_text, encoding="utf-8").read()
```

```

all_the_words = split_into_words(full_text)
meaningful_words = [word for word in all_the_words if word not in stopwords]
meaningful_words_tally = Counter(meaningful_words)
most_frequent_meaningful_words = meaningful_words_tally.most_common(number_of_desired_words)

most_frequent_meaningful_words

```

Out[72]:

```

[('god', 128),
 ('isabella', 114),
 ('time', 114),
 ('could', 104),
 ('master', 78),
 ('go', 69),
 ('good', 67),
 ('said', 67),
 ('mr', 67),
 ('mother', 65),
 ('see', 64),
 ('much', 57),
 ('found', 53),
 ('like', 51),
 ('never', 50),
 ('well', 50),
 ('place', 50),
 ('son', 49),
 ('little', 49),
 ('new', 48)]

```

In [74]:

```

import re
from collections import Counter

def split_into_words(any_chunk_of_text):
    lowercase_text = any_chunk_of_text.lower()
    split_words = re.split(r"\W+", lowercase_text)
    return split_words

filepath_of_text = "brown.txt"
number_of_desired_words = 20

stopwords = ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you',
'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'her', 'hers',
'herself', 'it', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves',
'what', 'which', 'who', 'whom', 'this', 'that', 'these', 'those', 'am', 'is',
'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'do',
'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as',
'while', 'of', 'at', 'by', 'for', 'with', 'about', 'against', 'between', 'into',
'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up',
'in', 'out', 'on', 'off', 'over', 'under', 'again', 'further', 'then', 'once',
'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few',
'most', 'other', 'some', 'such', 'no', 'nor', 'not', 'only', 'own', 'same',
'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', 'should', 'now',
'full_text = open(filepath_of_text, encoding="utf-8").read()

all_the_words = split_into_words(full_text)
meaningful_words = [word for word in all_the_words if word not in stopwords]
meaningful_words_tally = Counter(meaningful_words)
most_frequent_meaningful_words = meaningful_words_tally.most_common(number_of_desired_words)

most_frequent_meaningful_words

```

```
Out[74]: [('man', 86),  
          ('slave', 85),  
          ('slavery', 83),  
          ('master', 81),  
          ('upon', 80),  
          ('slaves', 75),  
          ('us', 62),  
          ('god', 52),  
          ('could', 52),  
          ('people', 49),  
          ('time', 48),  
          ('south', 43),  
          ('may', 43),  
          ('wife', 38),  
          ('men', 37),  
          ('government', 33),  
          ('yet', 31),  
          ('made', 31),  
          ('must', 31),  
          ('never', 30)]
```

Look at the 20 most frequent words for each narrative. In the Markdown cell below, write down **three observations you have about this data**. These might be similarities between the two narratives, differences between the two, or any other patterns or questions you notice based on their word frequency.

Brown's text uses 'slave' and 'slavery' frequently, but Truth not.

Both two text frequently use the word 'god'.

Both two text frequently use the word 'master'.

Bonus Questions

The text files you've used in this homework were not the original text files of these narratives. Instead, they've been cleaned by your instructor to make them shorter and easier to analyze. Your goal is to use Python to download the original .txt files from the website Project Gutenberg. Adapt the code from [these examples](#) and use Python's `urllib` package to download the narratives and save them as local files named `truth-original.txt` and `brown-original.txt`.

Here are the URL's for the two original text files on Project Gutenberg:

- Truth's narrative: <https://www.gutenberg.org/cache/epub/1674/pg1674.txt>
- Brown's narrative: <https://www.gutenberg.org/cache/epub/64992/pg64992.txt>

```
In [76]: import urllib.request  
  
urllib.request.urlretrieve('https://www.gutenberg.org/cache/epub/1674/pg1674.txt'  
import urllib.request  
  
urllib.request.urlretrieve('https://www.gutenberg.org/cache/epub/64992/pg64992.t
```

```
Out[76]: ('brown-original.txt', <http.client.HTTPMessage at 0x1da967a7c80>)
```

Write code for the following:

- Open and read each of the new text files you just downloaded.
- Print out the **number of lines** in each of the original (newly downloaded) text files.

```
In [78]: truth_original= open('truth-original.txt', mode='r', encoding='utf-8').read()
brown_original= open('brown-original.txt', mode='r', encoding='utf-8').read()
truth_original_lines=truth_original.split('\n')
brown_original_lines=brown_original.split('\n')
print(len(truth_original_lines))
print(len(brown_original_lines))
```

4087

2814

Compare the length of the original text files you just downloaded to the cleaned text files you used for the rest of the homework, measured by the number of lines.

Write two print() statements that calculate **how many lines were removed by the instructor for each narrative**.

```
In [80]: print(f"truth= {len(truth_original_lines)-len(truth_lines)}")
print(f"brown= {len(brown_original_lines)-len(brown_lines)}")
```

truth= 460

brown= 591

What sort of text did the instructor remove? Write Python code that allows you to compare the two versions Sojourner Truth's narrative. Then write a few sentences in the empty Markdown cell below explaining what you found.

```
In [82]: print(truth_original[0:1000])
```

The Project Gutenberg eBook of The Narrative of Sojourner Truth

This ebook is for the use of anyone anywhere in the United States and most other parts of the world at no cost and with almost no restrictions whatsoever. You may copy it, give it away or re-use it under the terms of the Project Gutenberg License included with this ebook or online at www.gutenberg.org. If you are not located in the United States, you will have to check the laws of the country where you are located before using this eBook.

Title: The Narrative of Sojourner Truth

Author: Olive Gilbert
Sojourner Truth

Release date: March 1, 1999 [eBook #1674]

Most recently updated: June 20, 2015

Language: English

Credits: This book is put on-line as part of the BUILD-A-BOOK Initiative at the Celebration of Women Writers through the combined work of Laura LeVine, Margaret Sylvia, and Mary Mark Ockerbloom

*** START OF THE PROJECT GUTENBERG EBOOK THE NARRATIVE OF SOJOURNER TRUTH ***

This

In [84]: `print(truth_fulltext[0:1000])`

NARRATIVE OF SOJOURNER TRUTH

HER BIRTH AND PARENTAGE.

THE subject of this biography, SOJOURNER TRUTH, as she now calls herself-but whose name, originally, was Isabella-was born, as near as she can now calculate, between the years 1797 and 1800. She was the daughter of James and Betsey, slaves of one Colonel Ardinburgh, Hurley, Ulster County, New York.

Colonel Ardinburgh belonged to that class of people called Low Dutch.

Of her first master, she can give no account, as she must have been a mere infant when he died; and she, with her parents and some ten or twelve other fellow human chattels, became the legal property of his son, Charles Ardinburgh. She distinctly remembers hearing her father and mother say, that their lot was a fortunate one, as Master Charles was the best of the family,-being, comparatively speaking, a kind master to his slaves.

James and Betsey having, by their faithfulness, docility, and respectful behavior, won his particular regard, received from him

The revised version has removed the copyright description from the original version, and the information about the book has been removed, and I presume some information about the non-article content at the end of the book has also been removed, leaving us with a pure text.

In []: