

DOKUZ EYLÜL UNIVERSITY
DEPARTMENT OF COMPUTER ENGINEERING

E-BOOK ANALYSIS AND REPRESENTATION

Assignment Report

by
R.Hakan Cankul

January 2020
İZMİR

Contents

1	Introduction	iii
2	Methodology	iii
2.1	Structure of Your Project	iii
2.2	Encountered Problems and Solutions	iii
2.3	Improvements	iii
3	Experimentation	iii
	Appendix A: Code.....	iii
	Appendix B: Screenshots of your use cases	iv

1 INTRODUCTION

In this project, we developed an e-book downloader with word frequency counters. The project basically, consists of 5 part: finding correct book from query, downloading pdf as txt file, reading txt with correct pre-processings, frequency calculation, and displaying. I personally believe, the most different topics that I implemented differently are using very smart and rare library for pdf to text editor, using regular expressions for scraping and pre-processing and using google infrastructure for finding correct books. About all 3 strategy, the details will be given at the below.

2 METHODOLOGY

2.1 Structure of Your Project

In this project following libraries is used, re, nltk, request, bs4, webbrowser, numpy, pandas and pdftotext. Except for pdftotext all of the python libraries is widely used. The request, bs4 libraries is used for finding correct pdf file within the web. The pandas and numpy is used for more aesthetic representation of words. The re library is (regular expressions) used for pre-processing of words. The nltk library is used for removing stopwords and counting frequencies of words. Lastly, the pdftotext library is used for saving pdf file as text file.

As data structures, python list, python string, nltk freqlist and pandas dataframe have been used. The lists and strings are convertible structures in natural language processing, especially in tokenization. After, saving all book as a single string, we used freqlist to store frequencies. Lastly, convert it to pandas for displaying results.

The finding correct pdf upon web and removing stopwords are the most critical algorithms that I developed. For clarity, simplicity and time performance freqdist and pdftotext.PDF method have been used, besides less important built-in methods like str.lower, request.get or dataframe.to_string.

The algorithm for webscraping:

To use google infrastructure, the following pseudo-code have been used which you can find in source code more detailed way.

SEARCH KEY WORDS = ['wikibooks', 'wikisource']

- Form a query for google like google.com/search?q+ keyword + bookname.
- Get the resulting pages (that page would be regular google search page)
- Get first link in the above page (that will give you, the first link that google suggest) (functions as I'm feeling Lucky button)
- Handle redirect links
- Get the page url, it supposed to be wikipedia website.
- Parse the above url to check if it is wikibooks or not by regex.
- If it is wikibooks, find the button that redirects PDF. And save that pdf url for downloading and saving as text.
- Otherwise, check for next keyword.

The algorithm for removing StopWords.

- download nltk stopwords for all language.
- open the book that you downloaded and saved, readlines.
- merge all of the sentences as a single string.
- lower the string and delete newlines (\n)
- create regex tokenizer to parse intended words
- check if these words are in stopword list or not.
- if not, save it as a new array, convert it into string again.
- put that string into nltk.FreqDist method.

2.2 Encountered Problems and Solutions

The most crucial 3 problem that I encounter were, finding correct book, converting pdf to text and displaying. To solve how to find correct book, I used basic google search method which I explained above. The rest of two solutions are below.

The PDF's are like solid structures which is very hard to parse, manipulate or extract information. Therefore, I spent most of my time to how save them as string. First, sub-problem was about the characters. Some string consist unreadable characters for the machine. Therefore, I used utf-8 encoding at many places in my code. The second problem is extracting only alphabetical information. I tried some libraries until the pdftotext that is developed by some student. I am aware of that library is very rare to implement therefore, I also put downloading links into source code. In case, it is forbidden to use, or tester could not download and test this library, I put some manual tests at the below of the code that you can easily check rest of the parts functions correctly. I took full responsibility, if any point reduction occurs about that part.

The second problem was displaying. The nltk.freqdist data structure needs some investigation to use by un-familiar student. Therefore, I imported two different libraries too, pandas and numpy which are essentials of data science. For getting common or distict words were easy with freqdist data structure because it makes us enable to just substract and add these lists. However, for better visualization I had to convert them into pandas.DataFrame at the end of the code.

2.3 Improvements

It is asked to do in assignment paper but I asked user's two extra different questions. How many book will you select 1 or 2? and How many words you would like to display at the end of the execution.

Also, for ease to tester, I direct the pdftotext github page for downloading. For getting many errors I also used try and except block. Lastly, after downloading english stopwords, I cached them into an array, so code execution will be better at time performance.

Last improvement was regular expression. Most of the words need many pre-process to get valuable information. In this project numbers and signs considered non-valuable information by my decision. Therefore, I used regex for parse and extract the one solid words. Thanks to the fact that regex are super fast, it could be considered as improvement.

3 EXPERIMENTATION

I personally experiment/test the provided two books Non-Programmer's Tutorial for Python 2.6 and Non-Programmer's Tutorial for Python 3 in wikibooks. Even if you could download, pdftotext or not, additional tests could give you proper results. The code will work on both one book or two book. The distinct words and common words also executes successfully. For downloading pdftotext, I also add some breach y/n to direct you to browser or built-in tests.

4 CONCLUSION

I believe I learned the web scraping in detail. Also, learning regex for data processing and nltk toolbox was an extra for me. In addition, most of the computer engineers use the tools that others develop. The pdftotext could be considered the hard library to implement therefore, to merge my project with it, was surely gave me experience on large project later.

APPENDIX A: CODE

Due to mass amonth of comments reading the source code from appendix must be hard. Please check .py file for readability.

```
1. import requests,bs4,webbrowser,sys,time
2. import re
3. import nltk
4.
5. from nltk.corpus import stopwords # For removing stopwords
6. nltk.download('stopwords') # If you did not download stopwords before, it w
   ill.
7. from nltk.tokenize import RegexpTokenizer # Just tokenizer for wide usage
8. import numpy as np
9. import pandas as pd
10.
11. TEST = False # This global variable exist for built-
   in tests or user based test.
12. # 0 for user based. 1 for built-in tests.
13.
14. # /$$ /$$ /$$ /$$ /$$$$$
15. # | $$ | $$ | $$ | $$ /$__ $$
16. # | $$ | $$ /$$$$$ | $$ /$$ /$$$$$ /$$$$$ | $$ \_ / $$$
   $$$ /$$$$$ | $$ /$$ /$$ /$ | $$
17. # | $$$$$$ |__ $ $ /$$ /$ |__ $ $ | $$__ $ $ | $$ |__
   $ $ | $$__ $ $ | $$ /$$ /$ | $ $ | $ $ | $ $
18. # | $$__ $ $ /$$$$$ | $$$$ / /$$$$$ | $ $ \ $ $ | $ $ /$$$
   $$$ | $ $ \ $ $ | $$$$ / | $ $ | $ $ | $ $
19. # | $$ | $$ /$__ $ $ | $ $ $ $ /$__ $ $ | $ $ | $ $ | $ $ $ $ /$__
   $ $ | $ $ | $ $ | $ $ $ $ | $ $ | $ $ | $ $ | $ $
20. # | $$ | $ $ | $$$$ $ $ \ $ $ | $$$$ $ $ | $ $ | $ $ | $$$$ /$ $ $
   $$$ | $ $ | $ $ | $ $ \ $ $ | $$$$ /$ $ $
21. # |_ / |_ \_ / |_ \_ / |_ \_ / |_ \_ / |_ \_ / |_ \_ / |_ \_ /
   \_ / |_ / |_ / |_ / \_ / \_ / |_ /
22.
23. numberOfBook = int(input("How many books will you download (1 or 2): "))
24. # To choose how many books will I display.
25. print()
26. bookname1 = input("Book1: ")
27.
28. # CAREFUL!!
29. # The best way to find e-
   books on internet is basically surfing on google with spefic keywords
30. # Below keywords is for your searching query. In this homework just wikiboo
   k and wikisource.
31. # Assume you want to download Hamlet in wikibook, then you need to search w
   ikibook hamlet in google.
32. searchKeywords = ['wikibook','wikisource']
33.
34. # Below you search in google respectively in wikibook and wikisource. If th
   e book is not find in wikibooks
35. # Then, it will try wikisource.
36. for keyword in searchKeywords:
37.     # var: searching guery in google however, it directs you to re-
       direct page.
38.     var = requests.get(r'http://www.google.com/search?q='+keyword+' '+bookn
       ame1+'&btnI')
39.     # This page is about where to direct.
```

```

40.     redirect = var.url
41.     raw = requests.get(redirect)
42.     soup = bs4.BeautifulSoup(raw.text, "html.parser")
43.     # You basically parse all 'a' symbols because it shows html urls.
44.     # We use find because, first result would be most relevant link.
45.     # Think of I'm feeling lucky button.
46.     lucky = soup.find('a').text
47.
48.     # You lucky would be like https://en.wikibooks.org/wiki/hamlet
49.     # To test if it is in wikibook or not we simply use regex.
50.     # \.[w]\. is finds all words between two dot. in this case only '.wiki
books.'
51.     # [1:-1] deletes first and last character .wikibooks. => wikibooks
52.     x = re.findall(r'\.[w]+\.', lucky)[0][1:-1]
53.     # If the link in wikibooks then simply break. Otherwise, loop again for
wikisource.
54.     if(x == 'wikibooks'):
55.         break
56.
57. # After getting lucky link we need to parse the wikibook webpage.
58. # To find pdf page, we need to find pdf url.
59.
60. raw_page = requests.get(lucky)
61. wiki_page = raw_page.content
62.
63. # urf-
    8 for just taking appoirate character. \x08 like chars is out of context
64. soup2 = bs4.BeautifulSoup(wiki_page.decode('utf-8'), "html.parser")
65.
66. # I manually inspect the website and pdf link is within the html statement
67. # <a class='internal' href='PDF_LINK' >PDF Version</a>
68. book_url = soup2.find_all('a', {"class": "internal"})[0]['href'][2:]
69. book_url = 'https://' + book_url
70.
71. # VOLAA ! We found pdf of book. Webscraping part is done.
72. print("The url of the book is:",book_url)
73.
74. # The next thing to do, save the pdf as txt. Therefore, I used external lib
rary called pdftotext.
75. # It will take pdf from above url and save it with the same name that you q
ueried above.
76. try:
77.     import pdftotext
78.     # Load your PDF
79.     with open(book_url, "rb") as f:
80.         pdf = pdftotext.PDF(f)
81.
82.     # Save all text to a txt file.
83.     with open(bookname1+'.txt', 'w') as f:
84.         f.write("\n\n".join(pdf))
85. except:
86.     # This is just an error message because I am aware of pdftotext is not
widely used.
87.     # Therefore, I put links to download, if may take a while.
88.     # If you do not want to install and trust me, I provide extra 2 .txt fi
le for built-in tests.
89.     # You can simply press 'n' to go manual tests otherwise, it will open d
ownload page.
90.     print()
91.     print("Looks like you did not download required library pdftotext. Plea
se install it from below link")
92.     print("https://github.com/jalan/pdftotext")

```



```

93.     print("Installing can take a while.")
94.     print("Or in the below, you can try the tests with pre-
    installed books which the developer provides")
95.     x = input("To install please enter 'y', to forward tests please enter '
n': ")
96.     if(x == "y"):
97.         try:
98.             import webbrowser
99.             webbrowser.open("https://github.com/jalan/pdftotext")
100.            exit()
101.        except:
102.            exit()
103.    elif(x == "n"):
104.        TEST = True
105.        print('\nTEST STARTED...\n')
106.    else:
107.        print("Incorrect button. Shutting Down.")
108.        exit()
109.
110.
111.    # If you choose 2 book and after first book all libraries are correc
tly installed.
112.    # We can take second book with simply the same strategy.
113.    if(numberOfBook == "2"):
114.        bookname2 = input("Book2: ")
115.        searchKeywords = ['wikibook', 'wikisource']
116.
117.        for keyword in searchKeywords:
118.            var = requests.get(r'http://www.google.com/search?q='+keywor
d+' '+bookname2+'&btnI')
119.            redirect = var.url
120.            raw = requests.get(redirect)
121.            soup = bs4.BeautifulSoup(raw.text, "html.parser")
122.            lucky = soup.find('a').text
123.            x = re.findall(r'\.([w]+\.)', lucky)[0][1:-1]
124.            if(x == 'wikibooks'):
125.                break
126.            raw_page = requests.get(lucky)
127.            wiki_page = raw_page.content
128.            soup2 = bs4.BeautifulSoup(wiki_page.decode('utf-
8'), "html.parser")
129.
130.            book_url = soup2.find_all('a', {"class": "internal"})[0]['href']
131.            [2:]
132.            book_url = 'https://' + book_url
133.
134.            # Load your PDF
135.            with open(book_url, "rb") as f:
136.                pdf = pdftotext.PDF(f)
137.
138.            # Save all text to a txt file.
139.            with open(bookname2 + '.txt', 'w') as f:
140.                f.write("\n\n".join(pdf))
141.
142.            # If you downloaded the books without trouble, in the current folder
143.            # there should be THEBOOKNAMEYOUENTERED.txt. Then we needed to parse
all words.
144.            if(TEST == False):
145.                if(numberOfBook == 1):
146.                    # If only one book is selected... Simply open it in read mod
e utf-8.
                    book1 = open(bookname1 + '.txt', 'r', encoding=('utf-8'))

```

```

147.         b1 = book1.readlines()
148.         book1.close()
149.         # First we stringfy all sentences in list array.
150.         # We need to lower all words and delete \n (new lines)
151.         # "Hello World!\n" => "hello world"
152.         bookstring1 = " ".join(b1).lower().replace("\n", "")
153.         text = bookstring1
154.
155.         # we are using Regular Expression Tokenizer for personally,
most accurate preprocesing
156.         # For comparing words we needed to eliminate signs (+, -
, ') or numbers (3.6, 5) etc.
157.         # Also, we need to count each of cat, cat's, non-
cat as one.
158.         # That is why we use regex with [A-Za-
z]+ which will take each word.
159.         tokenizer = RegexpTokenizer(r'[A-Za-z]+')
160.         text_tokens = tokenizer.tokenize(text)
161.
162.         # For performance, we cached stopwords one.
163.         cachedStopwords = stopwords.words('english')
164.         # Below, we shrink all statement block in below, into one li
ne.
165.         # tokens_without_sw = []
166.         # for word in text_tokens:
167.         #     if(word not in cachedStopwords):
168.         #         tokens_without_sw.append(word)
169.         tokens_without_sw = [word for word in text_tokens if not wor
d in cachedStopwords]
170.
171.         # Then we convert list to string and found frequency distrib
utions
172.         text1 = " ".join(tokens_without_sw)
173.         fdist1 = nltk.FreqDist(tokens_without_sw)
174.
175.         # we get how many words you want to display
176.         n = int(input("How many word you want to display: "))
177.
178.         # Some manipulations for more aestatic table
179.         df = pd.DataFrame(fdist1.most_common(n))
180.         df = df.rename(columns={0: "WORD", 1: "FREQ_1"})
181.
182.         df['NO'] = np.array(range(1,n+1))
183.         df = df[['NO', 'WORD', 'FREQ_1']]
184.
185.         # VOLAA !!
186.         print(df.to_string(index=False))
187.
188.         if(numberOfBook == 2):
189.             book1 = open(bookname1 + '.txt', 'r', encoding=('utf-8'))
190.             b1 = book1.readlines()
191.             book1.close()
192.             bookstring1 = " ".join(b1).lower().replace("\n", "")
193.             text = bookstring1
194.             tokenizer = RegexpTokenizer(r'[A-Za-z]+')
195.             text_tokens = tokenizer.tokenize(text)
196.
197.             cachedStopwords = stopwords.words('english')
198.             tokens_without_sw = [word for word in text_tokens if not wor
d in cachedStopwords]
199.
200.             text1 = " ".join(tokens_without_sw)
201.             fdist1 = nltk.FreqDist(tokens_without_sw)

```

```

202.         book2 = open(bookname2 + '.txt','r',encoding=('utf-8'))
203.         b2 = book2.readlines()
204.         book2.close()
205.
206.         bookstring2 = " ".join(b2).lower().replace("\n", "")
207.
208.         text2 = bookstring2
209.         text_tokens2 = tokenizer.tokenize(text2)
210.
211.         tokens_without_sw2 = [word for word in text_tokens2 if not word
ord in cachedStopwords]
212.         text2 = " ".join(tokens_without_sw2)
213.         fdist2 = nltk.FreqDist(tokens_without_sw2)
214.         n = int(input("How many word you want to display: "))
215.         df = pd.DataFrame(fdist1.most_common(n))
216.         df2 = pd.DataFrame(fdist2.most_common(n))
217.
218.         # Thanks to freqdist, If we have 2 booko to compare,
219.         # we can easily substract or add them to find common or disti
nc words.
220.         df3 = pd.DataFrame((fdist1 + fdist2).most_common(n))
221.         df4 = pd.DataFrame((fdist1 - fdist2).most_common(n))
222.         df5 = pd.DataFrame((fdist2 - fdist1).most_common(n))
223.
224.         # Rest of them is also table manipulation.
225.         df = pd.DataFrame(fdist1.most_common(n))
226.         df = df.rename(columns={0: "WORD", 1: "FREQ_1"})
227.
228.         df['NO'] = np.array(range(1,n+1))
229.         df = df[['NO', 'WORD', 'FREQ_1']]
230.         print(df.to_string(index=False))
231.
232.         print("\n")
233.         df2 = pd.DataFrame(fdist2.most_common(n))
234.         df2 = df2.rename(columns={0: "WORD", 1: "FREQ_2"})
235.
236.         df2['NO'] = np.array(range(1,n+1))
237.         df2 = df2[['NO', 'WORD', 'FREQ_2']]
238.         print(df2.to_string(index=False))
239.
240.         print("\n")
241.
242.         df = pd.DataFrame(fdist1.most_common(n))
243.         df2 = pd.DataFrame(fdist2.most_common(n))
244.         df3 = pd.DataFrame((fdist1 + fdist2).most_common(n))
245.
246.
247.         d = pd.DataFrame()
248.         d['NO'] = np.array(range(1,n+1))
249.         d['WORD'] = df3[0]
250.         d['FREQ_1'] = df[1]
251.         d['FREQ_2'] = df2[1]
252.         d['FREQ_SUM'] = df3[1]
253.         print(d.to_string(index=False))
254.         print("\n")
255.
256.         df4 = df4.rename(columns={0: "WORD", 1: "FREQ_1"})
257.         df4['NO'] = np.array(range(1,n+1))
258.         df4 = df4[['NO', 'WORD', 'FREQ_1']]
259.         print(df4.to_string(index=False))
260.         print("\n")
261.
262.         df5 = df5.rename(columns={0: "WORD", 1: "FREQ_2"})

```

```

263.         df5['NO'] = np.array(range(1,n+1))
264.         df5 = df5[['NO', 'WORD', 'FREQ_2']]
265.         print(df5.to_string(index=False))
266.         print("\n")
267.
268.         #####\ #####\ #####\ #####\
269.         #\_$$ _|$$ _|$$ _|$$ _|$$ _|
270.         #  $$ |  $$ |  $$ /  \_  $$ |
271.         #  $$ |  #####\  \#####\  $$ |
272.         #  $$ |  $$ _|  \_  $$\  $$ |
273.         #  $$ |  $$ |  $$\  $$  |  $$ |
274.         #  $$ |  #####\ \#####  |  $$ |
275.         #  \_  |  \_  |  \_  |  \_  |
276.
277.         ## WARNING!! ##
278.
279.         # These part is identically same with above if statement so any more
        comment is not necessary.
280.         # If you press n above, these statements will be executed.
281.         if(TEST == True):
282.             if(numberOfBook == 1):
283.                 print("BOOK 1: Non-Programmer's Tutorial for Python 2.6")
284.                 book1 = open('Non-
Programmer\'s_Tutorial_for_Python_2.6.txt', 'r', encoding=('utf-8'))
285.                 b1 = book1.readlines()
286.                 book1.close()
287.                 bookstring1 = " ".join(b1).lower().replace("\n", "")
288.                 text = bookstring1
289.                 tokenizer = RegexpTokenizer(r'[A-Za-z]+')
290.                 text_tokens = tokenizer.tokenize(text)
291.
292.                 cachedStopwords = stopwords.words('english')
293.                 tokens_without_sw = [word for word in text_tokens if not wor
d in cachedStopwords]
294.
295.                 text1 = " ".join(tokens_without_sw)
296.                 fdist1 = nltk.FreqDist(tokens_without_sw)
297.                 n = int(input("How many word you want to display: "))
298.
299.                 df = pd.DataFrame(fdist1.most_common(n))
300.                 df = df.rename(columns={0: "WORD", 1: "FREQ_1"})
301.
302.                 df['NO'] = np.array(range(1,n+1))
303.                 df = df[['NO', 'WORD', 'FREQ_1']]
304.                 print(df.to_string(index=False))
305.
306.             if(numberOfBook == 2):
307.                 print("BOOK 1: Non-Programmer's Tutorial for Python 2.6")
308.                 print("BOOK 2: Non-Programmer's Tutorial for Python 3")
309.                 book1 = open('Non-
Programmer\'s_Tutorial_for_Python_2.6.txt', 'r', encoding=('utf-8'))
310.                 b1 = book1.readlines()
311.                 book1.close()
312.                 bookstring1 = " ".join(b1).lower().replace("\n", "")
313.                 text = bookstring1
314.                 tokenizer = RegexpTokenizer(r'[A-Za-z]+')
315.                 text_tokens = tokenizer.tokenize(text)
316.
317.                 cachedStopwords = stopwords.words('english')
318.                 tokens_without_sw = [word for word in text_tokens if not wor
d in cachedStopwords]
319.
320.                 text1 = " ".join(tokens_without_sw)

```

```

321.         fdist1 = nltk.FreqDist(tokens_without_sw)
322.
323.         book2 = open('Non-
Programmer\'s_Tutorial_for_Python_3.txt', 'r', encoding=('utf-8'))
324.         b2 = book2.readlines()
325.         book2.close()
326.
327.         bookstring2 = " ".join(b2).lower().replace("\n", "")
328.
329.         text2 = bookstring2
330.         text_tokens2 = tokenizer.tokenize(text2)
331.
332.         tokens_without_sw2 = [word for word in text_tokens2 if not w
ord in cachedStopwords]
333.         text2 = " ".join(tokens_without_sw2)
334.         fdist2 = nltk.FreqDist(tokens_without_sw2)
335.         n = int(input("How many word you want to display: "))
336.         df = pd.DataFrame(fdist1.most_common(n))
337.         df2 = pd.DataFrame(fdist2.most_common(n))
338.         df3 = pd.DataFrame((fdist1 + fdist2).most_common(n))
339.         df4 = pd.DataFrame((fdist1 - fdist2).most_common(n))
340.         df5 = pd.DataFrame((fdist2 - fdist1).most_common(n))
341.
342.         df = pd.DataFrame(fdist1.most_common(n))
343.         df = df.rename(columns={0: "WORD", 1: "FREQ_1"})
344.
345.         df['NO'] = np.array(range(1,n+1))
346.         df = df[['NO', 'WORD', 'FREQ_1']]
347.         #print(df.to_string(index=False))
348.
349.         print("\n")
350.         df2 = pd.DataFrame(fdist2.most_common(n))
351.         df2 = df2.rename(columns={0: "WORD", 1: "FREQ_2"})
352.
353.         df2['NO'] = np.array(range(1,n+1))
354.         df2 = df2[['NO', 'WORD', 'FREQ_2']]
355.         #print(df2.to_string(index=False))
356.
357.         print("\n")
358.
359.         df = pd.DataFrame(fdist1.most_common(n))
360.         df2 = pd.DataFrame(fdist2.most_common(n))
361.         df3 = pd.DataFrame((fdist1 + fdist2).most_common(n))
362.
363.
364.         d = pd.DataFrame()
365.         d['NO'] = np.array(range(1,n+1))
366.         d['WORD'] = df3[0]
367.         d['FREQ_1'] = df[1]
368.         d['FREQ_2'] = df2[1]
369.         d['FREQ_SUM'] = df3[1]
370.         print("COMMON WORDS")
371.         print(d.to_string(index=False))
372.         print("\n")
373.
374.         print("BOOK 1: Non-Programmer's Tutorial for Python 2.6")
375.         print("DISTINCT WORDS")
376.         df4 = df4.rename(columns={0: "WORD", 1: "FREQ_1"})
377.         df4['NO'] = np.array(range(1,n+1))
378.         df4 = df4[['NO', 'WORD', 'FREQ_1']]
379.         print(df4.to_string(index=False))
380.         print("\n")
381.

```

```
382.         print("BOOK 2: Non-Programmer's Tutorial for Python 3")
383.         print("DISTINCT WORDS")
384.         df5 = df5.rename(columns={0: "WORD", 1: "FREQ_2"})
385.         df5['NO'] = np.array(range(1,n+1))
386.         df5 = df5[['NO', 'WORD', 'FREQ_2']]
387.         print(df5.to_string(index=False))
388.
389.
390.     print("\n\nProgram Finished.")
```

APPENDIX B: SCREENSHOTS OF YOUR USE CASES

```
TEST STARTED...

BOOK 1: Non-Programmer's Tutorial for Python 2.6
How many word you want to display: 20
NO      WORD      FREQ_1
1       print      517
2       number     282
3       program    231
4       license    228
5       name       222
6       b          198
7       first      190
8       python     176
9       line       155
10      list       151
11      work       149
12      input      147
13      http       141
14      index      139
15      value      134
16      org        130
17      user       127
18      function   120
19      version    111
20      use        107
```

```
BOOK 1: Non-Programmer's Tutorial for Python 3
How many word you want to display: 30
NO      WORD      FREQ_1
1       print      572
2       number     288
3       python     279
4       name       228
5       b          181
6       tutorial   177
7       program    172
8       list       159
9       non        156
10      programmer  155
11      line       153
12      wikibooks  152
13      index      144
14      value      133
15      input      132
16      numbers    129
17      function   121
18      first      112
19      org        104
20      true       99
21      menu       98
22      item       96
23      next       94
24      type       93
25      string     89
26      version    87
27      https      86
28      file       86
29      w          82
30      en         75
```

```
BOOK 1: Non-Programmer's Tutorial for Python 2.6
BOOK 2: Non-Programmer's Tutorial for Python 3
How many word you want to display: 10
```

COMMON WORDS

NO	WORD	FREQ_1	FREQ_2	FREQ_SUM
1	print	517	572	1089
2	number	282	288	570
3	python	231	279	455
4	name	228	228	450
5	program	222	181	403
6	b	198	177	379
7	list	190	172	310
8	line	176	159	308
9	first	155	156	302
10	index	151	155	283

```
BOOK 1: Non-Programmer's Tutorial for Python 2.6
```

DISTINCT WORDS

NO	WORD	FREQ_1
1	license	222
2	work	134
3	http	115
4	user	92
5	document	84
6	first	78
7	may	73
8	raw	66
9	copy	63
10	section	61

```
BOOK 2: Non-Programmer's Tutorial for Python 3
```

DISTINCT WORDS

NO	WORD	FREQ_2
1	tutorial	147
2	programmer	145
3	non	138
4	python	103
5	https	86
6	ope	74
7	sur	74
8	wikibooks	67
9	print	55
10	float	28

WARNING:

The difference in numbers could be result of the source of the pdf, pdf to text converter or preprocessing steps (RegExp Tokenizer in this case.) .

REFERENCES

1. Online Word Counter. (n.d.). CountWordsFree. Retrieved December 22, 2020, from <https://countwordsfree.com>
2. Wikibooks. (n.d.). Retrieved December 22, 2020, from https://en.wikibooks.org/wiki/Main_Page
3. Wikibooks contributors. (2020a). Non-Programmer's Tutorial for Python 3—Wikibooks, open books for an open world. Wikibooks, The Free Textbook Project. https://en.wikibooks.org/wiki/Non-Programmer%27s_Tutorial_for_Python_3
4. Wikibooks contributors. (2020b). Non-Programmer's Tutorial for Python 2.6—Wikibooks, open books for an open world. Wikibooks, The Free Textbook Project. https://en.wikibooks.org/wiki/Non-Programmer%27s_Tutorial_for_Python_2.6
5. *Jalan/pdftotext*. (2021, January 17). GitHub. <https://github.com/jalan/pdftotext>
6. Steven Bird, Ewan Klein, and Edward Loper (2009). Natural Language Processing with Python. O'Reilly Media Inc. <http://nltk.org/book>
7. *bs4*. (2021, January 17). PyPI. <https://pypi.org/project/bs4/>