

BIS 420 PROGRAMMING FOR DATA SCIENCE

PRAJAKTA POHARE

CHAPTER 13 EXERCISE 13.2

ILLINOIS STATE UNIVERSITY

Go to Project Gutenberg (<http://gutenberg.org>) and download your favorite out-of-copyright book in plain text format. Modify your program from the previous exercise to read the book you downloaded, skip over the header information at the beginning of the file, and process the rest of the words as before. Then modify the program to count the total number of words in the book, and the number of times each word is used.

Print the number of different words used in the book. Compare different books by different authors, written in different eras. Which author uses the most extensive vocabulary?

```
/usr/local/bin/python3 "/Users/prajaktapohare/Library/CloudStorage/OneDrive-ILStateUniversity/BIS420/Week 13/BIS420_PrajaktaPohare_Ch13_13.2.py"
• → ~ /usr/local/bin/python3 "/Users/prajaktapohare/Library/CloudStorage/OneDrive-ILStateUniversity/BIS420/Week 13/BIS420_PrajaktaPohare_Ch13_13.2.py"

Analysis of pride_and_prejudice.txt:
Total words: 4560
Unique words: 1091

Analysis of moby_dick.txt:
Total words: 4560
Unique words: 1091

Analysis of dracula.txt:
Total words: 4560
Unique words: 1091

Vocabulary comparison:
pride_and_prejudice.txt: 1091 unique words
moby_dick.txt: 1091 unique words
dracula.txt: 1091 unique words

Most extensive vocabulary: pride_and_prejudice.txt
○ → ~
```