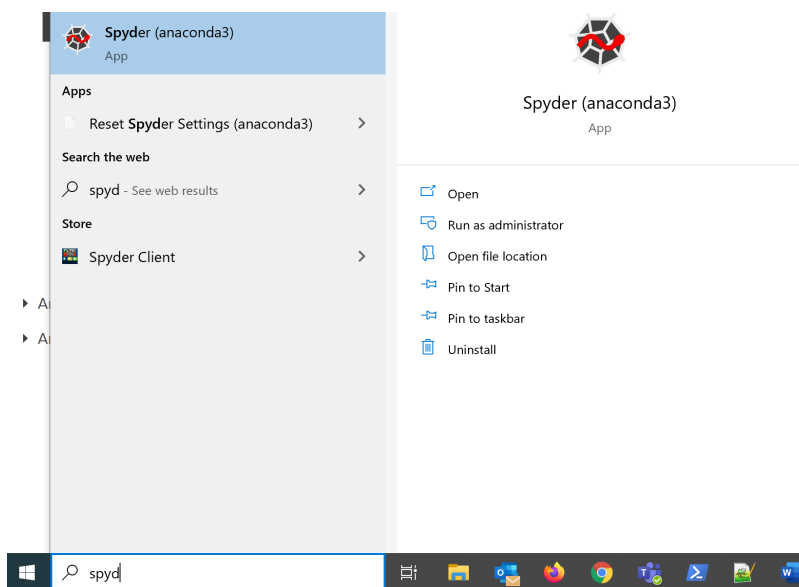# Explained Python text co-occurance calculator.

This is an explained example of using python to find word co-occurrences in text (example uses Jane Austen novels). For example, you can find that the word "pride" co-occurs with "darcy" more than with "elizabeth".
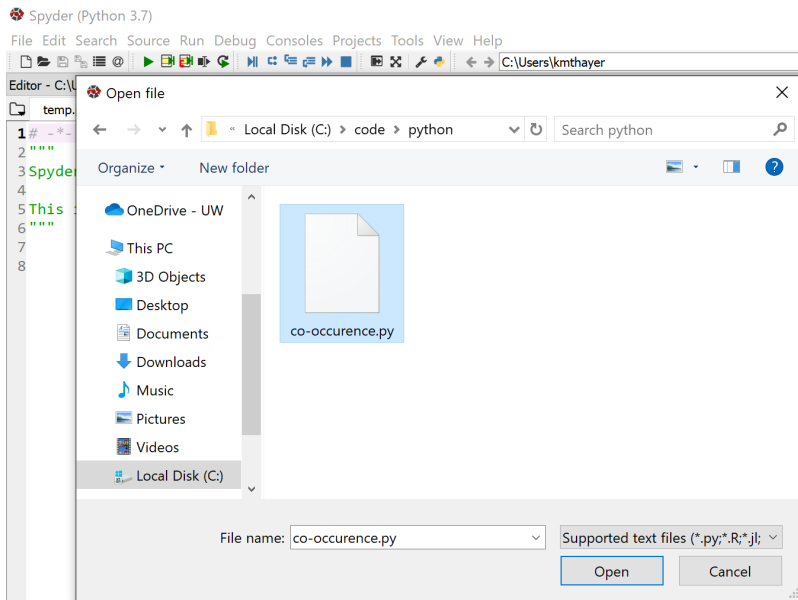
## Instructions

1. Download the files in this repository into a new directory (unzipping it if needed). Do this by using the green button with the down arrow that says "Code".

2. Download and install the Anaconda Distribution, which comes with Python and a code editor (called Spyder). Or find or use some other python code editor:
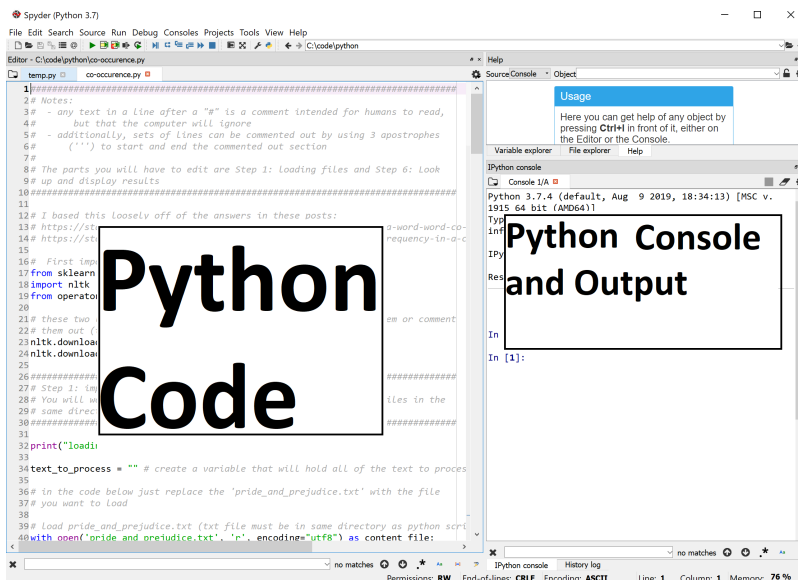
https://www.anaconda.com/distribution/

3. Open the Spyder code editor (screenshot shows opening Spyder on Windows



4. In Spyder, choose the "file" menu, and select "open." Then open the "co-occurence.py" in the new directory you made from the zip file
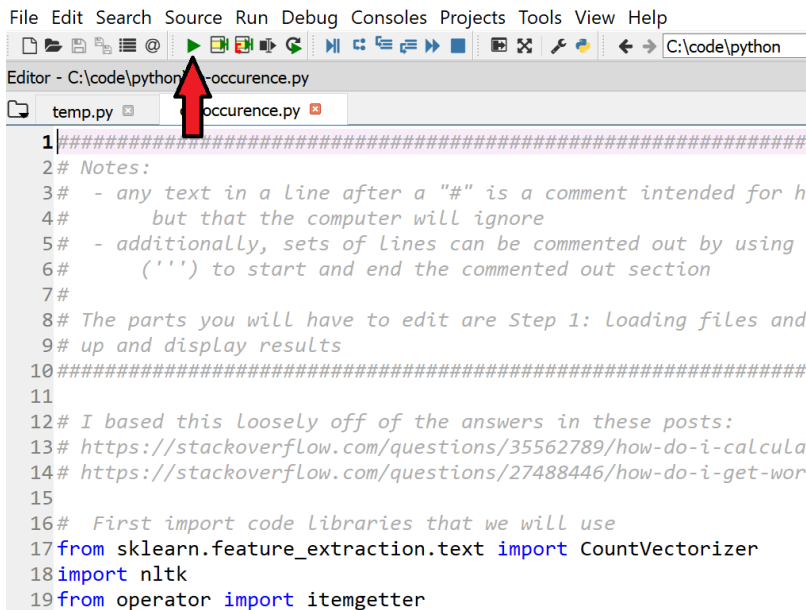
5. Now you should see the file open in Spyder like below (though without the giant labels I stuck on the screenshot). There is an area where the code file is on the left, and an area on the right called the "console" or "terminal" where the results of the code will appear, and where you can run additional commands if you want.



6. To run the Python file, click on the little green triangle at the top left.
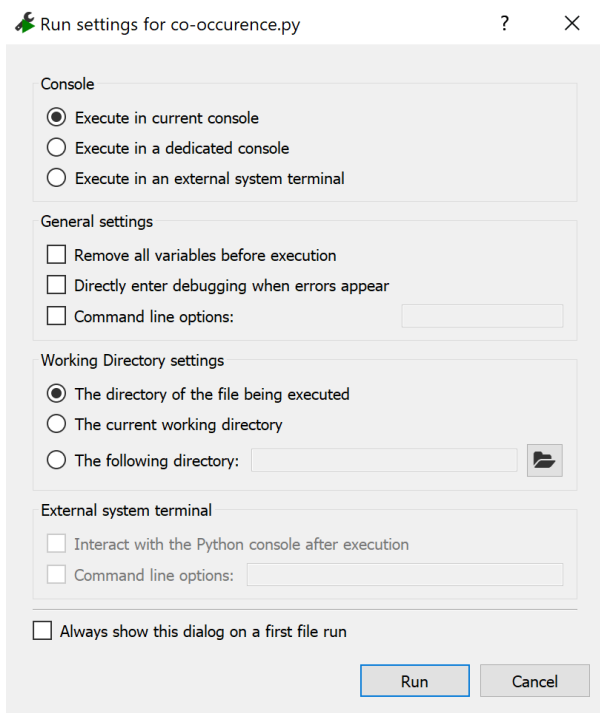
Spyder (Python 3.7)

File  Edit  Search  Source  Run  Debug  Consoles  Projects  Tools  View  Help

Editor - C:\code\python\...-occurence.py

temp.py ☒     ...occurence.py ☒

```
1 #############################################################
2 # Notes:
3 #  - any text in a line after a "#" is a comment intended for h
4 #       but that the computer will ignore
5 #  - additionally, sets of lines can be commented out by using
6 #       (''') to start and end the commented out section
7 #
8 # The parts you will have to edit are Step 1: Loading files and
9 # up and display results
10 #############################################################
11
12 # I based this loosely off of the answers in these posts:
13 # https://stackoverflow.com/questions/35562789/how-do-i-calcula
14 # https://stackoverflow.com/questions/27488446/how-do-i-get-wor
15
16 #  First import code libraries that we will use
17 from sklearn.feature_extraction.text import CountVectorizer
18 import nltk
19 from operator import itemgetter
```

7. The first time you run a program, it will ask you some options. The defaults should be fine, but double check that under "Working Directory settings" it has "The directory of the file being executed" selected.

Run settings for co-occurence.py          ?     X

Console
  ⦿ Execute in current console
  ◯ Execute in a dedicated console
  ◯ Execute in an external system terminal

General settings
  ☐ Remove all variables before execution
  ☐ Directly enter debugging when errors appear
  ☐ Command line options: [            ]

Working Directory settings
  ⦿ The directory of the file being executed
  ◯ The current working directory
  ◯ The following directory:  [            ] 📁

External system terminal
  ☐ Interact with the Python console after execution
  ☐ Command line options: [            ]

☐ Always show this dialog on a first file run

          [ Run ]    [ Cancel ]

8. After it runs, you should see the output of the script in the console panel on the right (you can scroll up and down). Note that for the top 40 terms co-occuring with pride, "pride" co-occurs with itself 69 times, with "and" 65 times, with "of" 64 times, etc.

9. Let me explain the file a little bit:

- Any text after a "#" in a line is a comment, meant for humans to read, and Python will ignore them. Additionally three apostrophes (''') can mark the start and end of a set of lines to be comments that Python will ignore

- After loading some code libraries at the top, I divided the code into six sections:
  - Step 1: import the text we want to process
  - Step 2: split words into sentences here
  - Step 3: clean up data (optional and currently disabled by commenting out all the lines)
  - Step 4: find the co-occurance matrix for all words in the text
  - Step 5: define custom functions to help retrieve and display information
  - Step 6: Look up and display results in the console

- The only sections you will be required to modify are Step 1 and Step 6.
  - Step 1 modifications: You will need to get whatever text you want to process, saved as txt files in the same directory as the python file. Then you will modify the line of code where it has "pride_and_prejudice.txt" and replace that with the name of your file. If you want to load multiple txt files, there is an example commented out for how to load emma.txt and add that one as well.
  - Step 2 modifications: I have several function calls to answer specific questions about words in Pride and Prejudice. You will want to modify and copy these lines of code to ask about whatever particular words you want to ask about, and you can change how many words to return in places where I asked for the top 40 words.

- There are several places where I have added optional code changes you can make:
  - Step 3: You can delete the three apostraphes at the start and end to enable that code cleaning step
  - Step 4: You can uncomment what is currently line 142 to only count each word once per sentence even if it occurs multiple times in a sentence.

- Step 5: In the definition of the get_cooccurances_sorted function, you can use an alternate line to get proportional co-occurance based on a term (number of co-occurances divided by number of times the other term appears). This finds more unique words that co-occur with what you are looking for.

- I left comments through most of the code, trying to explain what it does, and sometimes how, so you can look through it and try to understand and modify it yourself if you want to.

- If and when you try modifying things, back it up! It is easy to break a program and not know how to undo it. The easiest way I know to back things up is to just email myself the file as I work on it. Then if you ever get lost, you can download one of those working versions and start again from there.

10. If you want to run additional commands in the console, say like asking more questions about occurances and co-occurances, you can type (or paste) them directly into the console, like this:



11. Then press enter, and the results will display there (see below).



12. Note on the console vs. the code file. If you want to just look at values and try things out, that can be good to do on the console, but it is easy to lose which commands you ran. If you want to have lines of code that you remember and can run again, put those in the python file so they run every time you run the file.