

# MODULE 4.2 Performing Text Analysis: Basic Approaches with Python

## KEY TOOLS/PLATFORMS

### Python

A programming language that is good for working with data.

### pip

A package manager for Python

### pyplot

Visualization function in the Python data science package, Pandas

### HTRC Extracted Features

A downloadable dataset of text data and metadata extracted and abstracted from volumes in the HathiTrust Digital Library.

### HTRC Feature Reader

A Python library for working with HTRC Extracted Features.

## ACTIVITY: Identify the method

 Slide M4.2 - 6


What are the broad areas and methods used for the research examples we read earlier?

Project summaries: <http://go.illinois.edu/ddrf-research-examples>

	Broad area	Specific method
<i>Rowling and “Galbraith”: an authorial analysis</i>		
<i>Significant Themes in 19th Century Literature</i>		
<i>The Emergence of Literary Diction</i>		

## ACTIVITY: View adjectives in Extracted Features files

 Slide M4.2 - 26

1. Go to your PythonAnywhere dashboard and click on the “Browse files” button to check if all files and directories are in place 
1. On your “Files” page, you will see your directories on the left and your files on the right.
2. You should have two directories “1930/” and “1970/” listed on the left.

- Click on each directory.
- There should be 16 json.bz2 files in the “1970/” directory that correspond to the volumes in the poli\_science\_DDRF workset atn 5 json.bz2 files in the “1930/” directory of Presidential Papers volumes from the 1930s.

2. Open the Bash console.

3. Install the Feature Reader Library by typing:

```
pip install --user htrc-feature-reader
```

4. Hit enter.

5. From the Bash shell, run the script. Enter command (remember to use the tab key to help you automatically complete the file names):

```
python top_adjectives.py 1970
```

6. Hit Enter. The results will be printed out directly in the window.

```
Bash console 4179031
19:56 ~ $ pip install --user htrc-feature-reader
```

```
13:45 ~ $ python top_adjectives.py 1970
```

```
13:45 ~ $ python top_adjectives.py 1970
count
token                21248
other                 17824
own                   11552
new                   10400
good                  9984
great                 9440
American              7936
many                  7840
major                 6480
last                  6288
public                6224
important             6096
first                 5808
such                  5760
economic              5504
human                 5424
international         5216
national              4864
same                  4688
next                  3920
nuclear               3840
local                 3744
foreign               3696
political             3648
comprehensive         3488
few                   3472
sure                  3440
possible              3408
```

**ACTIVITY: Now you try!**

 **Slide M4.2 - 32**

- Compare the adjectives used by presidents in the 1970s with those used in the 1930s.
- Work and discuss with your neighbor.
  - How do you need to change your command?
  - What differences do you see? Similarities?

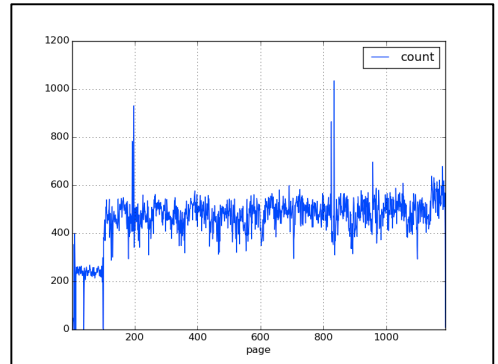
<b>Challenge</b>	<p>Try modifying the script to search verbs or another part of speech.          (Hint: Don't forget the Penn Tree Bank – you can find it here:  <a href="https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html">https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html</a> )</p>
------------------	---

### ACTIVITY: Visualize word trends

Slide M4.2 - 36

1. Open the Bash console
2. Enter the command “python word\_count.py”
3. Click the “PythonAnywhere” logo to go back to the home screen and go to your Files.
4. The “words.png” file should be in your Files. Click on the download icon to the right of the file name to open the file.

```
12:28 ~ $ python word_count.py
```



### ACTIVITY: Now you try!

Slide M4.2 - 40

- Can you modify the script to look at another volume?