

MCIS6273 Data Mining (Prof. Maull) / Fall 2017 / HW0

This assignment is worth up to 20 POINTS to your grade total if you complete it on time.

Points Possible	Due Date	Time Commitment (estimated)
20	Wednesday, Sep 06 @ Midnight	4 hours or less

- **GRADING:** Grading will be aligned with the completeness of the objectives.
- **INDEPENDENT WORK:** Copying, cheating, plagiarism and academic dishonesty *are not tolerated* by University or course policy. Please see the syllabus for the full departmental and University statement on the academic code of honor.

OBJECTIVES

- Familiarize yourself with Github
- Work with exploring unstructured data with Python and text
- Work with a large dataset and understand data munging in Python Pandas

WHAT TO TURN IN

You are being encouraged to turn the assignment in using the provided Jupyter Notebook. To do so, clone the course repository and modify the `hw0.ipynb` file in the `homework/` directory. If you do not know how to do this, please ask, or visit one of the many tutorials out there on the basics of using Github and cloning repositories.

Turn in a copy of a `.ipynb` file, a PDF or Word Document to Blackboard with the answers to the questions labeled with the § sign.

ASSIGNMENT TASKS

(0%) Familiarize yourself with Github

Github is the *de facto* home of open source software on the web. While there are other platforms that provide similar (and in some cases better) services, it is still *the place* for some of the largest open source software projects in the world. You will be creating an account on Github for several reasons:

1. we will be using Github to store, transmit and share homework and lecture notes,
2. we will use Github for all assignments,
3. at some point in your future, you will *very likely* be using Github, if for private work for a company or public work on an open source project.

§ **Find ONE Github repository of interest and explore it.** There is **nothing to turn in** for this task – just begin to explore Github.

(25%) Work with exploring unstructured data with Python and text

In this part, you will explore some text data and get a little familiar with Python's parsing and text capabilities. You will grab data from the free books provided online from Project Gutenberg and use the provided code to compare these documents. Turn in the answers to the given tasks after studying the code provided in `gutenberg_get_words()` which takes a url for a book on Project Gutenberg and returns a list of the words in the book. Notice the `stopwords=` parameter is used to eliminate words that relay low or no information. This is a common technique use in text processing.

These five books will be used for the tasks for this question:

Book	URL
The Prince, Machiavelli	http://www.gutenberg.org/cache/epub/1232/pg1232.txt
Frankenstein; Or, The Modern Prometheus by Mary Wollstonecraft Shelley	http://www.gutenberg.org/cache/epub/84/pg84.txt
Siddhartha by Hermann Hesse	http://www.gutenberg.org/cache/epub/2500/pg2500.txt
The Republic by Plato	http://www.gutenberg.org/cache/epub/1497/pg1497.txt
The Federalist Papers by Alexander Hamilton, John Jay, and James Madison	http://www.gutenberg.org/cache/epub/1404/pg1404.txt

```
import requests
import re
```

```
US_STOPWORDS = ["a", "about", "above", "above", "across", "after", "afterwards", "again", "against", "a
```

```
def gutenberg_get_words(url="http://www.gutenberg.org/cache/epub/1232/pg1232.txt",
                        range=slice(0,None), stopwords=[]):
```

```
    r = requests.get(url)
    data = re.sub(r"[^\w\s]", "", str(r.text)).lower()

    return \
        [w for w in data.split() if w not in stopwords]
```

```
## EXAMPLE USAGE
```

```
words = gutenberg_get_words(
    "http://www.gutenberg.org/cache/epub/84/pg84.txt",
    stopwords=US_STOPWORDS)
```

```
print(words[100:115])
```

§ submit the Python code that does the following:

- using the code and the 5 books provided above, explore and apply the very nice Python library called `collections`. Use the `Counter` class to load the word frequencies of each book into a Python dictionary.
- **NOTE:** you will need to be online with an internet connection for this to work, since it loads the data directly from the URLs of the books.

§ turn in *at least 2 sentences* and any code if you used code to answering the following:

- there are similarities and differences in the top 30 words of the five provided documents – be specific about describing what they are? How similar or different are each of the top 30 words list? You can

compare them by hand (look at them) or you are encouraged to write Python code to compare them automatically.

(75%) Work with a large dataset and understand data munging in Python Pandas

As we have learned in class, a great deal of time doing data mining involves understanding the data in a dataset and preprocessing it in preparation for working with it in a real analysis of some sort. We will develop an understanding of:

- importing CSV data into Pandas DataFrames,
- exporting JSON data from a Pandas DataFrame,
- understand how to compute data from DataFrame elements

You might find the Pandas IO library helpful in completing this task.

We will be using an interesting data source: **US Baseball Statistics Archive** by Sean Lahman (CCBY-SA 3.0), which can be found in these two locations. You can download the entire archive to your local machine.

- <http://seanlahman.com/baseball-archive/statistics/>
- <https://github.com/chadwickbureau/baseballdatabank>

§ submit the code to load the batting table from the Baseball Archive directly from the Github repository into a Pandas DataFrame.

- You can use the CSV file at this url. You will most certainly be using the `read_csv()` method.

§ submit the code to write the JSON file that contains all the batting records for the Cincinnati Reds (CIN) in 1981.

- This will require a simple query into the dataset (loaded as a DataFrame) – please read the documentation on boolean indexing.
- You will need to explore the `to_json()` method to then convert the DataFrame that results from the query above into a JSON object.

§ submit the code that calculates the following batting rank score for the 1981 Reds (CIN)

$$\varphi = G \frac{H}{AB} + \frac{RBI}{\sqrt{R}}$$

where G, H, AB, RBI, R are *games, hits, at bats, runs batted in* and *runs* respectively.

- Your answer should ONLY include the batting for the 1981 Cincinnati Reds batting roster.

§ for the 1981 data submit the top 4 batter IDs by their batting rank score, φ , previously calculated.