

Nathaniel Velarde  
W205-2, Fall 2017  
Exercise 2

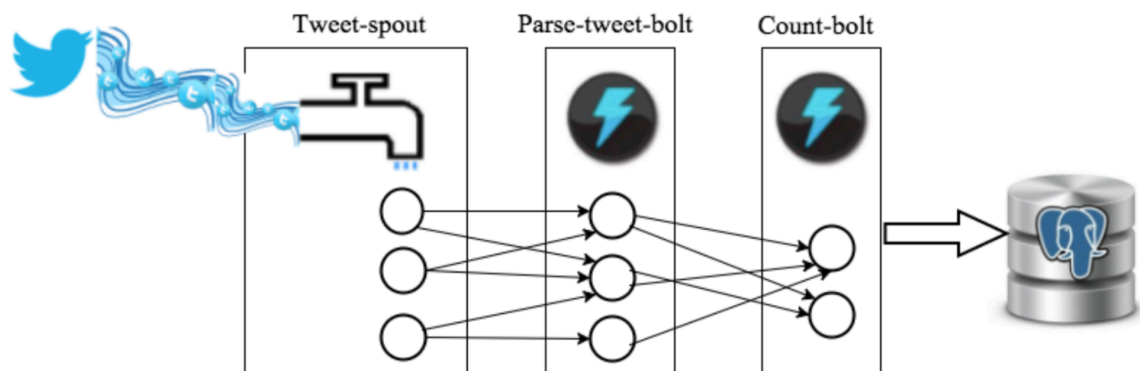
## Twitter Application Architecture

### Executive Summary

We have implemented an end-to-end streaming application that reads the stream of tweets from the Twitter streaming API using Apache Storm/streamparse, parses the tweets into words, keeps track of the number of the occurrences each word appears in the stream (word counts) and then aggregates the words and word counts into a Postgres database. The application includes two python scripts – `finalresults.py` and `histogram.py` – to perform basic analytics on the data collected and stored in Postgres.

### Architecture

The figure below shows the overall architecture of our application as well as the Storm topology implemented.



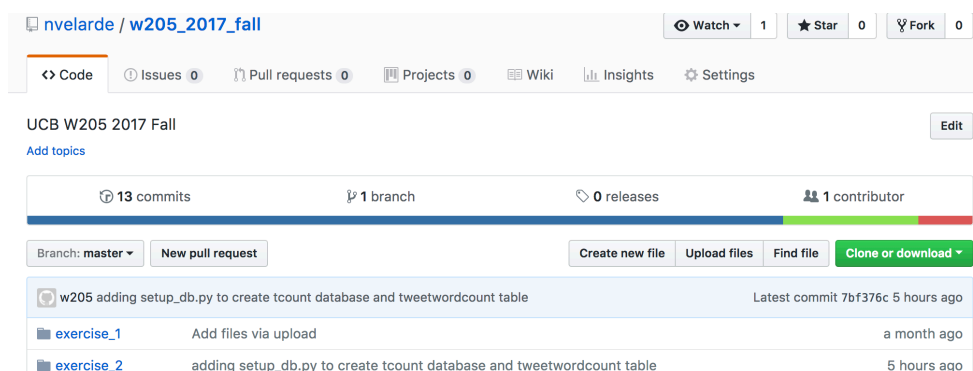
The application uses the Tweepy library (a python module used to interact with the Twitter API) to read the live stream of tweets from Twitter in the “Tweet-spout” component. Then, the “Parse-tweet-bolt” parses the tweets (cleaning them up by filtering out hash tags, user mentions, urls as well as leading and leading and lagging punctuations), extracts individual words from each parsed tweet and emits the words to the “Count-bolt” which counts the number of each word and updates the counts associated with each word in the “tweetwordcount” table in the “tcount” Postgres database. The “tcount” database and “tweetwordcount” table are created in advance of executing the streamparse code.

We implemented the topology as shown in the above figure where each circle represents an instance. Our topology has 3 instances of “Tweet-spout”, 3 instances of “Parse-tweet-bolt” and 2 instances of “Count-bolt.”

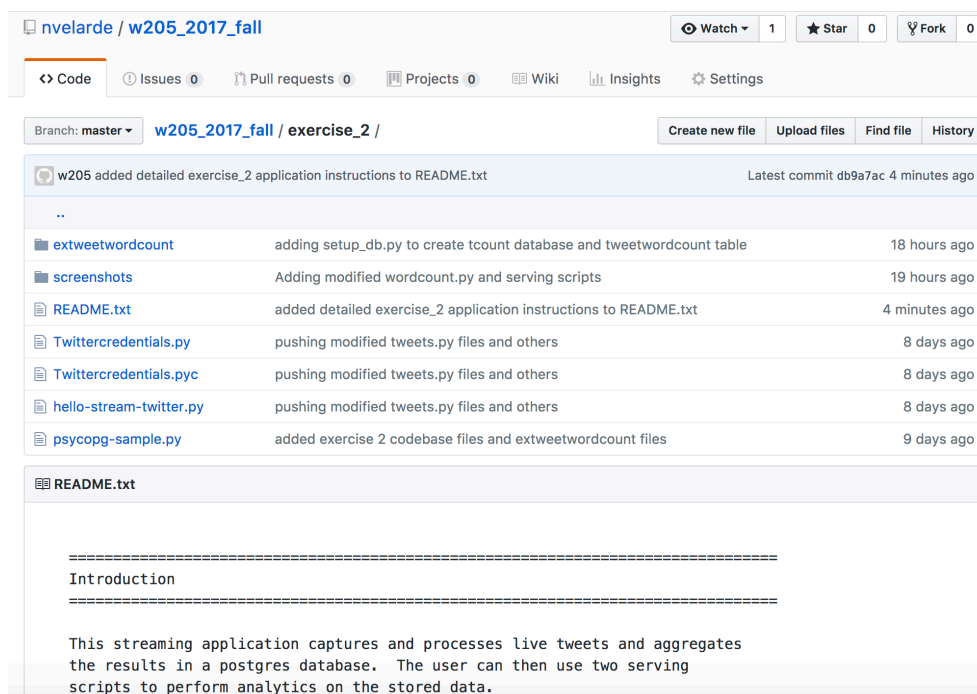
## Directory and File Structure

The application files can be accessed through my **w205\_fall\_2017** GitHub repository under the **exercise\_2/** directory.

[https://github.com/nvelarde/w205\\_2017\\_fall](https://github.com/nvelarde/w205_2017_fall)

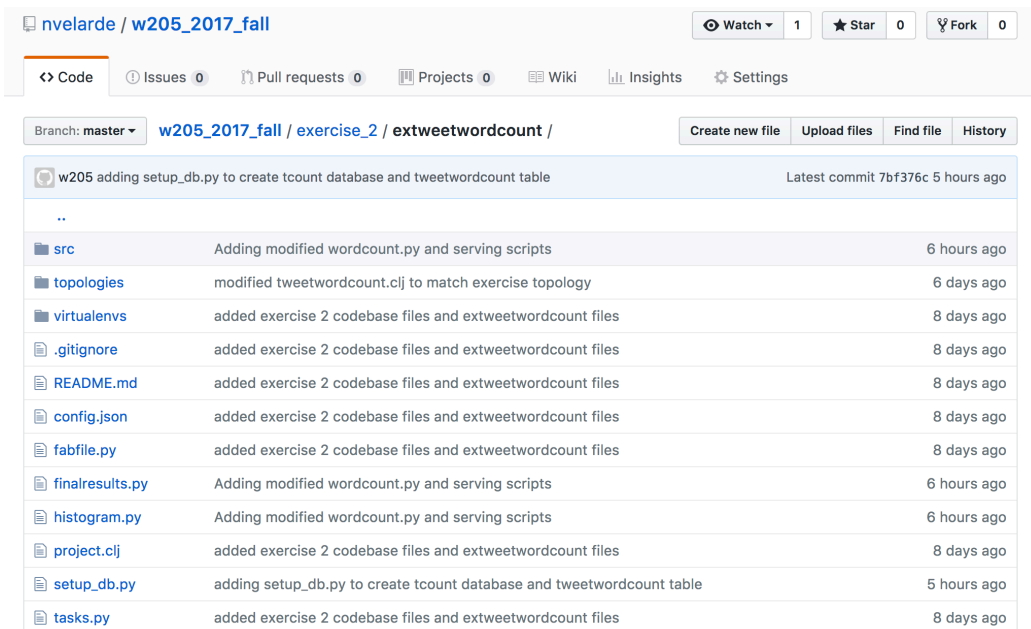


Within the **exercise\_2** directory, all of the files necessary to run the application are in the **exttweetwordcount** directory. The **exercise\_2** directory **README.txt** file has detailed instructions on how to run the application.



Key **exttweetwordcount** subdirectories and associated files:

- **src** (directory) – python source code for topology spouts and bolts
  - spouts (directory)
    - tweets.py
  - bolts (directory)
    - parse.py
    - wordcount.py
- **topologies** (directory)
  - tweetwordcount.clj



The **exttweetwordcount** directory also contains three python serving scripts:

- **setup\_db.py** – creates **tcount** database and **tweetwordcount** table – this script has to be run first as the rest of the application code is dependent on the existence of the **tcount** database and **tweetwordcount** table
- **finalresults.py** – when passed a single word as an argument, finalresults.py returns the total number of occurrences of the word in the captured Twitter stream. The user can run the program without passing an argument which will return all of the words in the stream and their total count of occurrences, sorted alphabetically.
- **histogram.py** – takes two arguments, k1 and k2, which must be positive integers separated by spaces and returns all the words with a total number of occurrences greater than or equal to k1 and less than or equal to k2.

## Application Execution Instructions

As mentioned on page 2 (see lower screenshot), detailed instructions and information required to setup and run the application can be found in the **README.txt** in the **exercise\_2** directory.

## Screenshots of Running Application

### Continuous log of incoming parsed tweets – words and word counts

```
keys — w205@ip-172-31-29-146:~/w205_2017_fall/exercise_2/extweetwordcount — ssh -i UCB_Lab1.pem
15444 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt are: 6
15446 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt of: 8
15447 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt Honduras: 1
15449 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt 5: 2
15451 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt so: 6
15453 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt seconds: 1
15459 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt is: 11
15466 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt of: 9
15473 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt just: 9
15474 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt summer&w: 1
15476 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt please: 1
15478 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt started: 1
15480 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt two: 1
15484 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt mummies: 1
15506 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt bare: 1
15508 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt voices*: 1
15511 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt with: 9
15514 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt on: 10
15520 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt me: 9
15533 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt the: 21
15536 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt me: 10
15539 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt picture: 1
15542 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt for: 5
15554 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt this: 7
15556 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt music: 2
15559 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt final: 1
15563 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt for: 6
15564 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt chaos: 1
15570 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt me: 11
15574 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt you: 6
15578 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt searching: 1
15585 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt are: 7
15591 [Thread-41] INFO backtype.storm.task.ShellBolt - ShellLog pid:6398, name:count-bolt best: 1
15598 [Thread-43] INFO backtype.storm.task.ShellBolt - ShellLog pid:6408, name:count-bolt Doesnt: 1
```

### Selected output from finalresults.py (case where argument is passed)

```
keys — w205@ip-172-31-29-146:~/w205_2017_fall/exercise_2/extweetwordcount —
[[w205@ip-172-31-29-146 extweetwordcount]$ python finalresults.py trump
Total number of occurrences of "trump" is: 2
[[w205@ip-172-31-29-146 extweetwordcount]$ python finalresults.py the
Total number of occurrences of "the" is: 359
[[w205@ip-172-31-29-146 extweetwordcount]$ python finalresults.py dsdfafdfdsf
There were no occurrences of "dsdfafdfdsf" in the captured Twitter stream
[[w205@ip-172-31-29-146 extweetwordcount]$
```

### Selected output from histogram.py

```
keys — w205@ip-172-31-29-146:~/w205_2017_fall/exercise_2/extweetwordcount
[[w205@ip-172-31-29-146 extweetwordcount]$ python histogram.py 45 55
get: 55
they: 55
how: 53
just: 52
we: 51
i: 49
your: 49
life: 46
[[w205@ip-172-31-29-146 extweetwordcount]$
```