

W205 Exercise 2

Name: Manish Shah

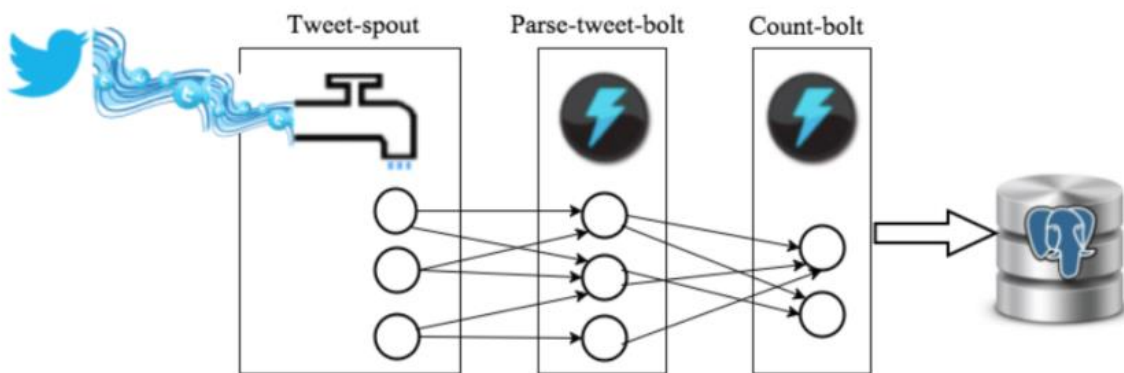
Date: 12/3/2017

The Application

This application subscribes to a tweeter feed and writes the number of occurrences of these words in the database. The tweets containing the words, ["a", "the", "i", "you", "u"] are tracked and only the tweets in English language are selected.

The application uses Storm technology to subscribe to the tweets and process them. The Tweet spout program subscribes to the tweet, queues them and emits it to the Parse bolt. The parse bolt processes the tweet, excludes retweets and tweets with non-ascii characters and creates a list of valid words and emits it to the wordcount bolt. The wordcount bolt updates the database. The application is implemented in python 2.7.

The application topology is illustrated in the diagram below.



Database

Postgres database is used to store the words and the counts of the words. The name of the database is tcount and it has one table, tweetwordcount. The tweetwordcount table has the following structure.

```
tcount=# \d+ tweetwordcount;
          Table "public.tweetwordcount"
  Column | Type          | Modifiers | Storage | Description
-----+-----+-----+-----+-----
 word   | text          | not null  | extended |
 count  | integer       | not null  | plain    |
Indexes:
    "tweetwordcount_pkey" PRIMARY KEY, btree (word)
Has OIDs: no
```

Topology

The application has one spout, tweet, with 3 instances, one parse bolt with 3 instances and one count bolt with 2 instances. The topology is defined as.

```
33 lines (31 sloc) 676 Bytes
1 (ns tweetwordcount
2   (:use [streamparse.specs])
3   (:gen-class))
4
5 (defn tweetwordcount [options]
6   [
7     ;; spout configuration
8     {"tweet-spout" (python-spout-spec
9       options
10      "spouts.tweets.Tweets"
11      ["tweet"]
12      :p 3
13      )
14    }
15     ;; bolt configuration
16     {"parse-tweet-bolt" (python-bolt-spec
17       options
18       {"tweet-spout" :shuffle}
19       "bolts.parse.ParseTweet"
20       ["word"]
21       :p 3
22       )
23      "count-bolt" (python-bolt-spec
24        options
25        {"parse-tweet-bolt" ["word"]}
26        "bolts.wordcount.WordCounter"
27        ["word" "count*"]
28        :p 2
29        )
30    }
31  ]
32 )
```

Directory Structure

The directory structure is as follows.

Files/Folders	Description
Create_database.py	Creates the database tcount and table tweetwordcount
Extweetwordcount	The streamparse project folder used by streamparse run
Final_results.py	Queries the database for words
Histogram.py	Queries the database for words with counts $\geq k$ and $\leq l$
Main.sh	Main program to run the script
Top20.py	Python program to query top 20 words based on count
Top20_plot.py	Python program to plot the bar chart (tested on windows)
Plot.png	Top 20 bar graph.
Output	Folder for the output files
Screen shots	Screen shots of the project

The topology, spout and bolts are defined in the Extweetwordcount folder.

Run Instructions:

The analysis has been tested on AWS hosted EC2 instance with 100 GB of disk space, 8 GB of Ram and 2 processors and using Linux operating system.

The AMI used is "UCB W205 Spring 2016 (ami-be0d5fd4)" .

The application needs tweepy and matplotlib libraries to be installed for python.

To run the program, on the command prompt, type,

`./main.sh`

This script will create the database, run the storm application for 5 mins and then run the result scripts.

The main.sh script is shown below.

```
Executable File 17 lines (16 sloc) 410 Bytes
1 echo "Creating the tcount database"
2 python create_database.py
3 cd extweetwordcount
4 sparse run > sparseoutput.txt&
5 sparse_pid=$!
6 date +%T"
7 sleep 5m
8 date +%T"
9 kill -9 $sparse_pid
10 cd ..
11 rm -rf output
12 mkdir output
13 python final_results.py this > output/final_results.txt
14 python final_results.py > output/final_results_all.txt
15 python histogram.py 10 300 > output/histogram.txt
16 python top_20.py > output/top_20.txt
```

Output:

1. Output of number of occurrences of word "this".

```
2 lines (1 sloc) 41 Bytes
1 Total number of occurrences of this :672
```

2. Output of all words sorted alphabetically with counts.

```
7996 lines (7995 sloc) | 87.6 KB
1  ! : 26
2  !! : 10
3  !!! : 4
4  !!!! : 2
5  !Proudly : 2
6  !SHILL! : 2
7  $$ : 2
8  $$$ : 2
```

3. Histogram of all words ≥ 20 and ≤ 50

```
2091 lines (2090 sloc) | 24.4 KB
1  Oregon's : 50
2  difficult : 50
3  Trust : 50
4  Harry : 50
5  middle-class : 50
6  Economics : 50
7  code : 50
8  won't : 50
9  Ocean : 50
10 Dreams : 50
```

4. Histogram bar graph.

