

Programming Assignment 2

Detecting the Most Popular Topics from Live Twitter Message Streams using the Lossy Counting Algorithm with Apache Storm

Due: March, 08 Tuesday 5:00PM

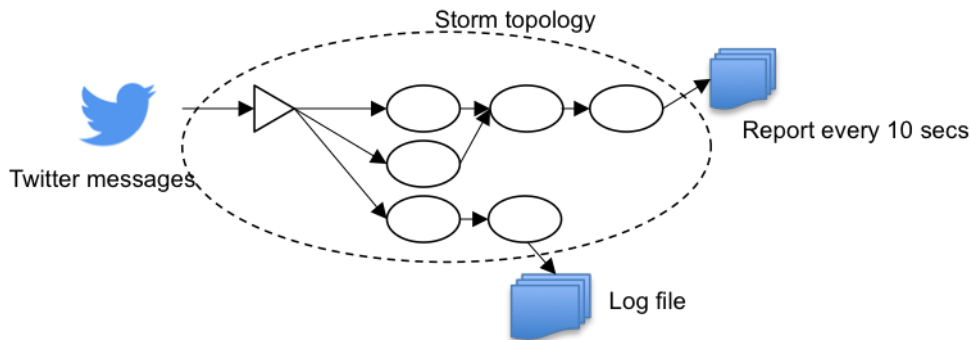
Submission: via Canvas, group submission

Instructor: Sangmi Lee Pallickara

Objectives

The goal of this programming assignment is to enable you to gain experience in:

- Implementing approximate on-line algorithms using a real-time streaming data processing framework
- Understanding and implementing parallelism over a real-time streaming data processing framework



1. Overview

In this assignment, you will design and implement a real-time streaming data analytics system using Apache Storm. The goal of your system is to detect the most frequently occurring hash tags from the live Twitter data stream in real-time.

A hashtag is a type of label or metadata tag used in social networks that makes it easier for users to find messages with a specific theme or content¹. Users create and use hashtags by placing the hash character (or number sign) # in front of a word or un-spaced phrase. Searching for that hashtag will then present each message that has been tagged with it. For example, *#springbreak* and *#zidane* were popular tags for the US on March 11, 2019.

Finding popular and trendy topics (via hashtags and named entities) in real-time marketing implies that you include something topical in your social media posts to help

¹ <https://en.wikipedia.org/wiki/Hashtag>

increase the overall reach. In this assignment, we will target data from live Twitter message provided by Twitter developers².

In this assignment, you will:

- Implement the Lossy Counting algorithm³
- List the top 100 most popular hashtags every 10 seconds
- Parallelize the analysis of your system

To perform above tasks, you are required to use Apache Storm, and Twitter Stream APIs.

2. Requirements of Programing Assignment 2

Your submission should include **the source codes** of two storm topologies specified in the section 5. Do not submit any data or result file.

To count the occurrences of hashtags, you should use the online algorithms included in this description. You are **not allowed to use any existing lossy counting algorithm implementations**.

Demonstration of your software should be on machines in CSB-120. This will include an interview discussing implementation and design details. Your submission should be via Canvas.

3. Install and setup your Storm cluster

You should create your own Storm cluster in CS120 with at least 5 nodes including the Nimbus node for this assignment. Here is the summary of the steps for setting up a Storm cluster⁴:

- a. Set up a Zookeeper cluster
- b. Install dependencies on Nimbus and worker machines
- c. Download and extract a Storm release to Nimbus and worker machines
- d. Fill in mandatory configurations into `storm.yaml`

Launch daemons under supervision using "storm" script and a supervisor of your choice

3.1. Set up your Zookeeper cluster

More information is available at:

<http://storm.apache.org/releases/2.1.1/Setting-up-a-Storm-cluster.html>

Storm uses Zookeeper for coordinating the cluster. Zookeeper is not used for message passing, so the load Storm places on Zookeeper is quite low. Single node Zookeeper clusters should be sufficient for most cases.

² <https://dev.twitter.com>

³ Gurmeet Singh Manku, and Rajeev Motwani, "Approximate Frequency Counts over Data Stream" 2002, VLDB

⁴ <https://storm.apache.org/documentation/Setting-up-a-Storm-cluster.html>

First, you should download a release of zookeeper:

```
$ wget https://archive.apache.org/dist/zookeeper/zookeeper-3.5.5/apache-  
zookeeper-3.5.5-bin.tar.gz  
$ tar xvf zookeeper-3.5.5-bin.tar.gz  
$ cd zookeeper-3.5.5-bin/conf
```

Now, in your /conf directory, save "zoo_sample.cfg" as "zoo.cfg" and modify with the following lines. This configuration does NOT process any environment variables such as \$HOSTNAME or \$USER. Please use an absolute path for this configuration. Change the PORT0 and PORT1 to a unique, unused port.

To avoid port conflict, please follow the port/server assignment used in PA1.

```
tickTime=2000  
dataDir=/s/*your_host_name*/a/tmp/zookeeper_*your_login_name*/data  
clientPort=PORT0  
admin.serverPort=PORT1
```

To run zookeeper, go to your zookeeper directory and run following command: (Make sure that you're already in a ssh session at *your_host_name* defined in zoo.cfg)

```
$ bin/zkServer.sh start
```

To check the status of your zookeeper instance, go to your zookeeper directory and run following command:

```
$ bin/zkServer.sh status
```

To stop your zookeeper instance, go to your zookeeper directory and run following command:

```
$ bin/zkServer.sh stop
```

You can also run zookeeper under supervision. See the section 3.3. By default, you will see the log file in your_zookeeper_directory/bin/zookeeper.out.

3.2. Set up a Storm cluster

Next, download a Storm release and extract the zip file somewhere on Nimbus and each of the worker machines.

```
$ wget http://mirrors.ocf.berkeley.edu/apache/storm/apache-storm-2.1.1/apache-storm-2.1.1.tar.gz  
$ tar xvf apache-storm-2.1.1.tar.gz  
$ cd apache-storm-2.1.1
```

The Storm release contains a file at conf/storm.yaml that configures the Storm daemons. You can see the default configuration values <https://github.com/apache/storm/blob/v1.1.1/conf/defaults.yaml>.

`storm.yaml` overrides anything in `defaults.yaml`. The mandatory configuration to get a working cluster includes:

- 1) `storm.zookeeper.servers`: This is a list of the hosts in the Zookeeper cluster for your Storm cluster. It should look something like:

```
storm.zookeeper.servers:  
  - "your_zookeeper_node.cs.colostate.edu"  
  
storm.zookeeper.port: <clientPort_in_zoo.cfg>
```

- 2) `nimbus.seeds`: The worker nodes need to know which machine the master is running on, in order to download topology jars and configurations. Specify the master node:

```
nimbus.seeds: ["your_nimbus_node.cs.colostate.edu"]
```

- 3) `storm.local.dir`: The Nimbus and Supervisor daemons require a directory on the local disk to store small amounts of state (like jars, confs, and things like that). You should create that directory on each machine, give it proper permissions, and then fill in the directory location using this configuration keyword. For example:

```
storm.local.dir: "/tmp/{Your-User-ID}-storm"
```

- 4) `supervisor.slots.ports`: For each worker machine, you can configure how many workers run on that machine with this configuration keyword. Each worker uses a single port for receiving messages, and this setting defines which ports are open for use. If you define five ports here, then Storm will allocate up to five workers to run on this machine. If you define three ports, Storm will only run up to three. By default, this setting is configured to run 4 workers on the ports 6700, 6701, 6702, and 6703. You can specify them with the ports in your port range. For example:

```
supervisor.slots.ports:  
  - your_port_number_1  
  - your_port_number_2  
  - your_port_number_3  
  - your_port_number_4
```

- 5) `nimbus.thrift.port`

```
nimbus.thrift.port: PORT5
```

- 6) `supervisor.thrift.port`:

```
supervisor.thrift.port: PORT6
```

- 7) Web UI: Specify the port number and host for your UI;

```
ui.port: PORT7
ui.host: "your_nimbus_node.cs.colostate.edu"
```

To avoid port conflict, please follow the port assignment posted on the assignment web page. To specify the port number for your UI, please see the default configuration file: <https://github.com/apache/storm/blob/master/conf/defaults.yaml>

3.3. Launch daemons under supervision

Next step is to launch all the Storm daemons. It is critical that you run each of these daemons under supervision. Storm is a fail-fast system, which means the processes will halt whenever an unexpected error is encountered. Storm is designed so that it can safely halt at any point and recover correctly when the process is restarted. This is why Storm keeps no state in-process -- if Nimbus or the Supervisors restart, the running topologies are unaffected. First, go to your storm directory and copy supervisor configuration. Supervisor is already installed in the CS120 cluster.

Go to your storm installation folder and create three configuration files.

```
more /etc/supervisord.conf > zk-supervisord.conf
more /etc/supervisord.conf > nimbus-supervisord.conf
more /etc/supervisord.conf > worker-supervisord.conf
```

Now, add zookeeper commands to your configuration file, zk-supervisord.conf.

```
**Modify following lines
[unix_http_server]
file=/tmp/{Your-User-ID}-storm/supervisor.sock ; (the path to the socket file)

[supervisord]
logfile=/tmp/{Your-User-ID}-storm/supervisord.log ; (main log file;default $CWD/supervisord.log)
pidfile=/tmp/{Your-User-ID}-storm/supervisord.pid ; (supervisord pidfile;default supervisord.pid)
;minfds=1024 ; (min. avail startup file descriptors; default 1024)
;minprocs=200 ; (min. avail process descriptors; default 200)

[rpcinterface:supervisor]
supervisor.rpcinterface_factory = supervisor.rpcinterface:make_main_rpcinterface

[supervisorctl]
serverurl=unix:///tmp/{Your-User-ID}-storm/supervisor.sock ; use a unix:// URL  for a unix socket

[program:zookeeper]
command={your_zookeeper_directory}/bin/zkServer.sh start-foreground

**comment below lines
;[include]
;files = supervisord.d/*.ini
```

Add nimbus commands to your configuration file, nimbus-supervisord.conf.

```
**Add/Modify following lines
[unix_http_server]
file=/tmp/{Your-User-ID}-storm/supervisor.sock ; (the path to the socket file)
```

```
[supervisord]
logfile=/tmp/{Your-User-ID}-storm/supervisord.log ; (main log file;default $CWD/supervisord.log)

[rpcinterface:supervisor]
supervisor.rpcinterface_factory = supervisor.rpcinterface:make_main_rpcinterface

[supervisorctl]
serverurl=unix:///tmp/{Your-User-ID}-storm/supervisor.sock ; use a unix:// URL  for a unix socket

[program:storm_ui]
command={your_storm_directory}/bin/storm ui

[program:storm_nimbus]
command={your_storm_directory}/bin/storm nimbus

**comment below lines
;[include]
;files = supervisord.d/*.ini
```

Add command for the worker nodes to your configuration file,
worker-supervisord.conf.

```
**Modify following lines
[unix_http_server]
file=/tmp/{Your-User-ID}-storm/supervisor.sock ; (the path to the socket file)

[supervisord]
logfile=/tmp/{Your-User-ID}-storm/supervisord.log ; (main log file;default $CWD/supervisord.log)

[rpcinterface:supervisor]
supervisor.rpcinterface_factory = supervisor.rpcinterface:make_main_rpcinterface

[supervisorctl]
serverurl=unix:///tmp/{Your-User-ID}-storm/supervisor.sock ; use a unix:// URL  for a unix socket

[program:storm_supervisor]
command={your_storm_directory}/bin/storm supervisor

**comment below lines
;[include]
;files = supervisord.d/*.ini
```

Now we have supervisor configuration to manage all the necessary processes to start storm cluster. On your zookeeper node, start your zookeeper cluster. First create directory for your log files and socket.

```
mkdir /tmp/{user_ID}-storm
supervisord -c {your_storm_directory}/zk-supervisord.conf
```

On your nimbus node,

```
mkdir /tmp/{user_ID}-storm
supervisord -c {your_storm_directory}/nimbus-supervisord.conf
```

On your worker nodes, start supervisors.

```
mkdir /tmp/{user_ID}-storm
supervisord -c {your_storm_directory}/worker-supervisord.conf
```

The UI can be accessed by navigating your web browser to
`http://{nimbus_host.cs.colostate.edu}:{your-UI-port}`

3.4. Launch your Storm job

To compile your file and launch your job, using Maven (maven is installed in CS120 cluster) is a good option.

To find examples, go to `{your_storm_directory}/examples/storm-starter` and follow the `README.markdown` in the same directory. You can start from the section, “## Build and install Storm jars locally”.

Open your `bashrc` (`vim ~/.bashrc`) and add this (update `PATH` variable):

```
Export PATH=$PATH:${HOME}/apache-storm-2.1.0/bin
```

Apply the change:

```
source ~/.bashrc
```

Once you have created your local Storm jar that also includes your topology, you can also use following command to launch your job.

```
storm jar target/storm-starter-1.1.1.jar  
org.apache.storm.starter.RollingTopWords production-topology remote
```

To kill a topology, simply run,

```
storm kill {topology name}
```

You should use the name that you have used when submitting the topology. For more information, please visit;

<https://storm.apache.org/documentation/Running-topologies-on-a-production-cluster.html>

3.5. Reading Live Twitter messages from your Storm Spout

1. Getting your application keys for OAuth

Create your Twitter account and log in. Next, go to

<https://developer.twitter.com/en/apps> and create your application. After creating the application on the application page, click “Keys and Access Tokens” tag and create your access token. You should use four keys (Consumer Key (API key), Consumer Secret (API secret), Access Token, and Access Token Secret) from this web page.

2. Downloading Twitter4J⁵

⁵ <http://twitter4j.org/en/>

Twitter4J is a Java library for the Twitter API. You can use pure HTTP GET to retrieve messages. However, Twitter4J will provide the simplest access to the Twitter messages from your Storm spout.

3. Using Twitter4J from your Storm Spout

```
TwitterStream twitterStream = new TwitterStreamFactory(
    new ConfigurationBuilder().setJSONStoreEnabled(true).build())
    .getInstance();

twitterStream.addListener(listener);
twitterStream.setOAuthConsumer(consumerKey, consumerSecret);
AccessToken token =
    new AccessToken(accessToken, accessTokenSecret);
twitterStream.setOAuthAccessToken(token);
```

4. Lossy Counting Algorithm

4.1. Definitions

The incoming stream is conceptually divided into *buckets* of width $w = \left\lceil \frac{1}{\epsilon} \right\rceil$ transactions each. Buckets are labeled with *bucket ids*, starting from 1. We denote the current bucket id by $b_{current}$, whose value is $\left\lceil \frac{N}{w} \right\rceil$. For an element e , we denote its true frequency in the stream seen so far by f_e . Note that ϵ and w are fixed while N , $b_{current}$ and f_e are running variables whose values change as the stream progresses.

Data structure D is a set of entries of the form (e, f, Δ) , where e is an element in the stream, f is an integer representing its estimated frequency, and Δ is a maximum possible error in f .

4.2. Algorithm

Initially D is empty. Whenever an element arrives, first lookup to see whether an entry for that element already exists or not. If the lookup succeeds, update the entry by incrementing its frequency f by one. Otherwise, create a new entry of the form $(e, 1, b_{current} - 1)$. We also prune D by deleting some of its entries at bucket boundaries, i.e., whenever $N \bmod w == 0$. The rule for deletion is: an entry (e, f, Δ) is deleted if $f + \Delta \leq b_{current}$. When a user requests a list of items with threshold s , we output those entries in D where $f \geq (s - \epsilon)N$.

For an entry (e, f, Δ) , f represents the exact frequency count of e ever since this entry was inserted into D . The value of Δ assigned to a new entry is the maximum number of times e could have occurred in the first $b_{current} - 1$ buckets. This value is exactly $b_{current} - 1$. Once an entry is inserted into D , its Δ value remains unchanged.

4.3. Example⁶

```

(The parameter  $s$  is not used until the end of the algorithm)
 $\epsilon = 0.2$ 

 $w = 1/\epsilon = 5$       (5 items per "bucket")

Input: 1 2 4 3 4 3 4 5 4 6 7 3 3 6 1 1 3 2 4 7
      |         | |         | |         | |         |
      +-----+ +-----+ +-----+ +-----+
      bucket 1  bucket 2  bucket 3  bucket 4

=====
 $b_{current} = 1$       inserted: 1 2 4 3 4
-----
Insert phase:
  D (before removing): (x=1;f=1; $\Delta$ =0) (x=2;f=1; $\Delta$ =0) (x=4;f=2; $\Delta$ =0) (x=3;f=1; $\Delta$ =0)

Delete phase: delete elements with  $f + \Delta \leq b_{current}$  (=1)
  D (after removing) : (x=4;f=2; $\Delta$ =0)

NOTE: elements with frequencies  $\leq 1$  are deleted
      New elements added has maximum count error of 0

=====
 $b_{current} = 2$       inserted: 3 4 5 4 6
-----
Insert phase:
  D (before removing): (x=4;f=4; $\Delta$ =0) (x=3;f=1; $\Delta$ =1) (x=5;f=1; $\Delta$ =1) (x=6;f=1; $\Delta$ =1)

Delete phase: delete elements with  $f + \Delta \leq b_{current}$  (=2)
  D (after removing) : (x=4;f=4; $\Delta$ =0)

      NOTE: elements with frequencies  $\leq 2$  are deleted
      Newly added elements have a maximum count error of 1

=====
 $b_{current} = 3$       inserted: 7 3 3 6 1
-----
Insert phase:
  D (before removing): (x=4;f=4; $\Delta$ =0) (x=7;f=1; $\Delta$ =2) (x=3;f=2; $\Delta$ =2) (x=6;f=1; $\Delta$ =2)
  (x=1;f=1; $\Delta$ =2)

Delete phase: delete elements with  $f + \Delta \leq b_{current}$  (=3)
  D (after removing) : (x=4;f=4; $\Delta$ =0) (x=3;f=2; $\Delta$ =2)

NOTE: elements with frequencies  $\leq 3$  are deleted
      Newly added elements have a maximum count error of 2

=====
 $b_{current} = 4$       inserted: 1 3 2 4 7
-----
Insert phase:
  D (before removing): (x=4;f=5; $\Delta$ =0) (x=3;f=3; $\Delta$ =2) (x=1;f=1; $\Delta$ =3) (x=2;f=1; $\Delta$ =3)

```

⁶ <http://www.mathcs.emory.edu/~cheung/Courses/584-StreamDB/Syllabus/07-Heavy/Manku.html>

```
(x=7;f=1;Δ=3)
```

Delete phase: delete elements with $f + \Delta \leq b_{\text{current}}$ ($=4$)

D (after removing) : (x=4;f=5;Δ=0) (x=3;f=3;Δ=2)

NOTE: elements with frequencies ≤ 4 are deleted

Newly added elements have a maximum count error of 3

Interpreting the content in D:

Item	f_{manku}	f_{actual}
4	5	5
3	3	5

5. Software Requirements

For this assignment, you should implement 2 topologies: (1) Non-parallel and (2) parallel topologies.

5.1. Non-parallel topology

For a non-parallel topology, your topology should:

- Retrieve data from the Twitter stream API in real-time from a **single** spout
- Provide a report of the top 100 most popular hashtags (sorted in descending order) every 10 seconds. If your software does not have 100 hashtags at the end of the 10 second interval, your report may contain hash tags those are available at that time.
- Here, the time stamp is the time stamp at the beginning of a **10-second window**.
- Generate a log file containing the top 100 hash tags and timestamps. Each log entry should have the following format:

```
Log file
<timestamp_in_epochtime><hashtagTop1><hashtagTop2><hashtagTop3>...<hashtagTop100>
```

Each of the messages should generate no more than 1 line of log entry. If there is no hashtag included in the message, please do not create a log entry. The bolts for counting and logging should be implemented as two different bolts.

5.2. Parallel topology

The parallel topology tracks the counts in parallel. Design and implement your topology including 4 different bolts that implement the lossy counting algorithm. These bolts should run on **different nodes**.

Your topology should,

- Retrieve data using the Twitter stream API in real-time from a single spout
- Perform counting without redundancy using the lossy counting algorithm on 4 bolts on different nodes
- Aggregate values and provide the top 100 most popular hashtags (sorted in descending order) every 10 seconds

- Generate log files that follow the format specified in section 6.1.

6. Grading

Your submissions will be graded based on the demonstration to the instructor in CSB120. The grading will be done on a 15-point scale. This assignment will account for 15% of your final course grade.

- 4 points: Set up the Storm cluster and other software
- 5 points: Implementation of the Non-Parallel topology and analysis
- 5 points: Implementation of the Parallel topology
- 1 point: Participation

7. Late Policy

Please check the late policy posted on the course web page.