Twitter Stream Crawler Project

UCR CS172 - Winter 2015

Nicolas Lawler and Henry Garcia

Collaboration Details

- Henry
 - o connected to the twitter stream, saved the tweets to json files, implemented the script for crawling the tweets, set up the documentation
- Nicolas
 - implemented the link identifying logic, implemented the title grabbing functionality, performed disk io as necessary to support the above, refined the documentation

Overview of System

- System Libraries
 - o json-simple 1.1.1
 - o isoup 1.8.1
 - o twitter4j 4.0.2
 - o JavaSE 1.7
- Architecture
 - The System is split into two main parts, the Crawler and the TitleFetcher.
 - The Crawler connects to the twitter streaming API and downloads a number of tweets (specified by the user), writing them to disk in batches of files. The size of the files are determined by the user
 - The TitleFetcher takes as input the directory that the tweets were saved in, and loads them a file at a time. It searches the text of the tweet for links, grabs the title of the html response when it finds one, and adds that title to the json of the tweet. All the tweets are then written back to the same file that they came from.
- Crawling/data collection strategy
 - o Crawler
 - The crawler first checks the directory for a file containing saved tweet ids. If the file exists the hashmap will be loaded with the tweet ids, otherwise the hashmap will start empty.

- The tweet crawler first creates a filter in order to grab more tweets that are in English and have a location. After creating the filter it establishes a connection with twitter's stream, using account and consumer tokens provided through https://dev.twitter.com/
- After establishing a connection with the twitter stream the system will start to grab tweets. First it will check if the tweet exists in the hashmap. If it does, then it will move onto the next tweet.
- If the tweet has not been crawled yet, then it will add the id to the hashmap, add the id to the hashmapque, and check if it has geolocation information.
- If the tweet has geolocation, then it will be added to the saveTweetQueue. otherwise it will move onto the next tweet.
- The hashmapque and saveTweetqueue are used to store the information into the files after the crawler has reached a certain amount of crawled tweets (100 by default).
- When saving the tweets the crawler will check if the file is larger than that specified by the user. The crawler will continue to increment the tweet file until it finds a file that is no longer larger than the file size provided by the user (useful if the user has already crawled tweets and doesn't want them to be placed in already crawled files).

TitleFetcher

- The TitleFetcher loads a File[] full of all the files in the provided directory, and then processes each file one at a time, skipping over files that don't have tweets.
- It parses the json into a JSONObject line by line and inserts that object into the input list. Then, for each item in the list, take the tweet and pass it off to a thread pool, where it will be serviced by the next available thread. Each worker begins by checking if the tweet has been processed before. If it hasn't, we get the text from the JSONObject, and attempt to find a url. If a url is found, we follow it and wait for an html response. If the response comes through, we parse it and get the title. Then, we add that title as a new field in the JSONObject. If the link turned out to be bad or the title was empty, we mark the tweet with a field "hasBadLink" so that it won't be searched again. Whether or not any fields have been added, we enqueue the tweet onto the output queue and the thread moves on to the next available tweet.
- When all the tweets have been processed, the contents of the output queue are written back to the same file that they came from. Then the process starts again for the next file.

- Data Structures
 - HashMap
 - For hashing crawled tweets so that they won't be crawled again
 - o ConcurrentLinkedQueue
 - For thread-safe writes when collecting processed tweets.
 - ArrayList
 - For loading and storing sets of tweets.
 - JSONObject
 - A simple json wrapper over a java HashMap, allows automatic encoding and decoding of JSON strings, used to represent tweets

Limitations

- If a file called "hashedtweets.txt" or a tweet file in the following format "tweets.*.json" exist and are not formated the way the crawler expects them to be, then the program will likely throw parsing errors and not function correctly.
- If the user cancels (^C) the system while the system is writing to the hashed tweet file or the tweetfile, then there is a possibility that the file will be corrupt and throw a parsing error on startup (if run in the same directory).
- The regular expression that matches to urls in the text string isn't fully robust. for example, it matches to urls ending in '.', and the resulting http request fails in a case where it may have succeeded.

Deployment

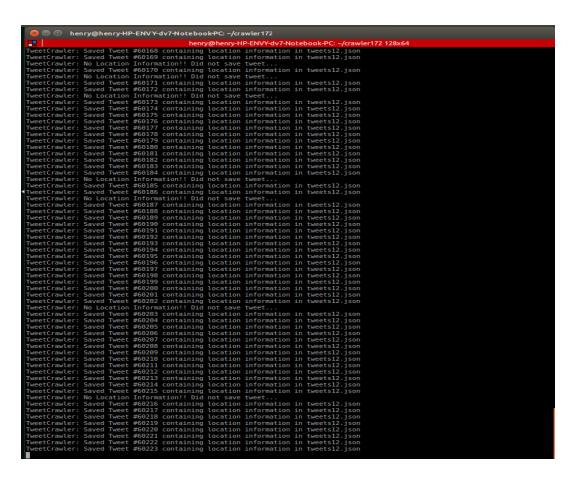
./runCrawlerExecutable.sh <Tweets> <File Sizes (mb)> <Threads> <outputdir> OR

./runCrawlerExecutable.sh <Tweets> <File Sizes (mb)> <Threads> <outputdir> <save Every # Tweets>

- <Tweets> refers to the number of tweets the user wants to grab.
- < File Sizes (mb) > refers to the size of files the user wants the tweets to be stored in
- <Threads> specifies the number of threads the user wants to create
- **<outputdir>** specifies the directory in which the user would like to store all of the data.
- <save Every # Tweets> specifies how many tweets the user would like to obtain before the crawler saves the tweetids and tweets into the corresponding files. 100 is the default setting.

Screenshots

```
The trailer: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets0.json
TweetCrawler: Saved Tweet #950 containing location information in tweets
```



```
0
                                                                                                                                                      USB9)
state: "Foday is going to be a good dayo<u>s IIIB</u>Malkingbead #BetterCaliSaul http://t.co/5Zs27UISH)" from tweet 13368
filmi
street 13372
Fuser 13372
                                                                                                                                          from tweet 13382
044U819J
from tweet 13384
UUX2qyX0
from tweet 13391
(5KXYcuT
    9
                                                               E. New London phases a com on Nutter: "Junior roller derby at gSDBollerBerby. Because. 11's junior roller derb
https://i.co/runMbask7
- "Alessa | Spanis Banks - Bog Area" from tweet 13415
http://i.co/dsHaffgFf
- "Alessa | Spanis Banks - Bog Area" from tweet 13415
http://i.co/dsHaffgFf
- "Instagram" from tweet 13419
http://i.co/dsBollerBerby | The Stage Banks - Ban
                                                                                                                                                      coups
es.com on Twitter: "Junior roller derby at @SDMollerDerby. Because, it's junior roller derby! http://t.co/75ce9CdGj4" from tweet 13409
BBLESY
Ø
· ko
  ube" from tweet 13459
HO6
                                                                                                               rier - YouTube" from tweet 13459
.co/OZFAmazHO6
le | Weiner Dog Derby" from tweet 13472
  0
                                                                                                                                    G3lek5ZY7
.com Earn 10s for every link visit by your friend* from tweet 13473
TDBECMKRW
from tweet 13476
                                                                                                                                                                        No nFultter: "i post so much on my story@i cant meny ic new burden with a much burden wit
                                                             nttp://t.co/Fm645MqmhV
2: "Careers Center - Hospital Account Manager" from tweet 13497
https://t.co/7Amjazd2sh
```

```
Deny@heny=HP-ENY_dyT-Notebook-PC:-/crawler172 | heny@heny=HP-ENY_dyT-Notebook-PC:-/crawler172 128x64 |

DRI: http://t.co/walbsUg38 | heny@heny=HP-ENY_dyT-Notebook-PC:-/crawler1
```