

Data Intensive Computing

Lab 2

Ashutosh Ahmad Alexandar

Harish Ganesan

This lab is broken down into many phases namely –

- Data Collection – Collect NYT URL and Tweets
- Data Processing – Scrape HTML Body from URL, take only text from tweet information
- Data Cleaning – Take only relevant data from scraped information, split data into several txt files
- Map Reduce – Feed multiple
- Sorting (Post Processing) – Sorting reduced output to get the highest occurrences
- Data Visualization – Creation of WordCloud



For this lab, we have considered three topics which are currently trending on the internet –

- **Trump+Russia**
- **NBA**
- **Blockchain**

As we can see from the above block diagram, we start off with the data collection.

We use the TwitterAPI and the NYTimes API to get tweets and article URLs based on keyword search. An example would be “trump+putin”.

Once we have the relevant tweets and URLs, we need to process these URLs using BeautifulSoup in python, to get the HTML body of the URLs we have. Then we must take only the paragraph tags, as these only contain the relevant information that we require. Once we have the paragraph tags, we take all the words in all the paragraphs and store it within a single file. We have one such file for each article.

We also take only the text portion of the tweet data and combine several tweet’s text into a single file, and we have multiple such files.

We can start Hadoop now on the VM, using start-hadoop.sh , which has already been provided.

Next, we need to input all these files into Hadoop, using the mapper and reducer code that we have provided. We will Then get an output file in the format ‘part-00000’. This output file contains the reduced output of input we have provided, and this is a set of words and their respective counts.

Once we have the output file, we run the text_sort.py script that we have provided, and this sorts the ‘part-00000’ file into a CSV file which is in descending order with respect to the word count.

As we can see from the screenshot below, an output folder is created in the HDFS. Then we need to copy that to our local file system using “hdfs dfs -get /output /home/Hadoop/op”.

Then we run the text_sort.py script to sort the file and obtain a new CSV.

After this step, we can stop Hadoop, using stop-hadoop.sh that was already provided.

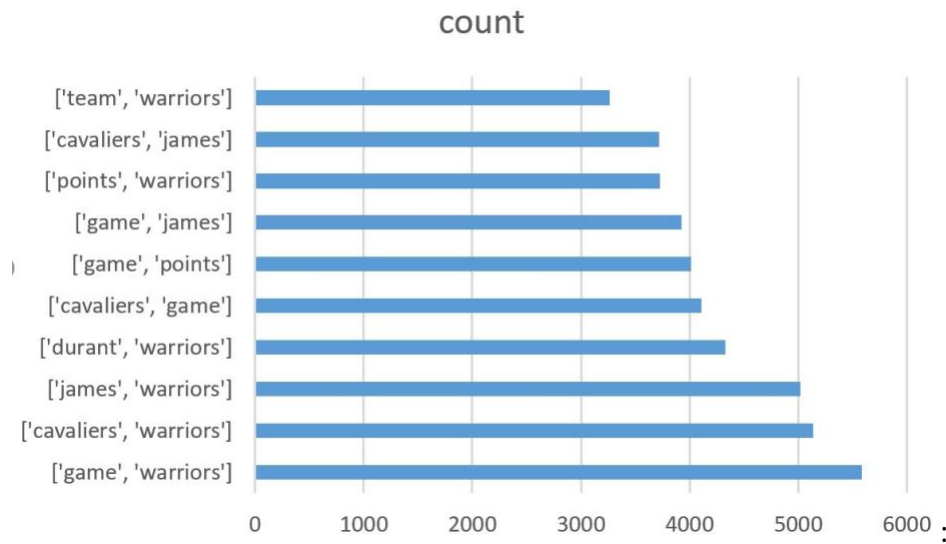
```
File Edit View Terminal Tabs Help Terminal - hadoop@hadoop-VirtualBox: ~/trump_nyt_co_oc

HDFS: Number of write operations=203
Map-Reduce Framework
  Map input records=1938
  Map output records=704624
  Map output bytes=16530412
  Map output materialized bytes=17940860
  Input split bytes=20692
  Combine input records=0
  Combine output records=0
  Reduce input groups=90
  Reduce shuffle bytes=17940860
  Reduce input records=704624
  Reduce output records=90
  Spilled Records=1409248
  Shuffled Maps =200
  Failed Shuffles=0
  Merged Map outputs=200
  GC time elapsed (ms)=55616
  CPU time spent (ms)=0
  Physical memory (bytes) snapshot=0
  Virtual memory (bytes) snapshot=0
  Total committed heap usage (bytes)=33352691712
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=291827
File Output Format Counters
  Bytes Written=2438
18/04/08 12:59:40 INFO streaming.StreamJob: Output directory: /user/hadoop/trump_tweet_cooc_op
hadoop at hadoop-VirtualBox in ~/hadoop using <> 18-04-08 - 12:59:40
  o hdfs dfs -get trump_tweet_cooc_op /home/hadoop/trump_tweet_cooc_op
18/04/08 13:00:22 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using b
uiltin-java classes where applicable
hadoop at hadoop-VirtualBox in ~/hadoop using <> 18-04-08 - 13:00:24
  o cd ../trump_nyt_co_oc
hadoop at hadoop-VirtualBox in ~/trump_nyt_co_oc using <> 18-04-08 - 13:01:00
  o python text_sort.py
hadoop at hadoop-VirtualBox in ~/trump_nyt_co_oc using <> 18-04-08 - 13:01:03
  o vi trump_nyt_co_oc.sorted.csv
hadoop at hadoop-VirtualBox in ~/trump_nyt_co_oc using <> 18-04-08 - 13:01:29
  o stop-hadoop.sh
18/04/08 13:06:56 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Stopping namenodes on [localhost]
localhost: stopping namenode
localhost: stopping datanode
Stopping secondary namenodes [0.0.0.0]
0.0.0.0: stopping secondarynamenode
18/04/08 13:07:19 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
hadoop at hadoop-VirtualBox in ~/trump_nyt_co_oc using <> 18-04-08 - 13:07:20
```

We then consider only the top 30 rows to create the WordCloud.

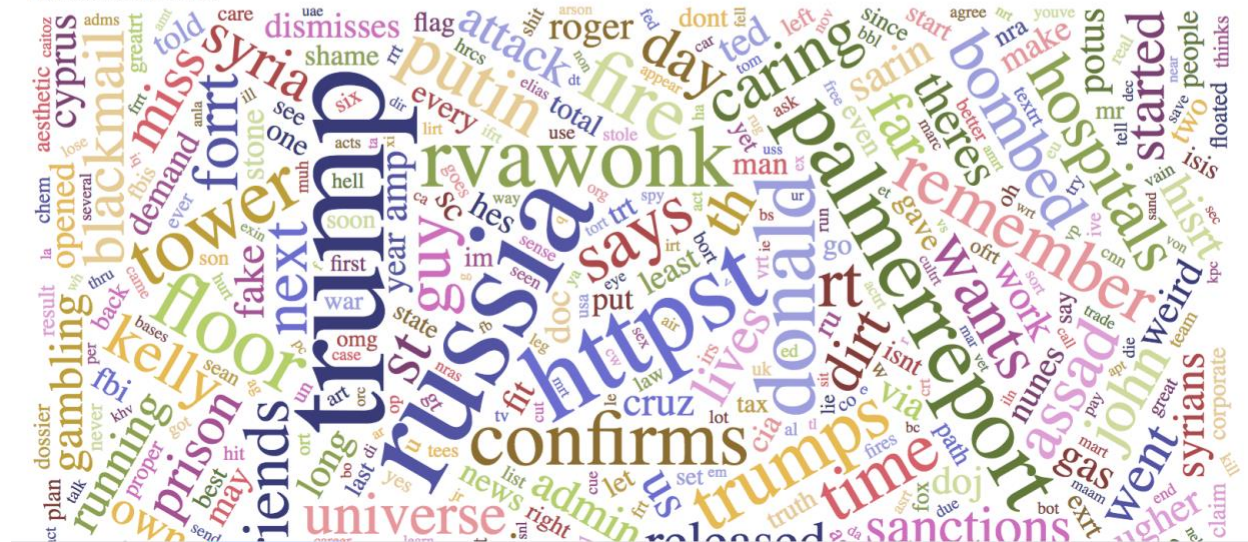
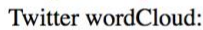
We use d3.js to create the word-clouds, by inputting these sorted CSV files into the HTML page, we then just take our data and draw a word-cloud with it. An example is shown below.

Upon removing these duplicates, we create a simple bar chart to visualize these co-occurrences. An example of these co-occurrences are is seen below:



This is the co-occurrence count for the pairs, which came from the NBA NY Times data.

Topic: ☒ trumprussia ☐ nba ☐ blockchain
NY Times wordCloud:



<https://buffalo.app.box.com/s/54kva33ijplicafiv23orccifsgot385>

