
Data Intensive Computing CSE 587- LAB 2

Apurba Mahanta-50288705

Sagnik Ghosh-50289151

University at Buffalo, The State University of New York

apurbama@buffalo.edu, sagnikgh@buffalo.edu

Abstract

The goal of this project is to demonstrate Data Aggregation, Big data analysis and Visualization on big data which is collected over multiple sources like Twitter, New York Times and Common crawl data using the API's offered by them. The outcomes of this project include evaluating the reliability of data from different sources like social media and news media. We worked on Hadoop 2.X and HDFS platforms and applied Big data algorithms like word count and word co occurrence. We used Python as our data processing language. We used Tableau for visualization and evaluate the results obtained from the above infrastructure and algorithms.

1 Data Analysis on Gun violence and control in the United States

We chose our topic as Gun violence and control in the United States. Our subtopic are as follows:

- a) Gunman
- b) Gun Control
- c) NRA
- d) Gun Sense
- e) Firearms and Second Amendment

Mass shooting and shooting at public places has been in the news in 2019. A lot of people and prominent figures have been expressing their views and concerns in various places like Twitter. A lot of criticism has also been seen in the public and hence this makes an interesting topic to understand the sentiments of people by collecting data from various sources and find the statistics about what people are thinking.

2 Big Data Infrastructure setup

- 1) We used the virtual image provided by our TA's to set up the environment.
- 2) We uploaded the image in virtual box and provided the ram memory size as 8 gb
- 3) We installed ubuntu by running the image.
- 3) After the infrastructure is set up, we ran sample commands to test the framework like start-dfs.sh and hadoop commands like `hadoop jar /home/cse587/hadoop-3.1.2/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.2.jar wordcount /test/exam.txt /test/output`.

3 Dataset size

Time Period : Jan 2019 to Apr 2019

- 1) Number of Tweets : 48390
- 2) New York Times articles : 757
- 3) Commoncrawl data articles : 722

34 4 Implementation

35 We collected data from 3 different sources Twitter, NYT and Commoncrawl. We then cleaned the
36 data, removed the stop words from the articles using NLTK word lemmatizer in the mapper pro-
37 cess. The reducer process gave us the results which includes top 10 most frequent words and their
38 co-occurences. The implementation flow of the data is as shown below:

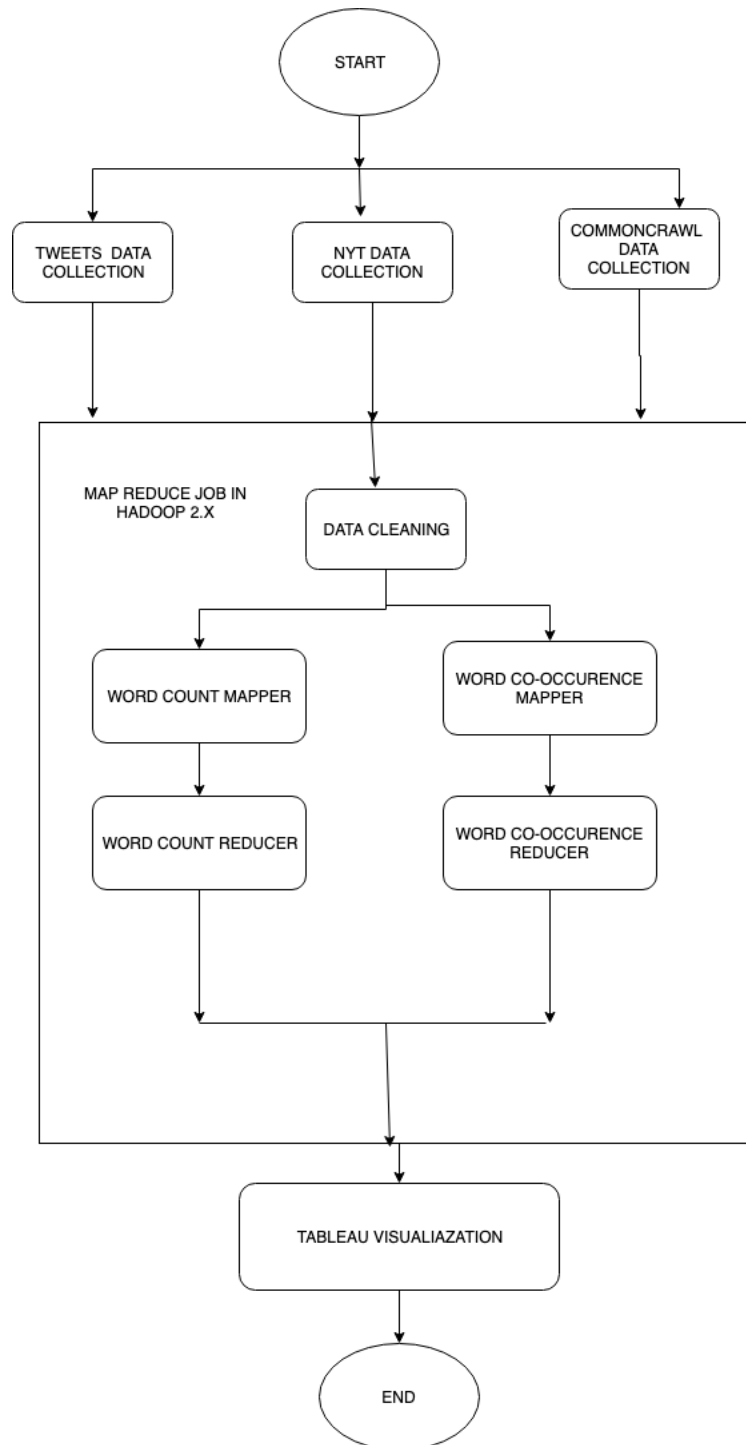


Figure 1: Implementation Flow

39 4.1 Twitter Data Collection

40 We used Tweepy api in python to collect data from twitter data.The script parses the tweets and
41 collects text and user id.We removed the re tweets and duplicated tweets and saved the tweets to be
42 used for map reduce jobs.The following hashtags was used for tweets collection:

43 **#gunman**
44 **#shooting**
45 **#guncontrol**
46 **#gunsense**
47 **#firearms**
48 **#secondamendment**
49 **#nra**
50 **#gunlaws**
51

52 4.2 NYT Data Collection

53 To gather New york times articles we used the api provided by NYT to collect the articles.After
54 getting the related html contents from the url's returned,we then used Beautiful soup api to scrape and
55 remove all the unnecessary tags in the html content .The api helped in getting the text from the body
56 in the html content.Each article is then merged eventually to be further processed by map reduce
57 framework. The following queries were used :

58 **"gun control","gunman","nra","gun sense","firearms","second amendment"**

59 4.3 Commoncrawl Data Collection

60 For collecting data from common crawl we used python to search the relevant domains in common
61 crawl archive directory.The following domains were used :

62 **guncontrolDomain = 'gun-control.procon.org/*'**
63 **guncontrolDomain = 'freerepublic.com'**
64 **gunsenseDomain = 'gunviolencearchive.org/*'**
65 **nraDomain = 'nra.org/articles/*'**
66 **firearmsDomain = 'usatoday.com/story/*'**

67 **The indices searched were "2019-04","2019-09","2019-11","2019-14","2019-16"**
68

69 After iterating over all the hits in the archive we get the warc files returned .We then unzip
70 them on the fly and get the html contents.We then parse the html contents using Beautiful soup to get
71 the text from the body of the html.We then have a check if the text is relevant to our keywords,we
72 then store it in our disk .The relevant text is later merged and forwarded to map reduce framework for
73 further processing

74 4.4 Word Count using map reduce framework

75 MapReduce is a processing technique and a program model for distributed computing based on java.
76 The MapReduce algorithm contains two important tasks, namely Map and Reduce. Map takes a set
77 of data and converts it into another set of data, where individual elements are broken down into tuples
78 (key/value pairs). Secondly, reduce task, which takes the output from a map as an input and combines
79 those data tuples into a smaller set of tuples

80 4.4.1 Mapper word count

81 1)We cleaned the data inside the mapper function using NLTK word lemmatizer python library.This
82 removed all the stop words to make the data more meaningful.The mapper then emits a key value
83 pair for each word and its count as 1.

84 4.4.2 Reducer word count

85 The reducer collects all the data and then calculates the word count for each word .We then sort all
86 the data based on count and emit only the top 10 most frequent words from the dataset.

87 4.4.3 Mapper word co-occurrence

88 1) We cleaned the data inside the mapper function using NLTK word lemmatizer python library. This
89 removed all the stop words to make the data more meaningful. The mapper then checks each word and
90 its neighbour if they are among the top 10 words and emits the pair of it and counts as 1

91 4.4.4 Reducer word co-occurrence

92 The reducer collects all the data and then calculates the count for each word pair co-occurrence. We
93 then emit the total count for each word co-occurrence.

94 4.5 Visualization using Tableau

95 To demonstrate the results we used Tableau as our visualization tool for word cloud visualization. There are in total 12 visualization results provided. The word cloud visualizations include the
96 following results:

98 **Twitter**

- 99 1) Word count - Top 10 words - All data
- 100 2) Word co-occurrence of the top 10 words - All data
- 101 3) Word count - Top 10 words - 1 day data
- 102 4) Word co-occurrence of the top 10 words - 1 day data

103 **NYT data**

- 104 1) Word count - Top 10 words - All data
- 105 2) Word co-occurrence of the top 10 words - All data
- 106 3) Word count - Top 10 words - 1 day data
- 107 4) Word co-occurrence of the top 10 words - 1 day data

108 **Commoncrawl data**

- 109 1) Word count - Top 10 words - All data
- 110 2) Word co-occurrence of the top 10 words - All data
- 111 3) Word count - Top 10 words - 1 day data
- 112 4) Word co-occurrence of the top 10 words - 1 day data

113

114 5 Results and Visualization

115 **These results have been published online and can be viewed by the following clickable hyper link:**

116

117 [https://public.tableau.com/profile/sagnik.ghosh#!/vizhome/CSE587-DataIntensiveComputing-](https://public.tableau.com/profile/sagnik.ghosh#!/vizhome/CSE587-DataIntensiveComputing-Lab2/Main-Dashboard?publish=yes)
118 [Lab2/Main-Dashboard?publish=yes](https://public.tableau.com/profile/sagnik.ghosh#!/vizhome/CSE587-DataIntensiveComputing-Lab2/Main-Dashboard?publish=yes)

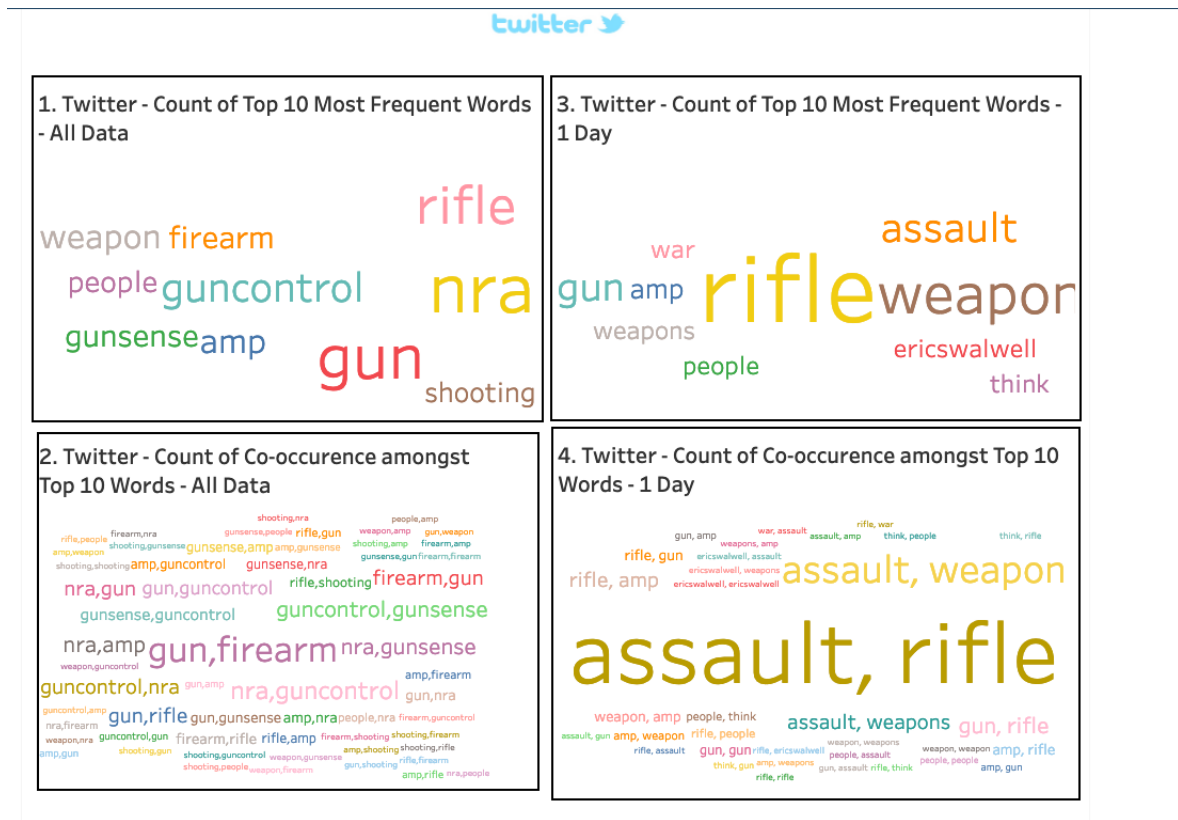


Figure 2: Twitter Word Cloud Visualization

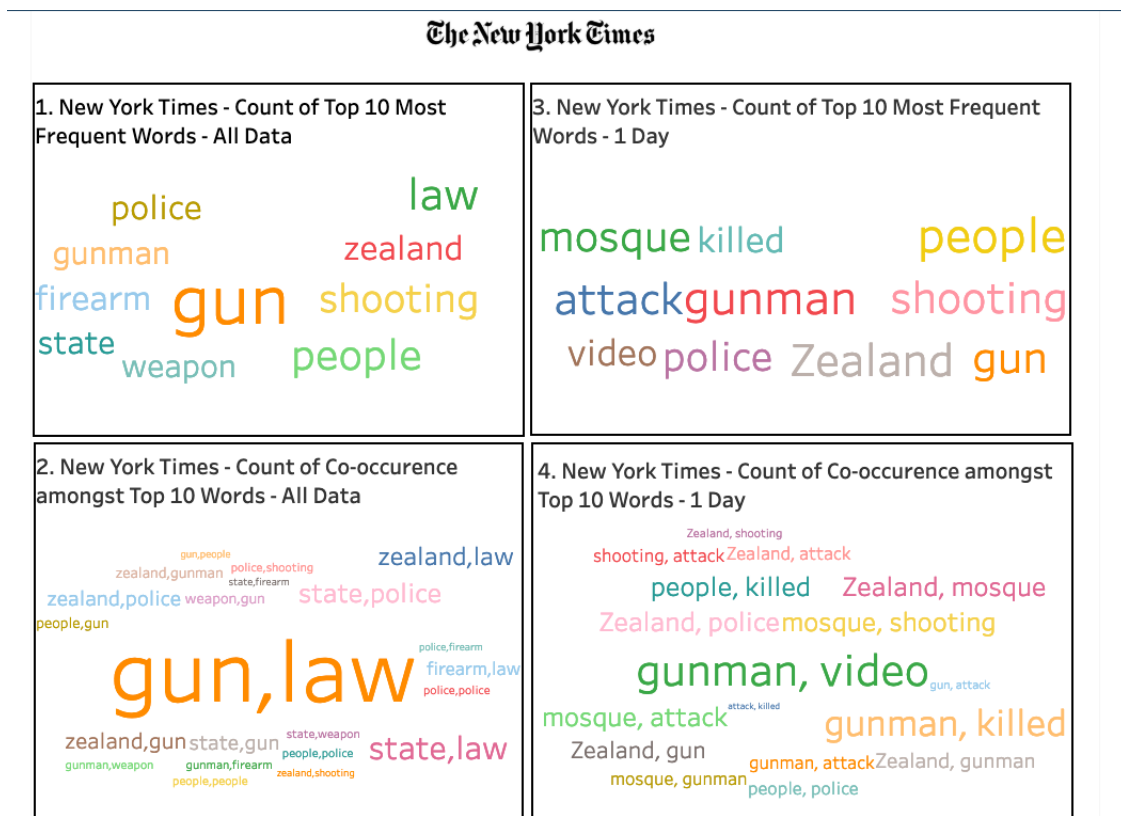


Figure 3: New York Times Cloud Visualization

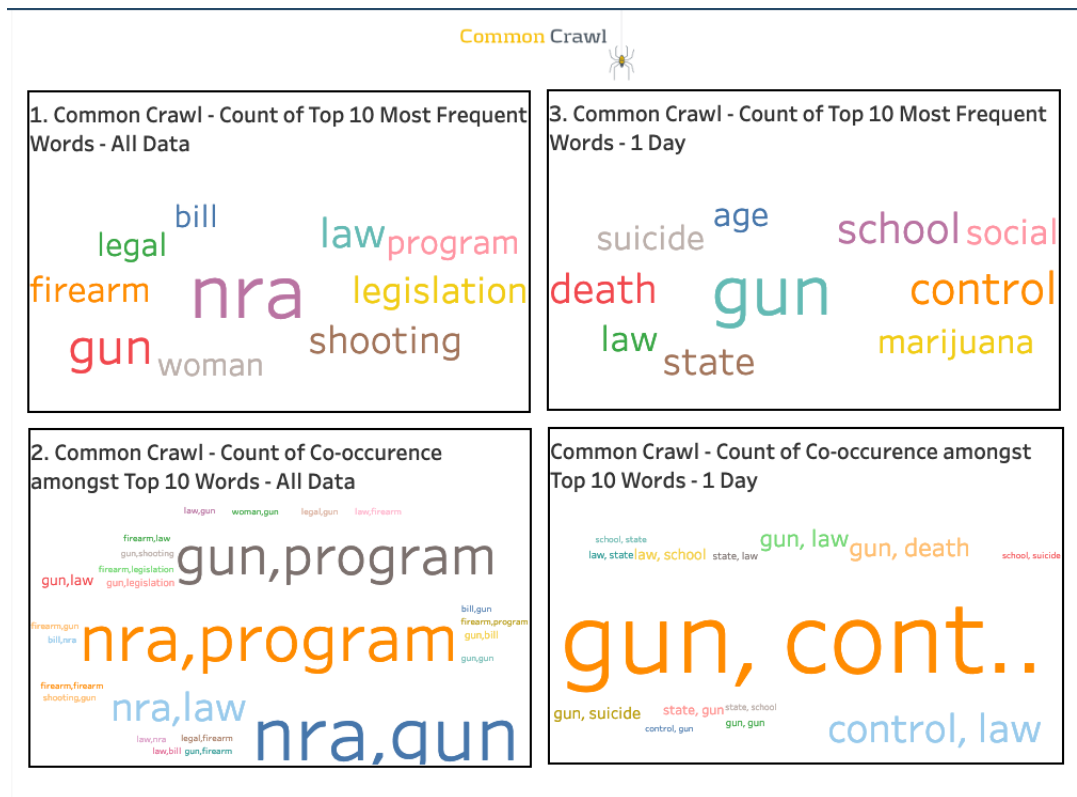


Figure 4: Commoncrawl word Cloud Visualization

6 Commands to execute

```

1) start-dfs.sh
2) hdfs dfs -put /part1/data/nyt/guns_nyt_all.txt
3) hadoop jar hadoop-3.1.2/share/hadoop/tools/lib/hadoop-streaming-3.1.2.jar -
file /home/cse587/mapper_nltk.py -mapper 'python3 mapper_nltk.py' -
file /home/cse587/reducer_nltk.py -reducer 'python3 reducer_nltk.py' -file
/home/cse587/nltkandyaml.mod -input /user/apurbama/MR/input/twitter/guns_twitter_All.txt
-output /user/apurbama/MR/output/twitter/wc/alldata/guns_all
4) hdfs dfs -getmerge /user/apurbama/MR/output/twitter/wc/alldata/guns_all guns_output.txt

```

References

```

1) https://cse.buffalo.edu/bina/cse487/spring2019/Lectures/Lab2/
2) https://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/
3) https://www.bellingcat.com/resources/2015/08/13/using-python-to-mine-common-crawl/
4) https://kb.tableau.com/articles/howto/creating-a-word-cloud

```