

LAB2: DATA AGGREGATION, BIG DATA ANALYSIS AND VISUALIZATION

DEBOSHREE MAZUMDAR(deboshe), VANDANA PRASAD GAALLEE(vandanap)

April 21, 2019

1 Introduction

The project involves data collection from three major sources: Twitter, NewYork Times and Common Crawl. We have selected the main topic as Politics with subtopics: Trump, Democrats, Republicans, Election and President.

2 Data Collection

2.1 Twitter

We have collected the tweets using tweepy streaming API. Total tweets collected for 14 MB in size. The tweets have been cleaned and pre-processed using regular expressions and python nltk package.

2.2 NewYork Times

NewYork times has the article API which has been used to collect 500 articles. The data is approximately 19 MB in size. All the HTML tags have been parsed using beautiful soup.

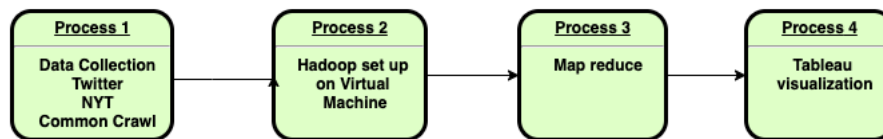
2.3 Common Crawl

We retrieved data from USA today filtering out with respect to the sub-topics we have. After collection, we processed the data using a python script. The data size is approximately 74 MB.

3 Hadoop setup

we installed oracle VM and added the image provided, created all folders in Hdfs file system and transferred the data and lastly, executed the map reduce program.

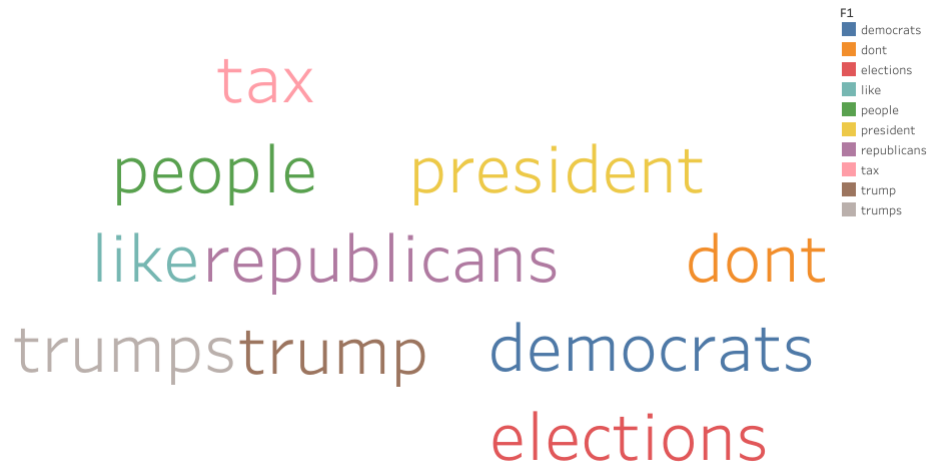
```
Activities Terminal Sun 13:03 cse587@cse587: ~
File Edit View Search Terminal Help
you." 13
you: 3
you; 5
you? 8
you?" 3
you_ 1
young 45
young, 3
young--here 1
young--like 1
young, 2
young; 1
younger 3
younger, 1
younger; 1
your 281
yours 3
yours! 2
yours, 5
yours, 3
yours." 1
yours; 1
yourself 8
yourself, 4
yourself," 1
yourself, 6
yourself?" 2
yourself?" 1
yourselves 2
youth 5
youthful 2
zeal; 1
zealous 1
zoöphagous 3
zoöphagous, 1
zoöphagy!" 1
{pg 8
{pg}184 1
£1 1
£10 1
at. 1
atat 1
cse587@cse587:~$ ^C
cse587@cse587:~$ hdfs dfs -copyToLocal /vandanap/MR/output/part-00000 /home/cse587/Documents/dic/output
cse587@cse587:~$ hadoop jar hadoop-3.1.2/share/hadoop/tools/lib/hadoop-streaming-3.1.2.jar -file /home/cse587/Documents/dic/part2/mapper.py -mapper /home/cse587/Documents/dic/part2/mapper.py -reducer /home/cse587/Documents/dic/part2/reducer.py -reducer /home/cse587/Documents/dic/part2/reducer.py -input /vandanap/MR/input/pg345.txt -output /vandanap/MR/output
```



4 Tableau visualization

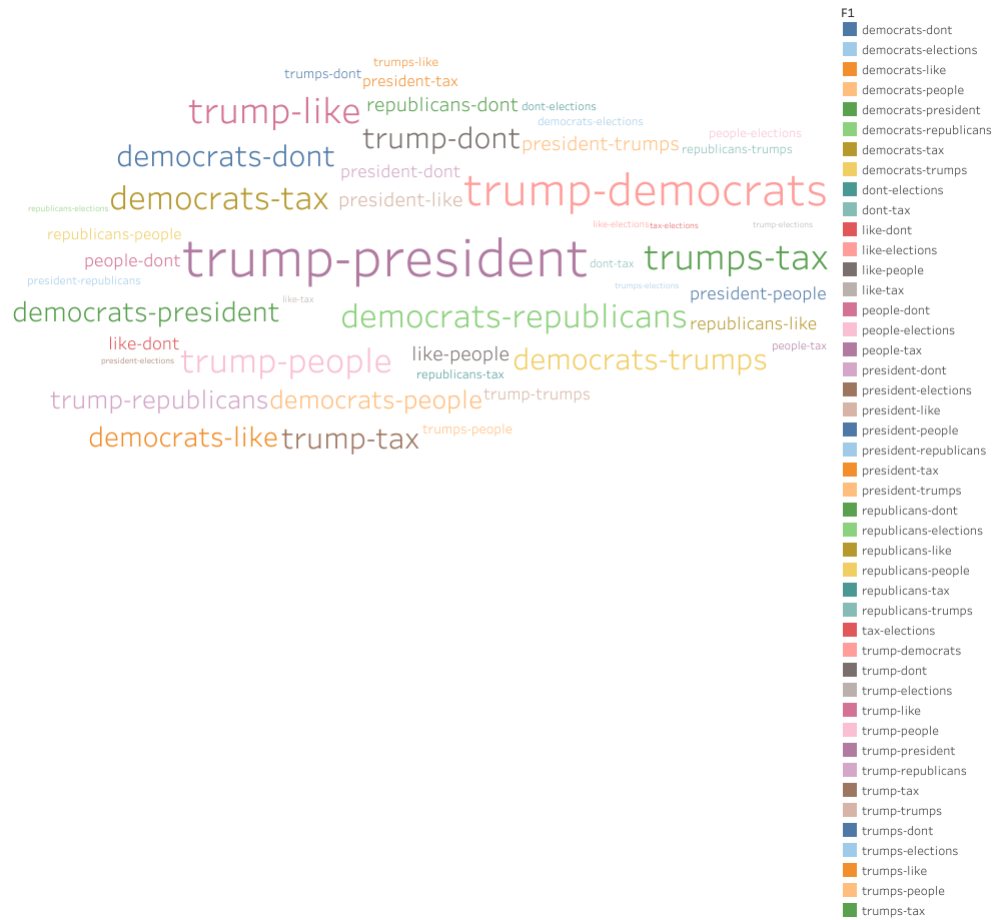
4.1 Twitter

word_count



4.2 NewYork Times

word_co-occurrence



4.3 Common Crawl

said people trump
 house president new
 party one advertisement
 democrats



democrats-house
democrats-new
democrats-one
democrats-party
democrats-people
house-one
house-party
house-people
new-advertisement
new-house
new-one
new-party
new-people
one-party
people-one
people-party
president-advertisement
president-house
president-new
president-one
president-party
president-people
said-advertisement
said-house
said-new
said-one
said-party
said-people
said-president
said-trump
trump-advertisement
trump-democrats
trump-house
trump-new
trump-one
trump-party
trump-people
trump-president



