

DIC LAB II REPORT

asishkak-50288695 | swapnika-50289464

Objective: The main goal of the project is to aggregate data from multiple sources and analyze the data using big data method of MapReduce and finally storing the data collected on WORM infrastructure Hadoop. A visualization data product is built from this analysis.

Procedure:

Development language used: Python

Topic of Interest: Politics

Sub-topics: Mueller Report, Border Wall, Trump, North Korea – USA Summit, Trump retreats on Healthcare

STEP-1:DATA COLLECTION

Data is collected from 3 sources which include twitter, NYTimes and Common Crawl.

Twitter:

- First, we get the API keys from the twitter developers website that are used to get the data.
- These keys are given OAuthHandler for authentication, just to make sure that the right person with right keys is accessing the data. The below lines explain it.

```
auth = tweepy.OAuthHandler(consumer_key,consumer_secret)
auth.set_access_token(access_token,access_secret)

api = tweepy.API(auth)
```

- Now, with the necessary keywords, we have collected tweets since 2019-01-01 which are originated within USA by using '39.8,-95.583068847656,2500km' as geocoordinates with radius.
- We wrote all the retrieved tweets to a text file by filtering out retweets by post fixing “-filter: retweets” to the keywords.
- Various arguments include query, since, geo-location etc.,

NYTimes:

- API keys are obtained from NYTimes and these keys are used to collect articles.
- The url's are collected by passing the arguments query, begin date, end date and pages range. This is done by using “articleAPI('key').search()”.
- From the list of URLs, the articles are obtained using BeautifulSoup.

Common Crawl:

- Common Crawl is a public archive which maintains snapshots of crawled data for various domains.
- In this assignment, we have to collect the most recent data since Jan 2019, and which also should be related to our keywords. So, the snapshots since 2019 are 2019-13, 2019-09, 2019-04.
- Now, the process of collecting common crawl is not by using API keys but a bit different.
- Common Crawl stores data as archives like WARC, WET (Text Archives), indexes etc.,
- Now, to retrieve data from WARC, first we should know where our data from a certain domain is located in the snapshots.
- So, we use <http://index.commoncrawl.org/> to get a JSON response of records where each record has its index, offsets, and length of data we want to retrieve.
- By using these indices, offset and length. We crawl the data using BeautifulSoup and filtering out all the <p> tags in the webpage.
- After retrieving the text, we now write the data to a text file.

STEP-2: PROCESSING OF DATA

- In this step, we remove all the unnecessary junk in our corpus.
- In data pre-processing we first make our entire corpus to lowercase. So, when we try to remove stop words in the later steps. We won't have any case-sensitive problems.
- In the next step links and usernames are removed with the help of regex library.
- Special characters and symbols are deleted with the help of a string matching and also, we've included one more step to remove all the left out special characters and numbers by removing all other than "\n" and alphabets by using regex library.
- Finally, stemming is done to reduce the words in our corpus to its stem words. So, when we reduce them to stem words. It becomes easy and the probability of achieving a higher word count is very high. This is done by using SnowballStemmer from nltk library.
- Code: *preprocessing.ipynb*

STEP-3: BIG DATA INFRASTRUCTURE

- Hadoop infrastructure is set up for storing and running MapReduce jobs on the collected corpus.
- This is done using the cloudera docker image which is provided to us.
- Now, we start the docker image using the below command.
 - `docker run --hostname=quickstart.cloudera --privileged=true -t -i -v /Users/Asish/Documents/DockerMR:/src --publish-all=true -p 8888 cloudera/quickstart /usr/bin/docker-quickstart`
 - I have provided /Users/Asish/Documents/DockerMR:/ as my src folder because dockers has a specific set of directories which are mounted to docker containers. In my case,
 - /Users
 - /Volumes

- /private
 - /tmp
- Now, let's create few folders in Hadoop file system where we can put our corpus. To create, we use the below commands.
- `hadoop fs -mkdir /user/asishkak`
 - `hadoop fs -mkdir /user/asishkak/MR`
 - `hadoop fs -mkdir /user/asishkak/MR/input`
 - `hadoop fs -mkdir /user/asishkak/MR/input/twitter`
 - `hadoop fs -mkdir /user/asishkak/MR/input/nyt`
 - `hadoop fs -mkdir /user/asishkak/MR/input/cc`
- Now, let's create few folders in Hadoop file system where we can get output. To create, we use the below commands.
- `hadoop fs -mkdir /user/asishkak/MR/output`
- Now, we copy the corpus from our local file system to hadoop file system by using the below commands.
- We move data to their specific folders.
 - For instance, twitter data is added to its specific folder in hadoop file system.
 - `hadoop fs -put *.txt /user/asishkak/MR/input/twitter/`
- Now, we've written MapReduce algorithms for both word count and word co-occurrence (Pair method).
- To execute them on our corpus, we give in mapper and reducer python files and input files on which mapping and reducing happens. The command to execute it is as below.
- `hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.6.0-cdh5.7.0.jar -file /src/mapper.py -mapper /src/mapper.py -file /src/reducer.py -reducer /src/reducer.py -input /user/asishkak/MR/input/twitter/* -output /user/asishkak/MR/output/twitter`

```

@quickstart/src/data/input/tw -- docker run --hostname=quickstart.cloudera --privileged=true -t -v ~/Documents/Docker/MR/src --pub...
~/wtf@Academia: Local/University/SEM2/CIC -- jupyter-notebook - python -- --publish-all=true -p 8888 cloudera/quickstart /usr/bin/docker-quickstart
[root@quickstart tw]# hadoop jar /usr/lib/hadoop-mapreduce/hadoop-streaming-2.6.0-cdh5.7.0.jar -file /src/mapper.py -mapper /src/mapper.py -file /src/reducer.py -reducer /src/reducer.py -input /user/asishkak/MR/input/* -output /user/asishkak/MR/output/twcount
19/04/21 03:42:21 WARN streaming.StreamJob: -file option is deprecated, please use generic option -files instead.
packageJobJar: [/src/mapper.py, /src/reducer.py] [/usr/jars/hadoop-streaming-2.6.0-cdh5.7.0.jar] /tmp/streamjob6604708394353196827.jar
tapDir=null
19/04/21 03:42:22 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
19/04/21 03:42:22 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
19/04/21 03:42:23 INFO mapred.FileInputFormat: Total input paths to process: 4
19/04/21 03:42:23 INFO mapreduce.JobSubmitter: number of splits:4
19/04/21 03:42:23 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1555803922186_0004
19/04/21 03:42:23 INFO impl.YarnClientImpl: Submitted application application_1555803922186_0004
19/04/21 03:42:23 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8080/proxy/application_1555803922186_0004/
19/04/21 03:42:23 INFO mapreduce.Job: Running job: job_1555803922186_0004
19/04/21 03:42:31 INFO mapreduce.Job: Job job_1555803922186_0004 running in uber mode : false
19/04/21 03:42:31 INFO mapreduce.Job:  map 0% reduce 0%
19/04/21 03:42:39 INFO mapreduce.Job:  map 25% reduce 0%
19/04/21 03:42:40 INFO mapreduce.Job:  map 50% reduce 0%
19/04/21 03:42:41 INFO mapreduce.Job:  map 100% reduce 0%
19/04/21 03:42:46 INFO mapreduce.Job:  map 100% reduce 100%
19/04/21 03:42:47 INFO mapreduce.Job: Job job_1555803922186_0004 completed successfully
19/04/21 03:42:47 INFO mapreduce.Job: Counters: 49
File System Counters
  FILE: Number of bytes read=4291136
  FILE: Number of bytes written=5171125
  FILE: Number of read operations=6
  FILE: Number of large read operations=6

```

```
Reduce input groups=1973
Reduce shuffle bytes=4283354
Reduce input records=406030
Reduce output records=1973
Spilled Records=912112
Shuffled Maps=4
Failed Shuffles=0
Merged Map outputs=4
GC time elapsed (ms)=313
CPU time spent (ms)=9356
Physical memory (bytes) snapshot=1454731264
Virtual memory (bytes) snapshot=6066159360
Total committed heap usage (bytes)=1405470400

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=2795577
File Output Format Counters
  Bytes Written=20258
13/04/21 03:42:47 INFO streaming.StreamJob: Output directory: /user/asishkak/MR/output/twcount
[root@quickstart tw]#
```

- A part_0000 is generated where we have our word count / word co-occurrence.
- To see the output. We execute the below command.
 - `hadoop fs -cat /user/asishkak/MR/output/part*`
- To view the output in sorted descending order of word count / co-occurrence.
 - `hadoop fs -cat /user/asishkak/MR/output/part* | sort -n k2 -r`

```
Bytes Written=20258
13/04/21 03:42:47 INFO streaming.StreamJob: Output directory: /user/asishkak/MR/output/twcount
[root@quickstart tw]# hadoop fs -cat /user/asishkak/MR/output/twcount/part* | sort -n -k2 -r | head -n1000
trump 39236
muller 19500
barack 17439
wall 14500
report 12782
presid 5140
class 4931
wang 3870
build 3500
relea 3335
investig 3000
doug 2954
shut 2815
vote 2554
subpoena 2184
house 2163
congress 2123
democrat 2061
mexico 1998
app 1800
republican 1750
imigr 1643
plan 1603
```

- To download the word count/ co-occurrence to our local file system. We can execute two commands which are
 - Hadoop fs -get /user/asishkak/MR/output/twitter/*
 - hadoop fs -cat /user/asishkak/MR/output/part* | sort -n k2 -r >> filename.txt

STEP-4: VISUALIZATION AND ANALYSIS

- In this step, we've used tableau for visualization.
- We have used the output file (text file) which contains the counts of words/ word co-occurrences respectively to generate a word cloud.
- For Instance, the word cloud generated for top 10 and 100 words for twitter corpus is below.

