

Quantitative Analyses for Customer Satisfaction and Brand Imaging

(ReadMe)

The following are the steps that we have followed for this project.

1. Data Scraping from different platforms
2. Code or steps for data ingest
3. Data cleaning for the scraped data and profiling code
4. Data analytics over the cleaned data.

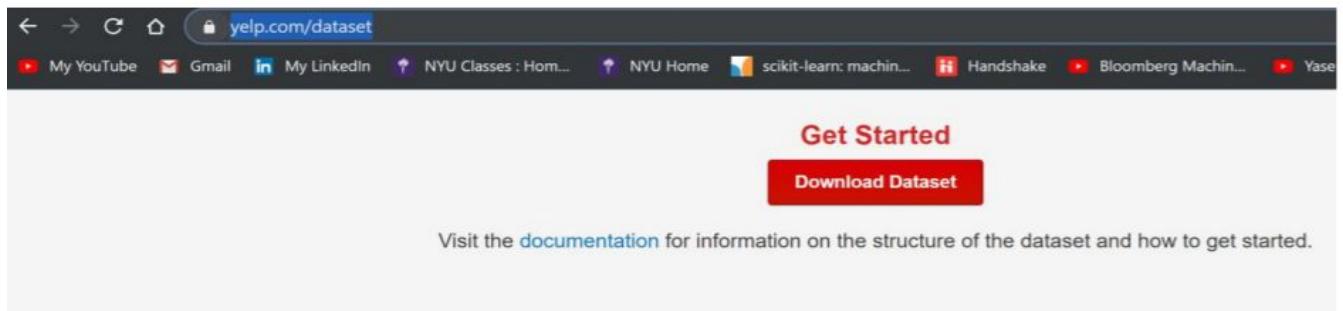
Additionally, the challenges that we faced are also mentioned in each step.

1. Data Scraping

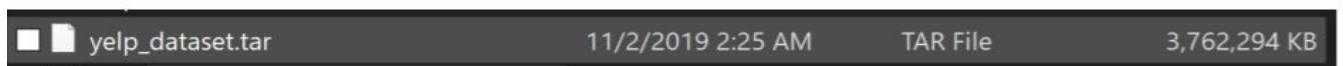
We have scrapped the data from 3 different sources: Yelp, Twitter and Reddit.

YELP :

For yelp, we got the data available from there website: → <https://www.yelp.com/dataset/download>. The total size of the data is **8GB**. Data had to be further cleaned for the analyses.



- The data is available in .tar format



The data files were extracted using the command on command prompt at the folder destination as:

tar -xvf yelp_dataset.tar

Yelp Extracted datafiles in JSON format:

Name	Date modified	Type	Size
business.json	11/15/2018 11:22 AM	JSON File	135,039 KB
checkin.json	11/15/2018 11:25 AM	JSON File	399,227 KB
Dataset_Challenge_Dataset_Agreement	1/14/2019 11:31 AM	Adobe Acrobat Do...	99 KB
photo.json	1/11/2019 7:06 PM	JSON File	25,060 KB
review.json	11/15/2018 11:35 AM	JSON File	5,222,145 KB
tip.json	11/15/2018 11:26 AM	JSON File	238,805 KB
user.json	11/15/2018 11:24 AM	JSON File	2,427,488 KB
Yelp_Dataset_Challenge_Round_13	1/14/2019 11:35 AM	Adobe Acrobat Do...	110 KB

TWITTER:

For this, the creation of a twitter development account is necessary which will provide *consumer_key*, *consumer_secret*, *access_key* and *access_secret*, which is necessary to establish a connection with Twitter. These values are unique for every user.

For data extraction, python script is used which is placed under this location: */user/ss13449/project/twitter_scrap.py*
Below is the screenshot for the code:

```
1 #!/usr/bin/env python3
2 #Author: Sonal Sharma
3 import tweepy
4 from tweepy import Stream
5 from tweepy import OAuthHandler
6 from tweepy.streaming import StreamListener
7
8 #consumer API key, consumer API secret key, access token and access token secret
9 consumer_key = "xx"
10 consumer_secret = "xx"
11 access_key = "xx"
12 access_secret = "xx"
13
14 class listener(StreamListener):
15     def on_data(self,tweets):
16         #creating .csv file to store extracted tweets
17         f = open("twitter_starbucks.json" , "a")
18         f.write(tweets + "\n")
19         f.close()
20         return(True)
21
22 #Accessing twitter with tweepy library for python with help of tokens and keys
23
24 auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
25 auth.set_access_token(access_key, access_secret)
26 stream = Stream(auth,listener())
27 stream.filter(track=["Starbucks"])
28 #stream.filter(track=["McDonalds"])
```

Twitter data for McDonalds and Starbucks was extracted using the commands as per the attached python file. (Reference: Twitter Assignment)

<input checked="" type="checkbox"/>	 twitter_mcdonalds.json	11/5/2019 2:24 PM	JSON File
<input checked="" type="checkbox"/>	 twitter_starbucks.json	11/5/2019 3:13 PM	JSON File

Through this script, we were able to scrap the data from twitter for different restaurants(of approx size **1.1 GB**)

REDDIT:

Setting up an account over Reddit is necessary for data scraping and which will provide unique values for the fields like *client_id*, *client_secret*, *user_agent*, *username*, and *password*. These values are required to successfully establish a connection with Reddit. Python script for data scraping from Reddit is placed under:

/user/ss13449/project/reddit_scrap.py

Below is the screenshot for the code:

```

1 #Reddit Data Extraction Code in Python:
2 import praw
3 import datetime as dt
4
5 reddit = praw.Reddit(client_id='xxxxxxxxxxxxxxxxxxxx', \
6                      client_secret='xxxxxxxxxxxxxxxxxxxx', \
7                      user_agent='xxxxxxxxxxxxxxxxxxxx', \
8                      username='xxxxxxxxxxxxxxxxxxxx', \
9                      password='xxxxxxxxxxxxxxxxxxxx')
10
11 subreddit = reddit.subreddit('SearchWords')
12 top_subreddit = subreddit.top() ## Similarly we have .top(), .new(), .hot(), .controversial() and .gilded()
13 top_subreddit = subreddit.top(limit=2000)
14
15 for submission in top_subreddit:
16     date_new = get_date(submission)
17     print("{}, {}, {}".format(submission.id, submission.title, date_new))
18
19 def get_date(submission):
20     time = submission.created
21     return datetime.datetime.fromtimestamp(time)
22
23 ## Ref: http://www.storybench.org/how-to-scrape-reddit-with-python/

```

The extracted two files are:

<input checked="" type="checkbox"/>  reddit_McDonalds	11/5/2019 2:14 PM	Microsoft Excel Co...
<input checked="" type="checkbox"/>  reddit_starbucks	11/5/2019 2:46 PM	Microsoft Excel Co...

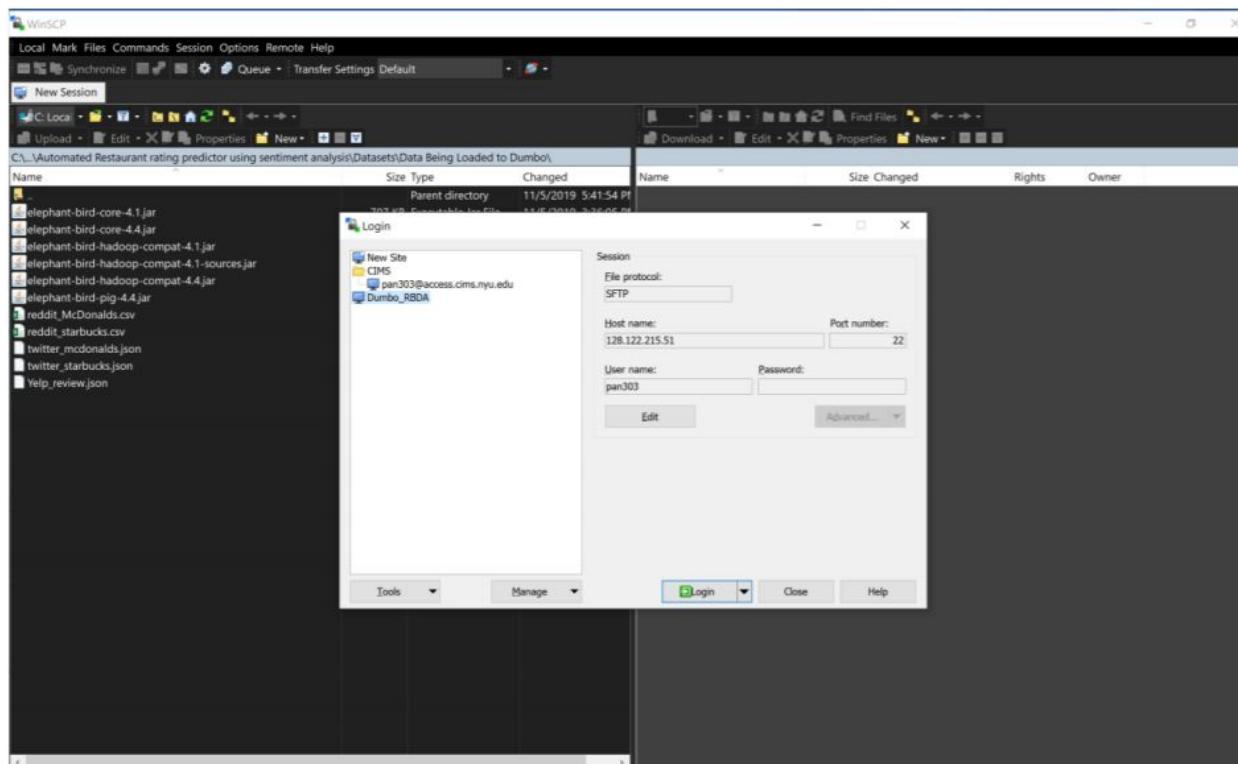
The data size of the data that we extracted from Reddit is around **1GB**.

2. Code or Steps for Data Ingest

WinSCP (Windows Secure Copy), a popular SFTP & FTP Client was used to copy the files from a local computer to the dumbo cluster.

Note: WinSCP can be used smoothly for file transfer if the NYU network is used. In case of non-NYU network connection, it is advisable to use the NYU VPN (i.e. Cisco AnyConnect Secure Mobility Client)

Step 1: Starting the WinSCP & connecting to NYU Dumbo Cluster as a remote host.




```
[kg2644@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/data_yelp/
Found 33 items
-rw-r--r-- 3 ss13449 users 0 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/_SUCCESS
-rw-r--r-- 3 ss13449 users 100175090 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00000
-rw-r--r-- 3 ss13449 users 10213925 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00001
-rw-r--r-- 3 ss13449 users 10256893 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00002
-rw-r--r-- 3 ss13449 users 10835035 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00003
-rw-r--r-- 3 ss13449 users 108467005 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00004
-rw-r--r-- 3 ss13449 users 108105833 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00005
-rw-r--r-- 3 ss13449 users 108251852 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00006
-rw-r--r-- 3 ss13449 users 108851757 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00007
-rw-r--r-- 3 ss13449 users 108305903 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00008
-rw-r--r-- 3 ss13449 users 108313385 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00009
-rw-r--r-- 3 ss13449 users 10836978 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00010
-rw-r--r-- 3 ss13449 users 10901995 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00011
-rw-r--r-- 3 ss13449 users 10813050 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00012
-rw-r--r-- 3 ss13449 users 10816761 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00013
-rw-r--r-- 3 ss13449 users 108238440 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00014
-rw-r--r-- 3 ss13449 users 108910274 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00015
-rw-r--r-- 3 ss13449 users 10805910 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00016
-rw-r--r-- 3 ss13449 users 10804058 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00017
-rw-r--r-- 3 ss13449 users 10816575 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00018
-rw-r--r-- 3 ss13449 users 108761439 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00019
-rw-r--r-- 3 ss13449 users 108161845 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00020
-rw-r--r-- 3 ss13449 users 10822208 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00021
-rw-r--r-- 3 ss13449 users 10829904 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00022
-rw-r--r-- 3 ss13449 users 108892359 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00023
-rw-r--r-- 3 ss13449 users 108349925 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00024
-rw-r--r-- 3 ss13449 users 10839292 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00025
-rw-r--r-- 3 ss13449 users 10853230 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00026
-rw-r--r-- 3 ss13449 users 10909898 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00027
-rw-r--r-- 3 ss13449 users 108157420 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00028
-rw-r--r-- 3 ss13449 users 108148653 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00029
-rw-r--r-- 3 ss13449 users 108309079 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00030
-rw-r--r-- 3 ss13449 users 98514129 2019-11-16 19:52 /user/ss13449/project/cleaned_data/data_yelp/part-00031
[kg2644@login-1-1 ~]$
```

From the above screenshot, we can see that we have successfully executed our file.

```
> cat /user/ss13449/project/mapper_clean_yelpbusiness.py
```

```
[kg2644@login-1-1 ~]$ hdfs dfs -cat /user/ss13449/project/mapper_clean_yelpbusiness.py
#!/usr/bin/env python
# -*- coding: utf-8 -*-
import sys
import json
import re
reload(sys)
sys.setdefaultencoding('utf-8')

for line in sys.stdin:
    try:
        d = json.loads(line)
        Business_id = str(d['business_id'])
        City = str(d['city']).encode('ascii','ignore')
        State = str(d['state'])
        Business_name=str(d['name']).encode('ascii','ignore')
        print '%s,%s,%s,%s' % (Business_id,Business_name,City,State)
    except:
        pass
[kg2644@login-1-1 ~]$
```

Use the below command for code execution:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files hdfs://dumbo/user/ss13449/project/ mapper_clean_yelpbusiness.py"
-mapper "python mapper_clean_yelpbusiness.py" -input
/usr/ss13449/project/data/business.json -output /user/ss13449/project/cleaned_data/data_yelp_business/
```

Output files are created at the below locations:

```
/user/ss13449/project/cleaned_data/data_yelp_business
```

```
[kg2644@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/data_yelp_business
Found 3 items
-rw-r--r-- 3 ss13449 users 0 2019-11-14 17:16 /user/ss13449/project/cleaned_data/data_yelp_business/_SUCCESS
-rw-r--r-- 3 ss13449 users 5388078 2019-11-14 17:16 /user/ss13449/project/cleaned_data/data_yelp_business/part-00000
-rw-r--r-- 3 ss13449 users 5394547 2019-11-14 17:16 /user/ss13449/project/cleaned_data/data_yelp_business/part-00001
[kg2644@login-1-1 ~]$
```

Datasource 2: Twitter (JSON Format)

Below are the cleaning scripts for Twitter.

```
/user/ss13449/project/mapper_clean_twitter_mcdonalds.py
/user/ss13449/project/mapper_clean_twitter_starbucks.py
```



```
[kg2644@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/data_twitter_starbucks/
Found 3 items
-rw-r--r--+ 3 kg2644 users          0 2019-11-17 21:59 /user/ss13449/project/cleaned_data/data_twitter_starbucks/_SUCCESS
-rw-r--r--+ 3 kg2644 users      114227 2019-11-17 21:59 /user/ss13449/project/cleaned_data/data_twitter_starbucks/part-00000
-rw-r--r--+ 3 kg2644 users      112960 2019-11-17 21:59 /user/ss13449/project/cleaned_data/data_twitter_starbucks/part-00001
[kg2644@login-1-1 ~]$
```

Datasource 3: Reddit (csv format)

Below are the cleaning scripts for Reddit:

```
/user/ss13449/project/mapper_clean_reddit_mcdonalds.py
/user/ss13449/project/mapper_clean_reddit_starbucks.py
```

cat /user/ss13449/project/mapper_clean_reddit_mcdonalds.py

Below is the code for the above file :

```
[kg2644@login-1-1 ~]$ hdfs dfs -cat /user/ss13449/project/mapper_clean_reddit_mcdonalds.py
#!/usr/bin/python 2.6.6
# -*- coding: utf-8 -*-

import re
import time
import csv
from csv import reader
import sys
from datetime import datetime
reload(sys)
sys.setdefaultencoding('utf-8')

for line in reader(sys.stdin) :
    d=line
    key = d[0]
    try:
        d[2]=d[2].rstrip('\n')
        timestamp=time.strftime('%m-%d-%Y',time.strptime(d[2],'%m/%d/%Y %H:%M'))
    except:
        if len(d[2])!= 0:
            continue
        timestamp=''

z = lambda x: re.compile("\#").sub("", re.compile("RT @[\w_]+:").sub("@", x).strip())
text_removedLinks= re.sub(r'''(?i)\b((?:https://|www\d{0,3}[.]|[a-z0-9.-]+[.][a-z]{2,4})|(?:[^s()>]+|\(\([^\s()>]+\)(\([^\s()>]+\))*)|(?:\(\([^\s()>]+\)(\([^\s()>]+\))*)\b)|[^\s`!()[]{}`.,>?<~^`^`])''',"",text_removedRT)
review=re.sub("[!@`~,\\\"~*$/%#@/\\r\\t\\n\\v\\l\\v\\l\\v\\l]",'', text_removedLinks.encode('ascii','ignore').lower())
print "%s : %s" : Starbucks : : : '%s' % (key,timestamp,review)
[kg2644@login-1-1 ~]$
```

Use the below command for code execution:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files hdfs://dumbo/user/ss13449/project/mapper_clean_reddit_mcdonalds.py -mapper "python mapper_clean_reddit_mcdonalds.py" -input /user/ss13449/project/data/reddit_McDonalds.csv -output /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/
```

Output files are created at the below locations:

```
/user/ss13449/project/cleaned_data/data_reddit_mcdonalds
```

```
[kg2644@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/
Found 3 items
-rw-r--r--+ 3 ss13449 users          0 2019-11-25 23:20 /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/_SUCCESS
-rw-r--r--+ 3 ss13449 users      315646 2019-11-25 23:20 /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/part-00000
-rw-r--r--+ 3 ss13449 users      318815 2019-11-25 23:20 /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/part-00001
[kg2644@login-1-1 ~]$
```

cat /user/ss13449/project/mapper_clean_reddit_starbucks.py

Below is the code for the above file :


```

[pan303@login-2-1:~]$ hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar \
> -D mapreduce.job.reduces=0 \
> -files hdfs://dumbo/user/ss13449/project/mapper_profile.py \
> -mapper "python mapper_profile.py" \
> -input /user/ss13449/project/output reddit_cleaned_data/part-00000 \
> -output /user/ss13449/project/output reddit_cleaned_data_profile
packageJobJar: [/opt/cloudera/parcels/CDH-5.15.2-1.cdh5.15.2.p0.3/jars/hadoop-streaming-2.6.0-cdh5.15.2.jar] /tmp/streamjob6521407989188451773.jar tmpDir=null
19/11/06 21:42:19 INFO mapred.FileInputFormat: Total input paths to process : 1
19/11/06 21:42:40 INFO hdfs.DFSClient: Exception in createBlockOutputStream
java.io.IOException: Bad connect ack with firstBadLink at 10.0.255.234:50010
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.createBlockOutputStream(DFSOutputStream.java:1785)
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.nextBlockOutputStream(DFSOutputStream.java:1683)
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:790)
19/11/06 21:42:40 WARN hdfs.DFSClient: Abandoning BP-1256044058-128.122.215.50-1440607644284blk_1148829710_75201059
19/11/06 21:42:40 WARN hdfs.DFSClient: Excluding datanode DatanodeInfoWithStorage[10.0.235.234:50010,DS-adeed62a-9e7d-4814-b6b5-68aa5683d6c2,DISK]
19/11/06 21:42:40 INFO mapreduce.JobSubmitter: number of splits:2
19/11/06 21:42:41 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1569350662793_36853
19/11/06 21:42:41 INFO impl.YarnClientImpl: Submitted application application_1569350662793_36853
19/11/06 21:42:41 INFO mapreduce.Job: url to track the job: http://babar.es.its.nyu.edu:8088/proxy/application_1569350662793_36853
19/11/06 21:42:41 INFO mapreduce.Job: Running job: job_1569350662793_36853
19/11/06 21:42:47 INFO mapreduce.Job: Job job_1569350662793_36853 running in uber mode : false
19/11/06 21:42:47 INFO mapreduce.Job: map 0% reduce 0%
19/11/06 21:42:52 INFO mapreduce.Job: map 100% reduce 0%
19/11/06 21:42:57 INFO mapreduce.Job: Job job_1569350662793_36853 completed successfully
19/11/06 21:42:57 INFO mapreduce.Job: Counters: 30
File System Counters
FILE: Number of bytes read=0
FILE: Number of bytes written=315172
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=210387
HDFS: Number of bytes written=1096
HDFS: Number of read operations=10
HDFS: Number of large read operations=0
HDFS: Number of write operations=4
Job Counters
Launched map tasks=2
Data-local map tasks=2
Total time spent by all maps in occupied slots (ms)=23392
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=5848
Total vcore-milliseconds taken by all map tasks=5848
Total megabyte-milliseconds taken by all map tasks=23953408
Map-Reduce Framework
Map input records=2490
Map output records=50
Input split bytes=246
Spilled Records=0
Failed Shuffles=0
Merged Map outputs=0
GC time elapsed (ms)=83
CPU time spent (ms)=1770
Physical memory (bytes) snapshot=712323072
Virtual memory (bytes) snapshot=7453806592
Total committed heap usage (bytes)=2379218944

```

```

[pan303@login-2-1 ~]$ hdfs dfs -cat /user/ss13449/project/output_reddit_cleaned_data_profile/part-00000
Total rows of Review Data: 1232
Column Profiles for Twitter Data:
1) ID
   Data Type: String
   Max String Length: 6
   Max String Size(bytes): 46
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: Yes
2) Date
   Data Type: String
   Max String Length: 10
   Max String Size(bytes): 50
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: No
3) Review
   Data Type: String
   Max String Length: 279
   Max String Size(bytes): 319
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: No
[pan303@login-2-1 ~]$
[pan303@login-2-1 ~]$ date
Wed Nov  6 22:35:55 EST 2019
[pan303@login-2-1 ~]$

```

Yelp Data Profiling using MapReduce:

Command:

```

hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar \
-D mapreduce.job.reduces=0 \
-files hdfs://dumbo/user/ss13449/project/mapper_profile.py \
-mapper "python mapper_profile.py" \
-input /user/ss13449/project/cleaned_data/yelp \
-output /user/ss13449/project/yelp_profile

```

```
[pan303@login-2-1 ~]$ hdfs dfs -cat /user/ss13449/project/yelp_out/part-00000
Total rows of Review Data: 157296
Column Profiles for Twitter Data:
1) ID
   Data Type: String
   Max String Length: 22
   Max String Size(bytes): 62
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: Yes
2) Date
   Data Type: String
   Max String Length: 10
   Max String Size(bytes): 50
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: No
3) Review
   Data Type: String
   Max String Length: 5419
   Max String Size(bytes): 6459
   Range: Doesn't Apply
   Contains Null: No
   is_primary_key: No
[pan303@login-2-1 ~]$
[pan303@login-2-1 ~]$
[pan303@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/yelp_out
Found 1 items
-rw-r--r-- 3 ss13449 users      526 2019-11-06 23:37 /user/ss13449/project/yelp_out/part-00000
[pan303@login-2-1 ~]$
[pan303@login-2-1 ~]$ █
```

Data Profiling for Twitter data using MapReduce:

mapper_profile.py has been used to for the Column Profiling of the data.

Command:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar \
-D mapreduce.job.reduces=0 \
-files hdfs://dumbo/user/ss13449/project/mapper_clean_yelp.py \
-mapper "python mapper_clean_py" \
-input /user/ss13449/project/cleaned_data/twitter* \
-output /user/ss13449/project/twitter_profile
```

```

pan303@login-2-1:~$ hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar \
> -D mapreduce.job.reduces=0 \
> -files hdfs://dumbo/user/ss13449/project/mapper_profile.py \
> -mapper "python mapper_profile.py" \
> -input /user/ss13449/project/output_twitter_cleaned_data/part-00000 \
> -output /user/ss13449/project/output_twitter_cleaned_data/part-00000
packageJobJar: [] [/opt/cloudera/parcels/CDH-5.15.2-1.cdh5.15.2.p0.3/jars/hadoop-streaming-2.6.0-cdh5.15.2.jar] /tmp/streamjob4509463206859604506.jar tmpDir=null
19/11/06 22:39:31 INFO mapred.FileInputFormat: Total input paths to process : 1
19/11/06 22:39:32 INFO mapreduce.JobSubmitter: number of splits:2
19/11/06 22:39:32 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1569350662793_36887
19/11/06 22:39:33 INFO impl.YarnClientImpl: Submitted application application_1569350662793_36887
19/11/06 22:39:33 INFO mapreduce.Job: The url to track the job: http://babar.es.its.nyu.edu:8088/proxy/application_1569350662793_36887/
19/11/06 22:39:33 INFO mapreduce.Job: Running job: job_1569350662793_36887
19/11/06 22:39:38 INFO mapreduce.Job: Job job_1569350662793_36887 running in uber mode : false
19/11/06 22:39:38 INFO mapreduce.Job: map 0% reduce 0%
19/11/06 22:39:45 INFO mapreduce.Job: map 100% reduce 0%
19/11/06 22:39:45 INFO mapreduce.Job: Job job_1569350662793_36887 completed successfully
19/11/06 22:39:45 INFO mapreduce.Job: Counters: 30
  File System Counters
    FILE: Number of bytes read=0
    FILE: Number of bytes written=315176
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=197689
    HDFS: Number of bytes written=1996
    HDFS: Number of read operations=10
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=4
  Job Counters
    Launched map tasks=2
    Data-local map tasks=2
    Total time spent by all maps in occupied slots (ms)=25500
    Total time spent by all reduces in occupied slots (ms)=0
    Total time spent by all map tasks (ms)=6375
    Total vcore-milliseconds taken by all map tasks=6375
    Total megabyte-milliseconds taken by all map tasks=26112000
  Map-Reduce Framework
    Map input records=1131
    Map output records=50
    Input split bytes=248
    Spilled Records=0
    Failed Shuffles=0
    Merged Map outputs=0
    GC time elapsed (ms)=82
    CPU time spent (ms)=1640
    Physical memory (bytes) snapshot=710389760
    Virtual memory (bytes) snapshot=7457062912
    Total committed heap usage (bytes)=2402811904
  File Input Format Counters
    Bytes Read=197441
  File Output Format Counters
    Bytes Written=1096
19/11/06 22:39:45 INFO streaming.StreamJob: Output directory: /user/ss13449/project/output_twitter_cleaned_data_profile
[pan303@login-2-1 ~]$ 

19/11/06 22:39:45 INFO streaming.StreamJob: Output directory: /user/ss13449/project/output_twitter_cleaned_data_profile
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/output_twitter_cleaned_data_profile
Found 3 items
-rw-r--r--+ 3 pan303 users          0 2019-11-06 22:39 /user/ss13449/project/output_twitter_cleaned_data_profile/_SUCCESS
-rw-r--r--+ 3 pan303 users      548 2019-11-06 22:39 /user/ss13449/project/output_twitter_cleaned_data_profile/part-00000
-rw-r--r--+ 3 pan303 users      548 2019-11-06 22:39 /user/ss13449/project/output_twitter_cleaned_data_profile/part-00001
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ hdfs dfs -cat /user/ss13449/project/output_twitter_cleaned_data_profile/part-00000
Total rows of Review Data: 566
Column Profiles for Twitter Data:
1) ID
  Data Type: String
  Max String Length: 19
  Max String Size(bytes): 59
  Range: Doesn't Apply
  Contains Null: No
  is_primary_key: Yes
2) Date
  Data Type: String
  Max String Length: 10
  Max String Size(bytes): 50
  Range: Doesn't Apply
  Contains Null: No
  is_primary_key: No
3) Review
  Data Type: String
  Max String Length: 138
  Max String Size(bytes): 178
  Range: Doesn't Apply
  Contains Null: No
  is_primary_key: No
[pan303@login-2-1 ~]$ 
[pan303@login-2-1 ~]$ 

```

4. Data Analysis

Below are some of the questions that we have answered through in Data Analysis part are:

1. Finding the polarity of the reviews for Starbucks and McDonalds (using sentiment analysis).
2. Finding the average polarity of a particular restaurant for any other restaurant (in our case- McDonalds, Starbucks) and comparing them with the actual ratings. This may help in verifying the ratings on social platforms with ratings from reviews.
3. Finding the frequent topics about which reviewers are talking about. This may help in restaurants to analyze their brand image and take appropriate actions and business decisions upon user frequent topics.

4. Try to compare the average polarity of any restaurant over the different platforms like Yelp, Twitter, Reddit and see on which platform they are popular and are most likely liked by people.
5. Compare the user ratings and average polarity ratings for different states and compare the results.
6. Compare the average user ratings and average polarity ratings weekly, monthly and quarterly, to see how the restaurant's performance.

Below are the key challenges that we encountered while dealing with this.

1. We will be using some of the *NLTK libraries like TextBlob, stopwords, ngrams* but these libraries were not present on the node. We were not aware of this and our code was getting failed every time. So, we contact HPC and Wensheng from the HPC team cooperated with us and getting the NLTK libraries installed over the node.
2. Though the libraries were installed, but we could not easily run it over the node. A python script was needed to run it and we created a *python wrapper script* and use it with the mapper code in the pipeline and were successfully able to execute our codes.
3. We also faced the issues with internal(*manage table*) table and we were not able to fetch the data from manage table to text file, whereas it becomes easier for us to fetch the data from *external tables to text files*.
 > But after a lot of research over the internet and consulting with TAs we found out that external tables are a better way to transfer data from hive to hdfs.

So, we have probably used the below technologies for our project:

1. Hive
2. Impala
3. Map-Reducer

Table Creation and Merging through Hive For YELP:

```
create external table Yelp_data(id string, date_stamp string, reviews string, business_id string, stars float) row format
delimited fields terminated by ',' LINES TERMINATED BY '\n' location '/user/ss13449/project/cleaned_data/data_yelp/';
-- Here an external table Yelp_data has been created on HIVE from the cleaned data which is present at the location :
/user/ss13449/project/cleaned_data/data_yelp/
```

```
0: jdbc:hive2://babar.es.its.nyu.edu:10000> describe yelp_data;
INFO : Compiling command(queryId=hive_20191130184848_c32c9f97-4485-4a9d-b2cf-56df9232f8f5): describe yelp_data
INFO : Semantic Analysis Completed
INFO : Building Hive schema [FieldSchemas:[FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer)]]
INFO : FieldSchemas[FieldSchema(name:comment, type:string, comment:from deserializer), FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer)]
INFO : Completed compiling command(queryId=hive_20191130184848_c32c9f97-4485-4a9d-b2cf-56df9232f8f5); Time taken: 0.159 seconds
INFO : Executing command(queryId=hive_20191130184848_c32c9f97-4485-4a9d-b2cf-56df9232f8f5): describe yelp_data
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20191130184848_c32c9f97-4485-4a9d-b2cf-56df9232f8f5); Time taken: 0.009 seconds
INFO : OK
+-----+
| col_name | data_type | comment |
+-----+
| id       | string    |          |
| date_stamp | string   |          |
| reviews   | string    |          |
| business_id | string   |          |
| stars     | float     |          |
+-----+
9 rows selected (0.256 seconds)
```

```
create external table Yelp_business_data(business_id string, business_name string, city string, state string) row format
delimited fields terminated by ',' LINES TERMINATED BY '\n' location
'/user/ss13449/project/cleaned_data/data_yelp_business/';
-- Here an external table Yelp_business_data has been created on HIVE from the cleaned data which is present at the
location : /user/ss13449/project/cleaned_data/data_yelp_business/
```

```
0: jdbc:hive2://babar.es.its.nyu.edu:10000> describe yelp_business_data;
INFO : Compiling command(queryId=hive_20191130184949_9cdab6e1-9908-41ab-bb71-b2d5a80dc411): describe yelp_business_data
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema[FieldSchemas:[FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer)]]
INFO : FieldSchemas[FieldSchema(name:comment, type:string, comment:from deserializer), FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer)]
INFO : Completed compiling command(queryId=hive_20191130184949_9cdab6e1-9908-41ab-bb71-b2d5a80dc411); Time taken: 0.277 seconds
INFO : Executing command(queryId=hive_20191130184949_9cdab6e1-9908-41ab-bb71-b2d5a80dc411): describe yelp_business_data
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId=hive_20191130184949_9cdab6e1-9908-41ab-bb71-b2d5a80dc411); Time taken: 0.008 seconds
INFO : OK
+-----+
| col_name | data_type | comment |
+-----+
| business_id | string    |          |
| business_name | string   |          |
| city       | string    |          |
| state      | string    |          |
+-----+
4 rows selected (0.403 seconds)
0: jdbc:hive2://babar.es.its.nyu.edu:10000>
```

```

create external table Yelp_review(id string, date_stamp date,reviews string,business_id string,business_name string,stars float,city string,state string)
STORED AS TEXTFILE
LOCATION 'hdfs://dumbo/user/ss13449/project/cleaned_data/tabledata/yelp';
-- Here an empty table Yelp_review has been created where the joined data from both the tables Yelp_data and Yelp_business_data will be stored and can also store them in Text File format

INSERT OVERWRITE TABLE Yelp_review select id,
cast(to_date(from_unixtime(unix_timestamp(date_stamp,'dd-MM-yyyy')))) as
date,reviews,yelp_data.business_id,yelp_business_data.business_name,stars,city,state
from yelp_data JOIN yelp_business_data ON yelp_data.business_id=yelp_business_data.business_id;
-- Here we will overwrite the empty table Yelp_review that we have created and will store the required data into it. One thing here is to note that we have made the JOIN condition on business_id, because it is the only unique identifier in both the tables.

```

Taking Data in Text File from Hive:

```

INSERT OVERWRITE DIRECTORY '/user/ss13449/project/cleaned_data/yelp' ROW FORMAT DELIMITED FIELDS TERMINATED BY ':' SELECT * FROM yelp_review;
-- We have finally overwritten our textfiles from the data which is present in the tables yelp_review, also see the table which is created in Hive.

```

```

0: jdbc:hive2://babar.es.its.nyu.edu:10000> describe yelp_review;
INFO : Compiling command(queryId:hive_20191130184444_bed60a5f-18a5-4ed0-863f-a3d7eb52a784): describe yelp_review
INFO : Semantic Analysis Completed
INFO : Returning Hive schema: Schema(fieldsSchemas:[FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer), FieldSchema(name:comment, type:string, comment:from deserializer), properties:null])
INFO : Completed compiling command(queryId:hive_20191130184444_bed60a5f-18a5-4ed0-863f-a3d7eb52a784); Time taken: 0.153 seconds
INFO : Executing command(queryId:hive_20191130184444_bed60a5f-18a5-4ed0-863f-a3d7eb52a784): describe yelp_review
INFO : Starting task [Stage-0:DDL] in serial mode
INFO : Completed executing command(queryId:hive_20191130184444_bed60a5f-18a5-4ed0-863f-a3d7eb52a784); Time taken: 0.008 seconds
INFO : OK
+-----+-----+-----+
| col_name | data_type | comment |
+-----+-----+-----+
| id       | string    |          |
| date_stamp | date     |          |
| reviews   | string    |          |
| business_id | string   |          |
| business_name | string  |          |
| stars     | float    |          |
| city      | string    |          |
| state     | string    |          |
+-----+-----+-----+
8 rows selected (0.309 seconds)
0: jdbc:hive2://babar.es.its.nyu.edu:10000> 

```

Table Creation in Impala for Twitter Data :

Similarly, we have created a table for Twitter just like we have created for Yelp, but this time we have created the tables in IMPALA.

```

create external table data_twitter_starbucks(id string, date_stamp string, business_id string, business_name string, user_rating string, city string, state string, polarity float, subjectivity float) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks';

```

```
[compute-1-1:21000] > describe data_twitter_starbucks;
Query: describe data_twitter_starbucks
+-----+-----+-----+
| name | type | comment |
+-----+-----+-----+
| id | string | |
| date_stamp | string | |
| business_id | string | |
| business_name | string | |
| user_rating | string | |
| city | string | |
| state | string | |
| polarity | float | |
| subjectivity | float | |
+-----+-----+-----+
Fetched 9 row(s) in 0.01s
[compute-1-1:21000] >
```

create external table data_twitter_mcdonalds(id string, date_stamp string, business_id string, business_name string, user_rating string, city string, state string, polarity float, subjectivity float) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds';

```
[compute-1-1:21000] > describe data_twitter_mcdonalds;
Query: describe data_twitter_mcdonalds
+-----+-----+-----+
| name | type | comment |
+-----+-----+-----+
| id | string | |
| date_stamp | string | |
| business_id | string | |
| business_name | string | |
| user_rating | string | |
| city | string | |
| state | string | |
| polarity | float | |
| subjectivity | float | |
+-----+-----+-----+
Fetched 9 row(s) in 0.01s
[compute-1-1:21000] >
```

Table Creation in Impala for Reddit Data :

Similarly, we have created a table for Reddit just like we have created for Yelp, but this time we have created the tables in IMPALA.

create external table data_reddit_starbucks(id string, date_stamp string, business_id string, business_name string, user_rating string, city string, state string, polarity float, subjectivity float) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks';

```
[compute-1-1:21000] > describe data_reddit_starbucks;
Query: describe data_reddit_starbucks
+-----+-----+-----+
| name | type | comment |
+-----+-----+-----+
| id | string | |
| date_stamp | string | |
| business_id | string | |
| business_name | string | |
| user_rating | string | |
| city | string | |
| state | string | |
| polarity | float | |
| subjectivity | float | |
+-----+-----+-----+
Fetched 9 row(s) in 0.01s
[compute-1-1:21000] >
```

create external table data_reddit_mcdonalds(id string, date_stamp string, business_id string, business_name string, user_rating string, city string, state string, polarity float, subjectivity float) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds';

```
[compute-1-1:21000] > describe data_reddit_mcdonalds;
Query: describe data_reddit_mcdonalds
+-----+-----+
| name | type | comment |
+-----+-----+
| id | string |          |
| date_stamp | string |          |
| business_id | string |          |
| business_name | string |          |
| user_rating | string |          |
| city | string |          |
| state | string |          |
| polarity | float |          |
| subjectivity | float |          |
+-----+-----+
Fetched 9 row(s) in 0.01s
[compute-1-1:21000] >
```

Now we will create an empty table data_for_analyses and will put the Reddit and Twitter data into this table. Below are the commands that we have used for this:

Insert into Table data_for_analyses select id, date_stamp, business_id, business_name, cast(user_rating as float), city, state, polarity, subjectivity from data_reddit_starbucks;

Insert into Table data_for_analyses select id, date_stamp, business_id, business_name, cast(user_rating as float), city, state, polarity, subjectivity from data_reddit_mcDonalds;

Insert into Table data_for_analyses select id, date_stamp, business_id, business_name, cast(user_rating as float), city, state, polarity, subjectivity from data_twitter_starbucks;

Insert into Table data_for_analyses select id, date_stamp, business_id, business_name, cast(user_rating as float), city, state, polarity, subjectivity from data_twitter_mcDonalds;

Till now we have created the tables and have placed all the cleaned data into it.

Now, let's run our analytics code, which will basically give us the polarity for each review. To run each analytics code using Hadoop pipeline, we will use the python Wrapper script which is present at the below location:

/user/ss13449/project/python_wrapper.sh

```
[kg2644@login-2-1 ~]$ hdfs dfs -cat /user/ss13449/project/python_wrapper.sh
#!/bin/bash
export NLTK_DATA=/share/apps/python/3.6.5/data/nltk
source /etc/profile.d/modules.sh
module purge
module load python/gnu/3.6.5
python $*
[kg2644@login-2-1 ~]$
```

Analytics Run for YELP :

hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics.py" -mapper "python_wrapper.sh analytics.py" -input /user/ss13449/project/cleaned_data/yelp -output /user/ss13449/project/cleaned_data/analysed_data

Analytics code:->

```
[kg2644@login-2-1 ~]$ hdfs dfs -cat /user/ss13449/project/analytics.py
#!/usr/bin/env python 3.6.5

import sys
import textblob
from textblob import TextBlob

try:
    for line in sys.stdin:
        lines = line.strip().split(':')
        blob = TextBlob(lines[2])
        print(lines[0],':',lines[1],':',lines[3],':',lines[4],':',lines[5],':',lines[6],':',lines[7],':',blob.sentiment[0],':',blob.sentiment[1])
except:
    pass
[kg2644@login-2-1 ~]$
```

Below is the screenshot, after successfully executing the above command:-->

```
[kg2644@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/analysed_data
Found 48 items
-rw-r--r-- 3 impala users 57968 2019-11-26 05:10 /user/ss13449/project/cleaned_data/analysed_data/5b4d3e88db967f8-a430dca00000000 1361128615_data.0.
-rw-r--r-- 3 impala users 96001 2019-11-26 05:10 /user/ss13449/project/cleaned_data/analysed_data/5b4d3e88db967f8-a430dca00000001_1970647877_data.0.
-rw-r--r-- 3 impala users 175469 2019-11-25 23:34 /user/ss13449/project/cleaned_data/analysed_data/67462653db2c7fb-7121313600000000 834729423_data.0.
-rw-r--r-- 3 impala users 174211 2019-11-25 23:34 /user/ss13449/project/cleaned_data/analysed_data/67462653db2c7fb-7121313600000001_2090129110_data.0.
-rw-r--r-- 3 impala users 180071 2019-11-25 05:05 /user/ss13449/project/cleaned_data/analysed_data/704a60f6d3e7340a-87e8c32800000000 1640151534_data.0.
-rw-r--r-- 3 impala users 186847 2019-11-25 05:05 /user/ss13449/project/cleaned_data/analysed_data/704a60f6d3e7340a-87e8c32800000001_1552240890_data.0.
-rw-r--r-- 3 impala users 175469 2019-11-25 05:08 /user/ss13449/project/cleaned_data/analysed_data/704a60f6d3e7340a-87e8c32800000000 1640151534_data.0.
-rw-r--r-- 3 impala users 174211 2019-11-25 05:08 /user/ss13449/project/cleaned_data/analysed_data/704a60f6d3e7340a-87e8c32800000001_1552240890_data.0.
-rw-r--r-- 3 ss13449 users 0 2019-11-17 05:05 /user/ss13449/project/cleaned_data/analysed_data/_SUCCESS
drwxrwxr-x+ 3 impala users 0 2019-11-25 05:10 /user/ss13449/project/cleaned_data/analysed_data/_impala_insert_staging
drwxrwxr-x+ 3 ss13449 users 0 2019-11-25 23:21 /user/ss13449/project/cleaned_data/analysed_data/reddit_mcdonalds
drwxrwxr-x+ 3 ss13449 users 0 2019-11-25 22:30 /user/ss13449/project/cleaned_data/analysed_data/reddit_starbucks
drwxrwxr-x+ 3 kg2644 users 0 2019-11-17 22:04 /user/ss13449/project/cleaned_data/analysed_data/twitter_mcdonalds
drwxrwxr-x+ 3 kg2644 users 0 2019-11-17 22:01 /user/ss13449/project/cleaned_data/analysed_data/twitter_starbucks
-rw-r--r-- 3 impala users 64620 2019-11-25 05:09 /user/ss13449/project/cleaned_data/analysed_data/704a60f6d3e7340a-87e8c32800000000 863814299_data.0.
-rw-r--r-- 3 impala users 63939 2019-11-25 05:09 /user/ss13449/project/cleaned_data/analysed_data/d3489dd702e242c9-60135deb00000000 2021977665_data.0.
-rw-r--r-- 3 ss13449 users 28816674 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00000
-rw-r--r-- 3 ss13449 users 28882116 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00001
-rw-r--r-- 3 ss13449 users 28644505 2019-11-17 00:50 /user/ss13449/project/cleaned_data/analysed_data/part-00002
-rw-r--r-- 3 ss13449 users 28376240 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00003
-rw-r--r-- 3 ss13449 users 28498548 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00004
-rw-r--r-- 3 ss13449 users 28287635 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00005
-rw-r--r-- 3 ss13449 users 27910164 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00006
-rw-r--r-- 3 ss13449 users 28569315 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00007
-rw-r--r-- 3 ss13449 users 28200017 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00008
-rw-r--r-- 3 ss13449 users 28025698 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00009
-rw-r--r-- 3 ss13449 users 28591428 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00010
-rw-r--r-- 3 ss13449 users 28338877 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00011
-rw-r--r-- 3 ss13449 users 27921502 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00012
-rw-r--r-- 3 ss13449 users 28674013 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00013
-rw-r--r-- 3 ss13449 users 28630894 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00014
-rw-r--r-- 3 ss13449 users 284857630 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00015
-rw-r--r-- 3 ss13449 users 27774275 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00016
-rw-r--r-- 3 ss13449 users 27608637 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00017
-rw-r--r-- 3 ss13449 users 28649691 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00018
-rw-r--r-- 3 ss13449 users 28640113 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00019
-rw-r--r-- 3 ss13449 users 23552580 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00020
-rw-r--r-- 3 ss13449 users 24214316 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00021
-rw-r--r-- 3 ss13449 users 244001106 2019-11-17 00:49 /user/ss13449/project/cleaned_data/analysed_data/part-00022
-rw-r--r-- 3 ss13449 users 19861963 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00023
-rw-r--r-- 3 ss13449 users 17110163 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00024
-rw-r--r-- 3 ss13449 users 15648601 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00025
-rw-r--r-- 3 ss13449 users 15701248 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00026
-rw-r--r-- 3 ss13449 users 15571149 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00027
-rw-r--r-- 3 ss13449 users 15625769 2019-11-17 00:48 /user/ss13449/project/cleaned_data/analysed_data/part-00028
-rw-r--r-- 3 ss13449 users 11432389 2019-11-17 00:47 /user/ss13449/project/cleaned_data/analysed_data/part-00029
-rw-r--r-- 3 ss13449 users 11295448 2019-11-17 00:47 /user/ss13449/project/cleaned_data/analysed_data/part-00030
-rw-r--r-- 3 ss13449 users 7183301 2019-11-17 00:47 /user/ss13449/project/cleaned_data/analysed_data/part-00031
[kg2644@login-2-1 ~]$
```

We have to make the `:` separated files, which will make it easier for us to make the processing going forward easier.

Analytics Run for Twitter:

Our analytics code will remain the same for all the analysis that we will be doing, only the input path gets changed in the command line.

Starbucks:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics.py" -mapper "python_wrapper.sh" analytics.py -input /user/ss13449/project/cleaned_data/data_twitter_starbucks -output /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks/
```

Below is the screenshot, after successfully executing the above command:

```
[kg2644@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks/
Found 3 items
-rw-r--r--+ 3 kg2644 users          0 2019-11-17 22:01 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks/_SUCCESS
-rw-r--r--+ 3 kg2644 users       65956 2019-11-17 22:01 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks/part-00000
-rw-r--r--+ 3 kg2644 users       66393 2019-11-17 22:01 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_starbucks/part-00001
[kg2644@login-2-1 ~]$
```

McDonald's:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics.py" -mapper "python_wrapper.sh" analytics.py" -input /user/ss13449/project/cleaned_data/data_twitter_mcdonalds -output /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds/
```

Below is the screenshot, after successfully executing the above command:

```
[kg2644@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds/
Found 3 items
-rw-r--r--+ 3 kg2644 users          0 2019-11-17 22:04 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds/_SUCCESS
-rw-r--r--+ 3 kg2644 users       105699 2019-11-17 22:04 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds/part-00000
-rw-r--r--+ 3 kg2644 users       108485 2019-11-17 22:04 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_twitter_mcdonalds/part-00001
[kg2644@login-2-1 ~]$
```

Analytics Run for Reddit:

Our analytics code will remain the same for all the analysis that we will be doing, only the input path gets changed in the command line.

Starbucks:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics.py" -mapper "python_wrapper.sh" analytics.py" -input /user/ss13449/project/cleaned_data/data_reddit_starbucks/ -output /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks/
```

Below is the screenshot, after successfully executing the above command:

```
[kg2644@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks/
Found 3 items
-rw-r--r--+ 3 ss13449 users          0 2019-11-25 22:30 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks/_SUCCESS
-rw-r--r--+ 3 ss13449 users      181946 2019-11-25 22:30 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks/part-00000
-rw-r--r--+ 3 ss13449 users      182671 2019-11-25 22:30 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_starbucks/part-00001
[kg2644@login-2-1 ~]$
```

McDonalds:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics.py" -mapper "python_wrapper.sh" analytics.py" -input /user/ss13449/project/cleaned_data/data_reddit_mcdonalds/ -output /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds/
```

Below is the screenshot, after successfully executing the above command:

```
[kg2644@login-2-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds/
Found 3 items
-rw-r--r--+ 3 ss13449 users          0 2019-11-25 23:21 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds/_SUCCESS
-rw-r--r--+ 3 ss13449 users      192705 2019-11-25 23:21 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds/part-00000
-rw-r--r--+ 3 ss13449 users      185461 2019-11-25 23:21 /user/ss13449/project/cleaned_data/analysed_data/analysed_data_reddit_mcdonalds/part-00001
[kg2644@login-2-1 ~]$
```

Analyses 1 & 2: Our first and second analyses is to find the average rating(that we will get through polarity) for Starbucks and McDonald's for all the data that we have gathered and compare them with the actual average rating.

Solution: We will be using Impala for faster processing of the data. Below are the commands that we will be using:

```
Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null;
```

Here we have used math to compute the rating. Polarity results after sentiment analysis comes in between [-1,+1] and so to push it in between [0,5], first multiply by 5, then add 5 and divide the ratings by 2. This will give us the rating in the range [0,5].

```
[compute-1-1:21000] > Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null;
Query: Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null
Query submitted at: 2019-11-30 18:16:48 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=b04bd27f2fdea224:19c51fa00000000
+-----+
| average_user_rating | average_polarity |
+-----+
| 3.084098344998759 | 2.827585428025279 |
+-----+
Fetched 1 row(s) in 0.73s
```

```
Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null;
```

```
[compute-1-1:21000] > Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null;
Query: Select avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null
Query submitted at: 2019-11-30 18:16:53 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=2c45d2f63922a90a:3d71f6ec00000000
+-----+
| average_user_rating | average_polarity |
+-----+
| 2.125 | 2.510612693991349 |
+-----+
Fetched 1 row(s) in 0.31s
```

From the above, we can see that for Starbucks the average user rating comes out to be 3.0 whereas the average polarity rating comes out to be 2.82, and for McDonald's average user rating is 2.1 whereas average polarity comes out to be 2.5. We can see that the average polarity results are close to the average user ratings, keeping the error margin of +/- 0.5.

Analyses 3: This involves making a word count program, to see which are those particular items that are most and least liked by the people (as in their reviews). This will basically help the restaurants in letting them know which are the things they are good at and in which they are not?

Solution: We have prepared the below code for the word count analysis. We have also made use of stopwords, as we don't want to count the common words.

STARBUCKS :

The code is present at location:

```
/user/ss13449/project/analytics_starbucks.py
```

```
[ss13449@login-1-1 ~]$ hdfs dfs -cat /user/ss13449/project/analytics_starbucks.py
#!/usr/bin/env python 3.6.5

import sys, re
from nltk.corpus import stopwords
from nltk.util import ngrams
businessName = 'starbucks'
twograms = dict()
stopword = set(stopwords.words('english'))
for line in sys.stdin:
    d = line.strip(" ").split(":")
    if businessName in d[3].lower() or businessName in d[4].lower():
        s = d[2].lower().replace("''", '')
        tokens = [token for token in s.split(" ") if token != ""]
        val = list(ngrams(tokens, 2))
        for i in val:
            x,y = i
            if x not in stopword and y not in stopword:
                tup = x.replace('\\\'','') + " " + y.replace('\\\'','')
                if tup in twograms:
                    twograms[tup] += 1
                else:
                    twograms[tup] = 1
print("\n".join("{!r}:{!r}".format(k, v) for k, v in twograms.items()))
#print(sorted(twograms.items(), key=lambda x: x[1], reverse=True))
```

Command used to calculate the word count program for Starbucks data is:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics_starbucks.py" -mapper "python_wrapper.sh analytics_starbucks.py" -input /user/ss13449/project/cleaned_data/all_data/ -output /user/ss13449/project/cleaned_data/word_count_starbucks
```

Below is the result for the same:

```
13/12/01 18:44:46 INFO Streaming.StreamJob: Output directory: /user/ss13449/project/cleaned_data/word_count_starbucks
[kg2644@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/word_count_starbucks
Found 41 items
-rw-r--r--+ 3 kg2644 users          0 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/_SUCCESS
-rw-r--r--+ 3 kg2644 users  251425 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00000
-rw-r--r--+ 3 kg2644 users  254969 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00001
-rw-r--r--+ 3 kg2644 users  198073 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00002
-rw-r--r--+ 3 kg2644 users  214348 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00003
-rw-r--r--+ 3 kg2644 users  190359 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00004
-rw-r--r--+ 3 kg2644 users  215360 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00005
-rw-r--r--+ 3 kg2644 users  196796 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00006
-rw-r--r--+ 3 kg2644 users  248437 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00007
-rw-r--r--+ 3 kg2644 users  225696 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00008
-rw-r--r--+ 3 kg2644 users  217065 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00009
-rw-r--r--+ 3 kg2644 users  210765 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00010
-rw-r--r--+ 3 kg2644 users  216429 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00011
-rw-r--r--+ 3 kg2644 users  194126 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00012
-rw-r--r--+ 3 kg2644 users  235470 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00013
-rw-r--r--+ 3 kg2644 users  213260 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00014
-rw-r--r--+ 3 kg2644 users  214739 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00015
-rw-r--r--+ 3 kg2644 users  222239 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00016
-rw-r--r--+ 3 kg2644 users  206330 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00017
-rw-r--r--+ 3 kg2644 users  210840 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00018
-rw-r--r--+ 3 kg2644 users  217729 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00019
-rw-r--r--+ 3 kg2644 users  192200 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00020
-rw-r--r--+ 3 kg2644 users  225289 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00021
-rw-r--r--+ 3 kg2644 users  216775 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00022
-rw-r--r--+ 3 kg2644 users  141690 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00023
-rw-r--r--+ 3 kg2644 users  117826 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00024
-rw-r--r--+ 3 kg2644 users  144724 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00025
-rw-r--r--+ 3 kg2644 users  105121 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00026
-rw-r--r--+ 3 kg2644 users  133878 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00027
-rw-r--r--+ 3 kg2644 users  141548 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00028
-rw-r--r--+ 3 kg2644 users  76930 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00029
-rw-r--r--+ 3 kg2644 users  105135 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00030
-rw-r--r--+ 3 kg2644 users  47248 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00031
-rw-r--r--+ 3 kg2644 users  88859 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00032
-rw-r--r--+ 3 kg2644 users  88561 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00033
-rw-r--r--+ 3 kg2644 users  84317 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00034
-rw-r--r--+ 3 kg2644 users  108932 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00035
-rw-r--r--+ 3 kg2644 users   4 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00036
-rw-r--r--+ 3 kg2644 users   4 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00037
-rw-r--r--+ 3 kg2644 users  65943 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00038
-rw-r--r--+ 3 kg2644 users  64256 2019-12-01 18:44 /user/ss13449/project/cleaned_data/word_count_starbucks/part-00039
[kg2644@login-1-1 ~]$
```

McDonalds:-

Code is present at location:

/user/ss13449/project/analytics_mcdonalds.py and below is the code for the same.

```
[ss13449@login-1-1 ~]$ hdfs dfs -cat /user/ss13449/project/analytics_mcdonalds.py
#!/usr/bin/env python 3.6.5

import sys, re
from nltk.corpus import stopwords
from nltk.util import ngrams
businessName = 'mcdonalds'
twoGrams = dict()
stopword = set(stopwords.words('english'))
for line in sys.stdin:
    d = line.strip().split(':')
    if businessName in d[3].lower() or businessName in d[4].lower():
        s = d[2].lower().replace("'", '')
        tokens = [token for token in s.split(" ") if token != ""]
        val = list(ngrams(tokens, 2))
        for i in val:
            x,y = i
            if x not in stopword and y not in stopword:
                tup = x.replace('\\','') + " " + y.replace('\\','')
                if tup in twoGrams:
                    twoGrams[tup] += 1
                else:
                    twoGrams[tup] = 1
print("\n".join("{}:{}\n".format(k, v) for k, v in twoGrams.items()))
#print(sorted(twoGrams.items(), key=lambda x: x[1], reverse=True))
```

Command used to calculate the word count program for McDonald's data is:

```
hadoop jar /opt/cloudera/parcels/CDH/lib/hadoop-mapreduce/hadoop-streaming.jar -D mapreduce.job.reduces=0 -files "python_wrapper.sh,analytics_mcdonalds.py" -mapper "python_wrapper.sh analytics_mcdonalds.py" -input /user/ss13449/project/cleaned_data/all_data/ -output /user/ss13449/project/cleaned_data/word_count_mcdonalds
```

Below is the result for the same:-->

```
[kg26440@login-1-1 ~]$ hdfs dfs -ls /user/ss13449/project/cleaned_data/word_count_mcdonalds
Found 41 items
-rw-r--r--+ 3 kg2644 users          0 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/_SUCCESS
-rw-r--r--+ 3 kg2644 users      10323 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00000
-rw-r--r--+ 3 kg2644 users       2351 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00001
-rw-r--r--+ 3 kg2644 users      7547 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00002
-rw-r--r--+ 3 kg2644 users      3144 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00003
-rw-r--r--+ 3 kg2644 users      5403 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00004
-rw-r--r--+ 3 kg2644 users      2297 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00005
-rw-r--r--+ 3 kg2644 users      2862 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00006
-rw-r--r--+ 3 kg2644 users     10123 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00007
-rw-r--r--+ 3 kg2644 users      7139 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00008
-rw-r--r--+ 3 kg2644 users      4740 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00009
-rw-r--r--+ 3 kg2644 users      6771 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00010
-rw-r--r--+ 3 kg2644 users      3679 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00011
-rw-r--r--+ 3 kg2644 users      2685 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00012
-rw-r--r--+ 3 kg2644 users      3403 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00013
-rw-r--r--+ 3 kg2644 users      3617 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00014
-rw-r--r--+ 3 kg2644 users      2937 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00015
-rw-r--r--+ 3 kg2644 users      1351 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00016
-rw-r--r--+ 3 kg2644 users      3068 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00017
-rw-r--r--+ 3 kg2644 users      3031 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00018
-rw-r--r--+ 3 kg2644 users      5227 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00019
-rw-r--r--+ 3 kg2644 users      2163 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00020
-rw-r--r--+ 3 kg2644 users      1308 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00021
-rw-r--r--+ 3 kg2644 users      356 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00022
-rw-r--r--+ 3 kg2644 users      1577 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00023
-rw-r--r--+ 3 kg2644 users      1858 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00024
-rw-r--r--+ 3 kg2644 users      895 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00025
-rw-r--r--+ 3 kg2644 users      2149 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00026
-rw-r--r--+ 3 kg2644 users      922 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00027
-rw-r--r--+ 3 kg2644 users      3677 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00028
-rw-r--r--+ 3 kg2644 users      2978 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00029
-rw-r--r--+ 3 kg2644 users      637 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00030
-rw-r--r--+ 3 kg2644 users      534 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00031
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00032
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00033
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00034
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00035
-rw-r--r--+ 3 kg2644 users     47229 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00036
-rw-r--r--+ 3 kg2644 users     46569 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00037
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00038
-rw-r--r--+ 3 kg2644 users      4 2019-12-01 18:46 /user/ss13449/project/cleaned_data/word_count_mcdonalds/part-00039
```

Analyses of Top topics in Impala:

Data is imported into the impala for sorting. The following command is used:

```
create external table fc_starbucks(feature string, count int) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/word_count_starbucks/';
```

```
create external table fc_mcdonalds(feature string, count int) row format delimited fields terminated by ':' location '/user/ss13449/project/cleaned_data/word_count_mcdonalds/';
```

To combine all topics and finding all top reviews:

```
select trim(feature) as features,sum(count) as frequency from fc_starbucks where count > 1 group by features order by sum(count) desc limit 10;
```

```
select trim(feature) as features,sum(count) as frequency from fc_mcdonalds where count > 1 group by features order by sum(count) desc limit 10;
```

Following results are obtained:

Starbucks:

```
[compute-1-1:21000] > select trim(feature) as features,sum(count) as frequency from fc_starbucks where count > 1 group by features order by sum(count) desc limit 10;
Query: select trim(feature) as features,sum(count) as frequency from fc_starbucks where count > 1 group by features order by sum(count) desc limit 10
Query submitted at: 2019-12-01 22:13:12 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=8140634f71f895d7:2e758a8b00000000
+-----+
| features | frequency |
+-----+
| 'drive thru' | 2948 |
| 'customer service' | 1868 |
| 'every time' | 882 |
| 'starbucks location' | 728 |
| 'parking lot' | 619 |
| 'iced coffee' | 568 |
| 'green tea' | 537 |
| 'worst starbucks' | 478 |
| 'ive ever' | 473 |
| 'dont know' | 447 |
+-----+
Fetched 10 row(s) in 0.42s
[compute-1-1:21000] >
```

McDonalds:

```
[compute-1-1:21000] > select trim(feature) as features,sum(count) as frequency from fc_mcdonalds where count > 1 group by features order by sum(count) desc limit 20;
Query: select trim(feature) as features,sum(count) as frequency from fc_mcdonalds where count > 1 group by features order by sum(count) desc limit 20
Query submitted at: 2019-12-01 22:12:29 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=6348e5f0c0ea0243:4c3a395700000000
+-----+
| features | frequency |
+-----+
| 'ice cream' | 74 |
| 'delicious ice' | 46 |
| 'drive thru' | 42 |
| 'mcdonalds ceo' | 38 |
| 'outta hand' | 25 |
| 'advertising getting' | 25 |
| 'mcdonalds advertising' | 25 |
| 'getting outta' | 25 |
| 'lol thats' | 24 |
| 'thats accurate' | 24 |
| 'cream machine' | 19 |
| 'mcdonalds sprite' | 19 |
| 'mcdonalds money' | 18 |
| 'become grandparents' | 18 |
| 'sudden got' | 18 |
| 'parents become' | 18 |
| 'ceo fired' | 18 |
| 'told zeka' | 17 |
| 'itz_nizdee told' | 17 |
| 'got food' | 17 |
+-----+
Fetched 20 row(s) in 0.42s
```

Analyses 4: Try to compare the average polarity of different restaurants over different platforms like Yelp, Reddit and Twitter and see on which platforms they are more popular/or are more likely to be liked by the people.

Solution:

For Starbucks: We will be running all our queries on IMPALA.

1. Over Reddit, average polarity comes out to be -->

```
[compute-1-1:21000] > select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_reddit_starbucks;
Query: select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_reddit_starbucks
Query submitted at: 2019-11-30 21:03:14 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=7d40a77a9317967d:f9ae772b00000000
+-----+
| count | average_polarity |
+-----+
| 4988 | 2.618489110807198 |
+-----+
Fetched 1 row(s) in 0.11s
[compute-1-1:21000] >
```

2. Over Twitter, average polarity comes out to be:

```
[compute-1-1:21000] > select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_twitter_starbucks;
Query: select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_twitter_starbucks
Query submitted at: 2019-11-30 21:03:53 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=df4e9584f226ff99:a94956ef00000000
+-----+
| count | average_polarity |
+-----+
| 1531 | 2.667556650827716 |
+-----+
Fetched 1 row(s) in 0.11s
[compute-1-1:21000] >
```

3. Over Yelp, average polarity comes out to be:

```
[compute-1-1:21000] > select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_yelp where business_name like '%Starbucks%';
Query: select count(polarity) as Count ,avg((polarity*5)+5)/2 as average_polarity from data_yelp where business_name like '%Starbucks%'
Query submitted at: 2019-11-30 21:05:04 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=3b4edbeb59f5daa2:9e2e73a400000000
+-----+
| count | average_polarity |
+-----+
| 20145 | 2.827585428025279 |
+-----+
Fetched 1 row(s) in 0.31s
[compute-1-1:21000] >
```

For McDonalds: We will be running all our queries on IMPALA.

1. Over Reddit, average polarity comes out to be:

```
[compute-1-1:21000] > select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_reddit_mcdonalds;
Query: select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_reddit_mcdonalds
Query submitted at: 2019-11-30 20:26:20 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=fc498e5ac30fab3a:9258616100000000
+-----+
| count | average_polarity |
+-----+
| 4973 | 2.676237589526467 |
+-----+
Fetched 1 row(s) in 0.22s
```

2. Over Twitter, average polarity comes out to be:

```
Received 1 row(s) in 0.11s
[compute-1-1:21000] > select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_twitter_mcdonalds;
Query: select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_twitter_mcdonalds
Query submitted at: 2019-11-30 20:26:28 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=ed42f7eb4884d126:6e66901b00000000
+-----+
| count | average_polarity |
+-----+
| 2588 | 2.582368420240764 |
+-----+
Fetched 1 row(s) in 0.11s
```

3. Over Yelp, average polarity comes out to be:

```
Received 1 row(s) in 0.11s
[compute-1-1:21000] > select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_yelp where business_name like '%McDonalds%';
Query: select count(polarity) as Count,avg((polarity*5)+5)/2 as average_polarity from data_yelp where business_name like '%McDonalds%'
Query submitted at: 2019-11-30 20:28:00 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=7043a5ed9f0031ad:32918a4e00000000
+-----+
| count | average_polarity |
+-----+
| 232 | 2.510612693991349 |
+-----+
Fetched 1 row(s) in 0.41s
[compute-1-1:21000] >
```

So, we have been able to compute the average polarity of different datasets.

1. If we talk about '**Starbucks**' dataset over 3 platforms, we can see that seeing the average polarity, Yelp has the highest average polarity. It also gives us an indication that people are more likely to post a review regarding '**Starbucks**' on **Yelp**.
2. If we talk about '**McDonalds**' dataset over 3 platforms, we can see that seeing the average polarity, Reddit has the highest average polarity. It also gives us an indication that people are more likely to post a review regarding '**McDonalds**' on **Reddit**.

NOTE: These results are based on the data which are given, the subject may likely change as data size increases.

Analyses 5: Compare the user ratings and average polarity ratings for different states and compare the results.

Solution: We have compared the average user ratings and average polarity ratings across different states for all the datasets. Let's see one by one

1. STARBUCKS

[compute-1-1:21000] > Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by state;		
Query: Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by state		
Query submitted at: 2019-12-01 16:34:29 (Coordinator: http://compute-1-1.local:25000)		
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=2e48f9b992eaa006:a419a67200000000		
state	average_user_rating	average_polarity
WI	3.129139072847682	2.845141008897388
ON	3.39685903191851	2.933296582941169
NV	2.980223848327137	2.768533948253982
SC	3.22680412371134	2.807625758284031
QC	3.640350877192982	2.943418216886481
NC	3.071868583162217	2.817971324967752
OH	3.42914979757085	2.877979269768498
IL	3	2.902597257988436
AZ	3.039394384830222	2.822624431696321
PA	3.449333333333333	2.938863927628388
AB	3.231770833333333	2.90656281360062

Fetched 11 row(s) in 1.40s

The query used is:

```
Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by state;
```

From the above screenshot, we can see that in the state 'QC' average user rating is 3.64 and the average polarity rating is 2.94, which is the highest among all. We can conclude that the STARBUCKS being the most popular brand in state 'QC' among all other states.

[compute-1-1:21000] > Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by state;		
Query: Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by state		
Query submitted at: 2019-12-01 16:34:44 (Coordinator: http://compute-1-1.local:25000)		
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=c949096e4f116ae3:b7b0a40000000000		
state	average_user_rating	average_polarity
ON	2.357142857142857	2.545783975995922
NV	2.382352941176471	2.748391650384292
NC	1.777777777777778	2.595914246824881
OH	1	1.883822185918689
IL	3.181818181818182	2.743339343896051
AZ	1.522388059701492	2.319336349915131
AB	1.833333333333333	2.376954749537011

Fetched 7 row(s) in 1.55s

The query used is:

```
Select state, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by state;
```

From the above screenshot, we can see that in state 'IL' average user rating is 3.18 and the average polarity rating is 2.74, which is the highest among all for McDonalds. We can conclude that McDonalds is the most popular brand in state 'IL' among all states.

Analyses 6: Compare the average user ratings and average polarity ratings yearly, to see how the restaurant's performance.

Solution:

Starbucks :

Below is the query and results for average user ratings and average polarity ratings for 'Starbucks':

```
Select year(to_timestamp(trim(date_stamp), 'yyyy-MM-dd')) as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by period order by period;
```

```
[compute-1-1:21000] > Select year(to_timestamp(trim(date_stamp), 'yyyy-MM-dd')) as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by period order by period;
Query: Select year(to_timestamp(trim(date_stamp), 'yyyy-MM-dd')) as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%Starbucks%" and user_rating is not null group by period order by period
Query submitted at: 2019-12-01 18:12:22 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=e04d6494caec02c1:64169dc00000000
+-----+-----+-----+
| period | average_user_rating | average_polarity |
+-----+-----+-----+
| 2006 | 2.25 | 2.036258010193706 |
| 2007 | 3.692307692307693 | 2.655653922746961 |
| 2008 | 3.615384615384615 | 2.931280138884456 |
| 2009 | 3.353658536585366 | 2.875694916931424 |
| 2010 | 3.572115384615385 | 2.950665247304781 |
| 2011 | 3.433613445378151 | 2.951887309108078 |
| 2012 | 3.473333333333333 | 2.964648868355774 |
| 2013 | 3.440397350993377 | 2.956344081909349 |
| 2014 | 3.265328874024526 | 2.89068241704015 |
| 2015 | 3.1389955630391 | 2.843670263465883 |
| 2016 | 3.002957607624055 | 2.814642328810233 |
| 2017 | 2.892805126711331 | 2.766470980075643 |
| 2018 | 2.97384886890803 | 2.779246716367921 |
| 2019 | 2.960170697012802 | 2.777272305092497 |
| 2020 | 3.072307692307692 | 2.822158169633342 |
+-----+-----+-----+
Fetched 15 row(s) in 0.62s
```

Seeing the above results, it's clear that Starbucks does the fine business between year 2010-2014 and after that there ratings declined between the year 2016-2019.

McDonalds :

Below is the query and results for average user ratings and average polarity ratings for 'McDonalds':

```
Select year(to_timestamp(trim(date_stamp), 'yyyy-MM-dd')) as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by period order by period;
```

```
[compute-1-1:21000] > Select year(trim(date_stamp), 'yyyy-MM-dd') as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by period order by period;
Query: Select year(trim(date_stamp), 'yyyy-MM-dd') as period, avg(user_rating) as average_user_rating, avg((polarity*5)+5)/2 as average_polarity from data_for_analyses where business_name like "%McDonalds%" and user_rating is not null group by period order by period
Query submitted at: 2019-12-01 18:12:26 (Coordinator: http://compute-1-1.local:25000)
Query progress can be monitored at: http://compute-1-1.local:25000/query_plan?query_id=554c4eabbbdb7c34:34968db00000000
+-----+
| period | average_user_rating | average_polarity |
+-----+
| 2010 | 3.5 | 2.83611118085682 |
| 2011 | 2 | 2.829422950744629 |
| 2012 | 3 | 2.622717749327421 |
| 2013 | 2.777777777777778 | 2.465335810557008 |
| 2014 | 1.65 | 2.349281987233553 |
| 2015 | 2.388888888888889 | 2.551084604864526 |
| 2016 | 2.326086956521739 | 2.545754286339101 |
| 2017 | 1.918367346938775 | 2.58898644701977 |
| 2018 | 2.022727272727273 | 2.531607875399376 |
| 2019 | 1.95 | 2.479516484410851 |
| 2020 | 1 | 2.028724737465382 |
+-----+
Fetched 11 row(s) in 0.62s
```

Seeing the above results, it can be seen that McDonald's average ratings declined between the years 2010-2014, but it shows progress in 2015 and after.

NOTE: All the analytics that we have made is based upon the data which is given. So, the results are subject to change as the data changes.