# BDT Final Project

## Facebook Post Analysis

Professor : Mrudula Mukadam
Student : Ling Sun 986265

# Facebook Post Analysis - inspirations

1. For a single account, What the most common words used in all of the posts

2. For a single account, how many posts posted each year?

3. For a single account, what are the most commented posts?

4. What are the most liked posts?

5. Who are my 'best friends'?

6. ......

# Development Environment

OS : mac OS high Sierra Version 10.13.4

Hadoop 2.6.5

Apache Spark 2.6

Hive 2.3.1

Python 3.6.4

Kafka 1.1.0

# Problems

1.  Failed to start Hadoop server
2.  Failed to start namenode server
3.  Failed to start hive metastore server

But till last Sunday, I got the dev environment set

4.  Spark didn't connect with Hive successfully
5.  Don't know anything about Python

......there will be more to be found

# Data Flow

1. Using [Facebook Graph API](#) and Http to request real-time streaming data

2. Receiving and parsing data

3. Sending data to Kafka

4. Using Spark Streaming to read data from Kafka

5. Saving data to HDFS via Hive

6. Using Spark SQL to query and analyse

7. Using Plotly to visualize data

# Code - Libs For Python

```python
1  import facebook
2  import requests
3  from kafka import KafkaProducer, KafkaClient
4  import json
5  import fpa_conf
6
```

```python
1   import os
2   import re
3   import json
4   import pandas as pd
5   import matplotlib.pyplot as plt
6   import plotly as py
7   import plotly.graph_objs as go
8
9   from pyspark import SparkContext
10  from pyspark.streaming import StreamingContext
11  from pyspark.streaming.kafka import KafkaUtils
12  from pyspark.sql import Row, SQLContext
13  from pyspark.sql import HiveContext
14  from pyspark.sql import SparkSession
15  from dateutil import parser
16  import numpy as np
17  import time
18  import fpa_conf
```

# Code - Producer

```python
60    def sendToKafka(data):
61        msg = json.dumps(parsePosts(data))
62        producer.send(fpa_conf.topic,msg.encode('utf-8'))
63
64    if __name__ == "__main__":
65        posts = graph.get_object(fpa_conf.user+"/posts?fields=message,created_time,comments{comment_count}")
66        sendToKafka(posts['data'])
67        while True:
68            try:
69                print("next page")
70                posts = requests.get(posts['paging']['next']).json()
71                if 'data' in posts:
72                    data2 = posts['data']
73                    sendToKafka(data2)
74            except KeyError:
75                break
76        print("done sending...")
```

# Code - Consumer

```python
148    if __name__ == "__main__":
149
150        sc = SparkContext(appName="CS523FinalProject")
151        sc.setLogLevel("ERROR")
152        sc.setSystemProperty("hive.metastore.uris", "")
153
154        ssc = StreamingContext(sc,20)
155        print("start reading data from kafka...")
156        kvs = KafkaUtils.createDirectStream(ssc,[fpa_conf.topic], {"metadata.broker.list": fpa_conf.brokers})
157        parsed = kvs.map(lambda v: json.loads(v[1])).flatMap(lambda post: post.values())
158
159        if parsed is not None:
160            print("start analysising...")
161            hc = getHiveContextInstance(sc)
162            hc.sql("drop table if exists t_posts")
163            posts = parsed.map(lambda r: (r['id'],r['message'],r['created_time'],r['likes'],r['comment_count']))
164
165            posts.foreachRDD(createTable)
166
167            #names = parsed.map(lambda r : r['like_names'])
168            #names.foreachRDD(countTop5LikedMostNames)
169            print('end analysising...')
170
171        else:
172            print('no data')
173        if ssc is not None:
174            ssc.start()
175            ssc.awaitTermination()
176
```
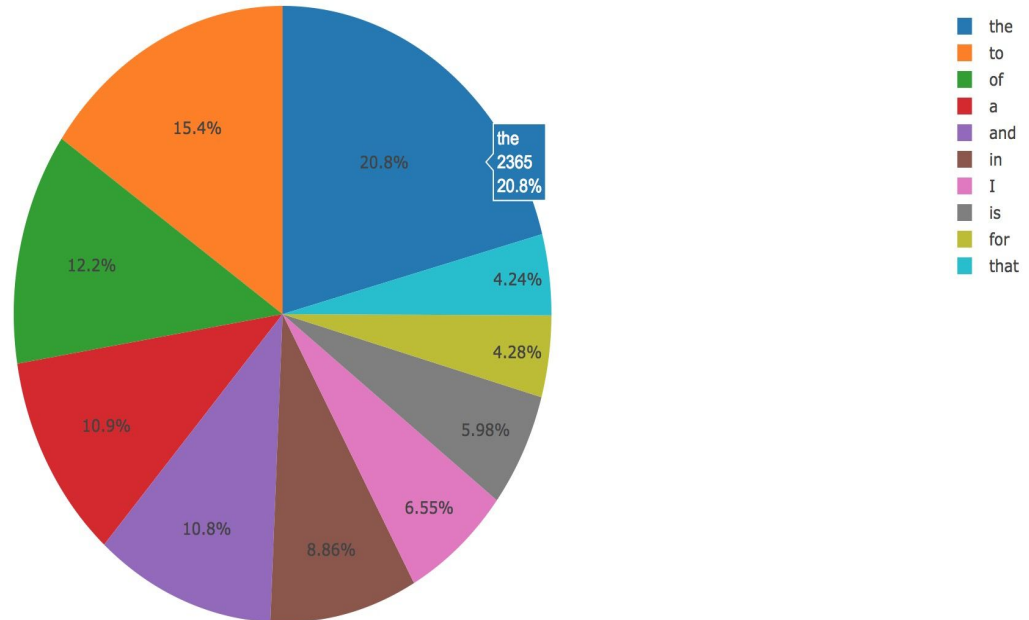
# Code - Consumer

```python
84    #fig2: count top 10 most used words
85    def countTop10MostUsedWords(hc):
86        messages = hc.sql("SELECT message FROM t_posts")
87        if messages is not None :
88            top10words = messages.rdd.flatMap(lambda p:p.message.split(' '))\
89                .map(lambda x: (x, 1)).reduceByKey(lambda x,y: x+y)\
90                .filter(lambda x: x[0]!='')\
91                .sortBy(lambda x: x[1],ascending=False)\
92                .take(10)
93
94            df = pd.DataFrame(top10words,columns=["word","count"])
95            data = go.Data([go.Pie(labels=df["word"],values=df["count"])])
96            layout = go.Layout(title=fpa_conf.user+': Top 10 Words Used')
97            fig = go.Figure(data=data, layout=layout)
98            py.offline.plot(fig, filename="/Users/sunling/MUM/BDT/project/BGFacebook/output/test/"+fpa_conf.user+
99
```
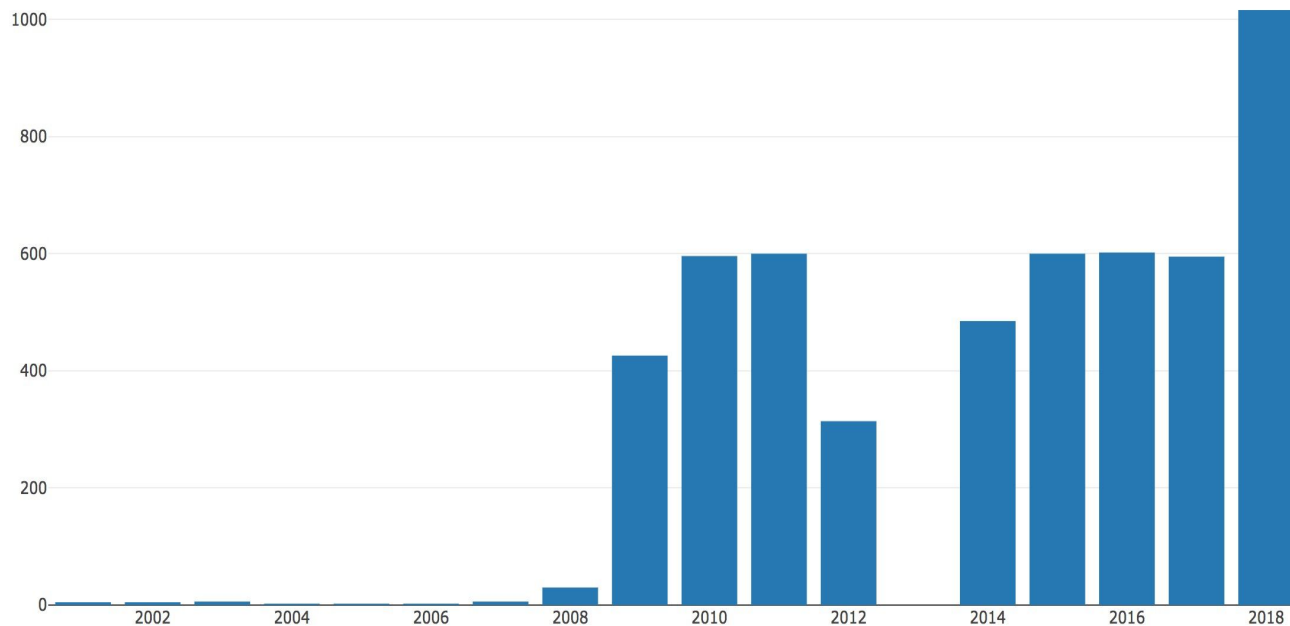
# Outputs - Pie



https://www.facebook.com/BillGates

BillGates: Top 10 Words Used

# Outputs - Bar

CNN Posts of Each Year | Total:5288



https://www.facebook.com/cnn/

# Outputs - Table

Top 10 most Commented posts

| Post | Comment_Count |
|------|---------------|
| Join me LIVE in Elkhart, Indiana! Great crowd for a #MAGA rally! | 24172 |
| I hereby demand, and will do so officially tomorrow, that the Department of Justice look into whether or not the FBI/DOJ infiltrated or surveilled the Trump Campaign for Political Purposes - and if any such demands or requests were made by people within the Obama Administration! | 22482 |
| Today I announced that the United States will withdraw from the Iran nuclear deal. | 18130 |
| Our great First Lady Melania Trump is doing really well. Will be leaving the hospital in 2 or 3 days. Thank you for so much love and support! | 14390 |
| To the students, families, teachers, and personnel at Santa Fe High School: We're with you in this tragic hour, and we will be with you forever. | 13147 |
| Can you believe that after ONE YEAR of a disgusting, illegal and unwarranted $10,000,000 Witch Hunt, we have had the most successful first 17 months of an administration in U.S. history - by far! | 13030 |
| Promises made, promises kept 🇺🇸 | 12844 |
| A great day for Israel! 🇺🇸🇮🇱 | 11504 |
| I want to deliver a message to the long-suffering people of Iran: The people of America stand with you. The future of Iran belongs to you. | 10639 |
| #Covfefe | 9820 |

https://www.facebook.com/DonaldTrump/

# Demo

**#start hadoop**

cd /usr/local/Cellar/hadoop-2.6.5/

sbin/start-dfs.sh

sbin/start-yarn.sh

**#start hive metastore**

hive --service metastore

**#start Kafka, mysql, spark**

brew service kafka start

......

**# run consumer**

spark-submit --jars
/usr/local/Cellar/kafka/1.1.0/libexec/libs/spark-streaming-kafka-assembly_2.11-1.6.1.jar
/Users/sunling/MUM/BDT/project/BGFacebook/BGConsumer.py

**# run producer**

python
/Users/sunling/MUM/BDT/project/BGFacebook/BGProducer.py

# Thank you!

Any questions ?