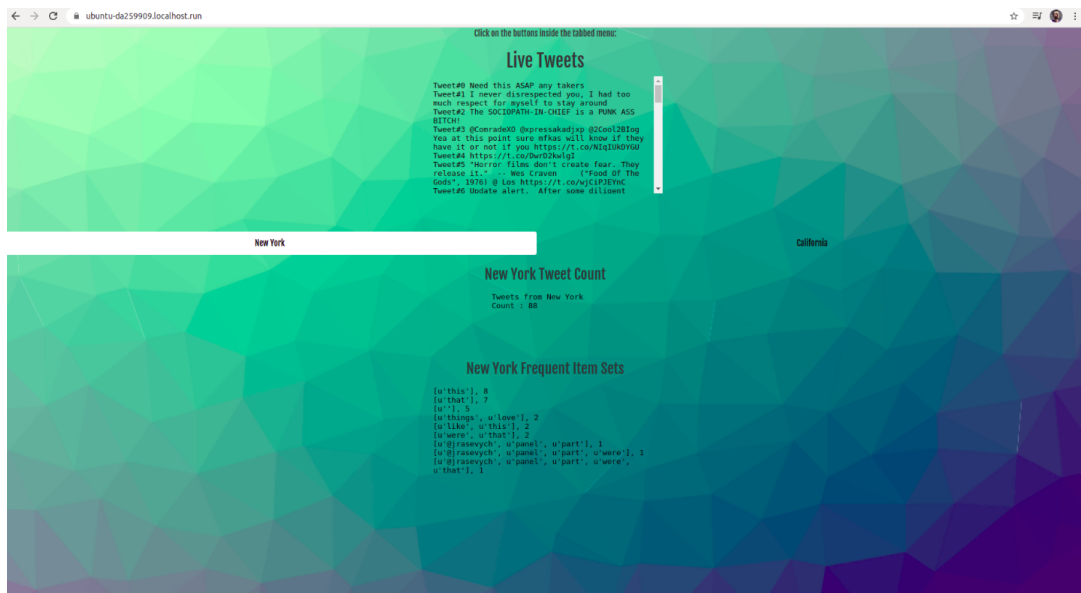# CS6847 Cloud Computing.
# Assignment - 5

Suhas Pai, CS17B116

## 1 Setup

I installed Apache Kafka, Apache Spark, and Zookeeper, initially on a single machine, and created its AMI Image. Zookeeper is initially started, and then the Kafka server is started. The config files are modified on both instances, to let the zookeeper server know the IP addresses of the Kafka brokers. The producer config file is changed as well. Zookeeper server is run only on one instance, while the Kafka server is started on both instances. Then, on one of the instances, a python script is run. The script uses Tweepy and PyKafka. Tweepy is used for collecting all tweets from New York and Los Angeles ,based on the location coordinates of the tweet. Two producers are created using PyKafka, and these tweets are then sent to the broker. Each instance holds one topic's data. Then, the Spark master is started. The Streaming Context is then started, and a stream is created. Tweets are then picked up from the appropriate topic, written to a file, and then the mllib library is used to run the frequent dataset algorithm. A http server is run in the folder of the html index file, and then using a python tool, the local website is made global.



Final website

# 2    Generating Tweets

```
producer1 = KafkaProducer(bootstrap_servers='localhost:9092',value_serializer=lambda m: json.dumps(m).encode('ascii'))
producer2 = KafkaProducer(bootstrap_servers='localhost:9092',value_serializer=lambda m: json.dumps(m).encode('ascii'))
topic_name1 = 'NY'
topic_name2 = 'CA'
NY_location = [-79.7619,40.4774,-71.7956,45.0159]
CA_location = [-124.48,32.53,-114.13,42.01]
access_token = "1255444235421179910-XQwfbYRU0YoY0qqsSQRXneGaN85Gza"
access_token_secret = "eBeNZrpK8c7iVuCUEtqRks3yieFmQLH5iXMj7rmPLZAUn"
consumer_key  =  "Yj0wHPWKToAkOIzbI1Ef4XAu6"
consumer_secret = "YZoyXWzTm8I2ldWghhi4s7jiKx7H3VqXnpmil030lp15IimE6C"


auth = OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_token_secret)

api = tweepy.API(auth)
```

Initial setup for tweepy and PyKafka

```python
class CustomStreamListener(tweepy.StreamListener):
    def on_status(self, status):
        try:

            text = status.text.encode('utf-8')
            loc = status.place.bounding_box.coordinates
            if loc[0][0][0]>=NY_location[0] and loc[0][0][0]<=NY_location[2]:

                producer1.send(topic_name1,{'city':"NY",'text':text})
            else:
                producer2.send(topic_name2,{'city':"CA",'text':text})

        except Exception as e:
            print(e)
            pass

    def on_error(self, status_code):
        print >> sys.stderr, 'Encountered error with status code:', status_code
        return True # Don't kill the stream

    def on_timeout(self):
        print >> sys.stderr, 'Timeout...'
        return True # Don't kill the stream

sapi = tweepy.streaming.Stream(auth, CustomStreamListener())
sapi.filter(locations=[-124.48,32.53,-114.13,42.01,-79.7619,40.4774,-71.7956,45.0159])
```

The class that filters tweets based on Box coordinates of New York and California

Two producers are started using PyKafka, one for each topic. The zookeeper ip address is the default address of localhost, and the tweets are dumped in json format, so that they can be filtered easily in the stream. Every tweet is read, and based on which box the coordinate falls in, it's location is identified.

# 3  Starting Kafka and Zookeeper

Zookeeper is started on one of the instances. Then two brokers, with different broker ids, are started on the two instances.



Zookeeper being started



Kafka broker with id -0, started on instance 1



Kafka broker with id - 1, started on instance 2

# 4 Spark Streaming

A stream is started, and tweets are pulled every 30 seconds. A window of 600 seconds is used, and all the tweets in this window are used for performing tweet count, and the frequent item set algorithm. The dstream cannot be used directly for the fp growth algorithm, and so, each rdd is written to a text file and this is used as input for the algorithm. The output of the algorithm is written to another text file, which is displayed in the website.

```python
ssc = StreamingContext(sc, 30)
kafkaStream = KafkaUtils.createStream(ssc, 'localhost:2181', 'spark-streaming', {'NY':1,'CA':1})
kafkaStream1 = kafkaStream.window(600,30)

parsed = kafkaStream1.map(lambda v: json.loads(v[1]))

parsed.foreachRDD(lambda x : print_all_tweets(x.collect()))
#parsed.count().map(lambda x:'Tweets in this batch: %s' % x).pprint()

NYstream = parsed.filter(lambda x : x['city']=='NY')
NYstream.count().map(lambda x : 'Tweets from New York Count : %s' %x).pprint()

CAstream = parsed.filter(lambda x : x['city']== 'CA')
CAstream.count().map(lambda x : 'Tweets from California Count : %s' %x).pprint()

NYstream.count().map(lambda x : 'Tweets from New York Count : %s' %x).foreachRDD(lambda x : f(x.collect()))
CAstream.count().map(lambda x : 'Tweets from California Count : %s' %x).foreachRDD(lambda x : g(x.collect()))

NYstream.foreachRDD(lambda x : print_list(x.collect(),'NY'))
CAstream.foreachRDD(lambda x : print_list(x.collect(),'CA'))

ssc.start()
ssc.awaitTermination(timeout=3600)
```
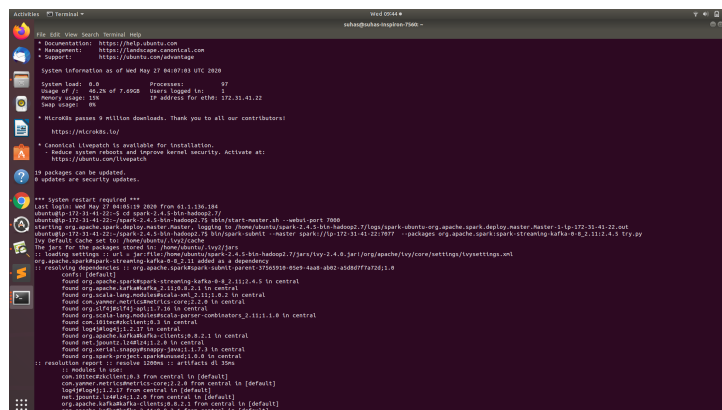
Relevant code for the Spark stream

```python
def FPGrowthAlgo(city):
    try:
        data = sc.textFile(home_dir+"/website/tweets_"+city+".txt")
        transactions = data.map(lambda line: line.strip().split(' '))
        model = FPGrowth.train(transactions, minSupport=0.001, numPartitions=1)
        result = model.freqItemsets().collect()
        print("################# FP Growth OP")
        result = sorted(result, key = lambda x : -x[1])
        ones = filter(lambda x : len(x[0])==1,result)
        twos = filter(lambda x : len(x[0])==2,result)
        threes = filter(lambda x : len(x[0])>=3,result)
        f = open(home_dir+"/website/freq_"+city+".txt","w")
        for fi in ones[:3]:
            f.write(str(fi[0]) + ", " + str(fi[1])+"\n")
        for fi in twos[:3]:
            f.write(str(fi[0]) + ", " + str(fi[1])+"\n")
        for fi in threes[:3]:
            f.write(str(fi[0]) + ", " + str(fi[1])+"\n")
    except Exception as e:
        print('Passed' +  str(e))
        pass
```
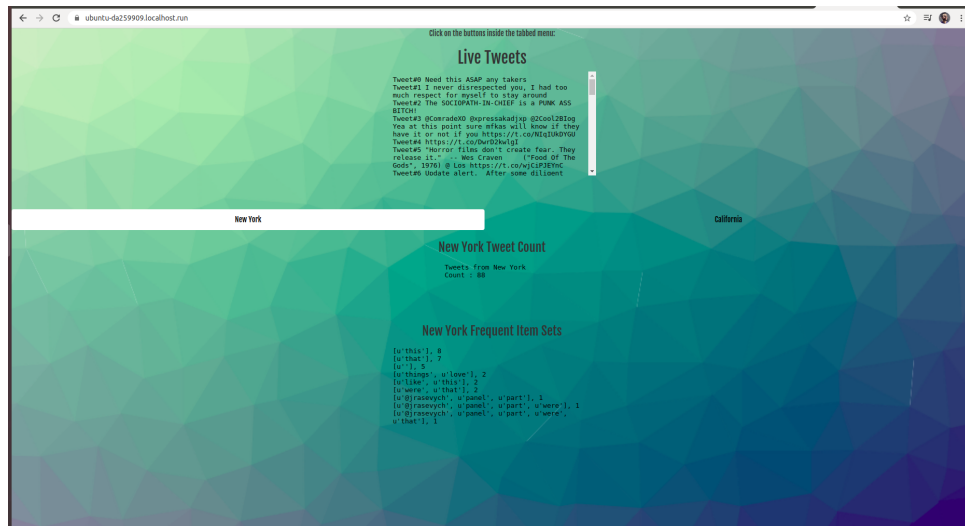
Fp Growth algorithm



The command for starting spark, and submitting the job

# 5 Relevant Screenshots

The website reloads every 15 seconds, and outputs the data of the Spark Fp Growth algorithm. More proof of working can be found in the video attached along with this report.



Initial New York Tweet count



New York Tweet count after 30 seconds