# Assignment 2

## Data Visualization with Kibana

# Introduction

In this project we will be working with NYC Open Data published by the city of New York pertaining to 311 service requests collected since 2010 with over 21 million rows with 41 columns.

https://nycopendata.socrata.com/Social-Services/311-Service-Requests-from-2010-toPresent/erm2-nwe9.

# Problem Background

We have been hired as a data analyst by the city of New York to gain valuable insights from their huge data set for 311 service requests. Our task is to use the ELK stack you successfully installed and configured in your GCP platform. Successful completion of this task includes creating a Logstash configuration file (a sample is given to you) as well as a geo-point template (for maps), creating a GCP instance and firing Logstash to ingest the NYC 311 service request data into Elasticsearch and using Kibana to analyze and visualize the results as per the questions given.

# Objectives

1. To further expose we are using ELK stack as an analytic tool to analyze streaming realistic big data.

2. To give us experience working with opened end problems, that are similar to problems that we will face in our career as a big data professional.

3. At the end of this project we should:

- Gain sufficient confidence in creating Logstash configuration files and creating Elasticsearch   indices, advanced queries, charts, maps and dashboards using Kibana i.e. fully using ELK stack in real big data scenarios.
- Gain an appetite for working with large streaming datasets.
- Be aware of the potential and benefits of analyzing large streaming data using big data tools.

# *Installation of Elasticsearch, Kibana, Logstash*

# *Logstash Config File*

```
input {
    file {
        path => " /home/amulpatel155/logstash-7.5.1/311_service.csv"
        start_position => "beginning"
        sincedb_path => "/dev/null"
    }
 }

filter {
        csv {
                separator =>","
                columns => ["Unique Key","Created Date","Closed Date","Agency","Agency
Name","Complaint Type","Descriptor","Location Type","Incident Zip","Incident
Address","Street Name","Cross Street 1","Cross Street 2","Intersection Street
1","Intersection Street 2","Address Type","City","Landmark","Facility Type","Status","Due
Date","Resolution Description","Resolution Action Updated Date","Community
Board","BBL","Borough","X Coordinate (State Plane)","Y Coordinate (State Plane)","Open
Data Channel Type","Park Facility Name","Park Borough","Vehicle Type","Taxi Company
Borough","Taxi Pick Up Location","Bridge Highway Name","Bridge Highway Direction","Road
Ramp","Bridge Highway Segment","Latitude","Longitude","Location"]
                }
date{  match => ["Created Date", "MM/dd/yyyy hh:mm:ss a"]
        target => "Created Date"
}
date{  match => ["Closed Date", "MM/dd/yyyy hh:mm:ss a"]
        target => "Closed Date"
}
date{  match => ["Due Date", "MM/dd/yyyy hh:mm:ss a"]
        target => "Due Date"
}
date{  match => ["Resolution Action Updated Date", "MM/dd/yyyy hh:mm:ss a"]
        target => "Resoultion Action Updated Date"
}
    mutate {convert => ["Incident Zip","integer"]}
    mutate {convert => ["BBL","integer"]}
    mutate {convert => ["X Coordinate (State Plane)","integer"]}
    mutate {convert => ["Y Coordinate (State Plane)","integer"]}
    mutate {convert => ["Latitude","float"]}
    mutate {convert => ["Longitude","float"]}
    mutate {copy =>
            { "Longitude" => "[location][lon]"
              "Latitude" => "[location][lat]" }
        }
    mutate {replace => { "Location" => "%{Longitude},%{Latitude}" }}
    }

output {
 elasticsearch {
 hosts => "localhost"
 index => "nycinfo"

  }
stdout {codec => dots}
```
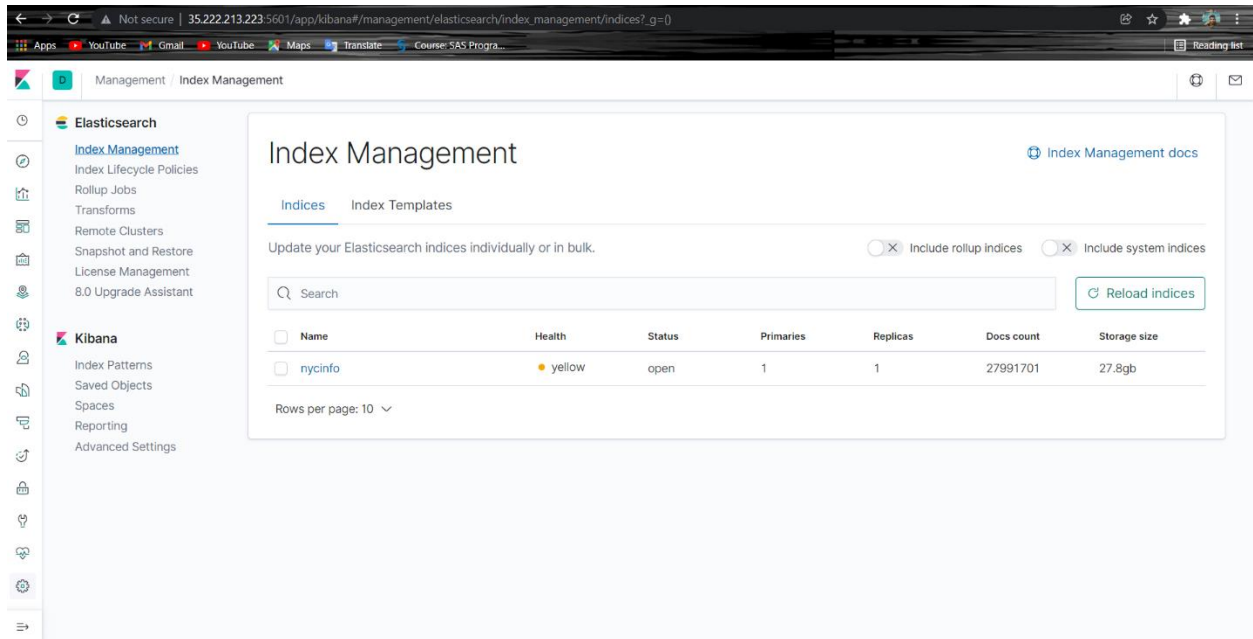
# Mapping

# Dataset is uploaded in Kibana
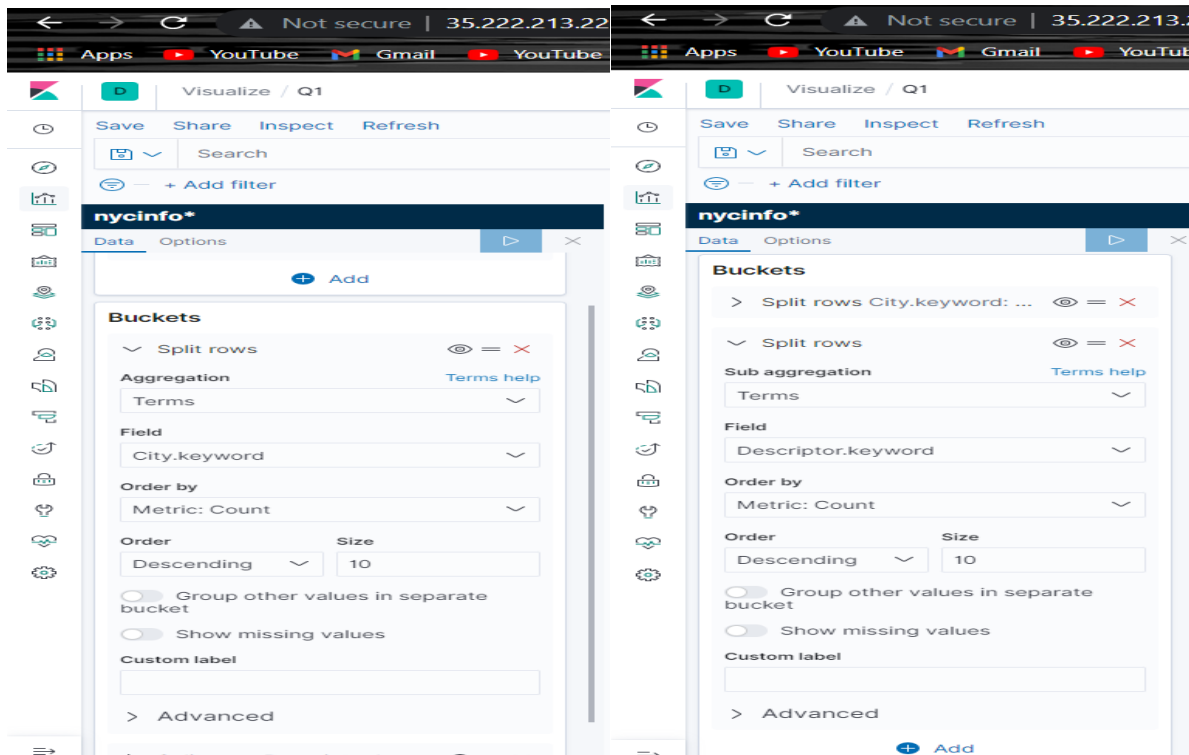


# Dataset

# *Analysis*

**Create a table showing the top 10 cities with the highest calls alongside the count of top 10 complaint calls (by Descriptor) in each city.**
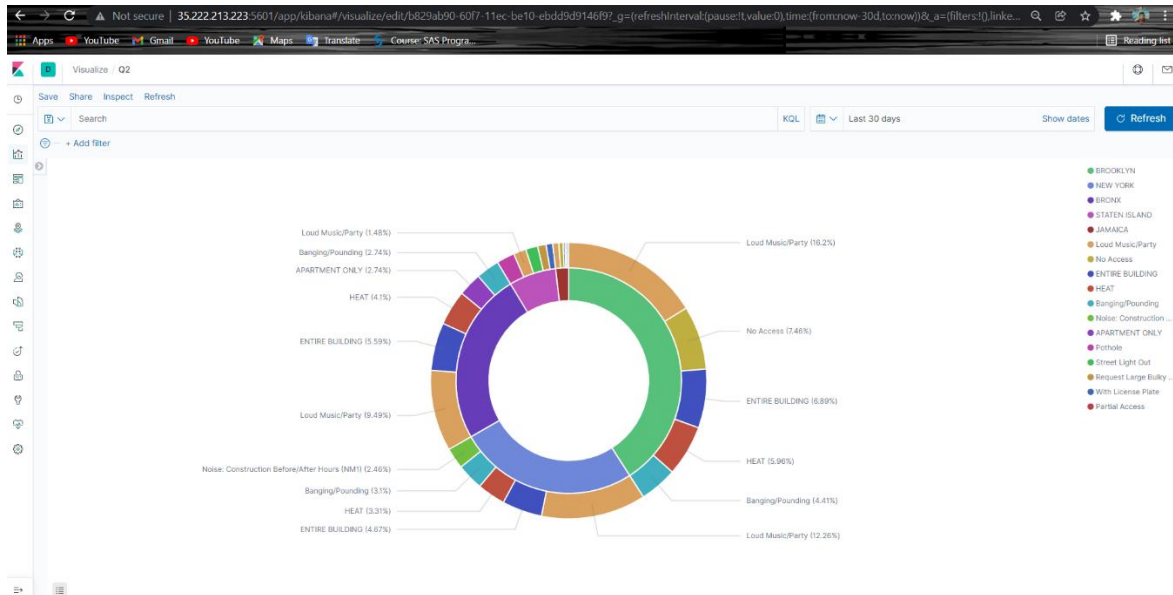
Visualization: -



Query: -

**Create a pie chart showing the top 5 cities with the highest calls alongside the top five calls (Descriptor) in each city.**
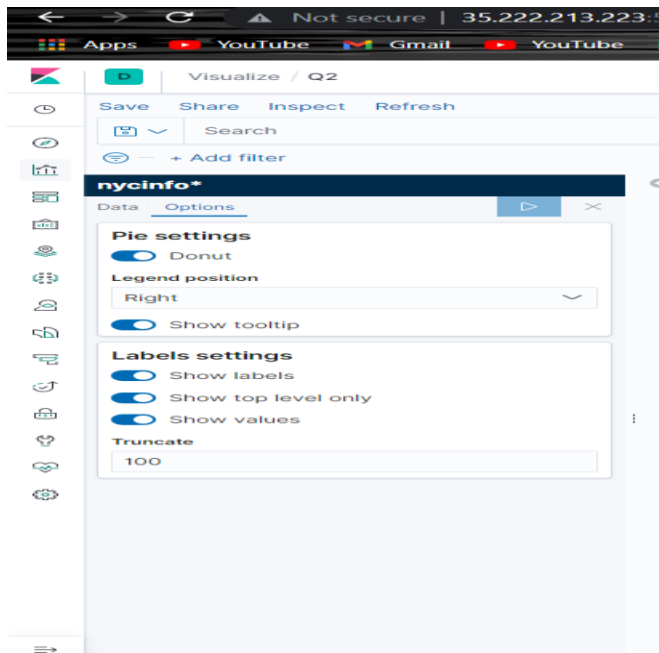
Visualization: -



Query: -

*Create a tag cloud representing the top 20 call descriptors.*

Visualization & Query: -



*Create a coordinated map of all the major call descriptors in each city*

*Create a dashboard for all visualizations of 1to 4 above.*

Visualization: -



# Teamwork: -

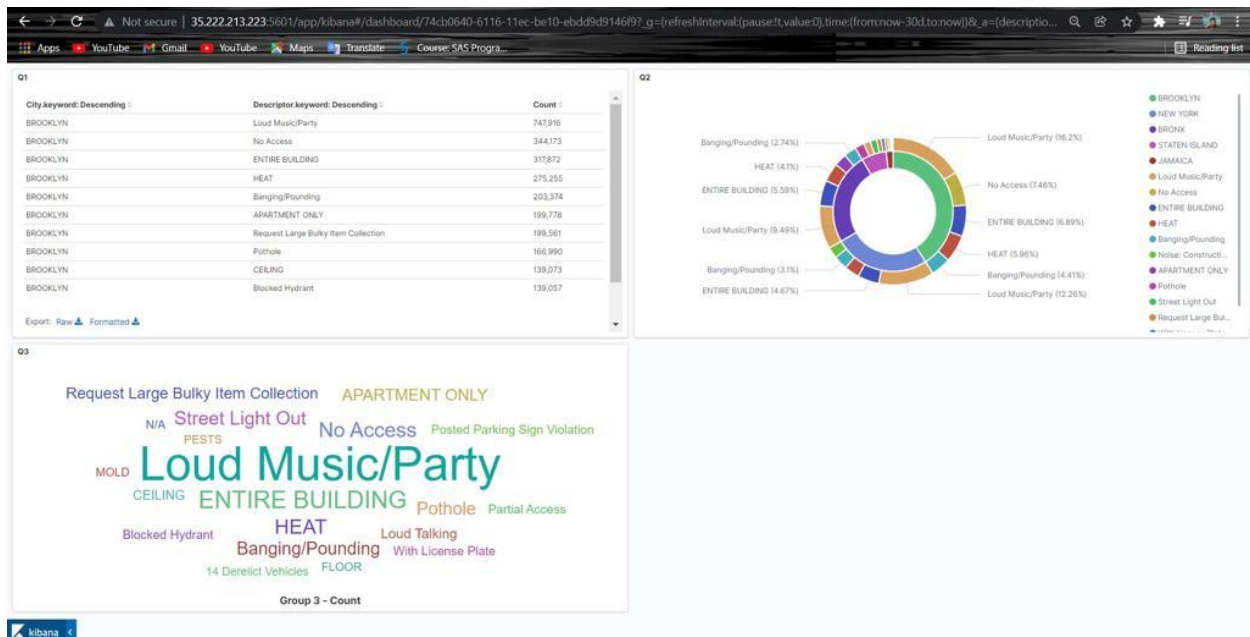As being team member, all the members have equally contribution in this project. Karun has installed the ELK, create firewall for Kibana and setup the Kibana for further steps. Main task of Amish was uploaded data set into Kibana. Amish has insert data into Logstash, he did some changes in Logstash file, create config file and uploaded dataset into Kibana. Amul has perform all the research question after competition of Karun and Amish task. As well as Amul has created this document file.

*Thank You*