

Project 4: Elasticsearch

Participants:

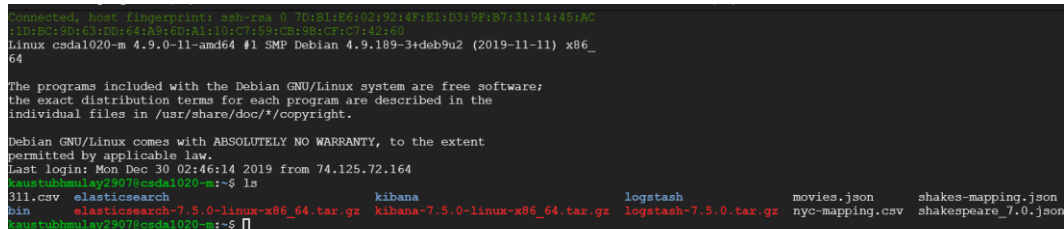
1. Deenu Yadav
2. Fanny Guevara
3. Kaustubh Mulay

Dataset used: <https://nycopendata.socrata.com/Social-Services/311-Service-Requests-from-2010-to-Present/erm2-nwe9>

Number of rows: 22.1 million

Initial steps for setup:

1. Upload dataset to VM instance (assuming Elasticsearch, Kibana and Logstash are all installed, the configuration files correctly modified as in Project 3). Checking contents of upload location:

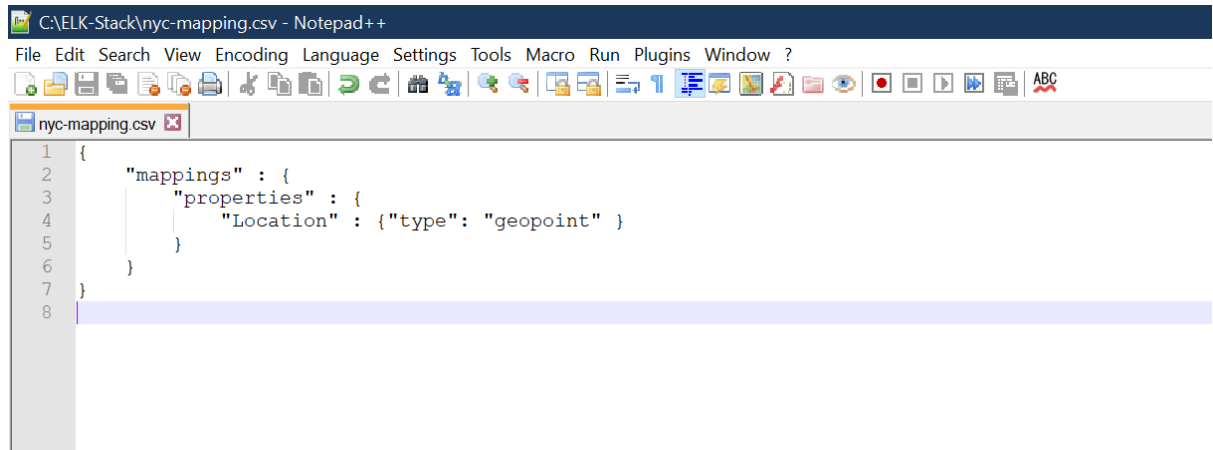


```
connected, host fingerprint: ssh-rsa 0 7b:81:e6:02:92:4f:e1:d3:9f:b7:31:14:45:ac
:1d:bc:90:43:dd:64:a9:6b:a1:10:c7:59:cb:9b:cf:c7:42:60
Linux csdal020-m 4.9.0-11-amd64 #1 SMP Debian 4.9.189-3+deb9u2 (2019-11-11) x86_
64

The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

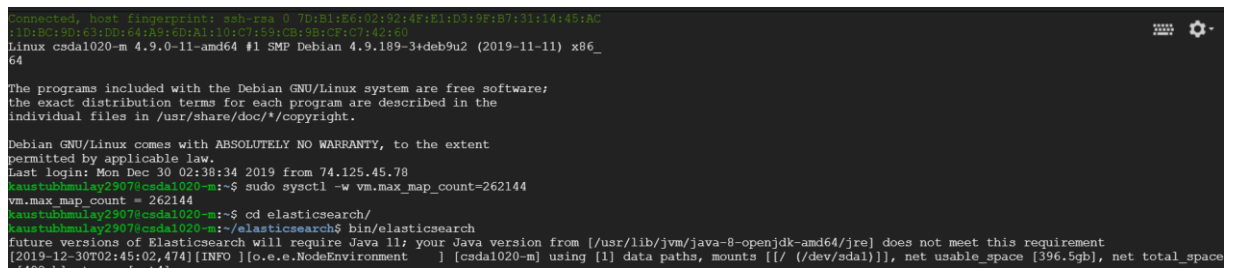
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Mon Dec 30 02:46:14 2019 from 74.125.72.164
kaustubhmulay2907@csdal020-m:~$ ls
311.csv  elasticsearch  kibana  logstash  movies.json  shakes-mapping.json
bin      elasticsearch-7.5.0-linux-x86_64.tar.gz  kibana-7.5.0-linux-x86_64.tar.gz  logstash-7.5.0.tar.gz  nyc-mapping.csv  shakespear_7.0.json
kaustubhmulay2907@csdal020-m:~$
```

2. Create mapping for coordinate map question.



```
C:\ELK-Stack\nyc-mapping.csv - Notepad++
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?
nyc-mapping.csv
1 {
2   "mappings" : {
3     "properties" : {
4       "Location" : {"type": "geopoint" }
5     }
6   }
7 }
8
```

3. Run commands for Elasticsearch after setting a higher virtual memory for the purpose:



```
connected, host fingerprint: ssh-rsa 0 7b:81:e6:02:92:4f:e1:d3:9f:b7:31:14:45:ac
:1d:bc:90:43:dd:64:a9:6b:a1:10:c7:59:cb:9b:cf:c7:42:60
Linux csdal020-m 4.9.0-11-amd64 #1 SMP Debian 4.9.189-3+deb9u2 (2019-11-11) x86_
64

The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Mon Dec 30 02:38:34 2019 from 74.125.45.78
kaustubhmulay2907@csdal020-m:~$ sudo sysctl -w vm.max_map_count=262144
vm.max_map_count = 262144
kaustubhmulay2907@csdal020-m:~$ cd elasticsearch/
kaustubhmulay2907@csdal020-m:~/elasticsearch$ bin/elasticsearch
future versions of Elasticsearch will require Java 11; your Java version from [/usr/lib/jvm/java-8-openjdk-amd64/jre] does not meet this requirement
[2019-12-30T02:45:02,474][INFO ][o.e.e.NodeEnvironment ] [csdal020-m] using [/] data paths, mounts [/ (/dev/sdal)], net usable_space [396.5gb], net total_space
[492gb], types [ext4]
```

4. Check if Elasticsearch is running or not:

```
← → ↻ ⌂ ⓘ Not secure | 35.185.18.78:9200
Apps W Wikipedia T Twitter SQL Mail Office 365 York U Stack Overflow - W...

{
  "name" : "csda1020-m",
  "cluster_name" : "elasticsearch",
  "cluster_uuid" : "ibe1yO4CRVuqj0G6L0xCCA",
  "version" : {
    "number" : "7.5.0",
    "build_flavor" : "default",
    "build_type" : "tar",
    "build_hash" : "e9ccaed468e2fac2275a3761849cbee64b39519f",
    "build_date" : "2019-11-26T01:06:52.518245Z",
    "build_snapshot" : false,
    "lucene_version" : "8.3.0",
    "minimum_wire_compatibility_version" : "6.8.0",
    "minimum_index_compatibility_version" : "6.0.0-beta1"
  },
  "tagline" : "You Know, for Search"
}
```

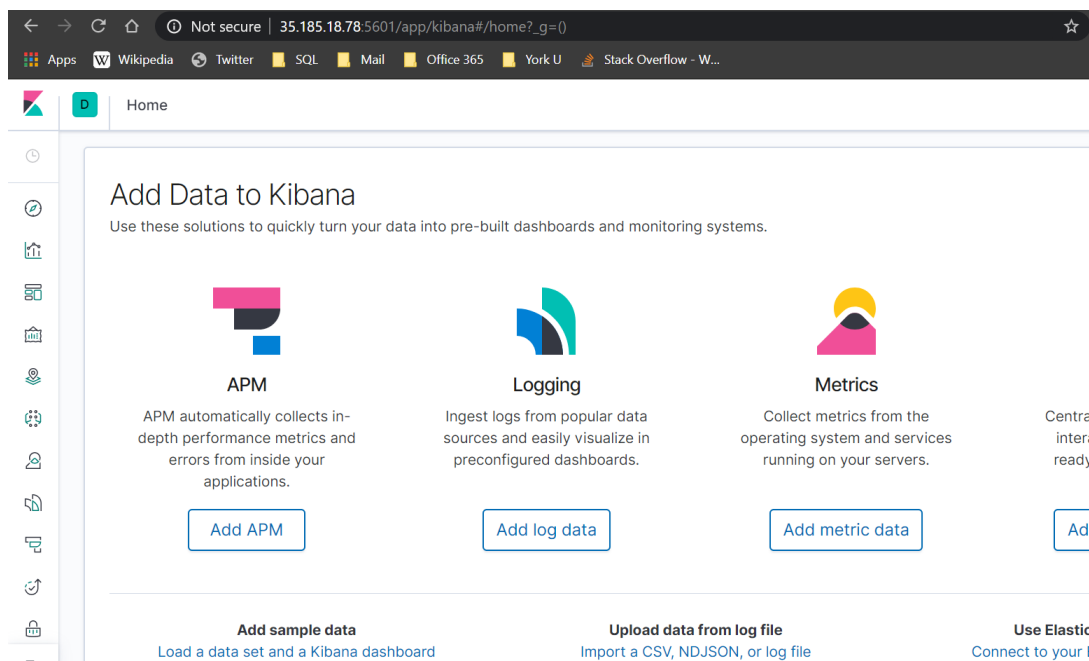
5. After opening another VM instance, running command to run Kibana in the background:

```
ssh.cloud.google.com/projects/csda1020-262300/zones/us-east1-c/instances/csda1020-m?authuser=1&hl=en_GB&projectNumber=1050248150075
Connected, host fingerprint: ssh-rsa 0 7D:B1:E6:02:92:4F:EI:D3:9F:B7:31:14:45:AC
1D:BC:90:63:DD:64:A9:6B:A1:10:C7:59:CB:9B:CF:C7:42:60
Linux csda1020-m 4.9.0-11-amd64 #1 SMP Debian 4.9.189-3+deb9u2 (2019-11-11) x86_
64

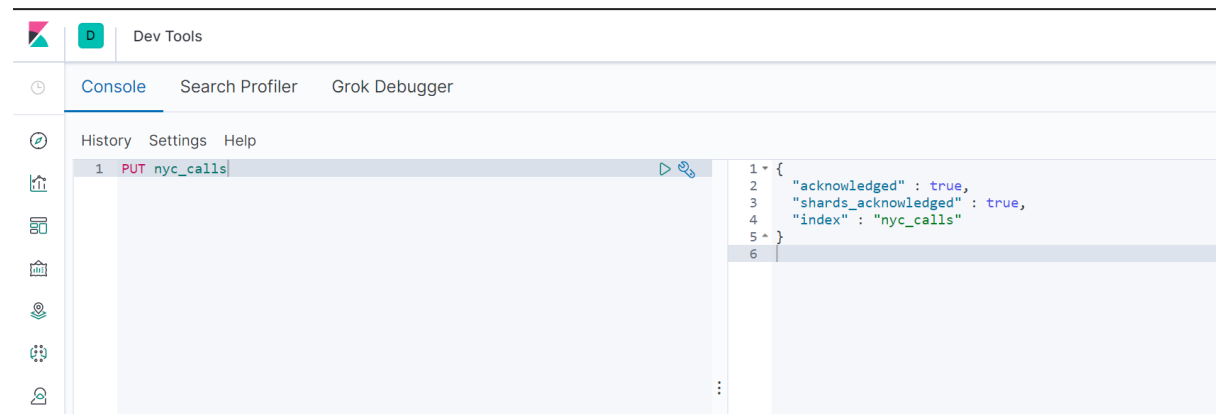
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Mon Dec 30 02:43:50 2019 from 74.125.189.97
kustubhulay2907@csda1020-m:~$ cd kibana/
kustubhulay2907@csda1020-m:~/kibana$ nohup bin/kibana
nohup: ignoring input and appending output to 'nohup.out'
]
```

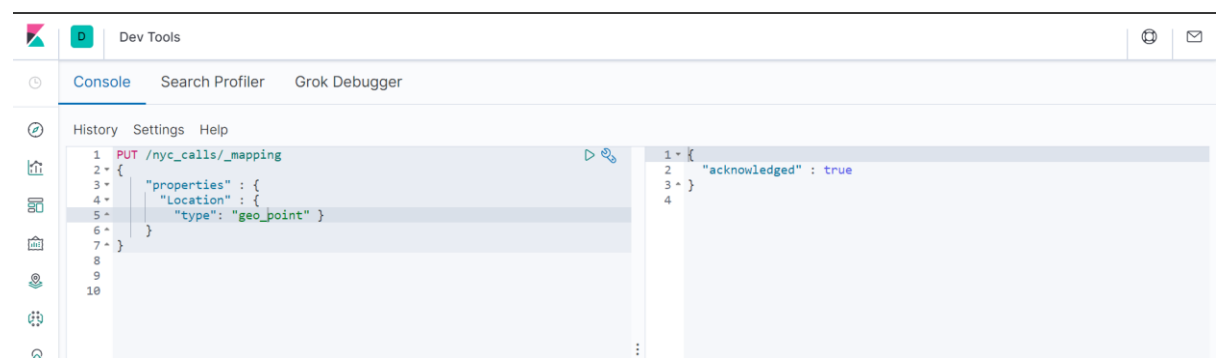
6. Checking if Kibana is successfully running or not:



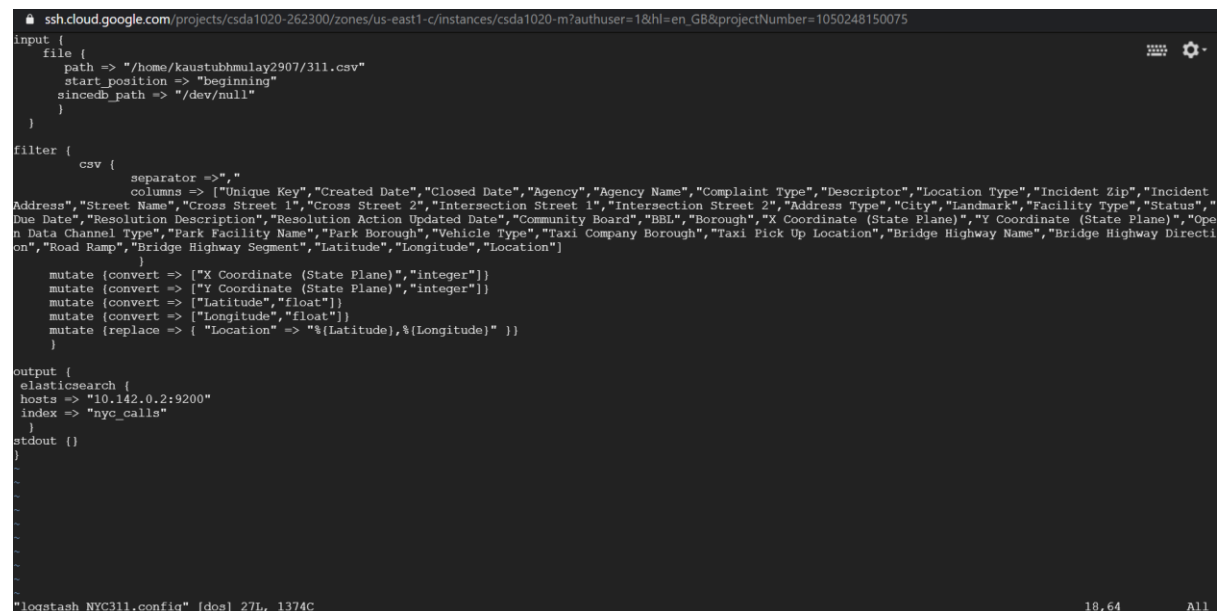
7. Creating index in Kibana:



8. Changing the mapping of the index in Kibana for the location columns:



9. Checking contents of logstash configuration file to ensure that the index name is correctly reflected:



10. Checking the contents of the Logstash file in VM instance and then running command to run Logstash in background:

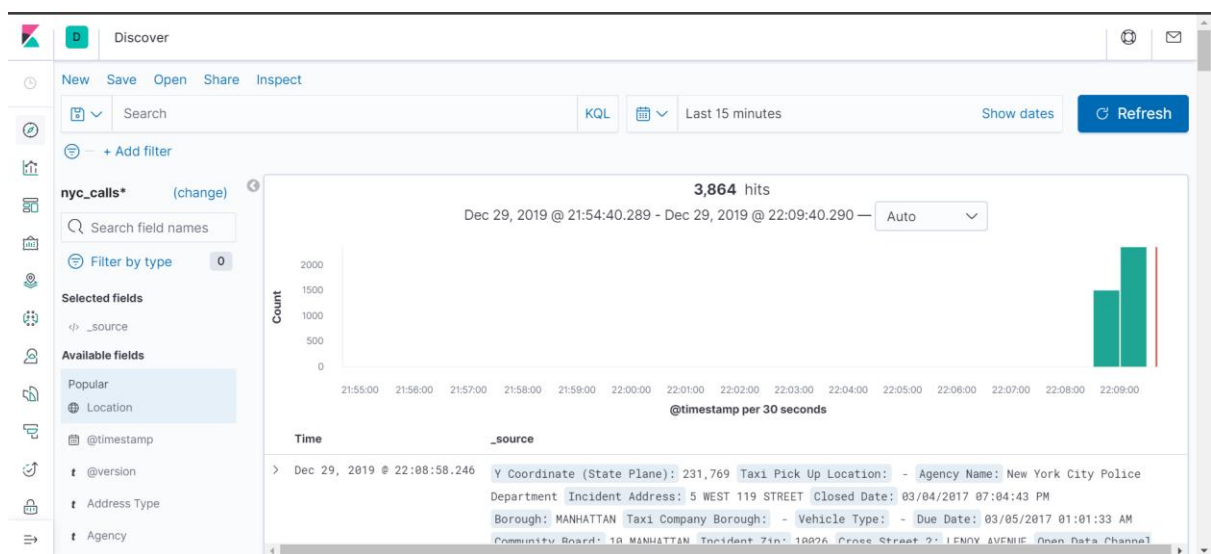
```
Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Mon Dec 30 02:46:14 2019 from 74.125.72.164
kaustubhmuly2907@csdal020-m:~$ ls
311.csv  elasticsearch-7.5.0-linux-x86_64.tar.gz  kibana  logstash  movies.json  shakes-mapping.json
bin      elasticsearch-7.5.0-linux-x86_64.tar.gz  kibana-7.5.0-linux-x86_64.tar.gz  logstash-7.5.0.tar.gz  nyc-mapping.csv  shakespear_7.0.json
kaustubhmuly2907@csdal020-m:~$ cd logstash/
kaustubhmuly2907@csdal020-m:~/logstash$ nohup bin/logstash -f /home/kaustubhmuly2907/logstash/logstash_NYC311.config
nohup: ignoring input and appending output to 'nohup.out'
^Ckaustubhmuly2907@csdal020-m:~/logstash$ ls
bin  CONTRIBUTORS  Gemfile  kibana.yml  LICENSE.txt  logstash-core  logstash_NYC311.config  nohup.out  tools  x-pack
config  data  Gemfile.lock  lib  logs  logstash-core-plugin-api  modules  NOTICE.TXT  vendor
kaustubhmuly2907@csdal020-m:~/logstash$ vi logstash_NYC311.config
kaustubhmuly2907@csdal020-m:~/logstash$ nohup bin/logstash -f /home/kaustubhmuly2907/logstash/logstash_NYC311.config
nohup: ignoring input and appending output to 'nohup.out'
^C
```

11. Checking if the data is successfully getting uploaded or not:

```
← → ↺ ⌂ ⓘ Not secure | 35.185.18.78:9200/nyc_calls/_count?
Apps W Wikipedia Twitter SQL Mail Office 365 York U Stack Overflow - W...

{"count":3191,"_shards":{"total":1,"successful":1,"skipped":0,"failed":0}}
```

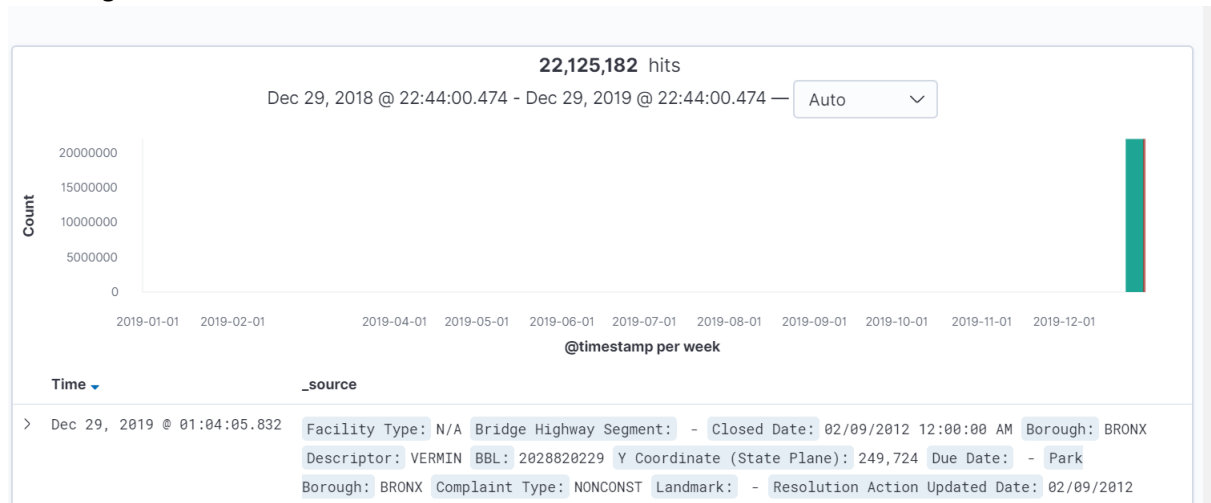
12. Checking if the documents are being created in Kibana or not:



13. Checking the progress of the data import into Kibana to see if the entire dataset was loaded or not. Full data successfully loaded:

```
{
  "count" : 22125182,
  "_shards" : {
    "total" : 1,
    "successful" : 1,
    "skipped" : 0,
    "failed" : 0
  }
}
```

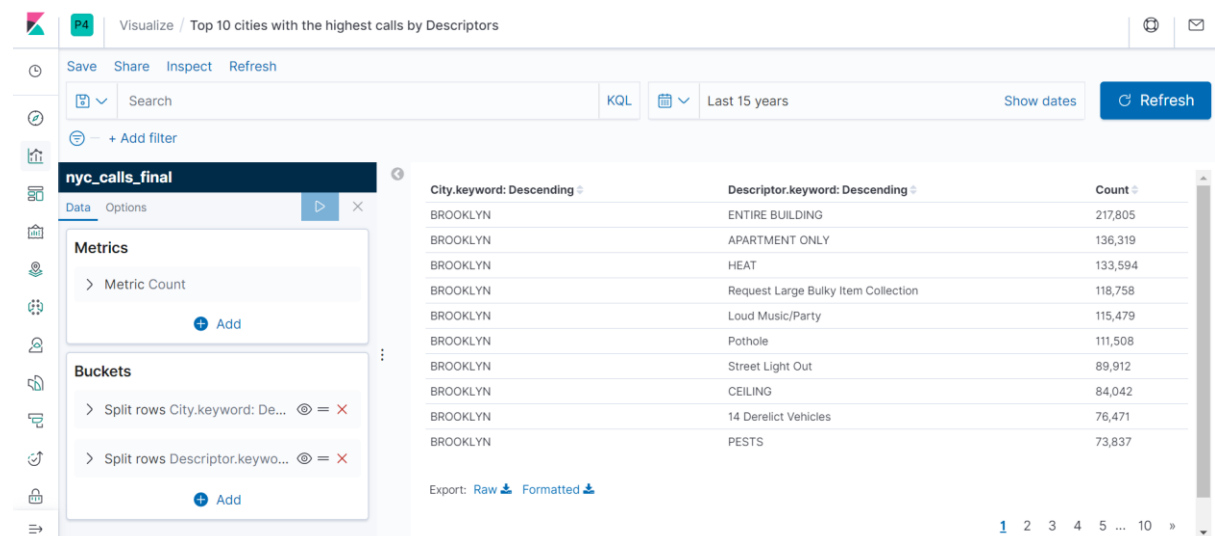
14. Checking if all the documents are visible in Kibana or not:

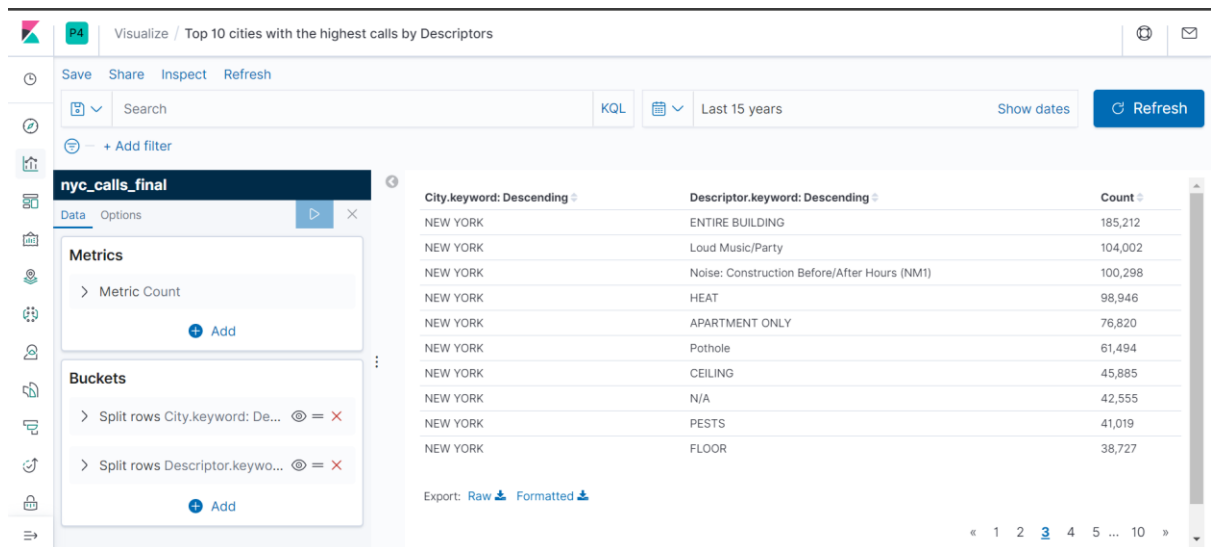
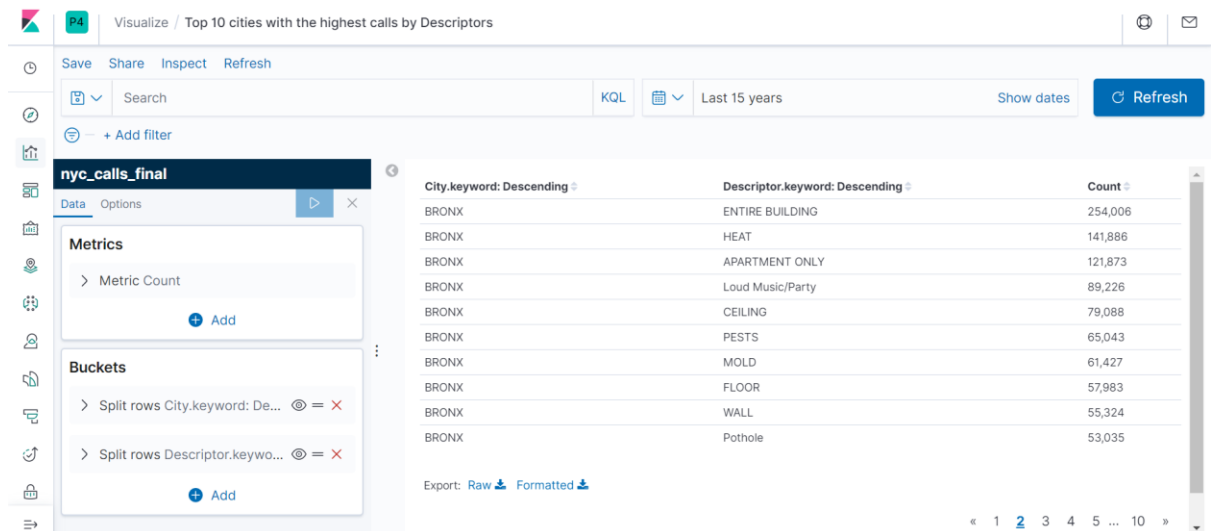


Analytical Questions:

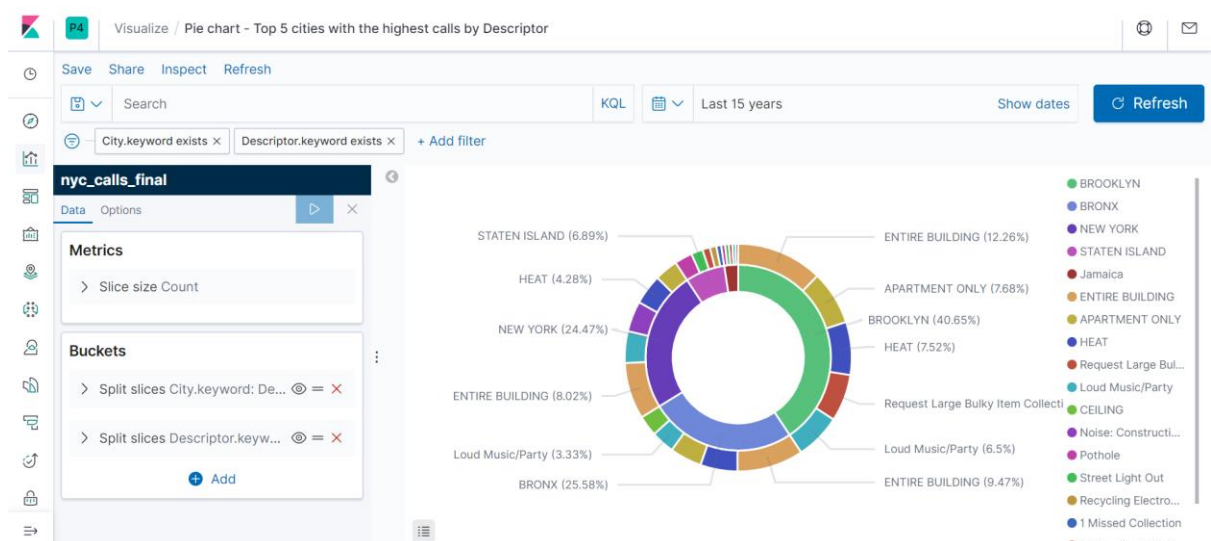
1. Create a table showing the top 10 cities with the highest calls alongside the count of top 10 complaint calls (by Descriptor) in each city.

Showing first few pages:

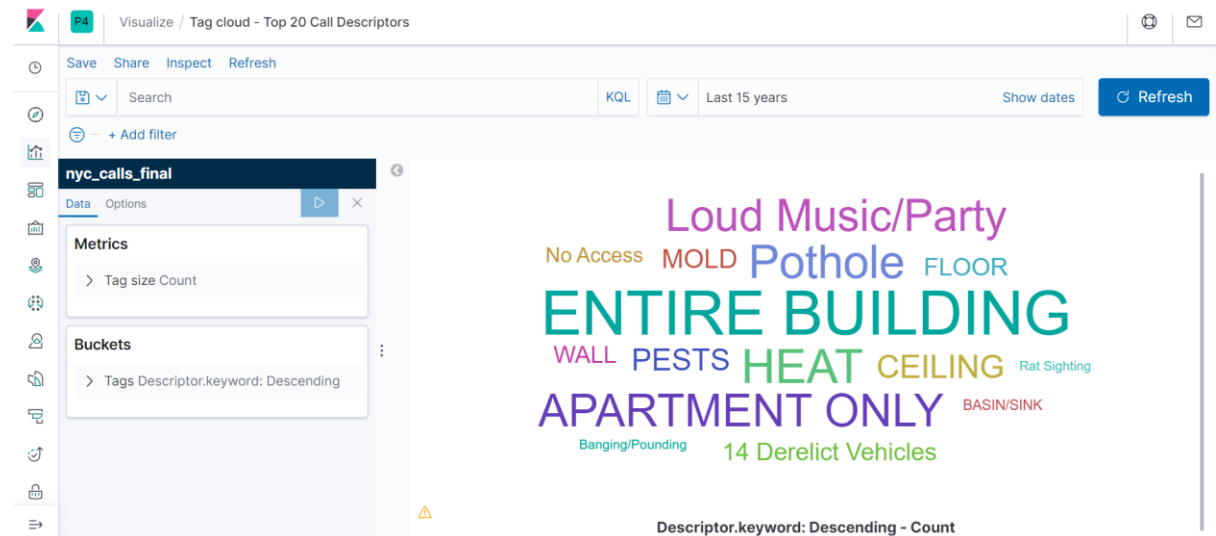




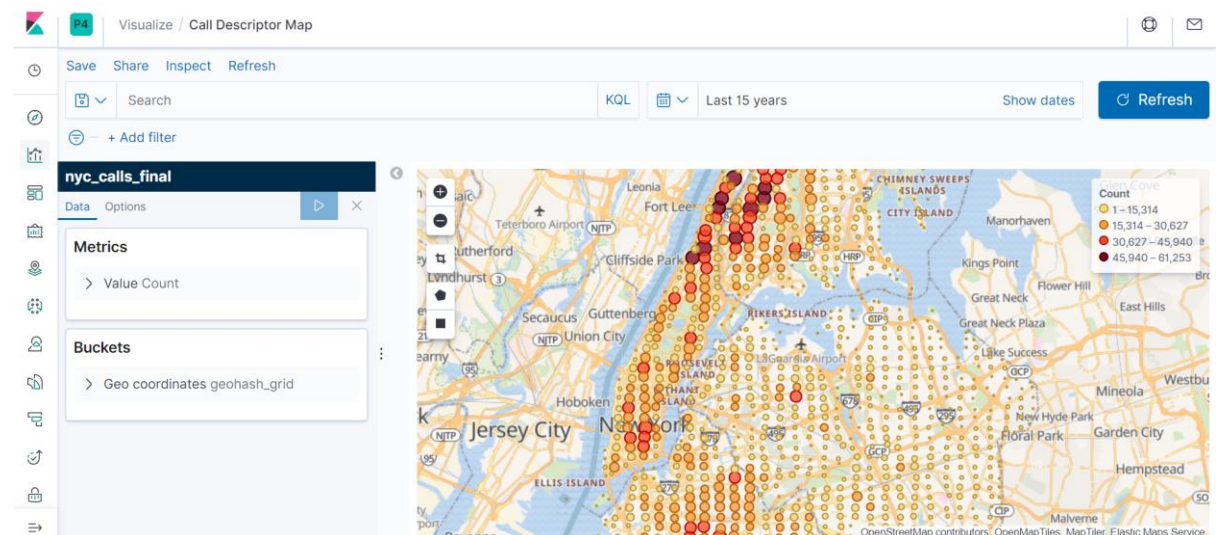
- Create a pie chart showing the top 5 cities with the highest calls alongside the top five calls (Descriptor) in each city.



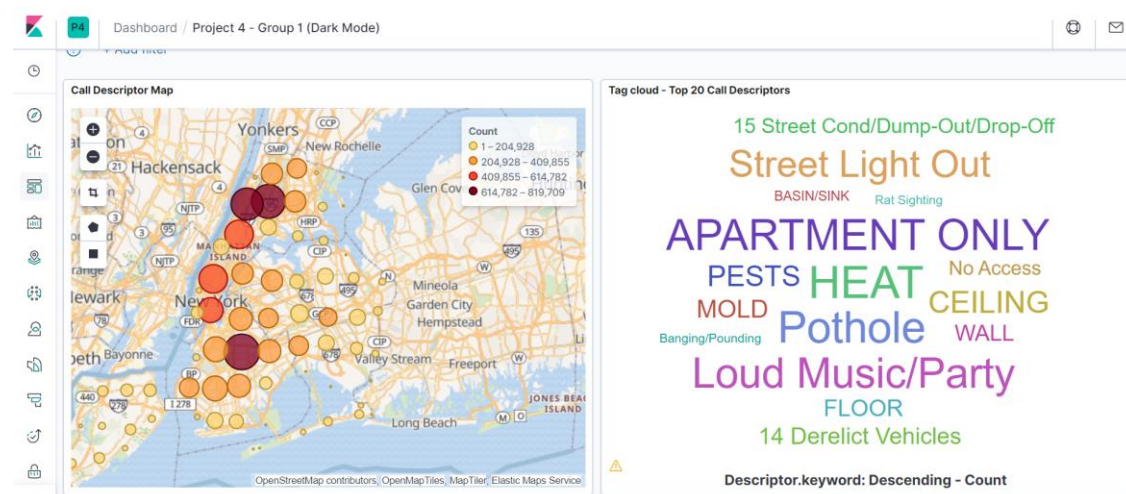
3. Create a tag cloud representing the top 20 call descriptors.

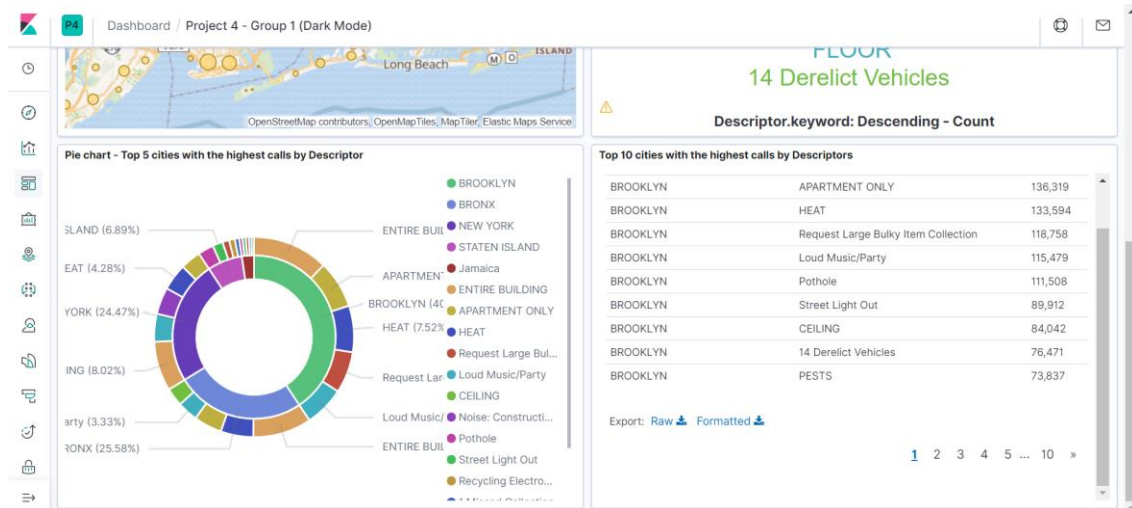


4. Create a coordinated map of all the major call descriptors in each city.



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





Thus, Project 4 is completed successfully.

---End---