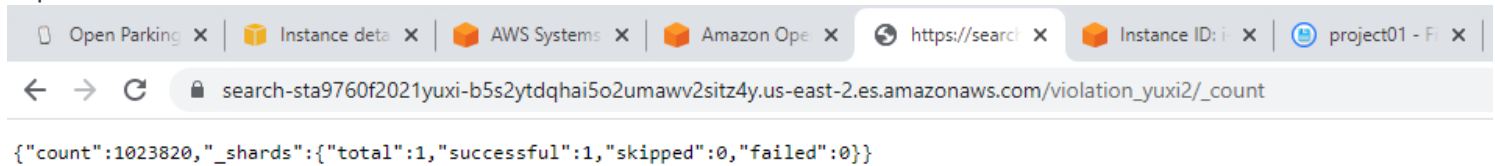


## CIS 9760 Fall 2021: Project 1 | Yuxi Song

I uploaded total 1m+ records into Elasticsearch:



```
{\"count\":1023820,\"_shards\":{\"total\":1,\"successful\":1,\"skipped\":0,\"failed\":0}}
```

### Steps to Build:

- Command prompt:
  - a. cd project01/
  - b. docker build -t project01:1.0 .
  - c. docker run project01:1.0 python src/main.py
  - d. docker run -v \$PWD:/app project01:1.0 python src/main.py
  - e. docker run \\\n-e DATASET\_ID=\"nc67-uf89\" \\\n-e APP\_TOKEN=\"Tcnb7vElccElonCJWpcHrQeNr\" \\\n-e ES\_HOST=\"https://search-sta9760f2021yuxi-b5s2ytdqhai5o2umawv2sitz4y.us-east-2.es.amazonaws.com\" \\\n-e INDEX\_NAME=\"violation\_yuxi2\" \\\n-e ES\_USERNAME=\"yuxi\" \\\n-e ES\_PASSWORD='Syxjzl1031##' \\\nproject01:1.0 --num\_pages=X --page\_size=Y
- o DATASET\_ID: the dataset identifier for Open Parking and Camera Violations
- o APP\_TOKEN: the unique key for accessing the API
- o ES\_HOST: the Elasticsearch host address
- o ES\_USERNAME: the Elasticsearch username
- o ES\_PASSWORD: the Elasticsearch password
- o page\_size: an optional argument to specify how many records to request from the API per call. If not provided, it defaults to 1,000 records per request
- o num\_pages: an optional argument to specify how many calls to make to the API. If not provided, it defaults to requesting all data in the dataset

### Command Line Arguments:

- o I used argparse library for handling of the command line argument, page\_size and num\_pages.

```
parser = argparse.ArgumentParser(description="Process data from violation.")
parser.add_argument("--page_size", type=int, help="How many rows to fetch per page", required=True)
parser.add_argument("--num_pages", type=int, help="How many pages to fetch")
args = parser.parse_args(sys.argv[1:])
print(args)
```
- o I used sodapy to import Socrata Open data API (To get dataset from Open Parking and Camera Violations).

### Elasticsearch Index and Mapping:

I wrote a function that sets the mapping for uploaded documents.

```
{
  "settings": {
    "number_of_shards": 1,
    "number_of_replicas": 1
  },
  "mappings": {
    "properties": {
      "plate": {"type": "text"},
      "state": {"type": "text"},
      "license_type": {"type": "text"},
      "summons_number": {"type": "text"},
      "issue_date": {"type": "date", "format": "MM/dd/yyyy"},
      "violation": {"type": "keyword"},
      "fine_amount": {"type": "float"},
      "penalty_amount": {"type": "float"},
      "interest_amount": {"type": "float"},
      "reduction_amount": {"type": "float"},
      "payment_amount": {"type": "float"},
      "amount_due": {"type": "float"},
      "precinct": {"type": "keyword"},
      "county": {"type": "keyword"}
    }
  }
}
```

## Usage of Elasticsearch:

I chose to use the Bulk API to upload a huge volume of data to Elasticsearch. I could upload to Elasticsearch the full page\_size of documents for each call to the API.

```
bulk_upload_data = ""
for i, line in enumerate(es_rows):
    print(f"Handling row {line['summons_number']} {i}")
    action = '{"index": {"_index": "' + INDEX_NAME + '", "_type": "_doc"}}'
    data = json.dumps(line)
    bulk_upload_data += f"{action}\n"
    bulk_upload_data += f"{data}\n"

print(bulk_upload_data)
```

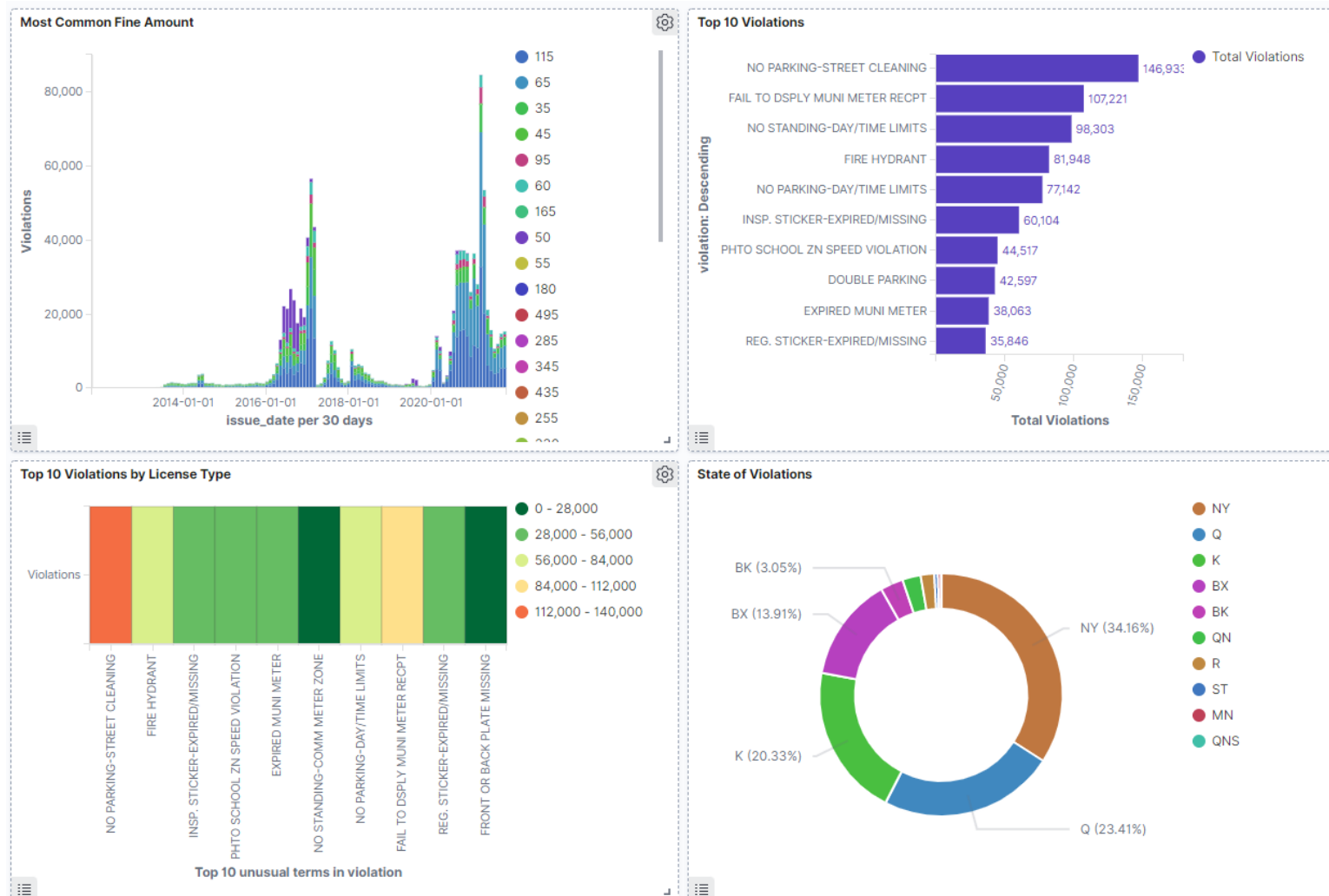
## Blog Post: Kibana Visualizations:

I built a dashboard. Timeframe: Last 10 years.

Total Number of Violations

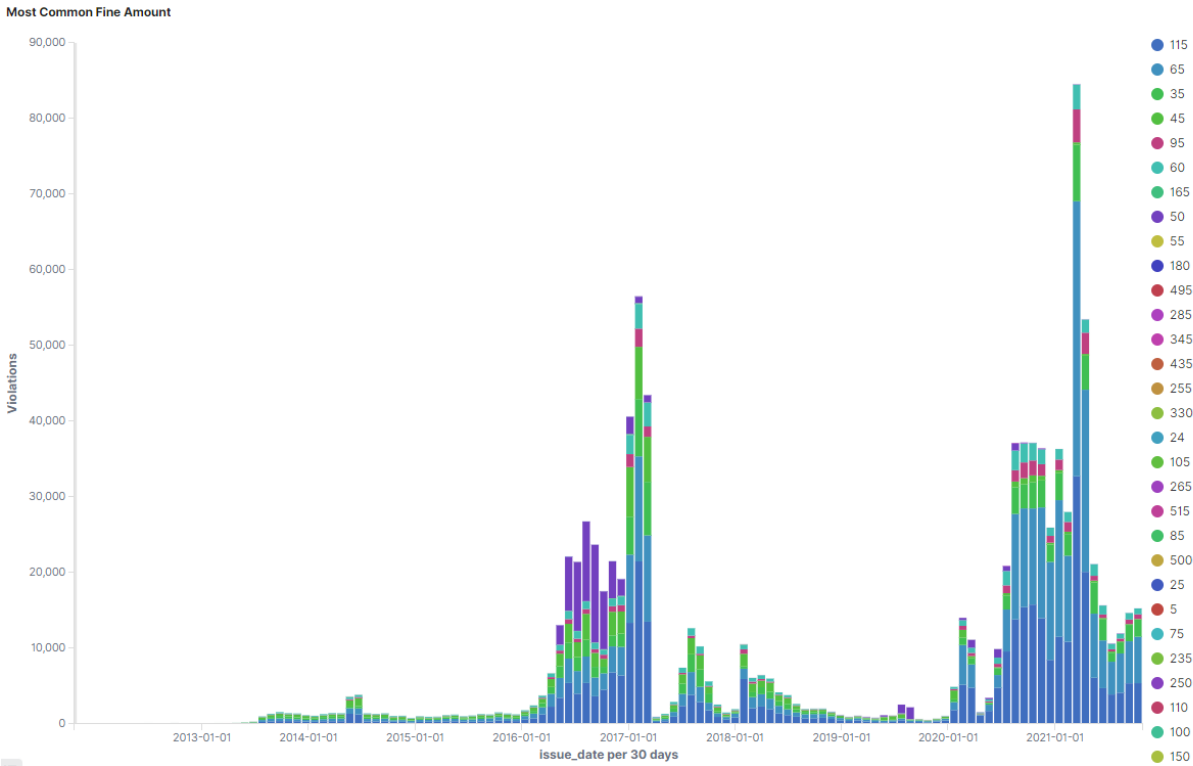
# 1,006,071

Violations

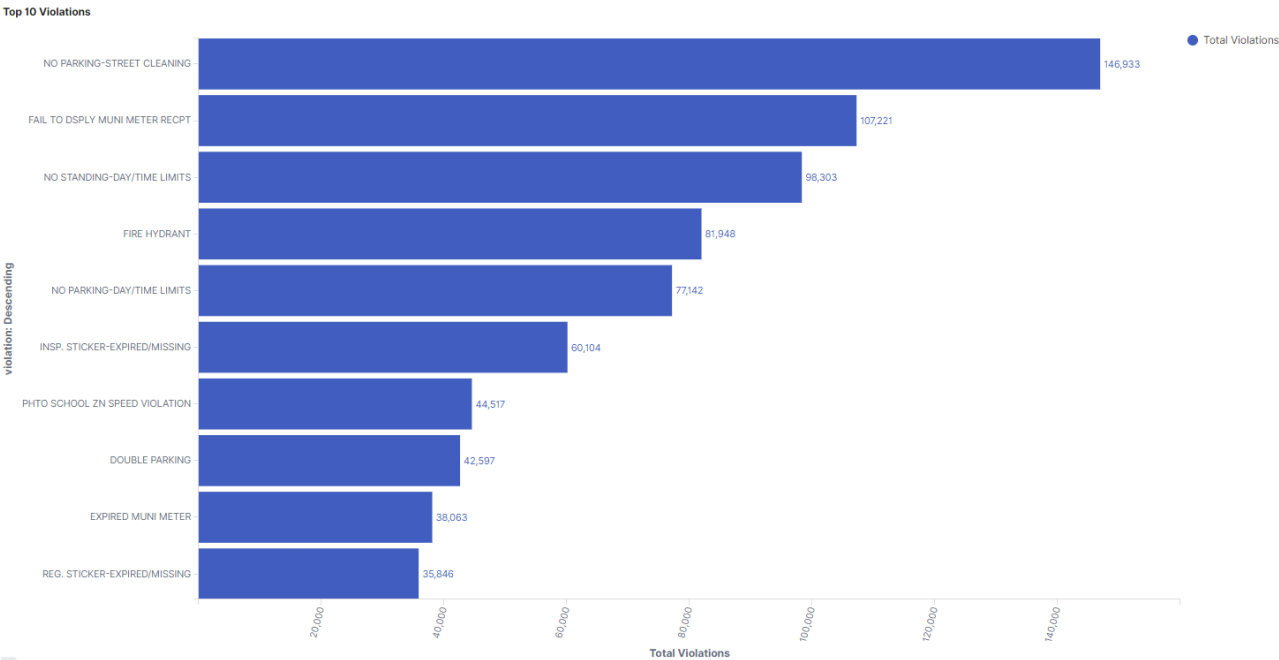


Visualizations:

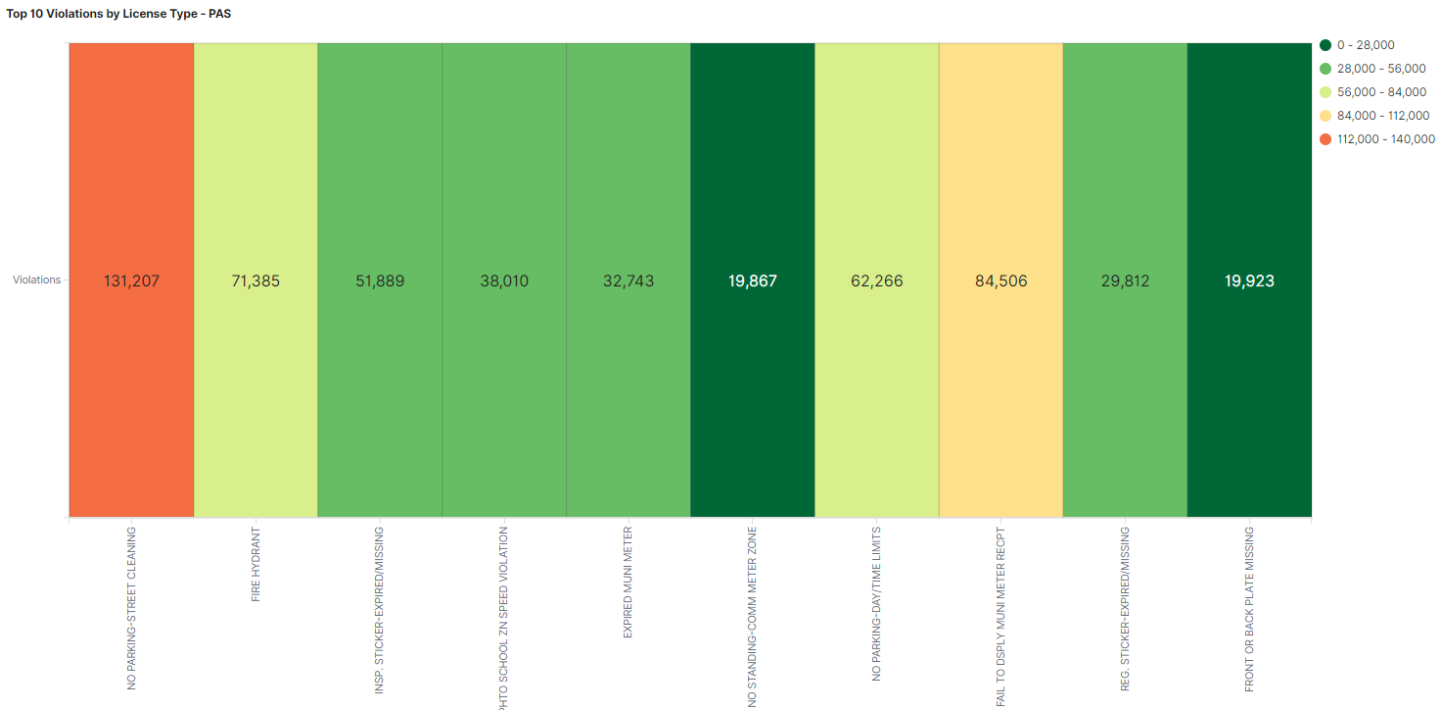
I had uploaded around 1 million records from the Open Parking and Camera Violation dataset to Elasticsearch. My first visualization displays fine\_amount within the last 10 years. According to the chart, the most common fine amount is \$115. It means most drivers have to pay more than \$100 dollars once they violate the rules. In other words, the fine amount is almost the same as some drivers’ daily earnings. If you want to save money, you would follow rules.



My second visualization shows the top 10 violations in the last 10 years. The most frequent violation is “No Parking-Street Cleaning” From this visualization, drivers need to check the date of street cleaning when they are parking, or they can move cars before street cleaning. If we are paying more attention, this type of violation could avoid.



More specifically, my third visualization shows the top 10 violations by license type “PAS”. License type “PAS” has the most violations due to No Parking – Street Cleaning. According to my second and third visualization, license type “PAS” violated 89%(131,207/146,933=89%) of Top 1 violation(“No Parking – Street Cleaning”) in the last 10 years.



My fourth visualization states that NY has the most violations. Everyone is in a hurry in NY, and a lot of people live and work in NY, so drivers can’t completely avoid parking violations. Therefore, we should be careful when we drive in NY.

State of Violations

