

Project: NYC Restaurant Inspection Data

Noé Breton et Reda Bourssouf

Installations de Elastic Search

On utilise docker compose et la commande `docker compose up -d` pour lancer elasticsearch et kibana.

```
version: '3.8'

services:
  elasticsearch:
    image: docker.elastic.co/elasticsearch/elasticsearch:8.4.2
    container_name: elasticsearch
    environment:
      - node.name=elasticsearch
      - cluster.name=es-docker-cluster
      - discovery.type=single-node
      - bootstrap.memory_lock=true
      - "ES_JAVA_OPTS=-Xms512m -Xmx512m"
      # Security settings for development
      - xpack.security.enabled=false
      - xpack.security.enrollment.enabled=false
      - xpack.security.http.ssl.enabled=false
      - xpack.security.transport.ssl.enabled=false
    ulimits:
      memlock:
        soft: -1
        hard: -1
    volumes:
      - es_data:/usr/share/elasticsearch/data
  ports:
    - "9200:9200"
    - "9300:9300"
  networks:
    - elastic
  healthcheck:
    test: ["CMD-SHELL", "curl -f http://localhost:9200/_cluster/health || exit 1"]
    interval: 30s
    timeout: 10s
    retries: 5

  kibana:
    image: docker.elastic.co/kibana/kibana:8.4.2
    container_name: kibana
    environment:
      - ELASTICSEARCH_HOSTS=http://elasticsearch:9200
```

—
PROF

```

- ELASTICSEARCH_USERNAME=kibana_system
- ELASTICSEARCH_PASSWORD=
ports:
- "5601:5601"
networks:
- elastic
depends_on:
  elasticsearch:
    condition: service_healthy
healthcheck:
  test: ["CMD-SHELL", "curl -f http://localhost:5601/api/status || exit 1"]
  interval: 30s
  timeout: 10s
  retries: 5

volumes:
  es_data:
    driver: local

networks:
  elastic:
    driver: bridge

```

Ingestion

On cherche à ingérer [NYC Restaurant Inspection Results](#).

La taille limite de fichier à ingérer est par défaut de 100b, malheureusement le csv à ingérer est de 123Mb, on change la limite à 150mb dans Advanced settings.

The top N most popular fields to show: 10

Maximum file upload size:
Sets the file size limit when importing files. The highest supported value for this setting is 1GB.
Default: 100MB
fileUpload:maxFileSize: 150MB

PROF

Filter editor suggest values:
Set this property to false to prevent the filter editor from suggesting values for fields.
 On

2. Questions

2.1 List all the neighborhoods in New York.

Par une aggregation sur la colonne BORO on récupère les neighborhoods unique en plus de leur nombre d'occurrence.

Request:

```

GET /restaurantny/_search
{
  "size": 0,
  "aggs": {

```

```

"unique_boro": {
    "terms": {
        "field": "BORO",
        "size": 100
    }
}
}

```

Response:

```

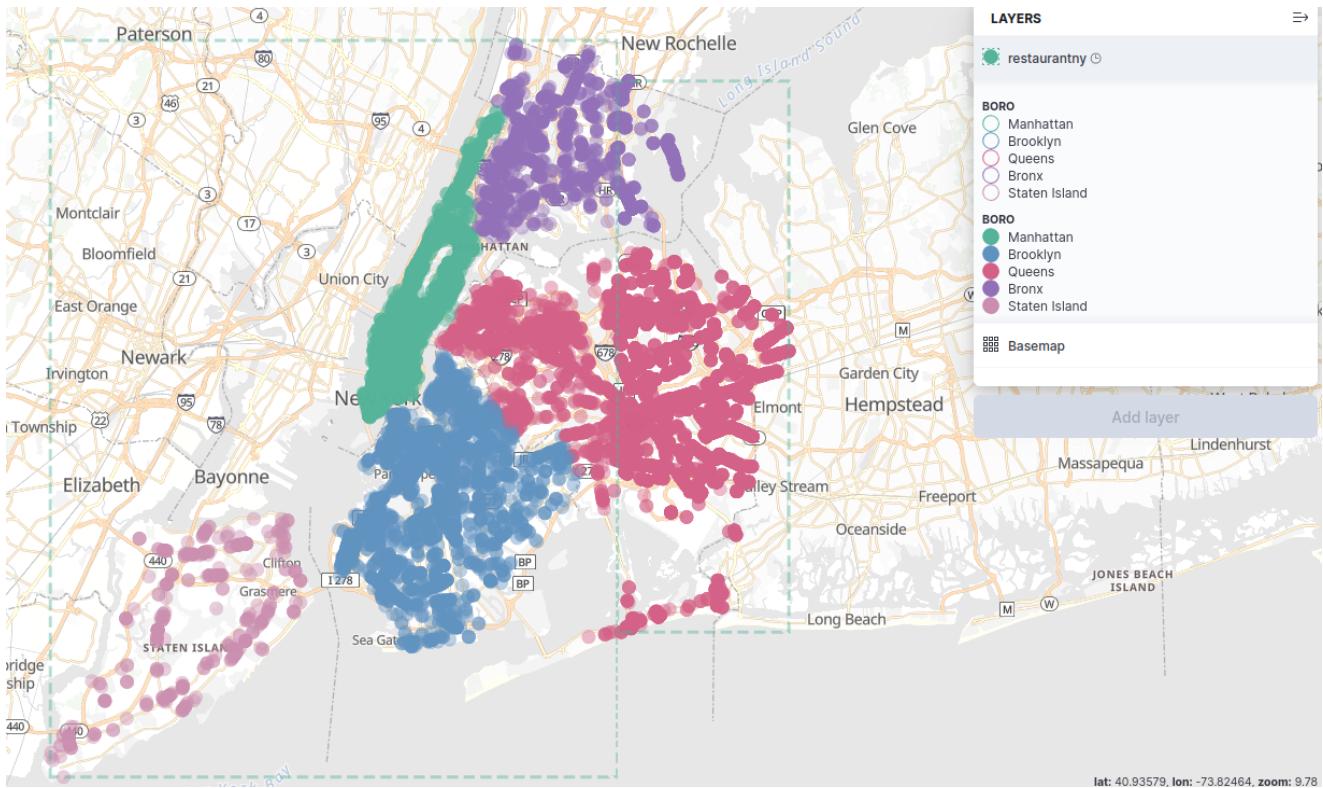
"aggregations": {
    "unique_boro": {
        "doc_count_error_upper_bound": 0,
        "sum_other_doc_count": 0,
        "buckets": [
            {
                "key": "Manhattan",
                "doc_count": 106988
            },
            {
                "key": "Brooklyn",
                "doc_count": 74803
            },
            {
                "key": "Queens",
                "doc_count": 71095
            },
            {
                "key": "Bronx",
                "doc_count": 26613
            },
            {
                "key": "Staten Island",
                "doc_count": 10071
            }
        ]
    }
}

```

—
PROF

Visalisation

Par cette visualisation on peut voir clairement tout les quartier de new york et les restauratant present a l'interieur.



2.2 Which neighborhood has the most restaurants ?

La commande precedente peut etre reprise pour recuperer le quartier avec le plus de restaurant en verifiant que l'ordonnement est bien decroissant et en ajoutant size=1 pour obtenir seulement le premier resultat. On utilise une deuxiemme aggregation de cardinalité sur l'id unique du restaurant pour obtenir le nombre de restaurant.

Request:

```
GET restaurantny/_search
{
  "size": 0,
  "aggs": {
    "restaurants_by_boro": {
      "terms": {
        "field": "BORO",
        "size": 1,
        "order": { "unique_restaurants.value": "desc" }
      },
      "aggs": {
        "unique_restaurants": {
          "cardinality": {
            "field": "CAMIS"
          }
        }
      }
    }
  }
}
```

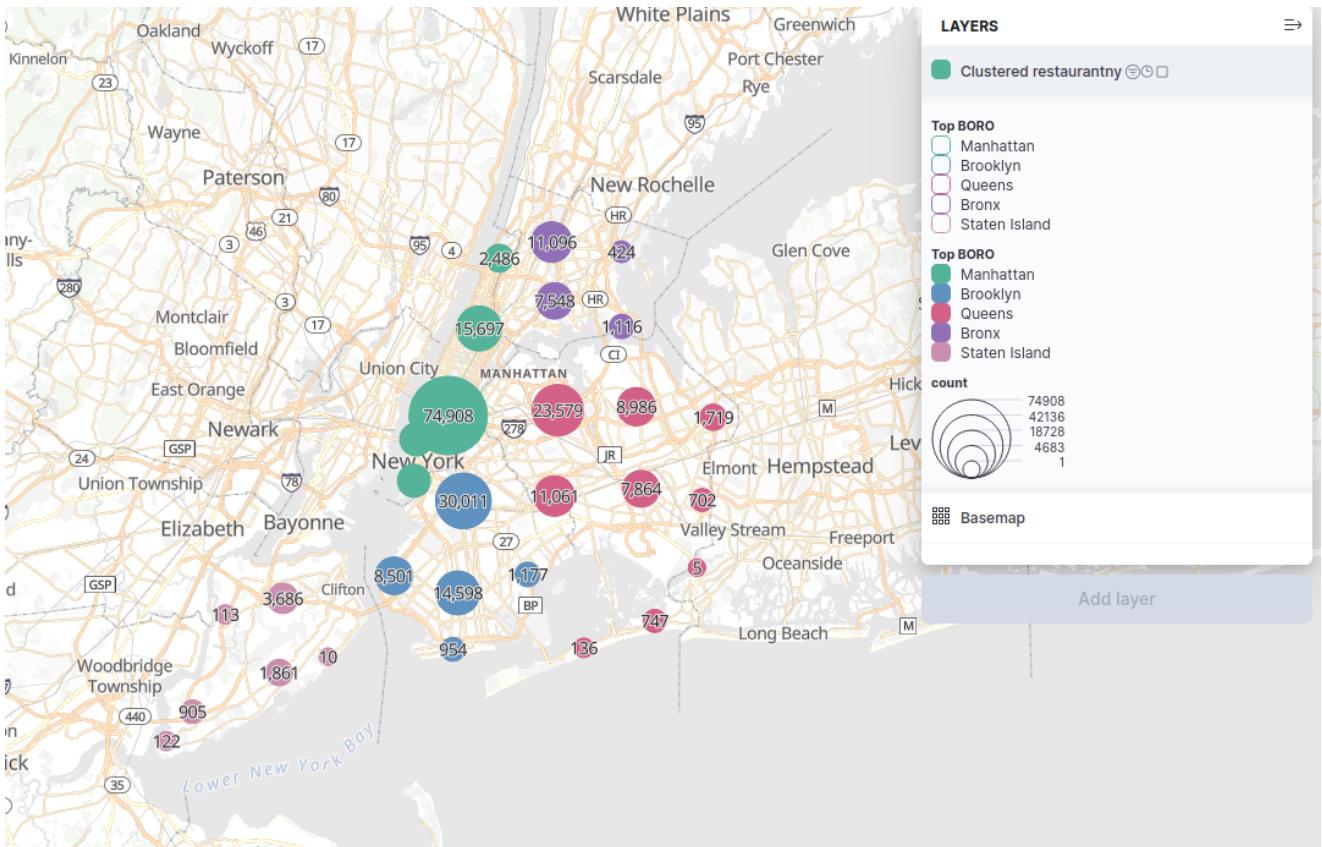
Resultat:

```
"aggregations": {  
    "restaurants_by_boro": {  
        "doc_count_error_upper_bound": 0,  
        "sum_other_doc_count": 182582,  
        "buckets": [  
            {  
                "key": "Manhattan",  
                "doc_count": 106988,  
                "unique_restaurants": {  
                    "value": 12113  
                }  
            }  
        ]  
    }  
}
```

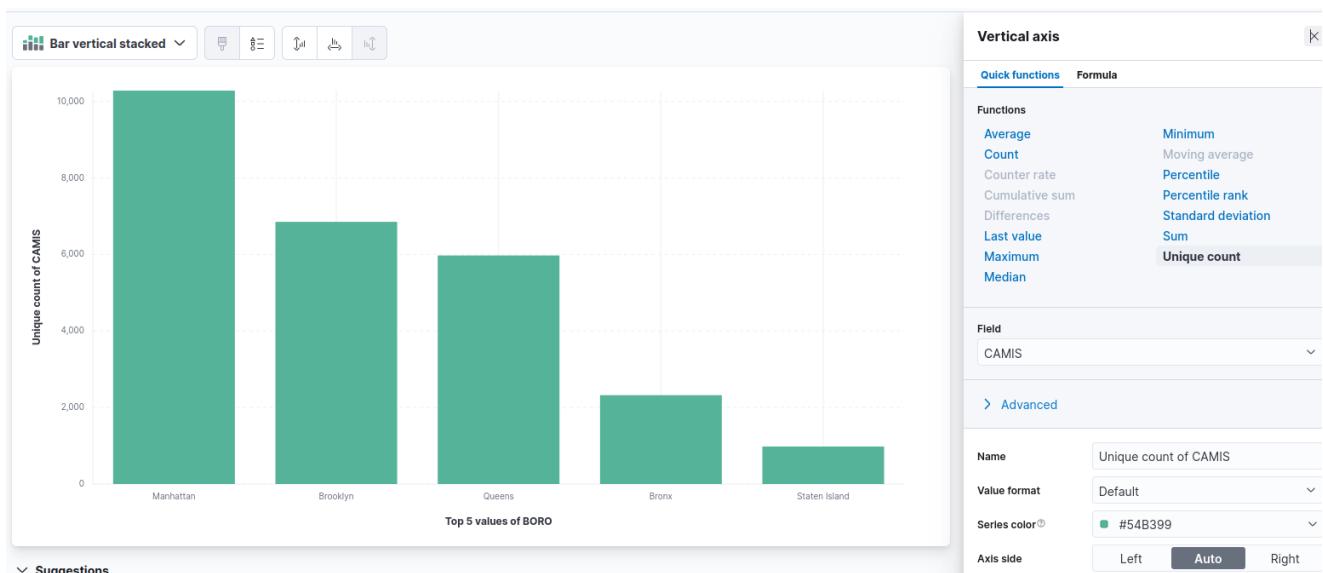
Manhattan possède 12113 restaurant, faisant de ce quartier celui avec la plus grande offre de restauration.

visualitation

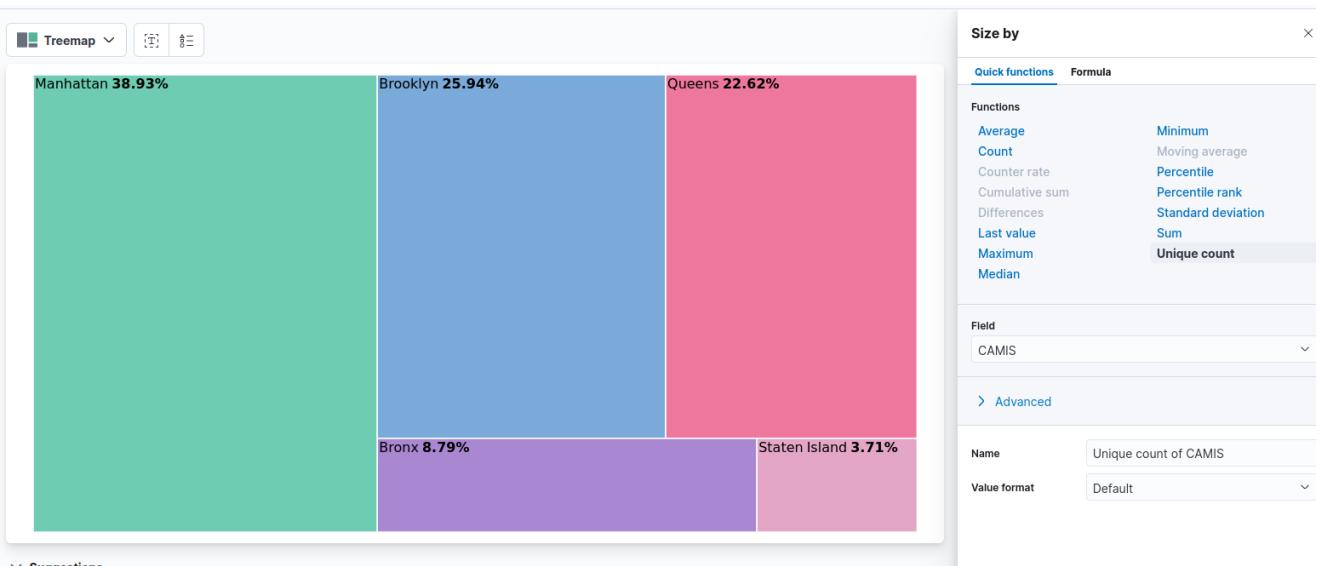
Resprendre la visualition precedente en ajoutant un label permet de se rendre compte approximativement du nombre de restaurant par quartier, mais cela manque de lisibilité



Comme la heatmap, on se rend compte de la plus grande densité de Manhattan mais ça reste moins évident qu'un barchart, on met en unique count CAMIS pour avoir seulement les restaurants.



On peut apprécier la répartition des restaurants de New York par BORO



▼ Suggestions

2.3 What does the violation code "04N" correspond to ?

On extrait l'attribut VIOLATION DESCRIPTION (et VIOLATION CODE histoire d'être sur) d'une ligne dont l'attribut VIOLATION CODE a comme valeur 04N.

Requete

```
GET restaurantny/_search
{
  "_source": ["VIOLATION DESCRIPTION", "VIOLATION CODE"],
  "query": {
    "term": {
      "VIOLATION CODE": "04N"
    }
  },
  "size": 1
}
```

PROF

Réultat

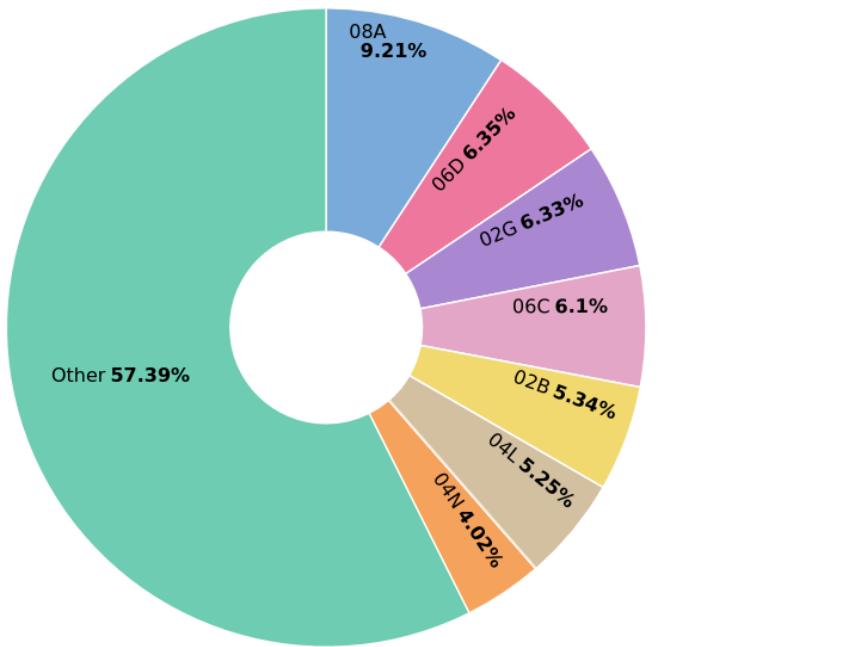
```
"hits": [
  {
    "_index": "restaurantny",
    "_id": "rBOH3JkByv84jpscZM78",
    "_score": 3.1635146,
    "_source": {
      "VIOLATION CODE": "04N",
      "VIOLATION DESCRIPTION": "Filth flies or food/refuse/sewage associated with (FRSA) flies or other nuisance pests in establishment's food and/or non-food areas. FRSA flies include house flies, blow flies, bottle flies, flesh flies, drain flies, Phorid flies and fruit flies."
    }
  }
]
```

```
        ]
    }
}
```

Le code 04N correspond à présence de mouches ou d'autres insectes dans les zones de préparation, de stockage ou de service des aliments.

viusalitation

Voici la proportion des codes de violation présent dans les restaurants enregistrés, on voit que 10F est le code de violation le plus présent, incluant les lignes qui ne présentent pas de code.



2.4 Where are the restaurants (name, address, neighborhood) that have a grade of A?

PROF

```
GET restaurantny/_search
{
  "_source": ["DBA", "BUILDING", "STREET", "ZIPCODE", "BORO", "GRADE"] ,
  "query": {
    "term": {
      "GRADE": "A"
    }
  }, "size": 100
}
```

response :

```
"hits": [
  {
```

```
        "_index": "restaurantny",
        "_id": "rROH3JkByv84jpscZM78",
        "_score": 0.3891273,
        "_source": {
            "DBA": "PRESSED JUICERY",
            "BUILDING": "2857",
            "BORO": "Manhattan",
            "ZIPCODE": 10025,
            "GRADE": "A",
            "STREET": "BROADWAY"
        }
    },
    {
        "_index": "restaurantny",
        "_id": "sBOH3JkByv84jpscZM78",
        "_score": 0.3891273,
        "_source": {
            "DBA": "KPOT KOREAN BBQ & HOT POT (ON 2ND FL, BAY PLAZA
MALL)",
            "BUILDING": "200",
            "BORO": "Bronx",
            "ZIPCODE": 10475,
            "GRADE": "A",
            "STREET": "BAYCHESTER AVENUE"
        }
    },
    {
        "_index": "restaurantny",
        "_id": "shOH3JkByv84jpscZM78",
        "_score": 0.3891273,
        "_source": {
            "DBA": "BLUE MOUNTAIN REST & BAKERY",
            "BUILDING": "959",
            "BORO": "Brooklyn",
            "ZIPCODE": 11226,
            "GRADE": "A",
            "STREET": "FLATBUSH AVENUE"
        }
    },
    {
        "_index": "restaurantny",
        "_id": "tBOH3JkByv84jpscZM78",
        "_score": 0.3891273,
        "_source": {
            "DBA": "Bronx Slice",
            "BUILDING": "37",
            "BORO": "Bronx",
            "ZIPCODE": 10454,
            "GRADE": "A",
            "STREET": "BRUCKNER BOULEVARD"
        }
    },
    {
```

—
PROF

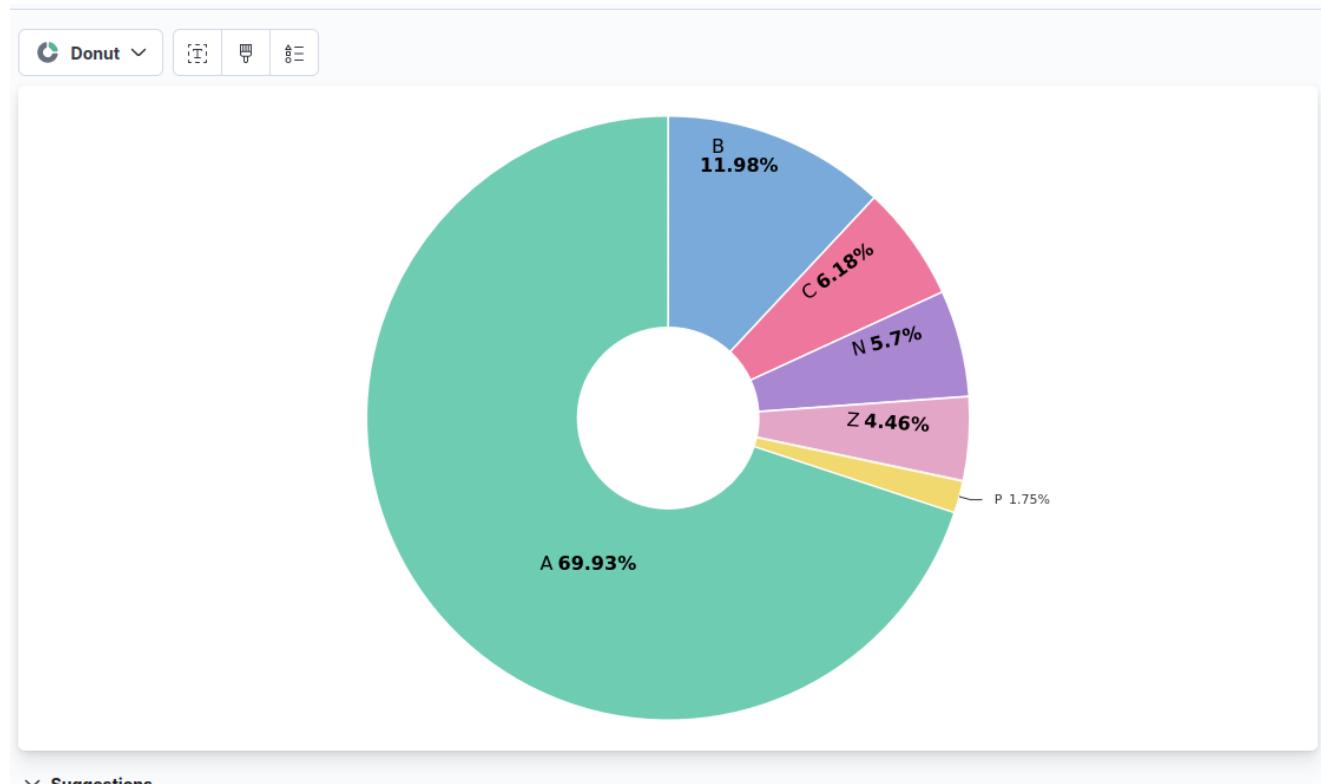
```

    "_index": "restaurantny",
    "_id": "txOH3JkByv84jpscZM78",
    "_score": 0.3891273,
    "_source": {
        "DBA": "CJ DIAMOND CAFE",
        "BUILDING": "4102",
        "BORO": "Queens",
        "ZIPCODE": 11355,
        "GRADE": "A",
        "STREET": "COLLEGE POINT BLVD"
    }
}, ...

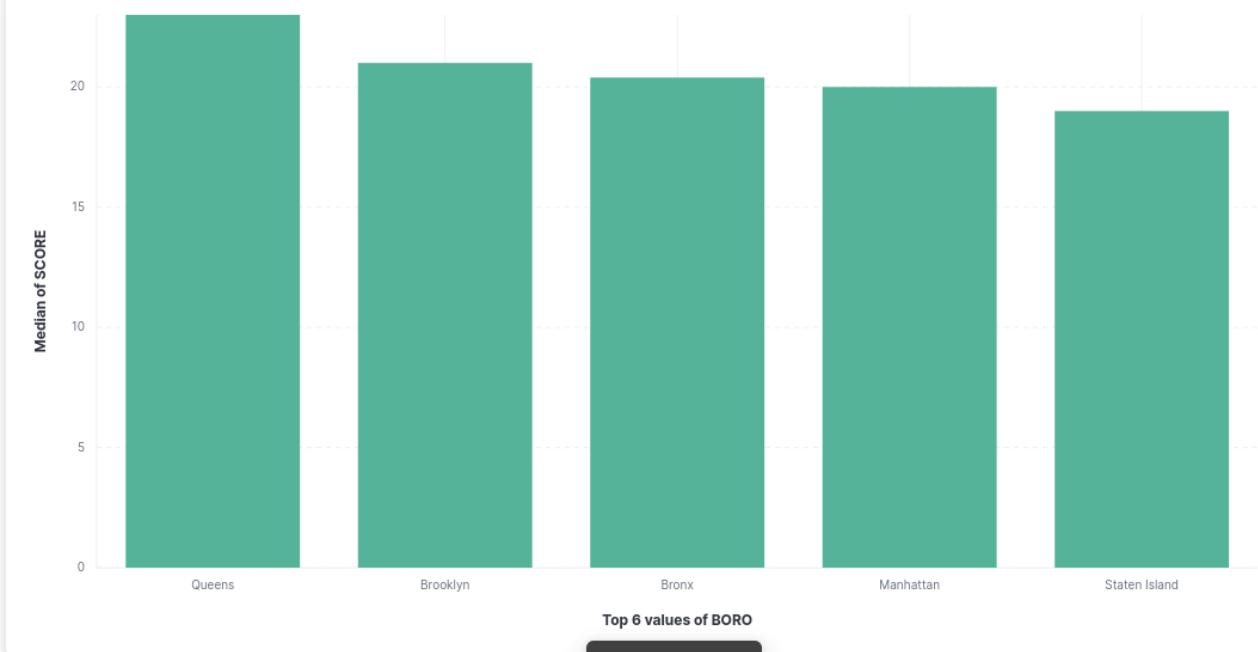
```

Visualisation

L'outil de visualisation n'autorise pas d'utiliser DBA comme une metrique, je choisis de visualisé la proportion des notes parmis tout les restaurant, les outils dynamique du dashbord permetteront d'ajuster la granulité.



On peut aussi ajouter la medianne des score par BORO, on decouvre que le Queens a la meilleure medianne parmis les BORO



2.5 What is the most popular cuisine? And by neighborhood?

```

GET restaurantny_final/_search
{
  "size": 0,
  "aggs": {
    "top3_Description": {
      "terms": {
        "field": "CUISINE DESCRIPTION",
        "size": 3,
        "order": { "unique_restaurants.value": "desc" }
      },
      "aggs": {
        "unique_restaurants": {
          "cardinality": {
            "field": "CAMIS"
          }
        }
      }
    }
  }
}

```

—
PROF

response:

```

"aggregations": {
  "top3_Description": {
    "doc_count_error_upper_bound": -1,
    "sum_other_doc_count": 192398,
    "buckets": [
      {
        "key": "American (New) (192398),"
        "doc_count": 192398
      }
    ]
  }
}

```

```

"buckets": [
    {
        "key": "American",
        "doc_count": 45070,
        "unique_restaurants": {
            "value": 4929
        }
    },
    {
        "key": "Chinese",
        "doc_count": 28142,
        "unique_restaurants": {
            "value": 2154
        }
    },
    {
        "key": "Coffee/Tea",
        "doc_count": 20167,
        "unique_restaurants": {
            "value": 2100
        }
    }
]
}
}

```

La cuisine la plus populaire de new york par le nombre de restaurant est la cuisine américaine, avec 4964 restaurant suivit par la cuisine chinoise et les coffee/tea shop.

Par BORO:

```

GET restaurantny_final/_search
{
    "size": 0,
    "aggs": {
        "by_boro": {
            "terms": {
                "field": "BORO",
                "size": 10
            },
        "aggs": {
            "top3_Description": {
                "terms": {
                    "field": "CUISINE DESCRIPTION",
                    "size": 1,
                    "order": { "unique_restaurants.value": "desc" }
                },
            "aggs": {
                "unique_restaurants": {

```

PROF

reponse:

```
"aggregations": {
  "by_boro": {
    "doc_count_error_upper_bound": 0,
    "sum_other_doc_count": 0,
    "buckets": [
      {
        "key": "Manhattan",
        "doc_count": 106988,
        "top3_Description": {
          "doc_count_error_upper_bound": -1,
          "sum_other_doc_count": 82973,
          "buckets": [
            {
              "key": "American",
              "doc_count": 22371,
              "unique_restaurants": {
                "value": 2476
              }
            }
          ]
        }
      }
    ]
  },
  {
    "key": "Brooklyn",
    "doc_count": 74803,
    "top3_Description": {
      "doc_count_error_upper_bound": -1,
      "sum_other_doc_count": 63742,
      "buckets": [
        {
          "key": "American",
          "doc_count": 10057,
          "unique_restaurants": {
            "value": 1094
          }
        }
      ]
    }
  }
}
```

```
},
{
  "key": "Queens",
  "doc_count": 71095,
  "top3_Description": {
    "doc_count_error_upper_bound": -1,
    "sum_other_doc_count": 62322,
    "buckets": [
      {
        "key": "American",
        "doc_count": 7988,
        "unique_restaurants": {
          "value": 854
        }
      }
    ]
  }
},
{
  "key": "Bronx",
  "doc_count": 26613,
  "top3_Description": {
    "doc_count_error_upper_bound": -1,
    "sum_other_doc_count": 23433,
    "buckets": [
      {
        "key": "American",
        "doc_count": 2948,
        "unique_restaurants": {
          "value": 317
        }
      }
    ]
  }
},
{
  "key": "Staten Island",
  "doc_count": 10071,
  "top3_Description": {
    "doc_count_error_upper_bound": -1,
    "sum_other_doc_count": 8237,
    "buckets": [
      {
        "key": "American",
        "doc_count": 1706,
        "unique_restaurants": {
          "value": 164
        }
      }
    ]
  }
}
```

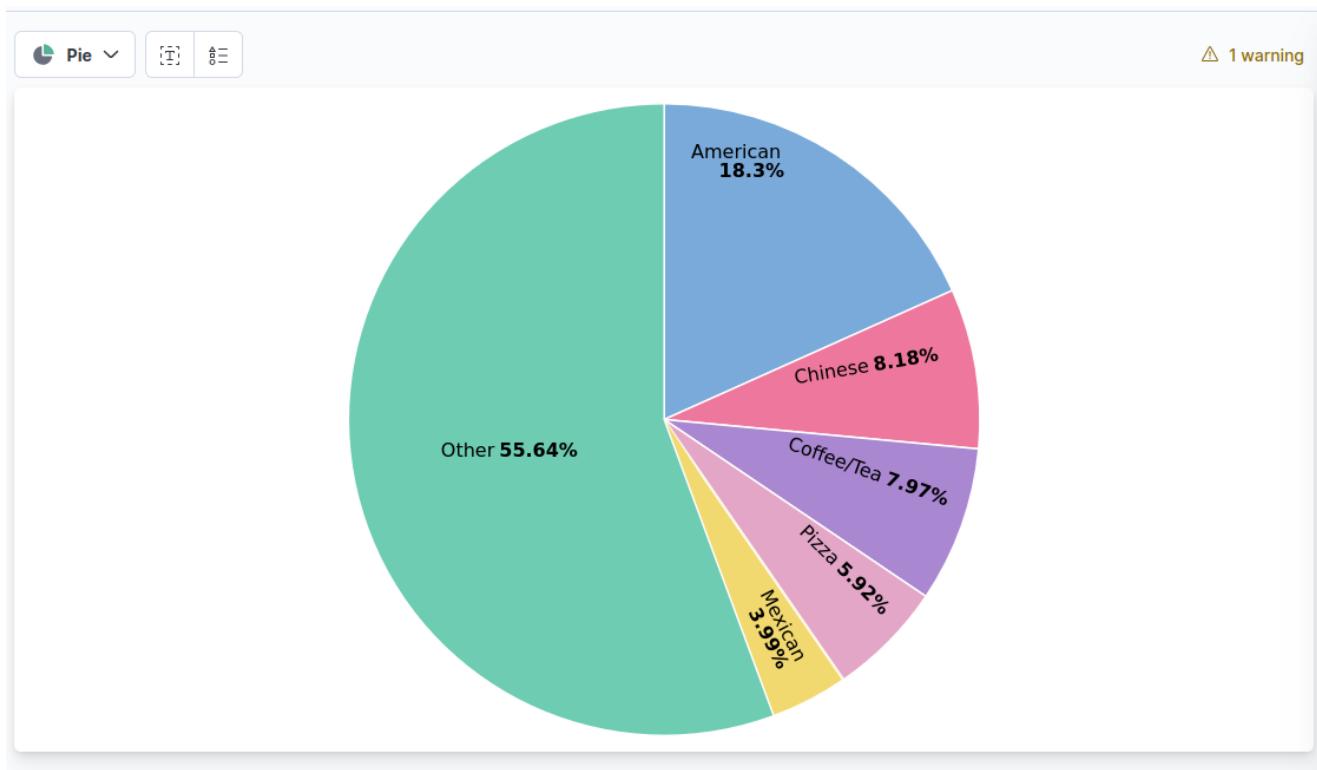
—
PROF

par BORO, la cuisine la plus populaire est

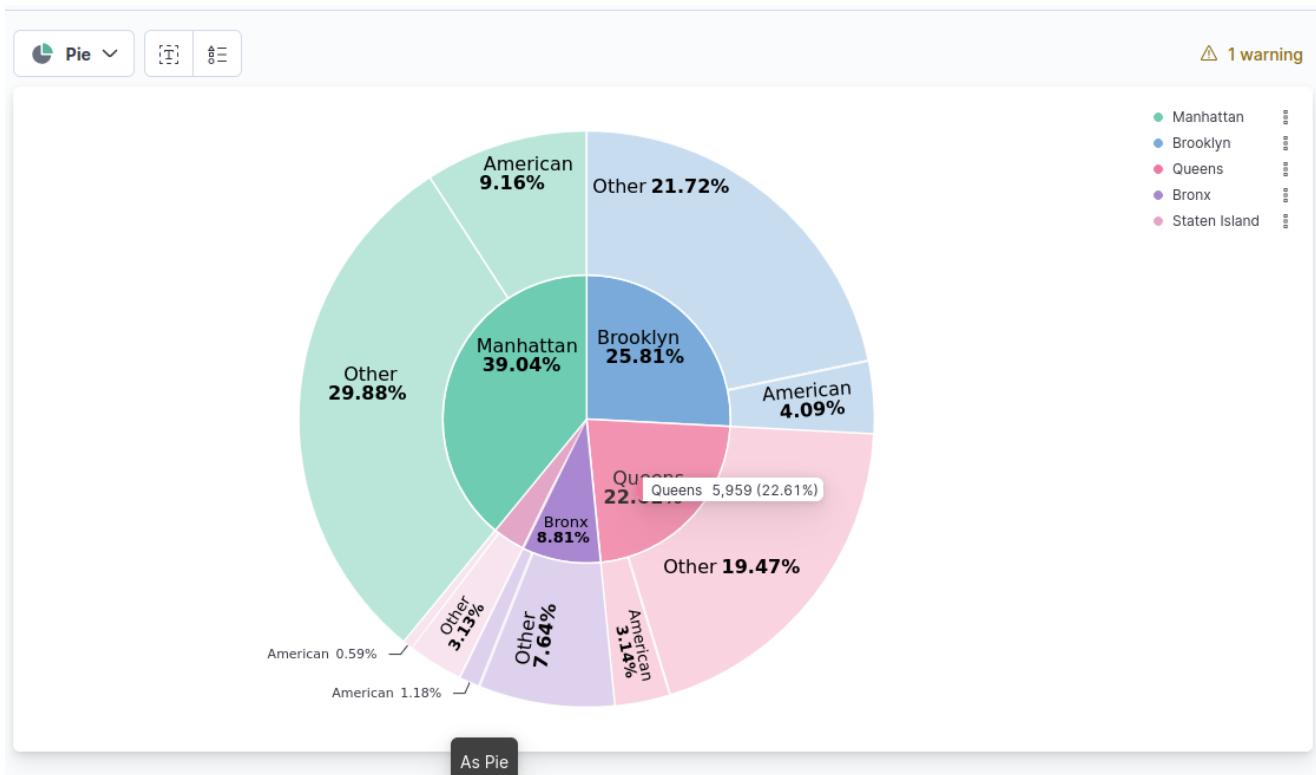
BORO	Cuisine la plus populaire	Nombre de restaurants
Manhattan	American	2476
Brooklyn	American	1094
Queens	American	854
Bronx	American	317
Staten Island	American	164

visualisation

Pour New-york, on peut utiliser un pie chart classique, ou un waffle chart (la visibilité reste moins bonne)

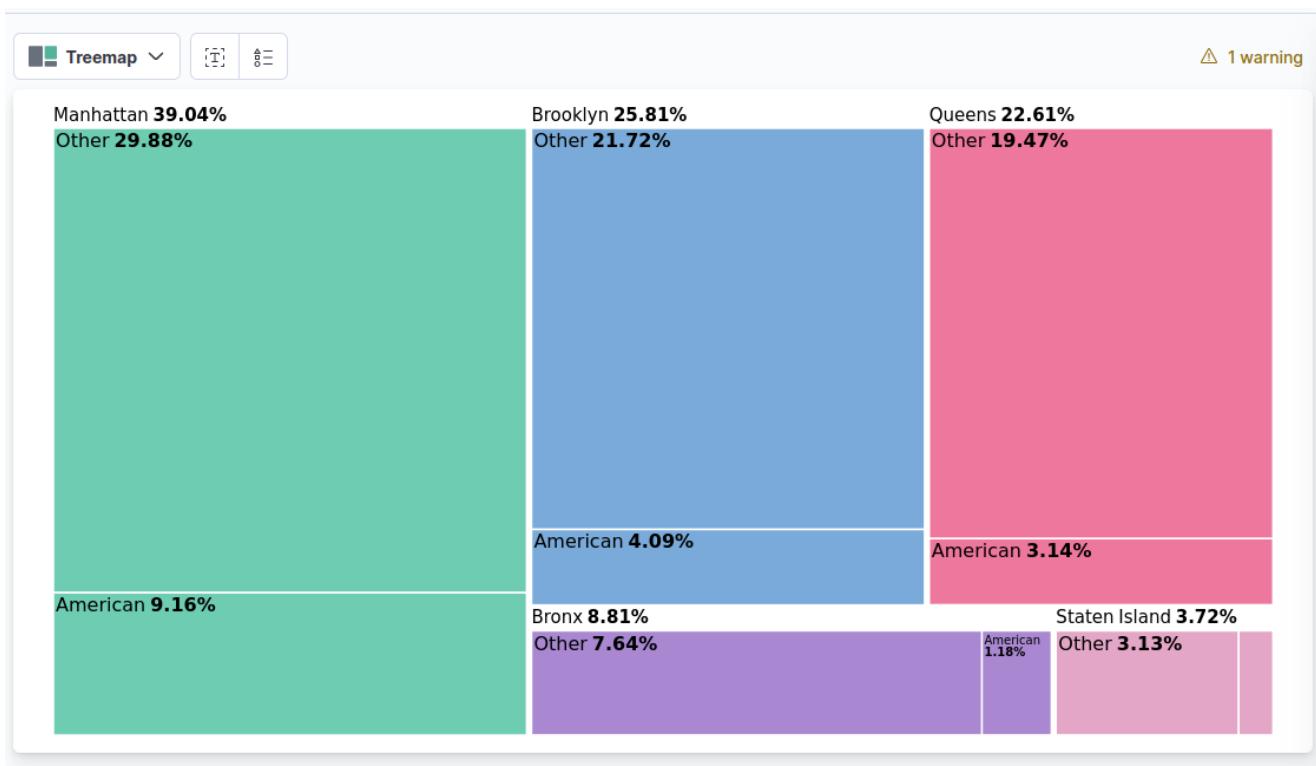


On peut utiliser un double pie chart ou une mosaique pour voir la cuisine la plus populaire par BORO



▼ Suggestions

⟳ Refresh



PROF

▼ Suggestions

2.6 What is the date of the last inspection?

```
GET restaurantny/_search
{
  "_source": ["INSPECTION DATE", "DBA"],
  "sort": [
    { "INSPECTION DATE": { "order": "desc" } }
```

```
],
  "size": 1
}
```

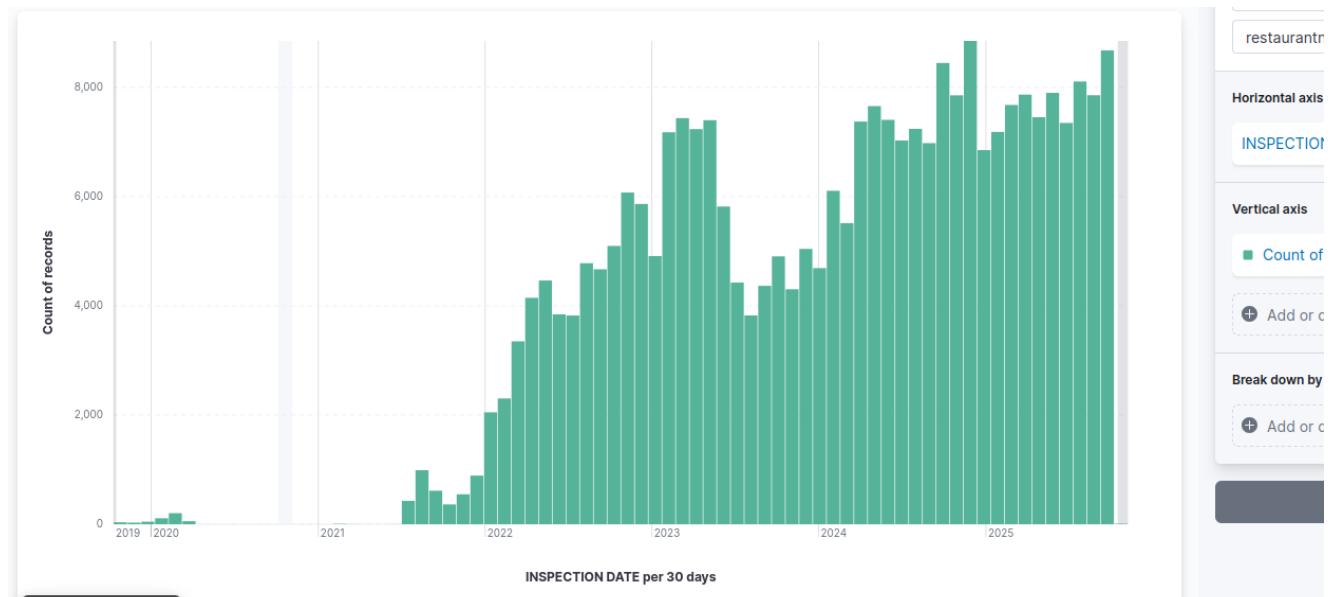
reponse :

```
"hits": [
  {
    "_index": "restaurantny",
    "_id": "UBOH3JkByv84jpscW101",
    "_score": null,
    "_source": {
      "DBA": "LE PETIT MONSTRE",
      "INSPECTION DATE": "10/09/2025"
    },
    "sort": [
      1759968000000
    ]
  }
]
```

la derniere inspection date du 10/09/2025 et concerne LE PETIT MONTRE.

visulation

On peut utiliser un diagramme temporel en barre ou en ligne pour connaitre l'évolution du nombre d'inspection en fonction du temps



2.7 Provide a list of Chinese restaurants with an A grade in Brooklyn.

On utilise la commande bool pour trouver les resaurant qui doivent (must) remplir les condition "CUISINE DESCRIPTION": "Chinese", "GRADE": "A" et "BORO": "Brooklyn".

```

GET restaurantny/_search
{
  "_source": ["DBA"],
  "query": {
    "bool": {
      "must": [
        { "term": { "CUISINE DESCRIPTION": "Chinese" } },
        { "term": { "GRADE": "A" } },
        { "term": { "BORO": "Brooklyn" } }
      ]
    }
  }
}

```

resultat:

```

"hits": [
  {
    "_index": "restaurantny",
    "_id": "rBOH3JkByv84jpscZM_8",
    "_score": 4.0605974,
    "_source": {
      "DBA": "NEW CENTURY CHINESE RESTAURANT"
    }
  },
  {
    "_index": "restaurantny",
    "_id": "HhOH3JkByv84jpscZND8",
    "_score": 4.0605974,
    "_source": {
      "DBA": "NEW GREAT WALL 1419"
    }
  },
  {
    "_index": "restaurantny",
    "_id": "nhOH3JkByv84jpscZND8",
    "_score": 4.0605974,
    "_source": {
      "DBA": "FOOD LOVER BAKERY"
    }
  },
  {
    "_index": "restaurantny",
    "_id": "UROH3JkByv84jpscZNH8",
    "_score": 4.0605974,
    "_source": {
      "DBA": "SUN GARDEN"
    }
  }
]

```

—
PROF

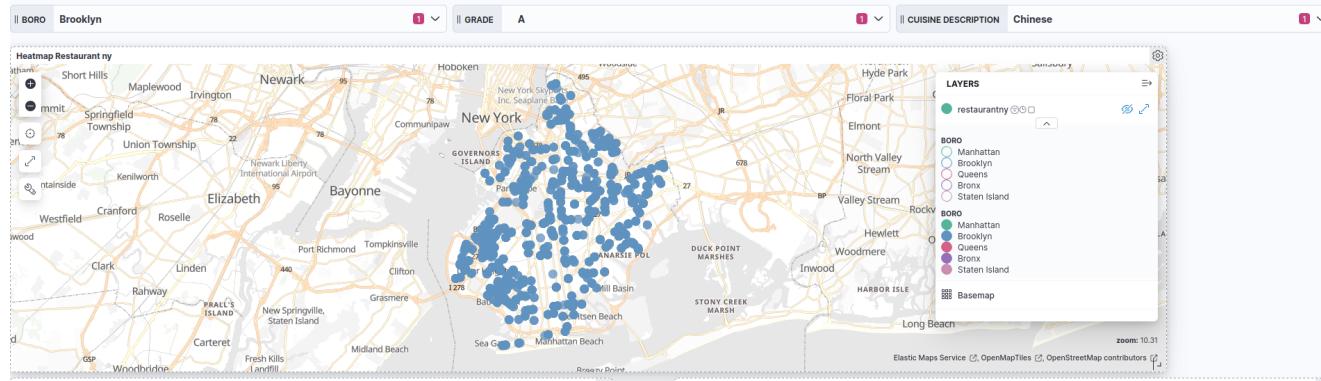
```

    "_index": "restaurantny",
    "_id": "WROH3JkByv84jpscZNH8",
    "_score": 4.0605974,
    "_source": {
        "DBA": "GRAND PANDA"
    }
},
{ ...

```

Visualisation

On peut utiliser les contrôles pour obtenir une liste des restaurants avec ces critères



2.8 Provide the address of the restaurant LADUREE

en faisant un match, on découvre qu'il y a plusieurs restaurants LADUREE (laduree, LADUREE soho et LADUREE)

```

History Settings Variables Help
1  r
2  GET restaurantny/_search
3  {
4      "query": {
5          "term": {
6              "field": "CUISINE DESCRIPTION",
7              "size": 1000,
8              "order": { "count": "desc" }
9          }
10     }
11  }
12  "size": 1
13  }
14  }
15  GET restaurantny/_search
16  {
17      "size": 0,
18      "query": {
19          "match": {
20              "DBA": "LADUREE"
21          }
22      },
23      "aggs": {
24          "unique_cuisines": {
25              "terms": {
26                  "field": "CUISE DESCRIPTION",
27                  "size": 1000
28              }
29          },
30          "aggs": {
31              "sample_doc": {
32                  "top_hits": {
33                      "size": 1
34                  }
35              }
36          }
37      }
38  }
39  "sort": [
40      {"INSPECTION DATE": { "order": "desc" } }
41  ],
42  "size": 1
43  }
44  }
45  GET restaurantny/_search
46  {
47      "size": 1000,
48      "query": {
49          "match": {
50              "DBA": "LADUREE"
51          }
52      },
53      "aggs": {
54          "unique_cuisines": {
55              "terms": {
56                  "field": "CUISE DESCRIPTION",
57                  "size": 1000
58              }
59          },
60          "aggs": {
61              "sample_doc": {
62                  "top_hits": {
63                      "size": 1
64                  }
65              }
66          }
67      }
68  }
69  "sort": [
70      {"INSPECTION DATE": { "order": "desc" } }
71  ],
72  "size": 1
73  }
74  }
75  }
76  }
77  }
78  }
79  }
80  }
81  }
82  }
83  }
84  }
85  }
86  }
87  }
88  }
89  }
90  }
91  }
92  }
93  }
94  }
95  }
96  }
97  }
98  }
99  }
100 }
101 }
102 }
103 }
104 }
105 }
106 }
107 }
108 }
109 }
110 }
111 }
112 }
113 }
114 }
115 }
116 }
117 }
118 }
119 }

200 - OK | 13 ms

```

On essaie une requête pour avoir le terme exact LADUREE. mais cela ne marche toujours pas, comme si DBA ne marchait pas avec une recherche exacte. En inspectant le mapping on s'aperçoit que DBA n'a pas d'option "keyword" qui permet une recherche exacte (contrairement à GRADE par exemple). On indexe keyword à DBA dans un nouveau index.

mapping de GRADE

```
{
  "restaurantny": {
    "mappings": {
      "GRADE": {
        "full_name": "GRADE",
        "mapping": {
          "GRADE": {
            "type": "keyword"
          }
        }
      }
    }
  }
}
```

mapping de DBA

```
{
  "restaurantny": {
    "mappings": {
      "DBA": {
        "full_name": "DBA",
        "mapping": {
          "DBA": {
            "type": "text"
          }
        }
      }
    }
  }
}
```

PROF

ajout de keyword a DBA a restaurant NY v1

```
PUT restaurantny_v1
{
  "mappings": {
    "properties": {
      "DBA": {
        "type": "text",
        "fields": {
          "keyword": { "type": "keyword" }
        }
      }
    }
  }
}
```

reindexation

```
POST _reindex
{
  "source": { "index": "restaurantny" },
  "dest": { "index": "restaurantny_v1" }
}
```

On ressaye la requete, on aggress sur le numero CAMIS pour être sur d'avoir la liste des restaurant unique en fonction de leur ID appelé LADUREE, on met la size de la seconde aggregation à 1 pour ne pas avoir des doublons en cas de multiple inspections.

```
{
  "size": 0,
  "query": {
    "term": {
      "DBA.keyword": "LADUREE"
    }
  },
  "aggs": {
    "unique_CAMIS": {
      "terms": {
        "field": "CAMIS",
        "size": 100
      },
      "aggs": {
        "sample_doc": {
          "top_hits": {
            "_source": ["DBA", "BUILDING", "STREET", "BORO"],
            "size": 1
          }
        }
      }
    }
  }
}
```

PROF

resultat

```
{
  "took": 0,
  "timed_out": false,
  "_shards": {
    "total": 1,
    "successful": 1,
    "skipped": 0,
```

```

    "failed": 0
  },
  "hits": {
    "total": {
      "value": 6,
      "relation": "eq"
    },
    "max_score": null,
    "hits": []
  },
  "aggregations": {
    "unique_CAMIS": {
      "doc_count_error_upper_bound": 0,
      "sum_other_doc_count": 0,
      "buckets": [
        {
          "key": 50054046,
          "doc_count": 6,
          "sample_doc": {
            "hits": {
              "total": {
                "value": 6,
                "relation": "eq"
              },
              "max_score": 10.704353,
              "hits": [
                {
                  "_index": "restaurantny_v1",
                  "_id": "jhWH3JkByv84jpscLOJK",
                  "_score": 10.704353,
                  "_source": {
                    "DBA": "LADUREE",
                    "BUILDING": "864",
                    "BORO": "Manhattan",
                    "STREET": "MADISON AVENUE"
                  }
                }
              ]
            }
          }
        }
      ]
    }
  }
}

```

PROF

l'adresse du restaurant LADUREE est 864 MADISON AVENUE, Manhattan.

visualisation

pour visualiser avec les nouveau filtre, on update notre dataview pour utiliser le nouvel index, mais je me suis rendu compte que le reindexage n'a pas pris en compte les dates ainsi que "location", je corrige le reindexage.

mapping de restaurantny

```
{  
  "restaurantny": {  
    "mappings": {  
      "_meta": {  
        "created_by": "file-data-visualizer"  
      },  
      "properties": {  
        "@timestamp": {  
          "type": "date"  
        },  
        "ACTION": {  
          "type": "text"  
        },  
        "BBL": {  
          "type": "long"  
        },  
        "BIN": {  
          "type": "long"  
        },  
        "BORO": {  
          "type": "keyword"  
        },  
        "BUILDING": {  
          "type": "keyword"  
        },  
        "CAMIS": {  
          "type": "long"  
        },  
        "CRITICAL FLAG": {  
          "type": "keyword"  
        },  
        "CUISINE DESCRIPTION": {  
          "type": "keyword"  
        },  
        "Census Tract": {  
          "type": "long"  
        },  
        "Community Board": {  
          "type": "long"  
        },  
        "Council District": {  
          "type": "long"  
        },  
        "DBA": {  
          "type": "text"  
        }  
      }  
    }  
  }  
}
```

PROF

```
},
"GRADE": {
    "type": "keyword"
},
"GRADE DATE": {
    "type": "date",
    "format": "MM/dd/yyyy"
},
"INSPECTION DATE": {
    "type": "date",
    "format": "MM/dd/yyyy"
},
"INSPECTION TYPE": {
    "type": "keyword"
},
"Latitude": {
    "type": "double"
},
"Longitude": {
    "type": "double"
},
"NTA": {
    "type": "keyword"
},
"PHONE": {
    "type": "keyword"
},
"RECORD DATE": {
    "type": "date",
    "format": "MM/dd/yyyy"
},
"SCORE": {
    "type": "long"
},
"STREET": {
    "type": "keyword"
},
"VIOLATION CODE": {
    "type": "keyword"
},
"VIOLATION DESCRIPTION": {
    "type": "text"
},
"ZIPCODE": {
    "type": "long"
},
"location": {
    "type": "geo_point"
}
}
}
}
```

—
PROF

Mapping de restaurantny_v1

```
{  
  "restaurantny_v1": {  
    "mappings": {  
      "properties": {  
        "@timestamp": {  
          "type": "date"  
        },  
        "ACTION": {  
          "type": "text",  
          "fields": {  
            "keyword": {  
              "type": "keyword",  
              "ignore_above": 256  
            }  
          }  
        },  
        "BBL": {  
          "type": "long"  
        },  
        "BIN": {  
          "type": "long"  
        },  
        "BORO": {  
          "type": "text",  
          "fields": {  
            "keyword": {  
              "type": "keyword",  
              "ignore_above": 256  
            }  
          }  
        },  
        "BUILDING": {  
          "type": "text",  
          "fields": {  
            "keyword": {  
              "type": "keyword",  
              "ignore_above": 256  
            }  
          }  
        },  
        "CAMIS": {  
          "type": "long"  
        },  
        "CRITICAL FLAG": {  
          "type": "text",  
          "fields": {  
            "keyword": {  
              "type": "keyword",  
              "ignore_above": 256  
            }  
          }  
        }  
      }  
    }  
  }  
}
```

PROF

```
        "ignore_above": 256
    }
}
},
"CUISINE DESCRIPTION": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"Census Tract": {
    "type": "long"
},
"Community Board": {
    "type": "long"
},
"Council District": {
    "type": "long"
},
"DBA": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword"
        }
    }
},
"GRADE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"GRADE DATE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"INSPECTION DATE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
}
```

—
PROF

```
        "ignore_above": 256
    }
}
},
"INSPECTION_TYPE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"Latitude": {
    "type": "float"
},
"Longitude": {
    "type": "float"
},
"NTA": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"PHONE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"RECORD_DATE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"SCORE": {
    "type": "long"
},
"STREET": {
    "type": "text",
    "fields": {
        "keyword": {

```

—
PROF

```

        "type": "keyword",
        "ignore_above": 256
    }
}
},
"VIOLATION CODE": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"VIOLATION DESCRIPTION": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
},
"ZIPCODE": {
    "type": "long"
},
"location": {
    "type": "text",
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    }
}
}
}
}
}

```

PROF

On recommande l'indexation

```

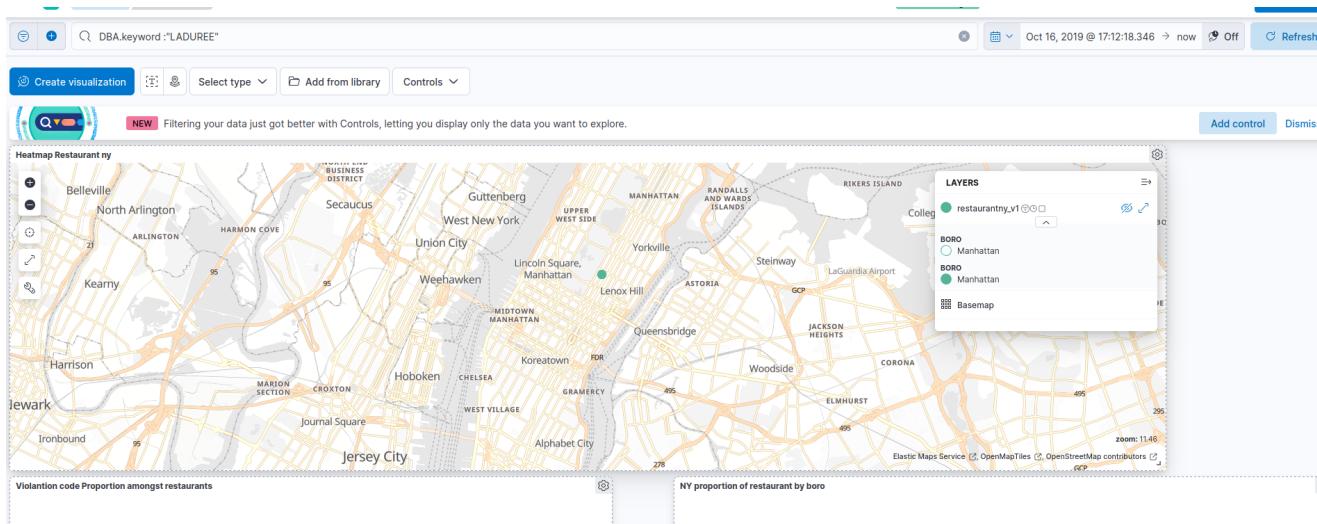
PUT restaurantny_v1
{
    "mappings": {
        "properties": {
            "DBA": {
                "type": "text",
                "fields": {
                    "keyword": { "type": "keyword" }
                }
            }
        }
    }
}
```

```

        },
        "INSPECTION DATE": {
            "type": "date",
            "format": "MM/dd/yyyy"
        },
        "GRADE DATE": {
            "type": "date",
            "format": "MM/dd/yyyy"
        },
        "RECORD DATE": {
            "type": "date",
            "format": "MM/dd/yyyy"
        },
        "location": {
            "type": "geo_point"
        }
    }
}

```

On peut finallement visualiser la localisation du restaurant LADUREE



Cette méthode, bien que optimisé, est assez longue et provoque des erreurs si l'on oublie des field dans le mapping ou si l'on ne réécrit pas le même mapping à la main (par exemple, tout mes type keyword ont été transformé en text avec une propriété keyword).

Une méthode plus rapide :

- pour les petit dataset, reupload la data et ajoutant keyword au field DBA
- ajouter un field DBA_keyword à notre index de base, mais ce seraient moins optimisé car les deux champs sont stocké et indexé différemment.

2.9 Identify the cuisine most affected by the violation “Hot food item not held at or above 140° F”

Pour cette requête, on peut utiliser match phrase pour trouver l'exacte phrase (non case insensitive) ou match, mais cela requiert d'ajuster la fuzziness pour ne pas avoir de résultat avec une description complètement différente.

requête:

```
GET restaurantny/_search
{
  "size": 0,
  "query": {
    "match_phrase": {
      "VIOLATION DESCRIPTION": "Hot food item not held at or above 140°
F"
    }
  },
  "aggs": {
    "most_affected_cuisine": {
      "terms": {
        "field": "CUISINE DESCRIPTION",
        "size": 10,
        "order": { "_count": "desc" }
      }
    }
  }
}
```

resultat

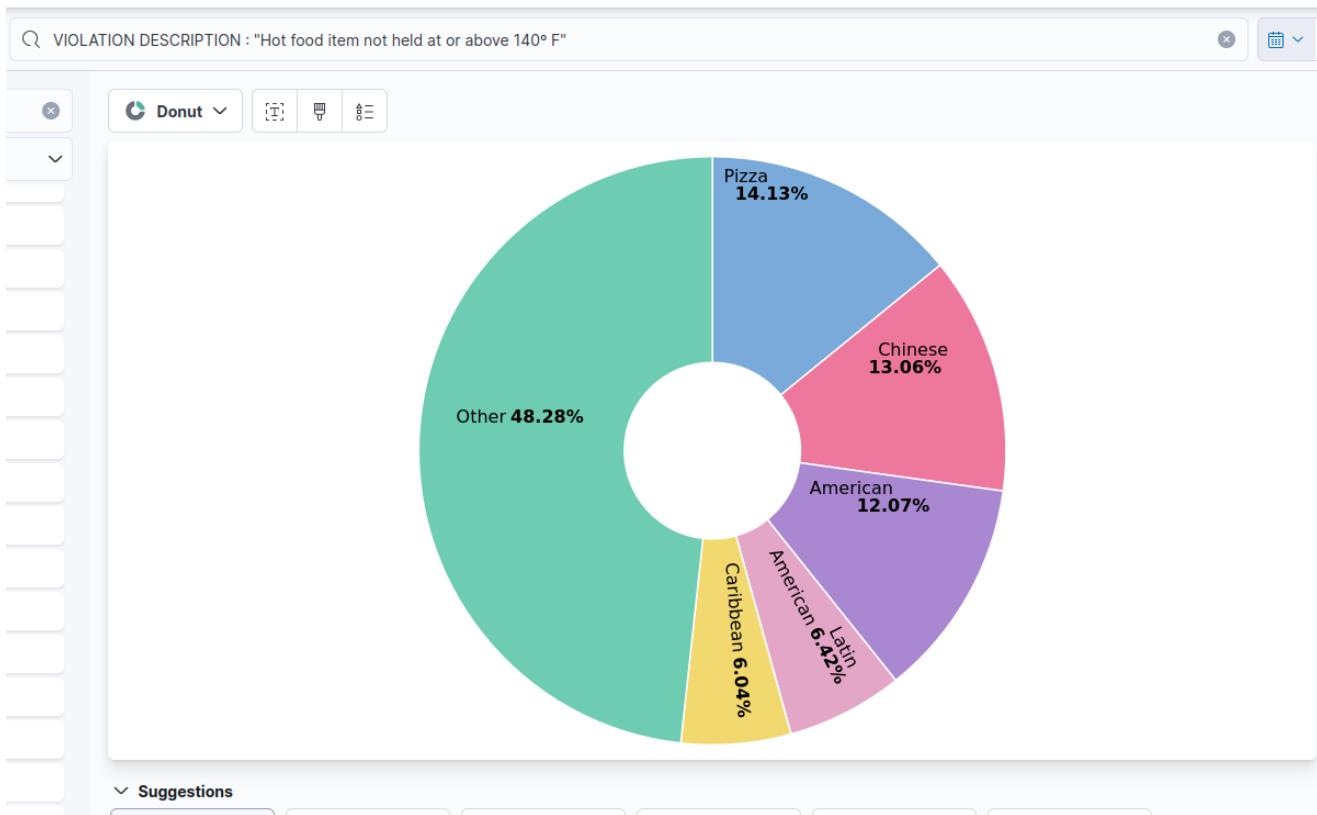
```
PROF
"aggregations": {
  "most_affected_cuisine": {
    "doc_count_error_upper_bound": 0,
    "sum_other_doc_count": 401,
    "buckets": [
      {
        "key": "Pizza",
        "doc_count": 187
      },
      {
        "key": "American",
        "doc_count": 179
      },
      {
        "key": "Chinese",
        "doc_count": 174
      },
      {
        "key": "Latin American",
        "doc_count": 84
      },
    ]
  }
}
```

```
{  
    "key": "Caribbean",  
    "doc_count": 80  
,  
{  
    "key": "Japanese",  
    "doc_count": 60  
,  
{  
    "key": "Mexican",  
    "doc_count": 49  
,  
{  
    "key": "Bakery Products/Desserts",  
    "doc_count": 48  
,  
{  
    "key": "Italian",  
    "doc_count": 47  
,  
{  
    "key": "Coffee/Tea",  
    "doc_count": 37  
}  
]  
}  
}
```

La cuisine la plus affecté par cette violation sont les pizzerias

visualisation

Pie Chart avec le filtre



2.10 Determine the most common violations (Top 5)

On aggrege sur les violation code et on affiche la violation description avec une size de 5

```
GET restaurantny/_search
{
  "size": 0,
  "aggs": {
    "most_affected_cuisine": {
      "terms": {
        "field": "VIOLATION CODE",
        "size": 5,
        "order": { "_count": "desc" }
      },
      "aggs": {
        "sample_description": {
          "top_hits": {
            "_source": ["VIOLATION DESCRIPTION"],
            "size": 1
          }
        }
      }
    }
  }
}
```

resultat

```
"aggregations": {
    "most_affected_cuisine": {
        "doc_count_error_upper_bound": 0,
        "sum_other_doc_count": 161778,
        "buckets": [
            {
                "key": "10F",
                "doc_count": 40108,
                "sample_description": {
                    "hits": {
                        "total": {
                            "value": 40108,
                            "relation": "eq"
                        },
                        "max_score": 1,
                        "hits": [
                            {
                                "_index": "restaurantny",
                                "_id": "qROH3JkByv84jpscZM78",
                                "_score": 1,
                                "_source": {
                                    "VIOLATION DESCRIPTION": "Non-food contact surface or equipment made of unacceptable material, not kept clean, or not properly sealed, raised, spaced or movable to allow accessibility for cleaning on all sides, above and underneath the unit."
                                }
                            }
                        ]
                    }
                }
            },
            {
                "key": "08A",
                "doc_count": 27330,
                "sample_description": {
                    "hits": {
                        "total": {
                            "value": 27330,
                            "relation": "eq"
                        },
                        "max_score": 1,
                        "hits": [
                            {
                                "_index": "restaurantny",
                                "_id": "sROH3JkByv84jpscZM78",
                                "_score": 1,
                                "_source": {
                                    "VIOLATION DESCRIPTION": "Establishment is not free of harborage or conditions conducive to rodents, insects or other pests."
                                }
                            }
                        ]
                    }
                }
            }
        ]
    }
},
```

—
PROF

```
        ]
    }
}
},
{
  "key": "06D",
  "doc_count": 18552,
  "sample_description": {
    "hits": {
      "total": {
        "value": 18552,
        "relation": "eq"
      },
      "max_score": 1,
      "hits": [
        {
          "_index": "restaurantny",
          "_id": "qxOH3JkByv84jpscZM78",
          "_score": 1,
          "_source": {
            "VIOLATION DESCRIPTION": "Food contact surface not properly washed, rinsed and sanitized after each use and following any activity when contamination may have occurred."
          }
        }
      ]
    }
  }
},
{
  "key": "02G",
  "doc_count": 18056,
  "sample_description": {
    "hits": {
      "total": {
        "value": 18056,
        "relation": "eq"
      },
      "max_score": 1,
      "hits": [
        {
          "_index": "restaurantny",
          "_id": "xhOH3JkByv84jpscZM78",
          "_score": 1,
          "_source": {
            "VIOLATION DESCRIPTION": "Cold TCS food item held above 41 °F; smoked or processed fish held above 38 °F; intact raw eggs held above 45 °F; or reduced oxygen packaged (ROP) TCS foods held above required temperatures except during active necessary preparation."
          }
        }
      ]
    }
  }
}
```

—
PROF

```

        }
    },
{
    "key": "10B",
    "doc_count": 17746,
    "sample_description": {
        "hits": {
            "total": {
                "value": 17746,
                "relation": "eq"
            },
            "max_score": 1,
            "hits": [
                {
                    "_index": "restaurantny",
                    "_id": "shOH3JkByv84jpscZM78",
                    "_score": 1,
                    "_source": {
                        "VIOLATION DESCRIPTION": "Anti-siphonage or back-flow prevention device not provided where required; equipment or floor not properly drained; sewage disposal system in disrepair or not functioning properly. Condensation or liquid waste improperly disposed of."
                    }
                }
            ]
        }
    }
}

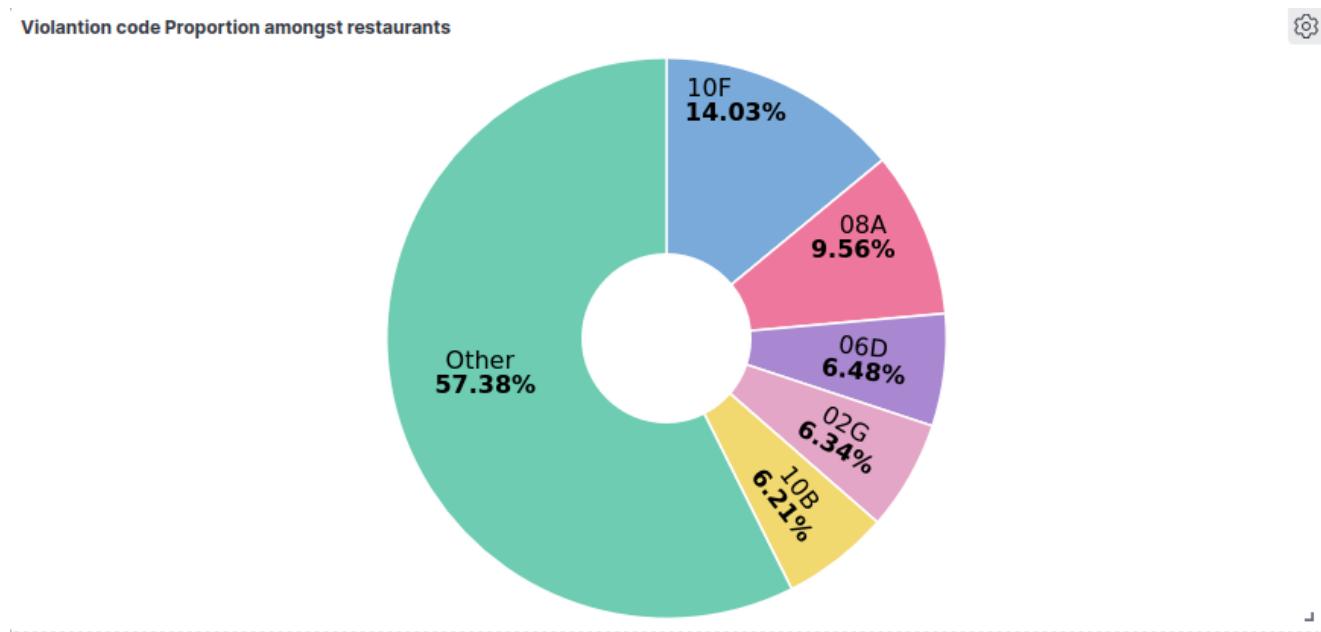
```

le top 5 des violations :

Rank	Violation Code	Description	Number of Occurrences
<hr/>			
1	10F	Non-food contact surface or equipment made of unacceptable material, not kept clean, or not properly sealed, raised, spaced, or movable to allow accessibility for cleaning on all sides, above and underneath the unit.	40 108
2	08A	Establishment is not free of harborage or conditions conducive to rodents, insects, or other pests.	27 330
3	06D	Food contact surface not properly washed, rinsed, and sanitized after each use and following any activity when contamination may have occurred.	18 552
4	02G	Cold TCS food item held above 41 °F; smoked or processed fish held above 38 °F; intact raw eggs held above 45 °F; or reduced-oxygen packaged (ROP) TCS foods held above required temperatures except during active necessary preparation.	18 056

Rank	Violation Code	Description	Number of Occurrences
5	10B	Anti-siphonage or back-flow prevention device not provided where required; equipment or floor not properly drained; sewage disposal system in disrepair or not functioning properly; condensation or liquid waste improperly disposed of.	17 746

visualisation



2.11 Identify the most popular restaurant chain

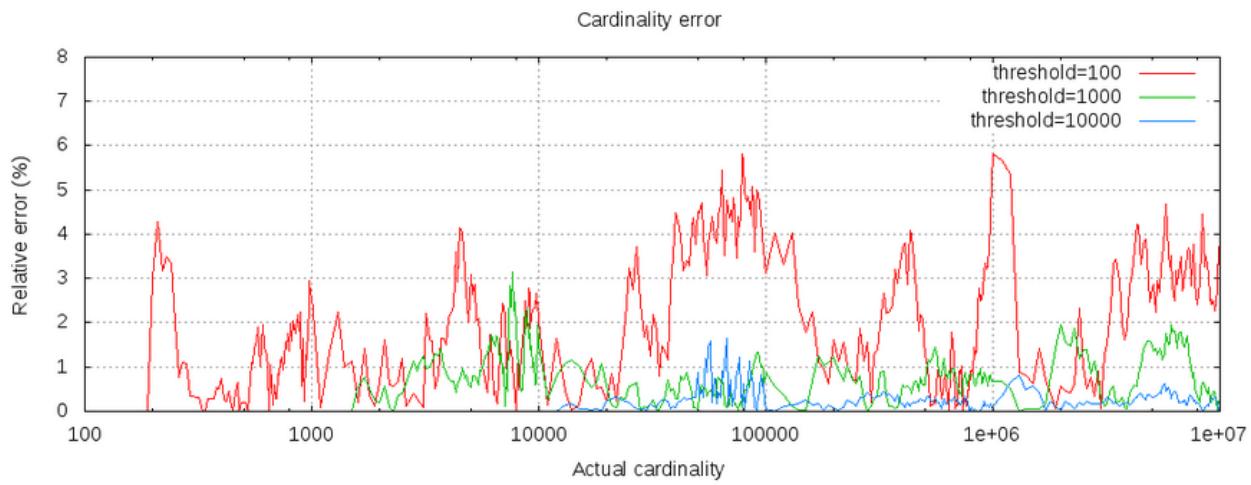
Je doit faire une aggregation sur les DBA des restaurant unique (cardinalité sur CAMIS).

pour la dernière requête, ma memoire n'est pas suffisante, je doit optimiser ma requête, plusieurs possibilité s'ouvre à moi

PROF

- Effectuer la requête sur un sample
- réduire la précision de Cardinality
- Effectuer la requête en batch (à la main ou via un script) pour avoir l'exact résultat.

Pour le projet, j'ai choisi de réduire le seuil de précision de la cardinalité de 3000 à 1000, en effet la cardinalité dans elastic search est basé sur l'algorithme Hyperloglog++ le comptage est approximatif. On peut réduire la précision pour récupérer de la puissance de calcul (<https://www.elastic.co/docs/reference/aggregations/search-aggregations-metrics-cardinality-aggregation>)



Pour vérifier la pertinence du résultat, On vérifie le nombre de document compté avec une aggregation classique

```
GET restaurantny_final/_search
{
  "size": 0,
  "aggs": {
    "restaurants_by_DBA": {
      "terms": {
        "field": "DBA.keyword",
        "size": 3,
        "order": { "_count": "desc" }
      } }
  }
}
```

réponse

PROF

```
"aggregations": {
  "restaurants_by_DBA": {
    "doc_count_error_upper_bound": 0,
    "sum_other_doc_count": 282779,
    "buckets": [
      {
        "key": "DUNKIN",
        "doc_count": 3241
      },
      {
        "key": "SUBWAY",
        "doc_count": 2003
      },
      {
        "key": "STARBUCKS",
        "doc_count": 1547
      }
    ]
  }
}
```

```
        ]
    }
}
```

On compte 3241 document pour DUNKIN, 2003 pour SUBWAY et 1547 pour STARBUCKS. Il est probable qu'on retrouve le même top 3 pour les restaurant unique en prenant en compte la correlation entre le nombre d'inspection et le nombre de restaurant.

requete avec cardinalité sur CAMIS et seuil de precision a 1500:

```
GET restaurantny_final/_search
{
  "size": 0,
  "aggs": {
    "restaurants_by_DBA": {
      "terms": {
        "field": "DBA.keyword",
        "size": 3,
        "order": { "unique_restaurants.value": "desc" }
      },
      "aggs": {
        "unique_restaurants": {
          "cardinality": {
            "field": "CAMIS",
            "precision_threshold": 1500
          }
        }
      }
    }
  }
}
```

PROF

reponse

```
"aggregations": {
  "restaurants_by_DBA": {
    "doc_count_error_upper_bound": -1,
    "sum_other_doc_count": 282779,
    "buckets": [
      {
        "key": "DUNKIN",
        "doc_count": 3241,
        "unique_restaurants": {
          "value": 348
        }
      },
      {
        "key": "SUBWAY",
        "doc_count": 2003,
        "unique_restaurants": {
          "value": 348
        }
      },
      {
        "key": "STARBUCKS",
        "doc_count": 1547,
        "unique_restaurants": {
          "value": 348
        }
      }
    ]
  }
}
```

```
"key": "STARBUCKS",
"doc_count": 1547,
"unique_restaurants": {
    "value": 209
},
{
    "key": "SUBWAY",
    "doc_count": 2003,
    "unique_restaurants": {
        "value": 176
    }
}
]
```

On obtient le même nombre de document, le seuil de precision à 1500 est donc assez précis, avec en top 3.

- Dukin : 348 restaurant
- STARBUCKS : 209 restaurant
- SUBWAY : 176 restaurant

donc DUKIN est le restaurant le plus populaire.

visualisation

On peut faire une table

Table ▾

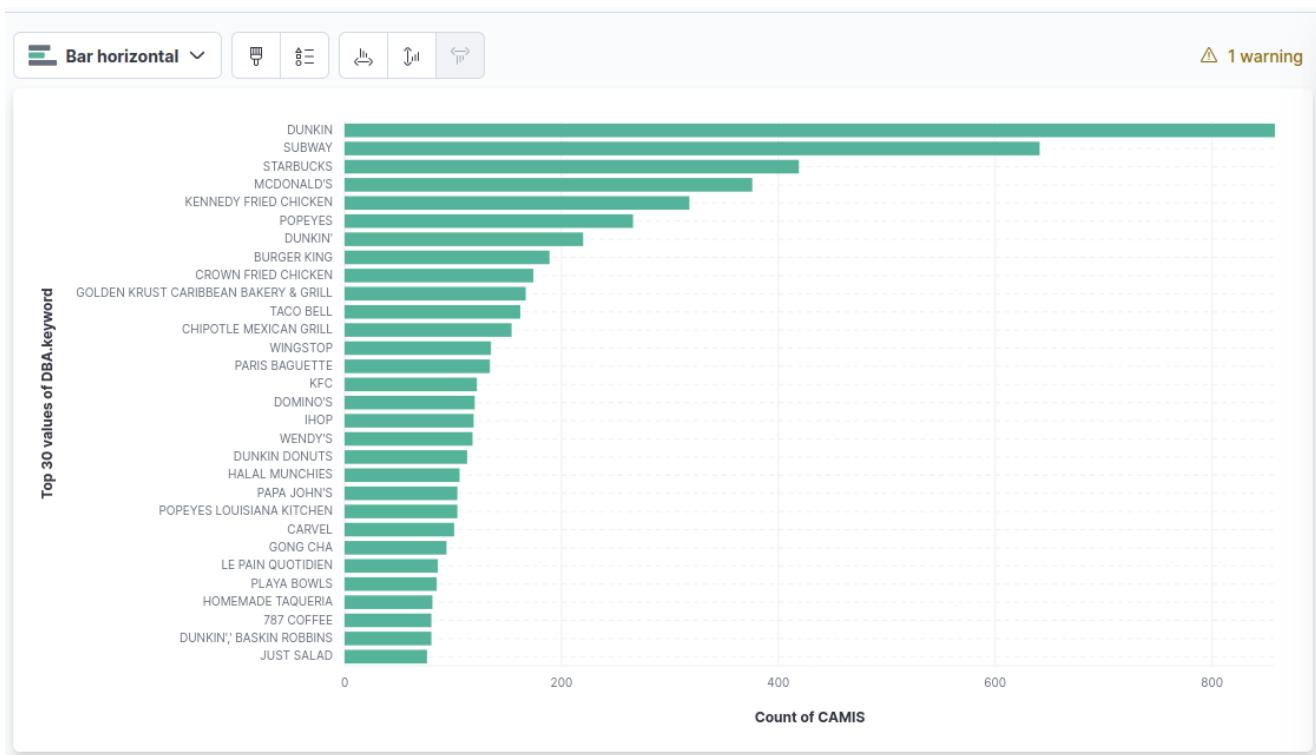
Count of CAMIS ▾

⚠ 1 warning

Top 100 values of DBA.keyword	Count of CAMIS
DUNKIN	8
SUBWAY	641
STARBUCKS	419
MCDONALD'S	376
KENNEDY FRIED CHICKEN	318
POPEYES	266
DUNKIN'	220
BURGER KING	189
CROWN FRIED CHICKEN	174
GOLDEN KRUST CARIBBEAN BAKERY & GRILL	167
TACO BELL	162
CHIPOTLE MEXICAN GRILL	154
WINGSTOP	135
PARIS BAGUETTE	134

▼ Suggestions

ou un bar chart (en enlevant "other" pour la lisibilité)

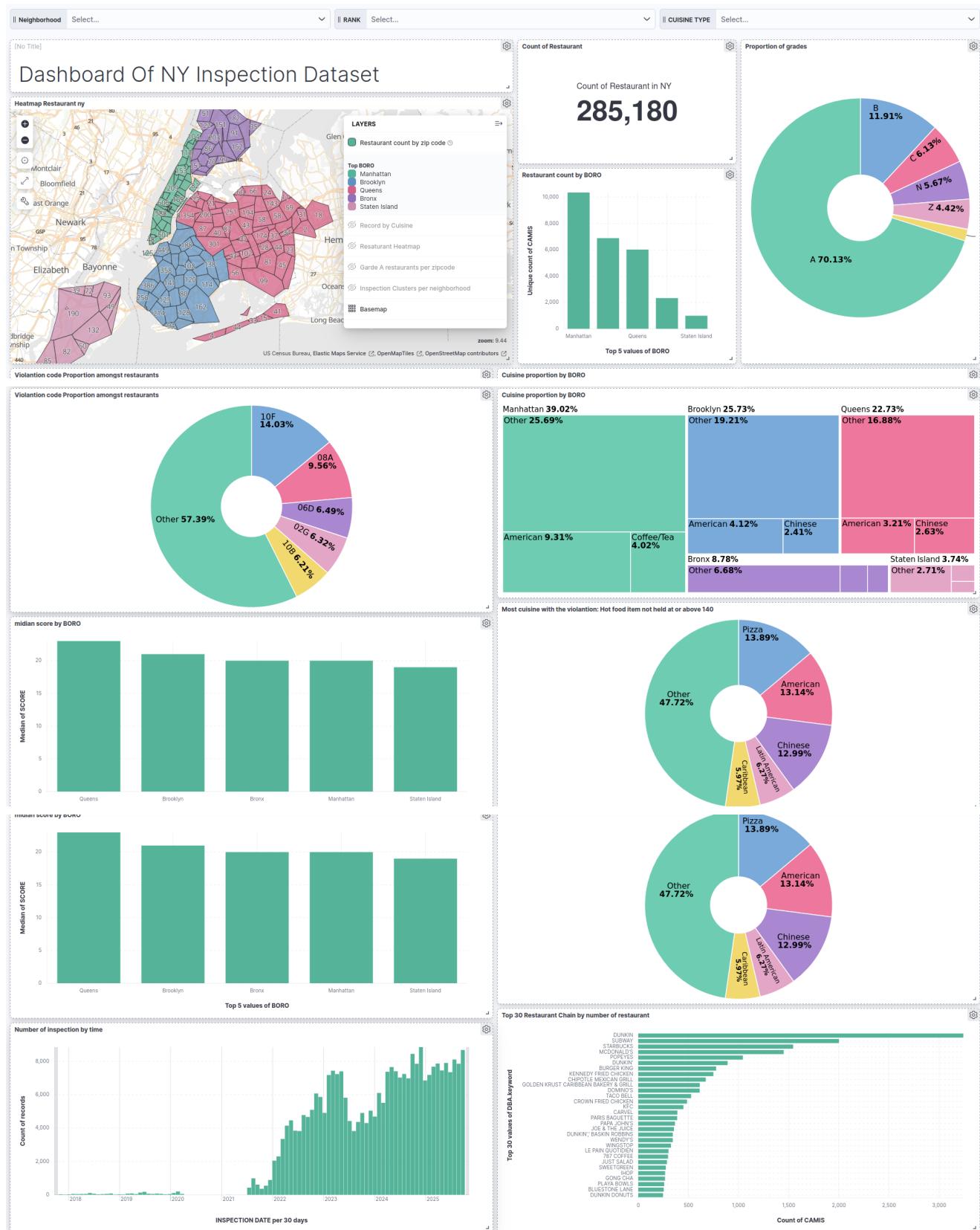


▼ Suggestions

On remarque DUNKIN et DUNKIN', on pourrait refaire l'index pour enlever les caractères spéciaux, mais cela pourrait aussi regrouper des restaurants différents, nous estimons que notre degré de précision est suffisant.

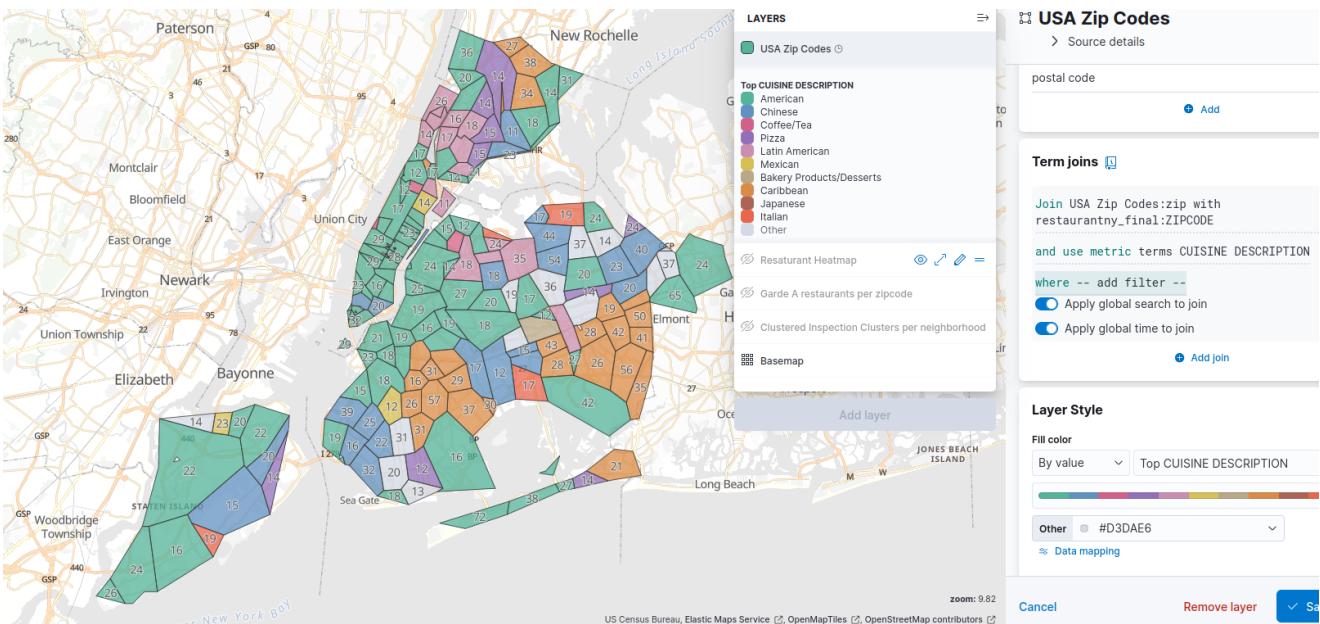
Dashboard

Visualisez le dashboard regroupant les visualisations ici :

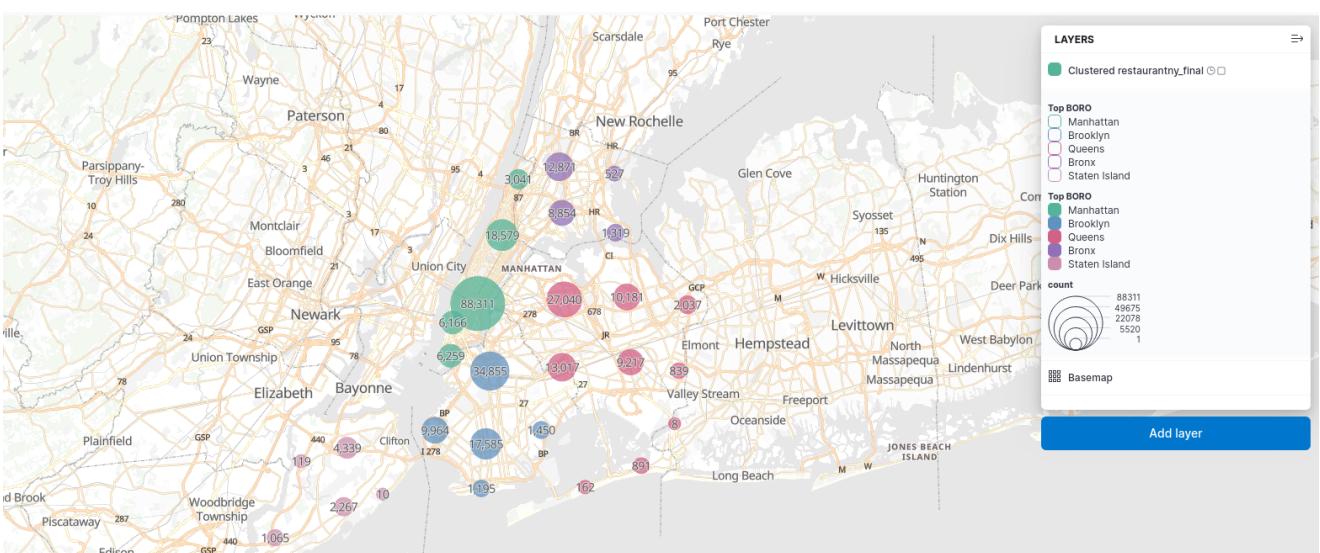


Visualisation with MAP

Nombre de document par zipcode en fonction du type de cuisine

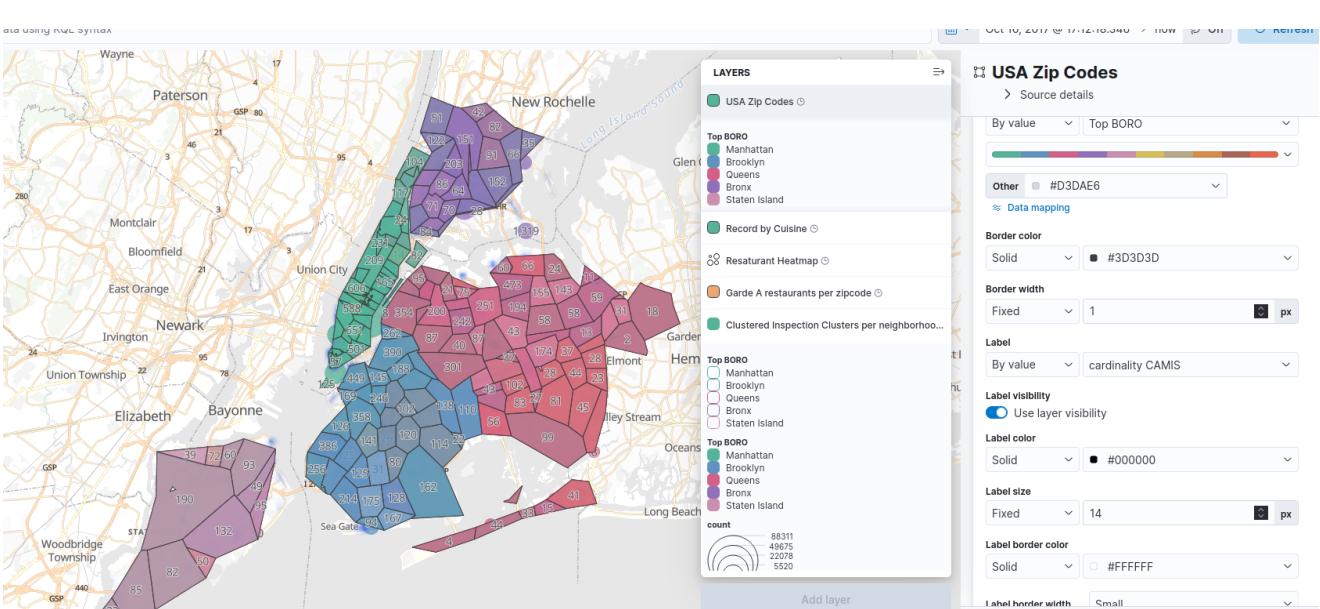


Cluster des inspection par Voisinage

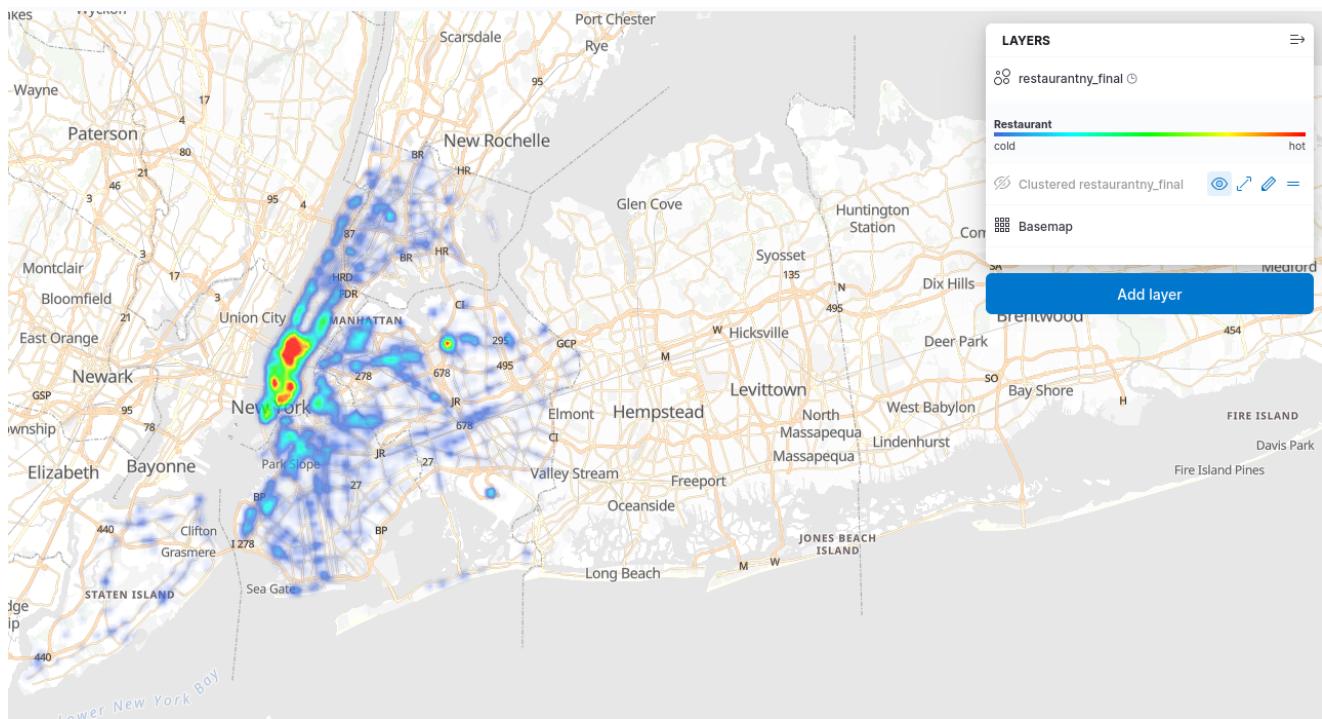


restaurant par voisinage

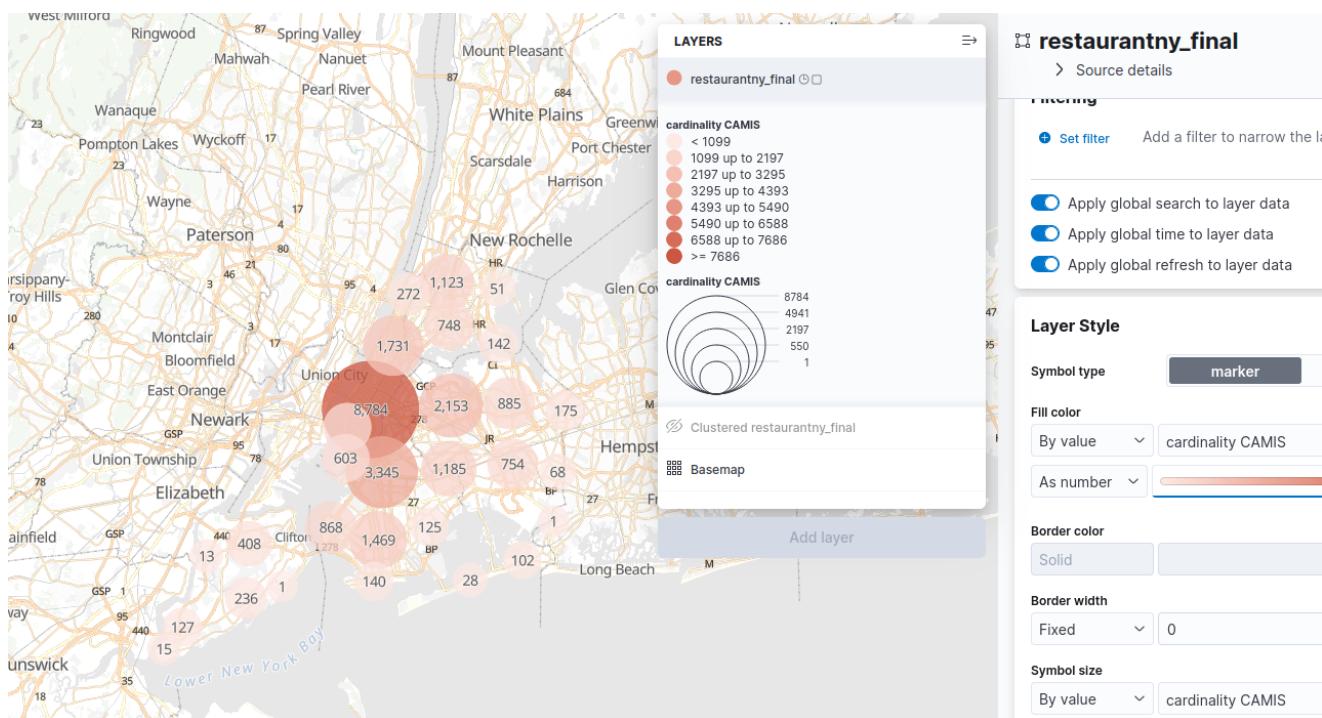
PROF



Heatmap of Unique Restaurant



Cluster of restaurant



Restaurant A la note de A par zipcode

