



Politecnico di Milano

Systems and methods for big and unstructured data

SMBUD project: third delivery

Group 24

Basso Paolo 10783951

Aiello Andrea 10863133

Borsatto Andrea 10628989

Cavalli Dario 10820532

Petriconi Emanuele 10577000

Academic Year 2021–2022

Specifications

The scope of this project is to import and analyze the following dataset on Elasticsearch, [somministrazioni-vaccini-latest](#). It contains information about the COVID-19 vaccine administrations in Italy.

Hypothesis

The only hypothesis made by our team and not by the authors of the dataset ([covid19-opendata-vaccini](#)) is that the numeric fields will not be bigger than $2^{31} - 1$. This hypothesis is reasonable knowing what each field represent.

Dataset Schema

After importing the dataset, an index called "**somministrazioni_vaccini**" was created with the following mappings:

Field name	Field type
data_somministrazione	date
fornitore	keyword
area	keyword
fascia_anagrafica	keyword
Sesso_maschile	integer
Sesso_femminile	integer
prima_dose	integer
seconda_dose	integer
pregressa_infezione	integer
dose_addizionale_booster	integer
codice_NUTS1	keyword
codice_NUTS2	keyword
codice_regione_ISTAT	keyword
nome_area	keyword

The following fields are of type *integer* because they represent an unbounded numeric value which is always less than $2^{31} - 1$ in the dataset, but some values are, or could be in the future, bigger than 32767 so they cannot be mapped to type *short*:

- Sesso_maschile
- Sesso_femminile
- prima_dose
- seconda_dose

- `pregressa_infezione`
- `dose_addizionale_booster`

The field "`data_somministrazione`" is of type *date* because it represents a date in the format YYYY-MM-DD (ISO8601).

The following fields are of type *keyword* because they are strings chosen from a limited number of possible values which are going to be used in [term-level queries](#):

- `fornitore`
- `area`
- `fascia_anagrafica`
- `codice_NUTS1`
- `codice_NUTS2`
- `nome_area`
- `codice_regione_ISTAT`

For "`codice_regione_ISTAT`" we decided to use the *keyword* type even though Kibana suggested the *numeric* type. We made this decision because, following the [documentation](#), the possible numeric values are from a limited set and this field will not be used for range queries of any kind while it could be used for [term-level queries](#).

Queries

1. Total number of somministrated doses in Lombardia after 01/12/2021.

```
1 GET /somministrazioni_vaccini/_search
2 {
3     "size":0,
4     "query":{
5         "bool":{
6             "must":[{
7                 "term":{
8                     "nome_area":{
9                         "value":"Lombardia"
10                    }
11                }},
12             {"range":{
13                 "data_somministrazione":{
14                     "gte": "2021-12-01"
15                 }
16             }}]
17         }
18     },
19     "aggs":{
20         "vaccinazioni":{
21             "sum":{
22                 "script": "doc['sesso_femminile'].value + doc['sesso_maschile'].value"
23             }
24         }
25     }
26 }
```

2. Number of injected doses in Italy during the last month to people under 19 years old.

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size":0,
4   "query":{
5     "bool":{
6       "must":[
7         {"range":{
8           "data_somministrazione":{
9             "gte": "now-1M/d"
10          }
11        }},
12       {"bool":{
13         "should":[
14           {"term":{"fascia_anagrafica": "05-11"}},
15           {"term":{"fascia_anagrafica": "12-19"}}
16         ]}
17       ]}
18     }
19   },
20   "aggs":{
21     "vaccinazioni":{
22       "sum":{
23         "script": "doc['sesso_femminile'].value + doc['sesso_maschile']
24                   .value"
25       }
26     }
27   }
```

3. Average number of vaccinated people with previous infection during the last 7 days.

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size":0,
4   "query":{
5     "bool":{
6       "must":{
7         "range":{
8           "data_somministrazione":{
9             "gte": "now-7d/d"
10          }
11        }
12      }
13    }
14  },
15  "aggs":{
16    "region":{
17      "terms":{
18        "field": "nome_area",
19        "size": 99
20      },
21      "aggs":{
22        "avg_pregressa_infezione":{
23          "avg": {
24            "field": "pregressa_infezione"
25          }
26        }
27      }
28    }
29  }
```

4. Average number of daily somministrazioni in Sardegna during the last month.

```
1 GET /sommministrazioni_vaccini/_search {
2   "size":0,
3   "query":{"
4     "bool":{"
5       "must":[{"
6         "range":{"
7           "data_somministrazione":{"
8             "gte": "now-1M/d"
9           }
10        }
11      }],
12      {"term":{"
13        "nome_area": "Sardegna"
14      }}}
15    }
16  },
17  "aggs":{"
18    "region":{"
19      "terms":{"
20        "field": "data_somministrazione"
21      },
22      "aggs":{"
23        "sommministrazioni":{"
24          "avg":{"
25            "script": "doc['sesso_femminile'].value + doc['sesso_maschile'].value"
26          }
27        }
28      }}
29    }
30  }
```

5. Number of vaccinated people for each kind of vaccine and for each age range from 22/12/2021.

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size":0,
4   "query":{"
5     "bool":{"
6       "must":{"
7         "range":{"
8           "data_somministrazione":{"
9             "gte": "2021-12-22"
10          }
11        }
12      }
13    }
14  },
15  "aggs": {
16    "fornitore&eta": {
17      "multi_terms":{"
18        "terms": [
19          {"field":"fornitore"},
20          {"field":"fascia_anagrafica"}
21        ],
22        "size": 99
23      },
24      "aggs":{"
25        "somministrazioni":{"
26          "sum":{"
27            "script": "doc['sesso_femminile'].value + doc['sesso_maschile'].value"
28          }
29        }
30      }
31    }
32  }
33 }
```


6. Number of booster doses given to male and female for each different supplier during the last 3 months.

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "bool": {
6       "must": [
7         {"range": {"data_somministrazione": {"gte": "now-3M/d"}}},
8         {"range": {"dose_addizionale_booster": {"gte": "1"}}}
9       ]
10    }
11  },
12  "aggs": {
13    "fornitore": {
14      "terms": {
15        "field": "fornitore"
16      },
17      "aggs": {
18        "somministrazioni_donne": {
19          "sum": {
20            "field": "sesso_femminile"
21          }
22        },
23        "somministrazioni_uomini": {
24          "sum": {
25            "field": "sesso_maschile"
26          }
27        }
28      }
29    }
30  }
31 }
```

7. People over 80 years old that have received the booster.

```
1 GET /somministrazione_vaccini/_search {
2   "size":0,
3   "query": {
4     "bool": {
5       "must": [
6         {"range":{"dose_addizionale_booster": {"gte": "1"}}},
7         {"bool":{"
8           "should":[
9             {"term":{"fascia_anagrafica":"80-89"}},
10            {"term":{"fascia_anagrafica":"90+"}}
11          ]}
12       ]}
13     ]
14   },
15   "aggs":{
16     "vaccinazioni": {
17       "sum": {
18         "script": "doc['sesso_femminile'].value + doc['sesso_maschile',
19           ].value"
20       }
21     }
22   }
23 }
```

8. Day with the highest number of injected doses during the last 30 days.

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "bool": {
6       "must": [
7         {
8           "range": {
9             "data_somministrazione": {
10              "gte": "now-30d/d"
11            }
12          }
13        }
14      ]
15    }
16  },
17  "aggs": {
18    "days": {
19      "terms": {
20        "field": "data_somministrazione"
21      },
22      "aggs": {
23        "vaccinations": {
24          "sum": {
25            "script": "doc['sesso_femminile'].value + doc['sesso_maschile'].value"
26          }
27        },
28        "vaccination_sort": {
29          "bucket_sort": {
30            "sort": [{
31              "vaccinations": {"order": "desc"}
32            }],
33            "size": 1
34          }
35        }
36      }
37    }
38  }
39 }
```

Commands

1. Add a new document to Elasticsearch. Use example values.

```

1 POST /somministrazioni_vaccini/_doc
2 {
3   "area": "BAS",
4   "codice_regione_ISTAT": "17",
5   "nome_area": "Basilicata",
6   "data_somministrazione": "2022-01-11",
7   "dose_addizionale_booster": 2,
8   "codice_NUTS1": "ITF",
9   "fascia_anagrafica": "12-19",
10  "prima_dose": 1,
11  "pregressa_infezione": 0,
12  "fornitore": "Janssen",
13  "seconda_dose": 0,
14  "sesso_maschile": 3,
15  "codice_NUTS2": "ITF5",
16  "sesso_femminile": 0
17 }
```

2. A document in the database has a small error that must be corrected. The document representing the vaccinations done in the “BAS” area with the “Janssen” vaccine on the "2021-12-22" (data_somministrazione) has reported less vaccinations than what actually happened. In particular sesso_maschile must be updated from 1 to 2 as well as prima_dose. An initial query found that the id of the document is “Btox530BZF4mmwtvcTdZ”.

```

1 POST /somministrazioni_vaccini/_update/Btox530BZF4mmwtvcTdZ
2 {
3   "doc": {
4     "prima_dose": 2,
5     "sesso_maschile": 2
6   }
7 }
```

3. Delete the document older than 6 months (since we do not want to analyse those).

```

1 POST /somministrazioni_vaccini/_delete {
2   "query": {
3     "range": {
4       "data_somministrazione": {
5         "lte": "now-6M/d"
6       }
7     }
8   }
9 }
```

Kibana Dashboard

To help with the understanding of the data, we built a **Kibana Dashboard**. The dashboard is built as follows:

Control panel

Area
Select...
Apply changes Cancel changes Clear form

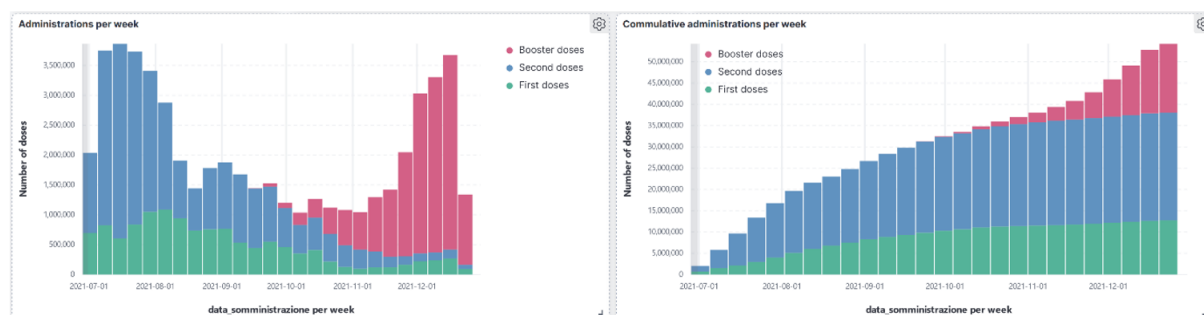
Supplier
Select...

Age group
Select...

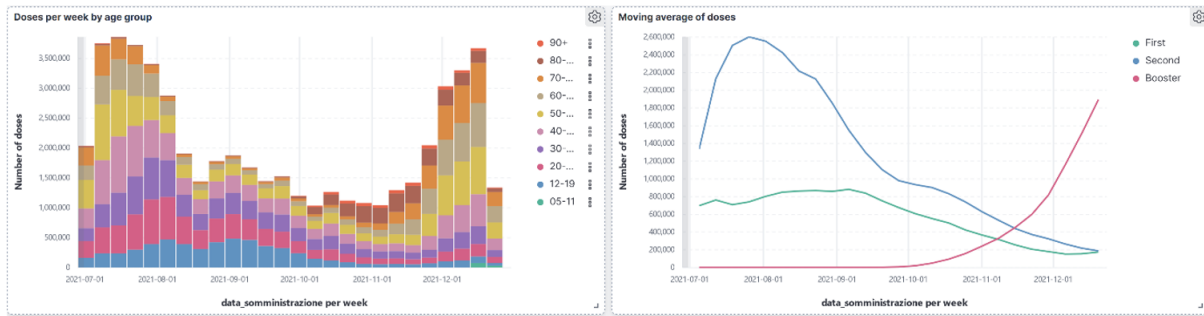
A control panel is used to filter the data according to the area, the supplier and the age group. In all three these fields multiple values can be chosen simultaneously.



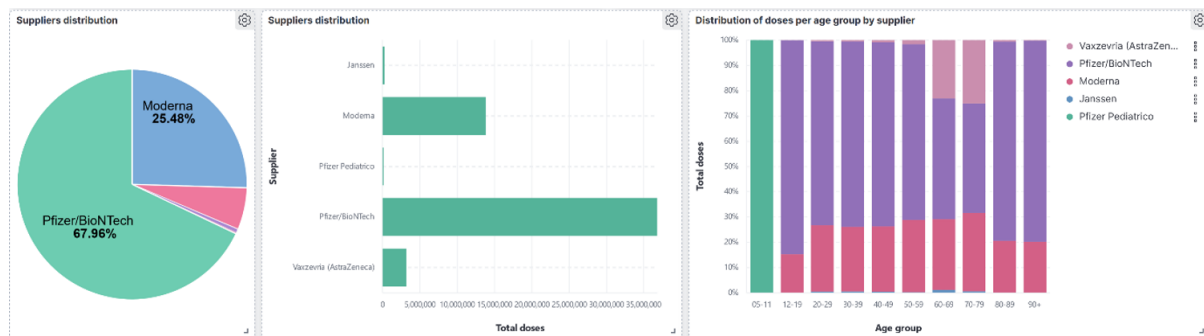
In the second row, we focus on the number of doses administered. The first graph shows the total number of doses, the second divides this number in first, second and third doses. The third graph shows how many doses have been administered to each age group and the graph is broken by dose type.



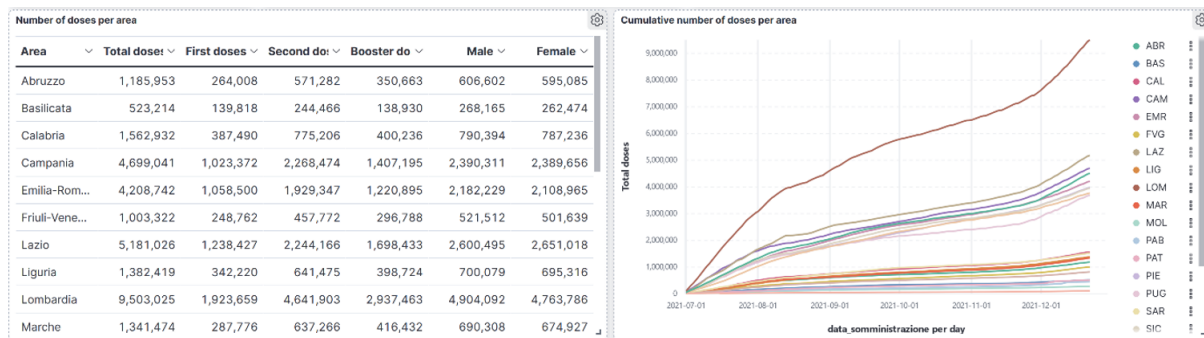
In the next row, we focus on the administrations. The first graph shows the number of administrations per week while the second graph shows the cumulative sum of administrations per week. Both graphs are broken by the dose type.



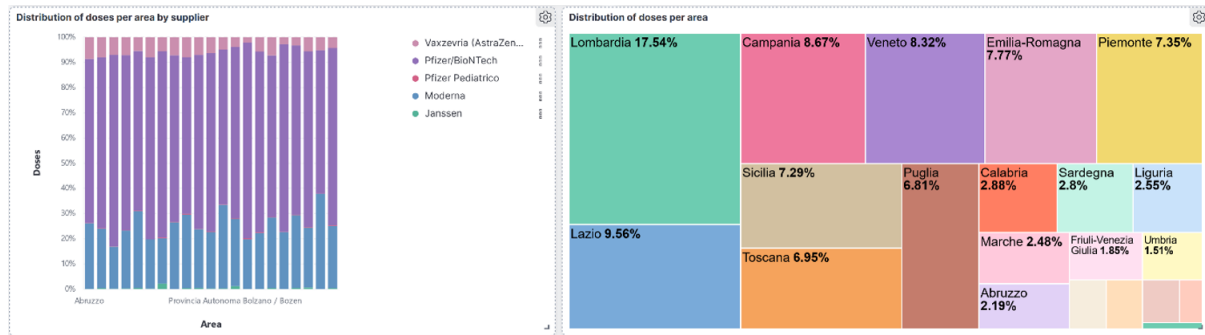
In the next row, we firstly show the number of doses administered per week but this time we break the graph by the age group. The second graph shows the moving average (window size = 7) of the administration per week broken by the type of dose.



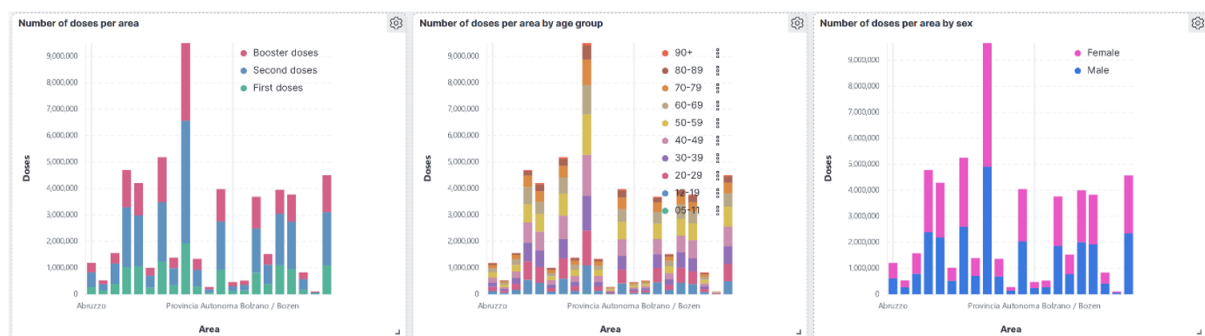
In the next row, we focus on the suppliers of the doses. The first graph is a pie chart that shows the suppliers distribution. The second graph is a horizontal bar chart which shows more clearly how many doses for each supplier have been administered. The last graph shows the distribution of the supplier used for each age group. This graph is interesting since it clearly shows that some suppliers are preferred for each age group.



We then focus on the administrations per region. We include a table with the more salient data for each region. On the left there is a graph which shows the cumulative number of doses administered in each area.



We then include a graph which for each region shows the suppliers distribution. On the left we included another graph which shows the distribution of doses per area in a different and maybe clearer way.



In the last row we included three graphs which shows some useful metrics about the number of doses administered for each region. The first graph focuses on the type of doses, the second on the age groups and the third on the sex.

Optional 1: Integrating other datasets

Specifications

Import in Elasticsearch the following dataset [dpc-covid19-ita-andamento-nazionale](#). It contains information about tests, hospitalizations, intensive care, home confinement, positive cases, recoveries, and deaths regarding the COVID-19 pandemic in Italy. A description of the fields is available [here](#).

Import only these fields (which are the ones still updated):

1. data
2. ricoverati_con_sintomi
3. terapia_intensiva
4. totale_ospedalizzati
5. isolamento_dominciliare
6. totale_positivi
7. variazione_totale_positivi
8. nuovi_positivi
9. dimessi_guariti
10. deceduti
11. totale_casi
12. tamponi

Importing the dataset to Elasticsearch

As specifications we need to whitelist some fields. A small python script is used to get only the wanted fields as a csv.

```

1 import pandas as pd
2
3 df = pd.read_csv("https://github.com/pcm-dpc/COVID-19/raw/master/dati-
    andamento-nazionale/dpc-covid19-ita-andamento-nazionale.csv")
4
5 df = df [[
6     "data",
7     "ricoverati_con_sintomi",
8     "terapia_intensiva",
9     "totale_ospedalizzati",
10    "isolamento_domiciliare",
11    "totale_positivi",
12    "variazione_totale_positivi",
13    "nuovi_positivi",
14    "dimessi_guariti",
15    "deceduti",
16    "totale_casi",
17    "tamponi",
18 ]]
19 print(df.to_csv(index=False))

```

Data Mappings

We imported the data and created an index **"covid19_italy"** with the following mappings:

Field name	Field type
data	date
ricoverati_con_sintomi	integer
terapia_intensiva	integer
totale_ospedalizzati	integer
isolamento_domiciliare	Integer
totale_positivi	Integer
variazione_totale_positivi	Integer
nuovi_positivi	Integer
dimessi_guariti	Integer
deceduti	Integer
totale_casi	Integer
tamponi	long

The following fields are of type *integer* because they represent an unbounded numeric value which is always less than $2^{31} - 1$ in the dataset, but some values are, or could be in the future, bigger than 32767 so they cannot be mapped to type *short*:

- ricoverati_con_sintomi
- terapia_intensiva
- totale_ospedalizzati
- isolamento_domiciliare
- totale_positivi
- variazione_totale_positivi
- nuovi_positivi
- dimessi_guariti
- deceduti
- totale_casi

Since "*tamponi*" is a cumulative value that could become big we decided to map this field to *long* to be sure that the upper bound will not be a problem in the future.

The field "*data*" is of type *date* because it represents a date in the format YYYY-MM-DDTHH:MM:SS (ISO 8601).

Queries

1. Daily average of new positive cases from 20/12/22.

```

1 GET /covid19_italy/_search
2 {
3   "size":0,
4   "query":{
5     "bool":{
6       "must":{
7         "range":{
8           "data":{
9             "gte": "2021-12-22"
10          }
11        }
12      }
13    }
14  },
15  "aggs": {
16    "media_positivi":{
17      "avg":{
18        "field": "nuovi_positivi"
19      }
20    }
21  }
22 }
```

2. Highest record of new positive cases.

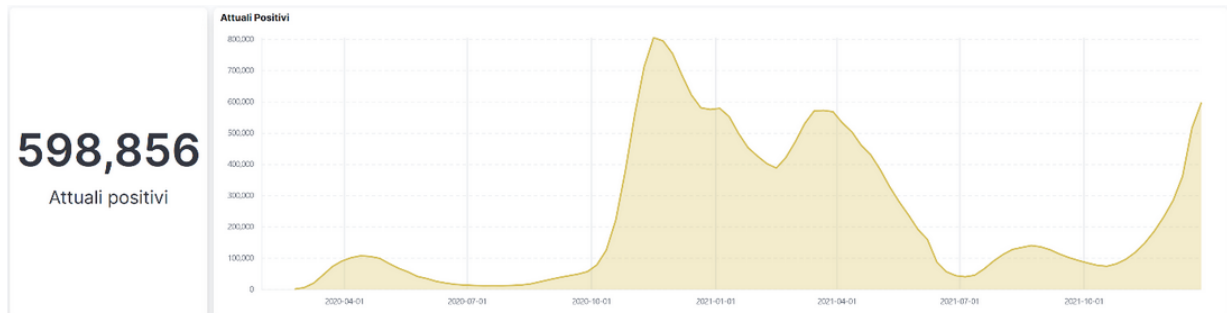
```
1 GET /covid19_italy/_search
2 {
3   "size": 0,
4   "aggs": {
5     "max_nuovi_positivi": {
6       "max": { "field": "nuovi_positivi" }
7     }
8   }
9 }
```

3. Day with the highest number of hospitalizations.

```
1 GET /covid19_italy/_search
2 {
3   "size": 5,
4   "sort": [
5     {"totale_ospedalizzati": {"order": "desc"}}
6   ],
7   "fields": [
8     "data",
9     "totale_ospedalizzati"
10  ],
11   "_source": false
12 }
```

Dashboard

We built a dashboard named “**Optional COVID Italy dashboard**” and we included it in the delivery.



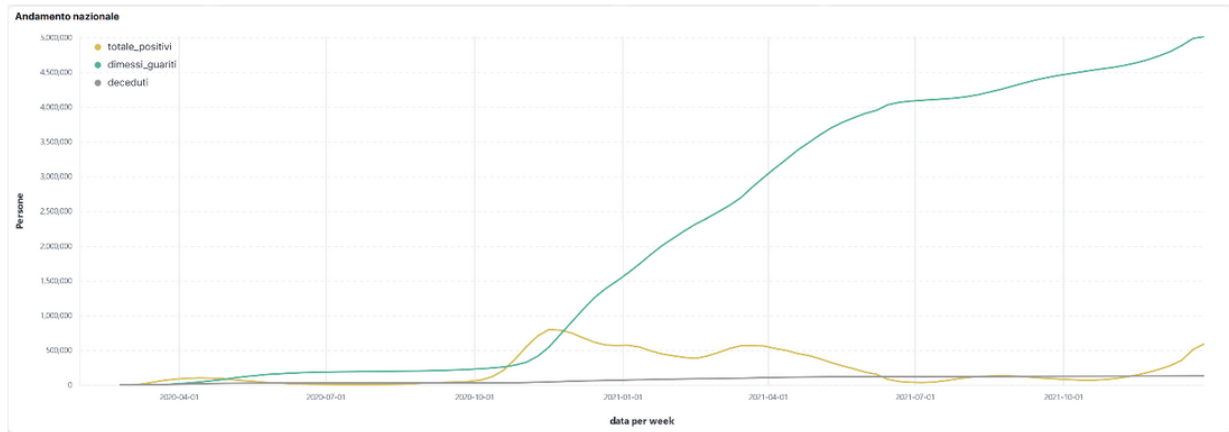
In the first row, we show the total number of positive people now and over time.



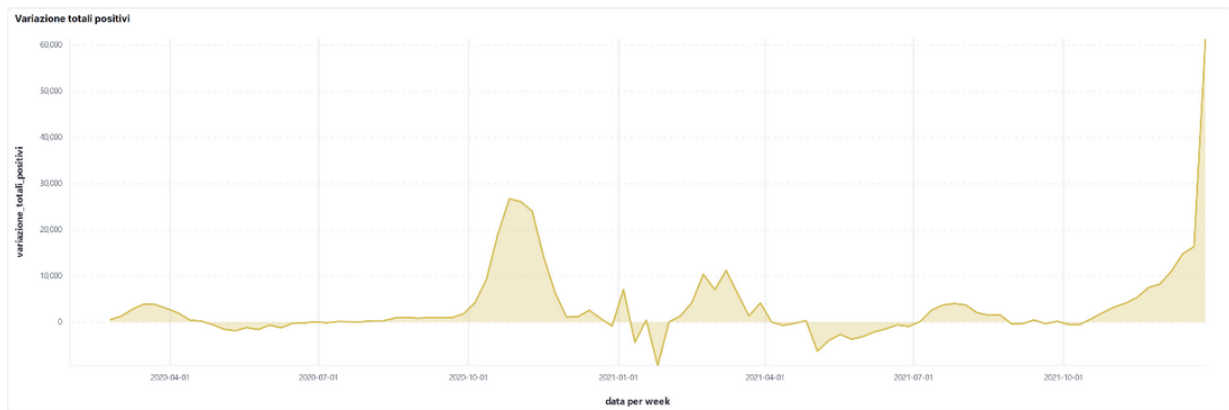
In the second row, we first show the cumulative number of healed people now and overtime. Then we show the number of new cases today and every day over time.



In the third row on the right, we show the total cumulative number of cases now and over time while on the left we show the number of cumulative number of deaths now and overtime.



We then include a graph which shows the number of positives, the cumulative number of healed and death over time.



In the last row we show the variation of the number of positives over time.

Optional 2: Implementing another NoSQL platform

We decided to implement the dataset in **Cassandra**.

Create Database structure

Firstly we create the keyspace:

```
1 CREATE KEYSPACE IF NOT EXISTS population WITH replication = {"class": "
   SimpleStrategy", "replication_factor": 3};
2
3 // Replace the double quotes with single ones.
```

Then we create the table:

```
1 USE population;
2
3 CREATE TABLE vaccinations (data_somministrazione date, fornitore text, area
   text, fascia_anagrafica text, sesso_maschile int, sesso_femminile int,
   prima_dose int, seconda_dose int, pregressa_infezione int,
   dose_addizionale_booster int, codice_NUTS1 text, codice_NUTS2 text,
   codice_regione_ISTAT int, nome_area text, PRIMARY KEY((
   data_somministrazione, fornitore, area, fascia_anagrafica)) );
```

Importing the data

We can now import the data:

```
1 // Import data from a csv file located in the specified folder (e.g: 'C:\
   OneDrive\Desktop\prova.csv')
2
3 // The same dataset used for ElasticSearch can be used here
4
5 COPY vaccinations (data_somministrazione, fornitore, area, fascia_anagrafica
   , sesso_maschile, sesso_femminile, prima_dose, seconda_dose,
   pregressa_infezione, dose_addizionale_booster, codice_NUTS1, codice_NUTS2
   , codice_regione_ISTAT, nome_area) FROM 'path-to-csv-file\file.csv' WITH
   DELIMITER=',' AND HEADER=TRUE;
```

Indexes

```
1 // for performance and avoiding ALLOW FILTERING in the queries
2
3 CREATE INDEX region_area ON vaccinations(area);
4
5 CREATE INDEX date_index ON vaccinations(data_somministrazione);
```

User defined functions

We define a user custom function that we will use in the queries:

```

1 // User defined functions must be enabled in path -to-cassandra-folder\
  cassandra\conf\cassandra.yaml, setting the variable
  enable_user_defined_functions to true
2
3 CREATE FUNCTION sum_different_columns ( arg1 int, arg2 int )
4
5     RETURNS NULL ON NULL INPUT
6
7     RETURNS int
8
9     LANGUAGE java
10
11     AS $$ return (arg1+arg2); $$;

```

Queries

1. Query data about vaccinations on 28 december 2021 using Moderna on people aged between 40-49 years old.

```

1     SELECT * FROM vaccinations WHERE data_somministrazione='2021-12-28'
      AND fornitore='Moderna' AND fascia_anagrafica='40-49' ALLOW
      FILTERING;

```

2. Find the total number of doses inoculated in Lombardia.

```

1     // Use ALLOW FILTERING if the index is not defined
2
3     SELECT area, sum_different_columns(SUM( Sesso_maschile), SUM(
      Sesso_femminile)) AS total_vaccination FROM vaccinations WHERE
      area='LOM';

```

3. Find the total number of doses inoculated in Lombardia on men and women.

```

1     // Use ALLOW FILTERING if the index is not defined
2
3     SELECT SUM(Sesso_maschile) AS man_total, SUM(Sesso_femminile) AS
      women_total FROM vaccinations WHERE area='LOM';

```

4. Find for which vaccine type and on which region and age group were inoculated the most doses of booster on 28 december 2021.

```

1     SELECT MAX(dose_addizionale_booster) AS max_booster, area,
      fascia_anagrafica, fornitore FROM vaccinations WHERE
      data_somministrazione = '2021-12-28';

```

Commands

1. Insert a new record on the database.

```
1      INSERT INTO vaccinations(data_somministrazione, fornitore, area,  
2      fascia_anagrafica, sesso_maschile, sesso_femminile, prima_dose,  
3      seconda_dose, pregressa_infezione, dose_addizionale_booster,  
      codice_NUTS1, codice_NUTS2, codice_regione_ISTAT, nome_area)  
  
      VALUES ('2021-12-29', 'Pfizer/BioNTech', 'EMR', '20-29', 1, 0, 0, 0, 0, 1, '  
      ITH', 'ITH5', 8, 'Emilia-Romagna');
```

2. Update a record on the database using "data__somministrazione", fornitore, area and "fascia__anagrafica".

```
1      UPDATE vaccinations  
2  
3      SET sesso_maschile= 2, dose_addizionale_booster=2  
4  
5      WHERE data_somministrazione = '2021-12-29' AND fornitore='Pfizer/  
      BioNTech' AND area = 'EMR' AND fascia_anagrafica='20-29';
```