

SMBUD

Vaccinations Analysis

*Specification, data model and Elasticsearch
analysis*



POLITECNICO
MILANO 1863

Filippo Lazzati (10629918) - Martina Magliani
(10682333) - Christian Grasso (10652464) - Sofia
Martellozzo (10623060) - Giacomo Lombardo
(10674987)
Year: 2021/2022

Contents

1	Problem specification	2
1.1	Dataset	2
2	Conceptual model	4
3	Queries and Commands	5
3.1	Queries	5
3.2	Commands	11
4	Kibana Dashboard	12
4.1	Dashboard initialization	13
4.2	Dashboard visualizations	13
5	Cassandra	13
5.1	Dataset implementation	17
5.2	Queries examples	18
6	References and sources	18

1 Problem specification

The purpose of this project is to design, store and query data on a NoSQL database in order to analyze data about COVID-19 vaccination statistics.

Starting from an open [dataset](#) containing raw data about the distribution and administration of COVID-19 vaccines in Italian regions, the main goal of the project is to obtain relevant information through queries and graphs.

To perform this task, the dataset has been loaded into **ElasticSearch**, a search engine with an internal database focused on high performances that interacts with external components through RESTful APIs. Once loaded into the search engine, information can be retrieved from the dataset not only through queries but also through Kibana, a powerful user interface that allows to visualize ElasticSearch data. More specifically, through **Kibana** it is possible to create Dashboards that generate and display graphs and other significant information in an intuitive way.

1.1 Dataset

The used dataset is provided by the Italian Government and contains the official data about vaccines distribution and administration. More specifically it provides a regional overview of the amount of doses administered every day in the time period between December 27th, 2020 and today. The dataset contains 15 fields, which are briefly summarized in the following:

- **index**, a unique identifier of the record;
- **area**, an acronym of the region at issue;
- **fornitore**, the brand provider of the vaccine;
- **data_somministrazione**, the date when the vaccine doses have been administered;
- **fascia_anagrafica**, the age range the people who received the vaccination belong to;
- **Sesso_maschile**, the number of males vaccinated;
- **Sesso_femminile**, the number of females vaccinated;
- **prima_dose**, the number of first doses;
- **seconda_dose**, the number of second doses;
- **pregressa_infezione**, the number of doses administered to people who came down with the disease between 3 and 6 months ago;
- **dose_addizionale_booster**, the number of third doses administered;

- **codice_NUTS1**, European classification of NUTS level 1;
- **codice_NUTS2**, European classification of NUTS level 2;
- **codice_regione_ISTAT**, ISTAT code of the region;
- **nome_regione**, the name of the region.

The way this dataset has been imported into **ElasticSearch** is shown in section 2.

According to the project's specifications, the possibility of integrating the project with an additional dataset is contemplated. Therefore, we have decided to add a secondary dataset, [consegne-vaccini-latest](#), that contains additional information about the amount of daily delivered doses in the same time period, to better contextualise the primary dataset. Such secondary dataset contains 8 field, briefly explained in the following:

- **index**, a unique identifier of the record;
- **area**, an acronym of the region at issue;
- **fornitore**, the brand provider of the vaccine;
- **data_consegna**, the date when the vaccine doses have been delivered;
- **numero_dosi**, the number of doses delivered;
- **codice_NUTS1**, European classification of NUTS level 1;
- **codice_NUTS2**, European classification of NUTS level 2;
- **codice_regione_ISTAT**, ISTAT code of the region;
- **nome_regione**, the name of the region.

2 Conceptual model

In `ElasticSearch`, there are two ways to create a schema (model) of the data we are importing in the database. We have to remember that `ElasticSearch` is an Information Retrieval system that contains an inner document-oriented database, which means that we can store data at the granularity level of documents, the same of MongoDB. In the `ElasticSearch` world, the schema of the data is called *mapping*, and can be of two types:

- [dynamic mapping](#), which is automatically carried out by `ElasticSearch`, which autonomously tries to infer the datatypes of the fields;
- [explicit mapping](#), in which we have to explicitly define the type of the fields when creating an index¹.

It is clear that dynamic mapping can fail, therefore we have to be cautious when using it.

In our project, we have decided to exploit dynamic mapping because the mapping inferred perfectly coincides with the one we would have defined. In fact, you can notice this in the following json code, which is provided by `ElasticSearch` to justify the dynamic mapping: The mapping requires some observations:

- the type **keyword** is used for structured content such as IDs, email addresses, hostnames, status codes, zip codes, or tags, therefore it perfectly fits for codes like "codice.NUTS1", or "area", or "fornitore", ... etc.; notice that it is ok for us also to have "fascia_anagrafica" as **keyword**, since we do not need to perform text search on it;
- the type **long** is a numeric datatype that represents a signed 64-bit integer with a minimum value of -2^{63} and a maximum value of $2^{63} - 1$; basically, it is type **integer** but enlarged with more bytes, therefore it is suited for all the numeric values we have in the dataset;
- the only date field, "data_somministrazione", has been properly inferred as a date, thanks to the fact that it was encoded in the standard format iso8601;
- it is better to spend some words about field "area". As a matter of fact, it has been inferred of type "keyword", which is ok for our purposes; however, it can be matched 1-to-1 to a region of Italy, therefore in order to provide a map view in the `Kibana` dashboard, later on we have provided some rows of explanation about how to match field "area" to "geo_shape" datatype, which allows to draw it as a polygon.

As far as the second dataset is concerned, there are few differences with what previously stated: As it can be easily seen by the json provided, the date "data_consegna" has been properly inferred, as well as all the other fields, therefore we can easily adopt dynamic mapping also for this dataset.

¹In the `ElasticSearch` world, an index can be seen as a table of RDBS.

3 Queries and Commands

In this section we present some queries and some commands that we consider useful to manage the considered data. It should be remarked that we can interact with `ElasticSearch` through REST APIs, since it is RESTful, and that there are various ways to interfacing with it; the most used are:

- a REST client application like [Postman](#);
- a client like [curl](#);
- Kibana, the graphical user interface of `ElasticSearch`. This is the way we adopt.

You have to notice also that you have to upload the datasets with names `"somministrazioni_vaccini"` and `"consegne_vaccini"` in order to use the following queries and commands.

3.1 Queries

We have identified the following queries:

1. **find the number of vaccines administered by brand in the last 30 days and also the brand with the highest number.**

In order to perform this query, we have to use the GET method (like all the other queries). Moreover, we have to filter the results through "query" to retrieve data that is not older than 30 days and then we have to exploit a script to sum the values of fields `"sesso_maschile"` and `"sesso_femminile"`. Then, we have to aggregate by brand (`"fornitore"`) to compute meaningful sums. Finally, to retrieve the brand with the highest number of vaccines administered we have to exploit also the `"max_bucket"` construct:

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "range": {
6       "@timestamp": {
7         "gte": "now-30d/d",
8         "lt": "now/d"
9       }
10    }
11  },
12  "runtime_mappings": {
13    "vaccini": {
14      "type": "long",
15      "script": ""
16        long total = doc['sesso_maschile'].value + doc['
17        sesso_femminile'].value;
18        emit(total);
19      ""
20    }
21  },
22 }
```

```

21 "aggs": {
22   "vaccinations_per_brand": {
23     "terms": {
24       "field": "fornitore"
25     },
26     "aggs": {
27       "total_vaccinations": {
28         "sum": {
29           "field": "vaccini"
30         }
31       }
32     }
33   }
34 },
35 "most-used_brand": {
36   "max_bucket": {
37     "buckets_path": "vaccinations_per_brand >
38     total_vaccinations"
39   }
40 }
41 }

```

2. find the region with the highest number of first vaccinations in the last 3 months.

The query is similar to the previous one. We have to filter the results by date through "query", and then exploit nested "aggs" to compute the sum of field "prima_dose" aggregated by region; finally we draw the maximum in the same way as before:

```

1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "range": {
6       "@timestamp": {
7         "gte": "now-90d/d",
8         "lt": "now/d"
9       }
10    }
11  },
12  "aggs": {
13    "first_vaccinations_per_region": {
14      "terms": {
15        "field": "area"
16      },
17      "aggs": {
18        "first_vaccinations": {
19          "sum": {
20            "field": "prima_dose"
21          }
22        }
23      }
24    },
25    "most-vaccinated_region": {
26      "max_bucket": {

```

```

27     "buckets_path": "first_vaccinations_per_region >
28         first_vaccinations"
29     }
30 }
31 }

```

3. find the total number of boosters in Italy in last 30 days.

This query is quite easy; indeed, we have to filter the results by date and then return the sum of field "dose_addizionale_booster":

```

1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "range": {
6       "@timestamp": {
7         "gte": "now-30d/d",
8         "lt": "now/d"
9       }
10    }
11  },
12  "aggs": {
13    "total_number_of_boosters": {
14      "sum": {
15        "field": "dose_addizionale_booster"
16      }
17    }
18  }
19 }

```

4. find the age range with the highest number of vaccinations in the last 30 days.

The pattern is similar to the first query; namely, we use a script after having filtered the results by date to compute the sum of male and female vaccine doses and finally we return the maximum value:

```

1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "range": {
6       "@timestamp": {
7         "gte": "now-30d/d",
8         "lt": "now/d"
9       }
10    }
11  },
12  "runtime_mappings": {
13    "vaccini": {
14      "type": "long",
15      "script": ""
16        long total = doc['sesso_maschile'].value + doc['sesso_femminile'].value;
17        emit(total);
18      ""
19    }
20  }
21 }

```



```

20 },
21 "aggs": {
22   "vaccinations_per_age_range": {
23     "terms": {
24       "field": "fascia_anagrafica"
25     },
26     "aggs": {
27       "total_vaccinations": {
28         "sum": {
29           "field": "vaccini"
30         }
31       }
32     }
33   }
34 },
35 "most-used_brand": {
36   "max_bucket": {
37     "buckets_path": "vaccinations_per_age_range >
38     total_vaccinations"
39   }
40 }
41 }

```

5. find the total number of Moderna doses received in the last 3 days in Italy.

We simply have to filter the data both by "fornitore" and date before computing the sum of field "numero_dosi":

```

1 GET /consegne_vaccini/_search
2 {
3   "size": 0,
4   "query": {
5     "bool": {
6       "filter": [
7         {"term": {
8           "fornitore": "Moderna"
9         }}
10      ],
11       "range": {
12         "@timestamp": {
13           "gte": "now-3d/d",
14           "lt": "now-1d/d"
15         }
16       }
17     }
18   }
19 },
20 "aggs": {
21   "vaccine_doses_received": {
22     "sum": {
23       "field": "numero_dosi"
24     }
25   }
26 }
27 }

```

6. **find the region with the highest disparity of vaccinations between males and females.**

This query is more tricky than the previous one because it involves a script in which we have to compute the absolute value of a difference (we do not know whether there are more males that undergo vaccination or females, and we do not want to return negative values). The rest of the query is similar to the others:

```
1 GET /somministrazioni_vaccini/_search
2 {
3   "size": 0,
4   "runtime_mappings": {
5     "vaccini": {
6       "type": "long",
7       "script": ""
8       long total = doc['sesso_maschile'].value - doc['
9       sesso_femminile'].value;
10      if(total < 0){
11        total = total * - 1;
12      }
13      emit(total);
14    }
15  },
16  "aggs": {
17    "vaccinations_per_region":{
18      "terms": {
19        "field": "nome_area"
20      },
21      "aggs": {
22        "max_dif_of_vaccinations_between_male_female": {
23          "max": {
24            "field": "vaccini"
25          }
26        }
27      }
28    },
29    "highest_disparity_of_vaccinations": {
30      "max_bucket": {
31        "buckets_path": "vaccinations_per_region >
32        max_dif_of_vaccinations_between_male_female"
33      }
34    }
35  }
```

7. **number of children between 5-11 vaccinated this week in Basilicata.**

DOES NOT WORK

```
1 GET /somministrazioni_vaccini/_count
2 {
3   "query": {
4     "bool": {
5       "must": [
6         { "match" :
7           "fascia_anagrafica": "5-11"
8         }
9       ]
10     }
11   }
```

```

8   },
9   { "range" : {
10    "@timestamp": {
11     "gte": "now-7d/d",
12     "lte": "now/d"
13    }
14   },
15   { "match" : {
16     "nome_area": "Basilicata"
17   }
18   }
19 }
20 }
21 }

```

8. to do.
DOES NOT WORK

9. brand of the less-used vaccine with young people (<30 years old).
DOES NOT WORK

```

1 GET /somministrazioni_vaccini/_search
2 {
3   "query": {
4     "bool": {
5       "should": [
6         { "match" :
7           "fascia_anagrafica": "5-11"
8         },
9         { "match" :
10           "fascia_anagrafica": "12-19"
11         },
12         { "match" :
13           "fascia_anagrafica": "20-29"
14         }
15       ]
16     },
17     "runtime_mappings": {
18       "vaccini": {
19         "type": "long",
20         "script": ""
21         long total = doc['sesso_maschile'].value + doc['sesso_femminile'].value;
22         emit(total);
23       }
24     },
25     "aggs": {
26       "vaccinations_per_region": {
27         "terms": {
28           "field": "nome_area"
29         },
30         "aggs": {
31           "total_vaccinations": {
32             "sum": {
33               "field": "vaccini"
34             }
35           }
36         }
37       }
38     }
39   }
40 }

```

```

35     }
36   },
37   "less-used_brand": {
38     "min_bucket": {
39       "buckets_path": "vaccinations_per_region >
total_vaccinations"
40     }
41   }
42 }

```

3.2 Commands

We have identified the following commands to show how the system works:

1. **Update one specific document (data of Abruzzo 2020/03/30 for age range 40-49 for Moderna) with specific data (100 males vaccinated and 2012 second doses).**

The construct we use is the "_update_by_query" one, which can be used with the POST method; we have to perform a query to find the document to update, and then we exploit a script (with parameters) to update the document:

```

1 POST /somministrazioni_vaccini/_update_by_query
2 {
3   "query": {
4     "bool": {
5       "filter": [
6         {"term": {"area": "ABR"}},
7         {"term": {"fascia_anagrafica": "40-49"}},
8         {"term": {"fornitore": "Moderna"}},
9         {
10          "range": {
11            "@timestamp": {
12              "time_zone": "+02:00",
13              "gte": "2020-03-30T00:00:00",
14              "lt": "2020-03-31T00:00:00"
15            }
16          }
17        }
18      ]
19    }
20  },
21  "script": {
22    "source": "ctx.sesso_maschile = params.sesso_maschile; ctx
.seconda_dose = params.seconda_dose",
23    "lang": "painless",
24    "params": {
25      "sesso_maschile": 100,
26      "seconda_dose": 2012
27    }
28  }
29 }

```

2. Create one new document.

This command is the easiest; we simply have to make a `POST` with `"_doc"`, and then specify a value for all the fields:

```
1 POST /somministrazioni_vaccini/_doc
2 {
3   "data_somministrazione": "2021-12-28",
4   "fornitore": "Moderna",
5   "area": "BAS",
6   "fascia_anagrafica": "30-39",
7   "sesso_maschile": 910,
8   "sesso_femminile": 1211,
9   "prima_dose": 93,
10  "seconda_dose": 563,
11  "pregressa_infezione": 33,
12  "dose_addizionale_booster": 1465,
13  "codice_NUTS1": "ITF",
14  "codice_NUTS2": "ITF5",
15  "codice_regione_ISTAT": 17
16 }
```

3. Add a runtime field.

Runtime fields are fields evaluated at query time. They allow us to add fields to existing documents without reindexing data and to define additional fields without modifying the schema. Here, a useful runtime field might be the sum of the total doses, since in each record the doses are grouped by type of administration. By creating the runtime field `dosi_totali`, we add an additional field to each record containing the desired information and to facilitate the creation of other queries. To create a runtime field, we have to use a script along with the `PUT` method, specifying the name of the field, its type and its source:

```
1 PUT somministrazioni_vaccini/_mapping
2 {
3   "runtime": {
4     "dosi_totali": {
5       "type": "long",
6       "script": {
7         "source": "emit(doc['prima_dose'].value + doc['seconda_dose'].value + doc['pregressa_infezione'].value + doc['dose_addizionale_booster'].value);"
8       }
9     }
10  }
11 }
```

4 Kibana Dashboard

By using the Kibana tool we realised a simple dashboard containing some relevant graphs to visualise the data provided by the dataset. The dashboard displays the following information and graphs:

- total number of administered vaccines;
- administered doses of vaccines per region;
- administered doses per supplier;
- administered doses per week, grouped by:
 - supplier;
 - age range;
 - gender.
- total number of delivered doses;
- delivered doses per region;
- delivered doses per week.

4.1 Dashboard initialization

To create the Kibana dashboard:

- Install the two datasets as it is explained in Section 3 (remember to name them `somministrazioni_vaccini` and `consegne_vaccini`);
- From the Dev Tools, run the command 3 (Add a runtime field), since some visualizations require the field `dosi_totali`;
- In the Maps tool, select "Add layer" and then "Upload GeoJSON". Then upload the `.geo.json` file provided in the exports section and name the index "italyregions";
- Go to Stack Management (you can leave the Maps page without saving) and then to Saved Objects. Then upload the `.ndjson` file provided in the exports section. It is not necessary to edit any import setting;
- Go to the Dashboard section and open the Italy Vaccines dashboard;
- Set a relevant time period (e.g. January 1st, 2021 - December 27th, 2021).

4.2 Dashboard visualizations

In this section are provided screenshots of the various visualizations created on the Kibana Dashboard.

5 Cassandra

An alternative way to implement the vaccine administrations dataset can be to store it in a columnar database such as Cassandra. Columnar databases are particularly suitable for very large datasets and for analytical purposes such as aggregation operations on entire columns, such as retrieving the total number of administered doses.

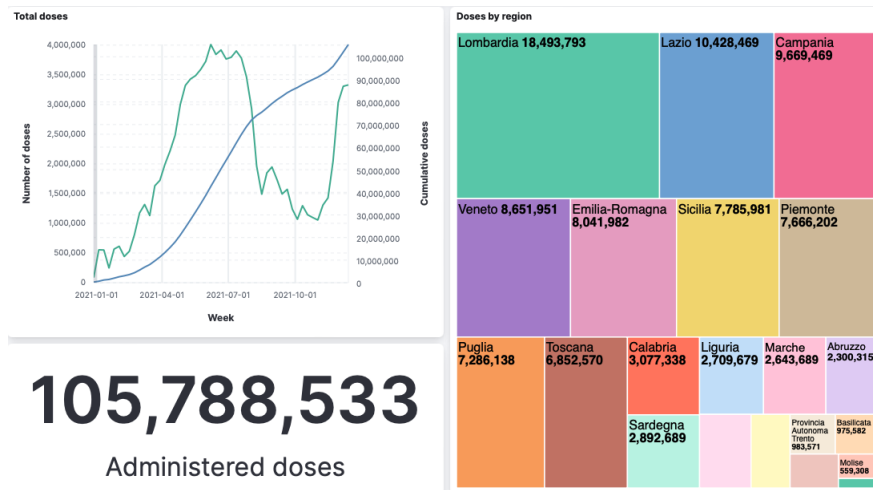


Figure 1: Overview of the vaccine doses administrations

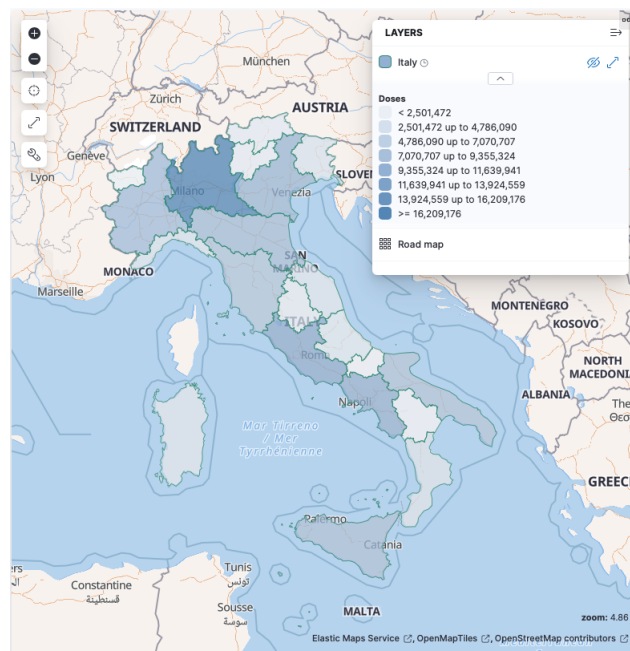


Figure 2: Doses administration by region, visualized on the Italian territory

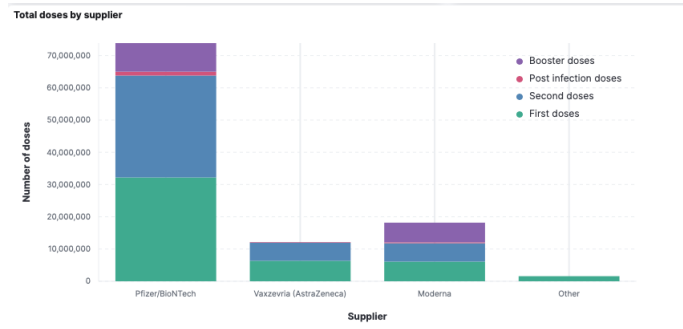


Figure 3: Administered doses per supplier

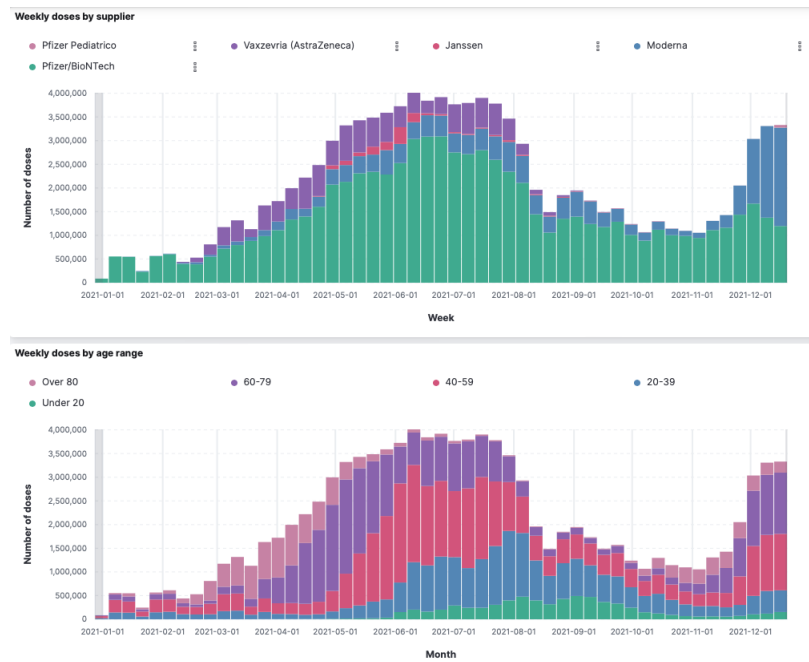


Figure 4: Weekly administered doses, grouped by supplier and age range

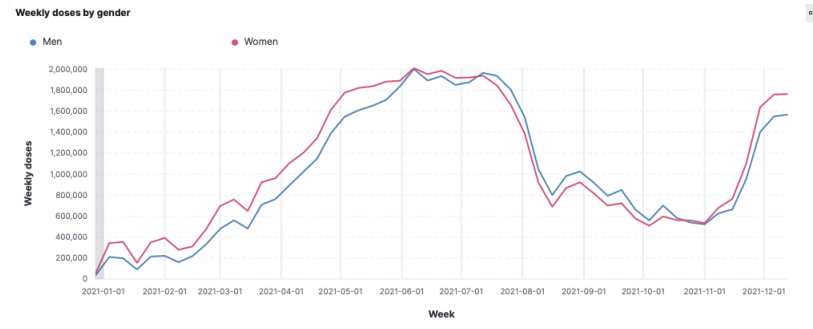


Figure 5: Weekly administered doses, grouped by gender

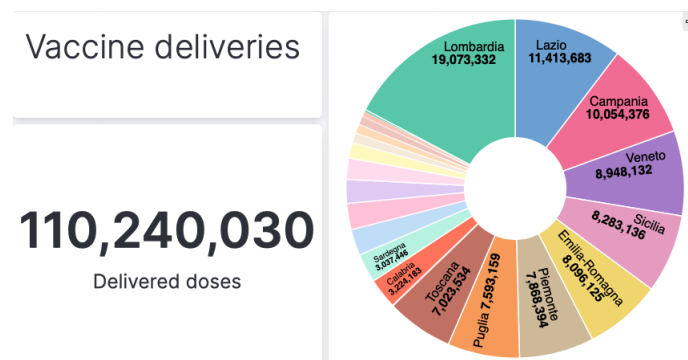


Figure 6: Overview of the vaccine deliveries

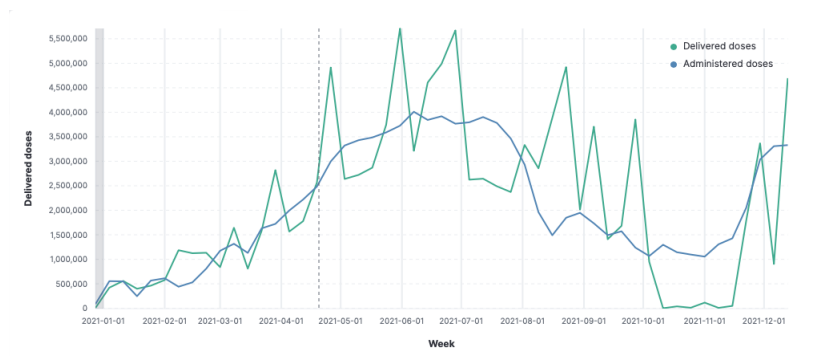


Figure 7: Weekly delivered/administered doses comparison

5.1 Dataset implementation

To implement the dataset in Cassandra, first of all start Cassandra, then move to the dataset folder and start the cqlsh console.

- Create a keyspace:

```
1 CREATE KEYSPACE somministrazioni WITH REPLICATION = {'  
  class': 'SimpleStrategy', 'replication_factor':3};  
2
```

- Select the keyspace:

```
1 USE somministrazioni;  
2
```

- Create a table:

```
1 CREATE TABLE somministrazioni (  
2     data_somministrazione date,  
3     fornitore text,  
4     area text,  
5     fascia_anagrafica text,  
6     sesso_maschile varint,  
7     sesso_femminile varint,  
8     prima_dose varint,  
9     seconda_dose varint,  
10    pregressa_infezione varint,  
11    dose_addizionale_booster varint,  
12    codice_NUTS1 text,  
13    codice_NUTS2 text,  
14    codice_regione_ISTAT varint,  
15    nome_area text,  
16    PRIMARY KEY (codice_regione_istat, fornitore,  
17    fascia_anagrafica, data_somministrazione)  
18    );
```

- Import the dataset:

```
1 COPY somministrazioni (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 '<.csv dataset>' WITH HEADER = false;  
2
```

- Create an index on nome_area:

```
1 CREATE INDEX regione ON somministrazioni (nome_area);  
2
```

Now the dataset is ready to be queried.

5.2 Queries examples

Here are some examples of possible queries in CQL. Note that in aggregation queries, it is not possible to order the results by the aggregated value since in Cassandra ordering is possible only on clustering columns.

- Total number of administered doses:

```
1      SELECT sum( Sesso_maschile)+sum( Sesso_femminile) AS
      total_doses FROM somministrazioni;
```

- Administered doses per supplier in June, 2021:

```
1      SELECT fornitore, sum( Sesso_maschile)+sum(
      Sesso_femminile) AS total_doses FROM somministrazioni
      WHERE data_somministrazione >= '2021-06-01' AND
      data_somministrazione <= '2021-06-30' GROUP BY
      codice_regione_ISTAT, fornitore ALLOW FILTERING;
```

- Total number of booster doses in Lombardia since December 1st, 2021:

```
1      SELECT sum(dose_addizionale_booster) AS booster_doses
      FROM somministrazioni WHERE data_somministrazione >=
      '2021-12-01' AND nome_area = 'Lombardia' ALLOW FILTERING;
```

6 References and sources

- [Elasticsearch](#)
- [Kibana](#)
- [Cassandra](#)
- [Dataset](#)
- [draw.io](#)