# Politecnico di Milano

Systems and methods for big and unstructured data

# SMBUD project: second 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 design and store the information needed to support a certification app for COVID-19.

The database includes detailed information about each person, vaccine and authorized body that administers the vaccine doses.

For each individual, their personal data must be recorded (first and last name, telephone number, email and birthdate) along with information about their residence (address, country, city, etc...). Every person stored is identified by a fiscal code, that can be deduced from the data listed above.

The authorized bodies, in charge of administering the vaccines, can be of different types, going from hospitals to vaccination hub and they are further divided into different departments. The data stored regards their name, type, departments ("Immunology", "Epidemiology", "Covid Emergency") and location.

Four vaccines are stored and for each of them the dataset maintains their name, minimum number of doses to be effective, type, website and even some information about the manufacturer (name, website and telephone number).

People can get vaccine shots or take tests to obtain a certificate that verifies their negativity to a Covid-19 infection.

When a person gets a vaccine shot, the database must record the vaccine inoculated, the dose number, the authorized body which delivered the vaccine, time and date the vaccine was delivered, the lot number, the vaccine production date, the staff members who performed the shot and a single emergency contact designated by the patient (name and telephone would suffice).
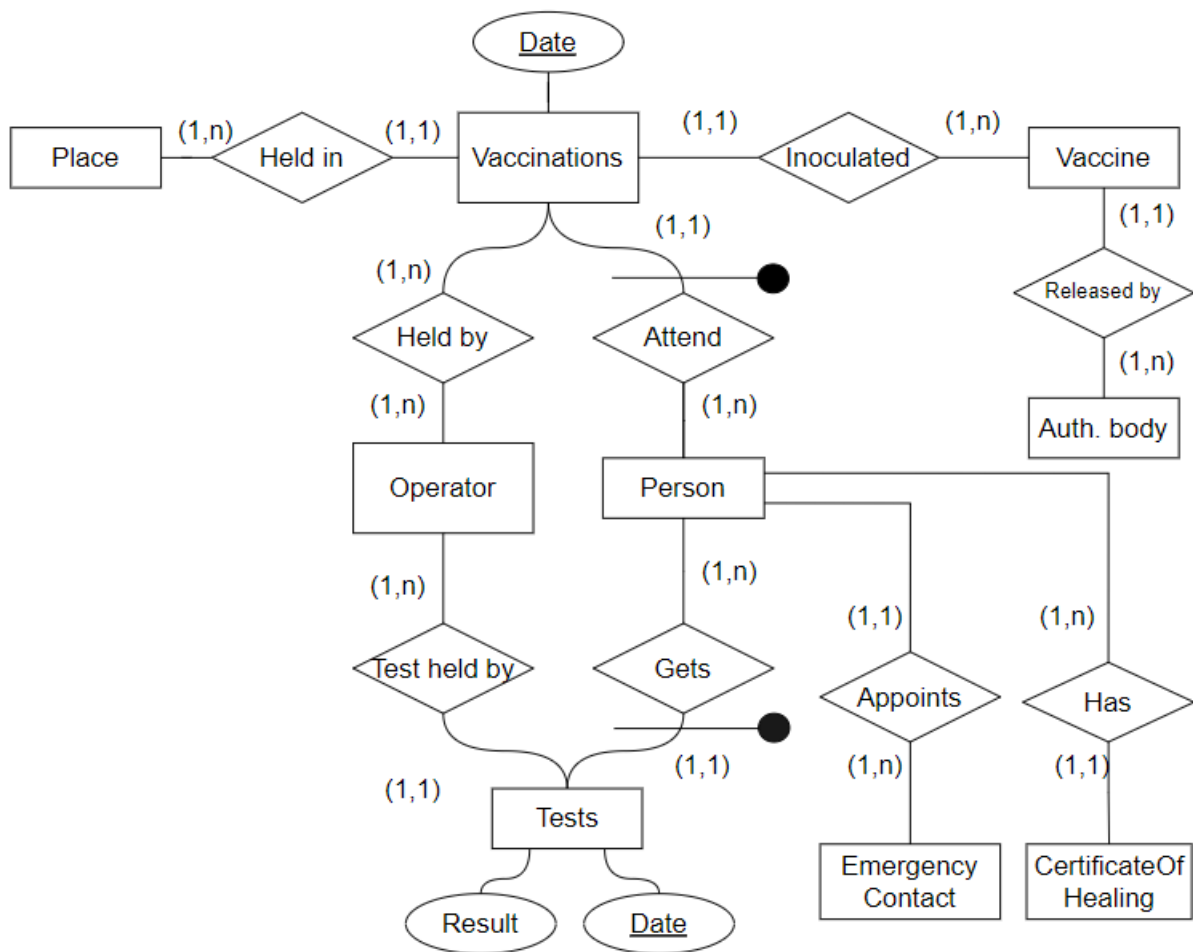
When a person does a test, it must record: the authorized body which performed the test, time and date test was performed, the operator that performed the test, the type (rapid, molecular, or serologic), the result and again a single emergency contact in case the test result is positive, and the person test is unreachable.

A certificate can be also obtained after a negative test result subsequently to a Covid-19 infection with a duration on 6 months.

## Check Validity

It is the duty of the verification apps to decide whether a certificate is valid or not. This is a mandatory requirement because the rules of when a certificate is valid or not will likely change in the future (as they did in the past) and it is not feasible to edit all certificates. For this reason, updating the certificate verification app must be forced. In this particular case the validity of a certificate is proved by checking the date a vaccine was administered, a test was taken or a negative result after a recovery from the infection has been delivered.

# Conceptual Model



The ER model translates into two collections in the MongoDB database, one for the certificates and the other for the authorized bodies. In APPENDIX A there are some tables which defines the structure of documents.

# Dataset Generation

We created a python script which uses a library that interacts with randomuser.me API to generate random people data.

The script connects directly to a database using PyMongo and performs queries against it. At the end the script also exports the created database as json files.

In the script we first generate first a fixed set of vaccines, places, and staff members and then we populate the authorized collection with random data. Lastly, we populate the certificates collection with new random data and the data created before (vaccines, places, staff members, authorized bodies).

## Assumptions

- We supposed to create the dataset for the full 2021, like if the dataset was created the 31st of December. This means that, for example, some people could not have yet completed the vaccination cycle because their second dose is scheduled for the 2022. No vaccinations nor tests performed in a year different from 2021 are present.

- If a person starts the vaccination cycle, then it will end it (no dropouts after the first dose).

- If a person results positive, it will take another test 10 days after.

- There are 3000 certificates.

- There are 100 places where people can perform tests or vaccinations.

- There are 500 staff members that can perform vaccinations and tests.

- There are 50 authorized bodies which issue the certificates.

- All the doses are made with the same vaccine.

- Some people get a booster dose.

In APPENDIX B all parameters used for the dataset generation are reported.

## Importing the dataset

In the delivery there is a folder called *dumps* where each json file corresponds to a collection. To import the dataset using MongoDB Compass, it is needed to create a collection in the database for each file and then using the import utility selecting the corresponding json as source.

# Queries

1. **Get top 10 places of vaccination (hospitals, vaccination centers, etc...) which delivered more vaccination doses**

```
db.certificates.aggregate([

// Filter people with a vaccination

{$match: {

'vaccinations.0': {$exists: 'True'}

}},

// Unwind vaccinations

{$unwind: {

path: '$vaccinations'

}},

// Group by place name

{$group: {

_id: '$vaccinations.place.name',

nVaccinations: {$sum: 1}

}},

// Sort by number of vaccinations DESC

{$sort: {

nVaccinations: -1

}},

// Get top 10

{$limit: 10}

])
```

## 2. Calculate the percentage of every vaccine used

```
1      db.certificates.aggregate([
2
3  // Match vaccinated people
4
5  {$match: {
6
7  'vaccinations.0': {$exists: 'True'}
8
9  }},
10
11  // Filter only the first dose
12
13  { $project: {
14
15  "vaccinations": { $first: "$vaccinations" }
16
17  }},
18
19  // Count by vaccine name
20
21  { $group: {
22
23  _id: "$vaccinations.vaccine.name",
24
25  nPeople: { $sum: 1 }
26
27  }},
28
29  // Find total and push counts to a vaccines array
30
31  {$group: {
32
33  _id:null,
34
35  total: {$sum:"$nPeople"},
36
37  vaccines: { $push: {name: "$_id", nPeople:"$nPeople"} }
38
39  }},
40
41  // Unwind vaccines array
42
43  {$unwind: {
44
45  path: "$vaccines"
46
47  }},
48
49  // Find percentage and project
50
51  {$project: {
52
53  _id:"$vaccines.name",
54
55  "percentage": {$multiply:[{$divide:["$vaccines.nPeople","$total"]},100]}
56
57  }},
```

```
58
59  ])
```

3. **Top 3 authorized bodies per country that delivered more doses**

```
1       db.certificates.aggregate([
2
3   // Filter people who got a vaccination
4
5   { $match: {
6
7   'vaccinations.authorized_body_id': {$exists: 'True'}
8
9   }},
10
11  // Unwind vaccinations
12
13  { $unwind: '$vaccinations' },
14
15  // Count vaccination doses by authorized_body_id
16
17  { $group: {
18
19  _id: '$vaccinations.authorized_body_id',
20
21  nVaccinationDoses: { $sum: 1 }
22
23  }},
24
25  // Sort by number of vaccination doses
26
27  { $sort: {
28
29  nVaccinationDoses: -1
30
31  }},
32
33  { $lookup: {
34
35  from: 'authorized_bodies',
36
37  localField: '_id',
38
39  foreignField: '_id',
40
41  as: 'authorized_body'
42
43  }},
44
45  // Unwind authorized bodies
46
47  { $unwind: '$authorized_body' },
48
49  // Group by country
50
51  { $group: {
52
53  _id: '$authorized_body.location.country',
54
```

```
55  authorized_bodies: {$push: {
56
57  _id: "$authorized_body._id",
58
59  name: "$authorized_body.name",
60
61  nVaccinationDoses: "$nVaccinationDoses"
62
63  }}
64
65  }},
66
67  // Keep only top 3
68
69  {$project: {
70
71  _id: "$_id",
72
73  "authorized_bodies": {$slice: ["$authorized_bodies", 3]}
74
75  }}
76
77  ])
```

4. **Determines if Leroy Payne (fiscalCode: "LRYPYN98A30R790Z") has a valid certificate, supposing that a certificate is valid if he is not currently positive and one of the following conditions is true:**

   - He got the first dose of a vaccine more than 15 days ago (before 2021-12-25)
   - He got a full vaccination cycle less than 9 months ago (270 days) and more than 15 days ago (from 2021-04-05 to 2021-12-16)
   - He took a negative test less than 48h ago (after 2021-12-30T00:00:00)
   - Healed from COVID less than 6 months ago (180 days, after 2021-06-04)

```
1     db.certificates.aggregate([
2
3   // Match the right person
4
5   {
6
7     $match: {
8
9       fiscal_code: "LLYLMA84D02P040R", // LRYPYN98A30R790Z
10
11    },
12
13   },
14
15   // Add useful fields
16
17   { $addFields: { n_doses: { $size: "$vaccinations" } } },
18
19   {
20
21     $addFields: { total_doses: { $max: "$vaccinations.vaccine.
         total_doses" } },
22
```

```
},

{ $addFields: { last_dose: { $first: "$vaccinations" } } },

{ $addFields: { last_test: { $first: "$tests" } } },

// Project

{

  $project: {

    fiscal_code: 1,

    n_doses: 1,

    total_doses: 1,

    last_dose: 1,

    last_test: 1,

    tests: 1,

    healings: 1,

  },

},

// Match not currently positive

{

  $match: {

    $or: [

      { "last_test.result": "negative" },

      { last_test: { $exists: false } },

    ],

  },

},

// Match conditions

{

  $match: {

    $or: [

      // 1. He got the first dose of a vaccine more than 15 days ago (
         before 2021-12-25)

      {
```

```
 82
 83              "last_dose.datetime": {
 84
 85                $lt: ISODate("2021-12-25T00:00:00.000Z"),
 86
 87              },
 88
 89          },
 90
 91          // 2. He got a full vaccination cycle less than 9 months ago
                  (270 days) and more than 15 days ago
 92
 93          {
 94
 95            $and: [
 96
 97              { $expr: { $gte: ["$n_doses", "$total_doses"] } },
 98
 99              {
100
101                "last_dose.datetime": {
102
103                  $gte: ISODate("2021-04-05T00:00:00.000Z"),
104
105                },
106
107              },
108
109              {
110
111                "last_dose.datetime": {
112
113                  $lte: ISODate("2021-12-17T00:00:00.000Z"),
114
115                },
116
117              },
118
119            ],
120
121          },
122
123          // 3. He took a negative test less than 48h ago
124
125          {
126
127            "last_test.datetime": {
128
129              $gte: ISODate("2021-12-30T00:00:00.000Z"),
130
131            },
132
133          },
134
135          // 4. Healed from COVID less than 6 months ago (180 days)
136
137          {
138
139            healings: {
140
```

```
141            $elemMatch: {
142
143               end: {
144
145                  $gte: ISODate("2021-06-04T00:00:00.000Z"),
146
147               },
148
149            },
150
151         },
152
153       },
154
155     ],
156
157   },
158
159  },
160
161 ]);
```

5. **Find fiscal code, first name and last name of the people who got a booster dose** (an additional dose after the normal vaccination cycle)

```
 1     db.certificates.aggregate([
 2
 3 // Match people who got a vaccination
 4
 5 {$match: {
 6
 7 'vaccinations.0': {$exists: 'True'}
 8
 9 }},
10
11 // Project interesting fields
12
13 {$project: {
14
15 "first_name": 1,
16
17 "last_name": 1,
18
19 "fiscal_code": 1,
20
21 "vaccinations": 1
22
23 }},
24
25 // Match people who did more doses than the vaccine.total_doses
26
27 {$match: {
28
29 $expr: { $gt: [
30
31 {$size: '$vaccinations'},
32
33 {$max: '$vaccinations.vaccine.total_doses'}
34
```

```
35  ]}

36
37  }},

38
39  // Project required fields

40
41  {$project: {

42
43  "first_name": 1,

44
45  "last_name": 1,

46
47  "fiscal_code": 1

48
49  }},

50
51  ])
```

# Commands

1. **Add a certificate to the database**

```
db.certificates.insertOne({

    first_name: "Mario",

    last_name: "Rossi",

    telephone: "051-806-7094",

    email: "mario.rossi@example.com",

    birth_datetime: ISODate("1984-10-22T06:10:00.000Z"),

    location: {

        address: "via delle primule, 56",

        country: "Italy",

        city: "Rovereto",

        zip_code: 26010,

        coordinates: {

            latitude: "45.8905",

            longitude: "11.0397",

        },

    },

});
```

2. **Add an authorized body to the database**

```
1      db.authorized_bodies.insertOne({
2
3    name: "provincial health office 1",
4
5    type: "provincial health office",
6
7    telephone: "10-988-7777",
8
9    location: {
10
11       address: "via delle azalee 88",
12
13       country: "Italy",
14
15       city: "Sermoneta",
16
17       zip_code: 04013,
18
19    },
20
21 });
```

3. **Add a test to Mario Rossi** (uuid: "61b63f7f98b9ace842754bf0")

```
1      db.authorized_bodies.insertOne({
2
3    name: "provincial health office 1",
4
5    type: "provincial health office",
6
7    telephone: "10-988-7777",
8
9    location: {
10
11       address: "via delle azalee 88",
12
13       country: "Italy",
14
15       city: "Sermoneta",
16
17       zip_code: 04013,
18
19    },
20
21 });
```

4. **Add a vaccination to Mario Rossi** (uuid: "61b63f7f98b9ace842754bf0")

```
db.certificates.updateOne(

{ _id: ObjectId("61b63f7f98b9ace842754bf0") },

{

    $push: {

        vaccinations: {

            datetime: ISODate("2021-09-29T08:04:45.000Z"),

            vaccine: {

                name: "COVID-19 Vaccine Janssen",

                total_doses: 1,

                type: "viral_vector",

                ema_url:

                    "https://www.ema.europa.eu/en/medicines/human/
    EPAR/covid-19-vaccine-janssen",

                manufacturers: [

                    {

                        name: "Johnson & Johnson",

                        website: "https://www.jnj.com/",

                        telephone: "(732) 524-0400",

                    },

                ],

                lot: 385,

                production_date: ISODate("2021-03-11T11:55:09.000Z")
    ,

            },

            staff: [

                {

                    first_name: "Roseneide",

                    last_name: "Porto",

                    role: "nurse",

                    telephone: "(06) 7086-2238",
```

```
56
57                        email: "roseneide.porto@example.com",
58
59                        birth_date: ISODate("1975-09-15T00:00:00.000Z"),
60
61                        location: {
62
63                            address: "2397 Rua Sete de Setembro ",
64
65                            country: "Brazil",
66
67                            city: "Salvador",
68
69                            zip_code: 43748,
70
71                            latitude: 80.7739,
72
73                            longitude: -170.5816,
74
75                        },
76
77                    },
78
79                ],
80
81                authorized_body_id: ObjectId("61b638eead9e15bece773fb7")
   ,
82
83                place: {
84
85                    name: "hospital 43",
86
87                    type: "hospital",
88
89                    telephone: "0179-3597946",
90
91                    location: {
92
93                        address: "440 Gr ner Weg",
94
95                        country: "Germany",
96
97                        city: "Crimmitschau",
98
99                        zip_code: 26168,
100
101                        latitude: 18.2513,
102
103                        longitude: -34.8376,
104
105                    },
106
107                },
108
109        },
110
111    },
112
113  }
114
```

```
115   );
```

# UI

| Home | Authorized Bodies | Add Certificate | Statistics | Fiscal Code |
|------|-------------------|-----------------|------------|-------------|

## Authorized Bodies

Click on the name of one authorized body to see data about all people vaccinated there.

| name | type | telephone |   | first_name | last_name | telephone |
|------|------|-----------|---|------------|-----------|-----------|
| regional health office 1 | regional health office | (061)-766-7561 |   | Theodore | Smith | (007)-196-7399 |
| provincial health office 1 | provincial health office | 0464-439-370 |   | Simon | Johansen | 22392556 |
| provincial health office 2 | provincial health office | 0908-498-5901 |   | Marielle | Konink | (248)-899-7130 |
| national health office 1 | national health office | (28) 3491-4241 |   | Lily | Roger | 06-68-25-02-20 |
| provincial health office 3 | provincial health office | 0449-598-514 |   | Rick | Ortiz | (826)-832-9131 |
| regional health office 2 | regional health office | 076 893 99 32 |   | Silas | Moraes | (77) 1951-5314 |
| provincial health office 4 | provincial health office | (581)-722-6245 |   | Eddie | Williams | 081-978-4555 |
| regional health office 3 | regional health office | 615-149-3754 |   | Verônica | da Paz | (39) 3102-4975 |
| regional health office 4 | regional health office | 43681683 |   | Michal | Lied | 43544704 |
| national health office 2 | national health office | 698-219-734 |   | Carla | Moreno | 691-272-391 |
| provincial health office 5 | provincial health office | (950)-194-7172 |   | Malik | Reistad | 40056823 |

| Home | Authorized Bodies | Add Certificate | Statistics | Fiscal Code |
|------|-------------------|-----------------|------------|-------------|

## Covid reports

**Report to generate**

Top places of vaccination  ⌄

Generate



Most used places to get a vaccination

# Appendix A

| Field | Type | Example |
|---|---|---|
| first_name | String | Mae |
| last_name | String | Meunier |
| telephone | String | 06-24-89-52-61 |
| email | String | mae.meunier@example.com |
| birth_date | Date | 1969-05-23T00:00:00.000+00 |
| fiscal_code | String | MAEMNR69M23M826B |
| Location | LOCATION | Location object (see correspon |
| emergency_contact | Object | |
| emergency_contact.first_name | String | Arthur |
| emergency_contact.last_name | String | Kowalski |
| emergency_contact.telephone | String | 311-427-2628 |
| vaccinations | Array | |
| vaccinations.datetime | Timestamp | 2021-10-28T17:31:37.000+00: |
| vaccinations.staff | Array[STAFF] | |
| vaccinations.authorized_body_id | ObjectId | |
| vaccinations.place | PLACE | |
| vaccinations.vaccine | Object | |
| vaccinations.vaccine.name | String | Comirnaty |
| vaccinations.vaccine.total_doses | Int32 | 2 |
| vaccinations.vaccine.type | String | mRNA |
| vaccinations.vaccine.ema_url | String | https://www.ema.europa.eu/en |
| vaccinations.vaccine.lot | Int32 | 226 |
| vaccinations.vaccine.production_date | Date | 2021-08-10T10:39:36.000+00: |
| vaccinations.vaccine.manufacturers | Array | |
| vaccinations.vaccine.manufacturers.name | String | Pfizer Inc. |
| vaccinations.vaccine.manufacturers.website | String | https://www.pfizer.com/ |
| vaccinations.vaccine.manufacturers.telepho | String | 1-800-879-3477 |
| healings | Array | |
| healings.start | Timestamp | 2021-08-21T17:12:29.000+00: |
| healings.end | Timestamp | 2021-08-31T09:33:18.000+00: |
| healings.authorized_body_id | ObjectId | 61b5cbc4318ef70857a7bddb |
| tests | Array | |
| tests.datetime | Timestamp | 2021-08-31T09:33:18.000+00: |
| tests.result | String | negative |
| tests.staff | STAFF | |
| tests.authorized_body_id | ObjectId | 61b5cbc4318ef70857a7bddb |
| tests.place | PLACE | |

| Field | Type | Example |
|---|---|---|
| location.address | String | 1686 Skibhusvej |
| location.country | String | Denmark |
| location.city | String | Brondby |
| location.zip_code | Int32 | 83109 |
| location.latitude | Double | 33.9805 |
| location.longitude | Double | 104.3353 |

| Field | Type | Example |
|---|---|---|
| place.name | String | hospital 28 |
| place.type | String | hospital |
| place.telephone | String | 46633601 |
| place.location | LOCATION | |

| Field | Type | Example |
|---|---|---|
| staff.first_name | String | Matilda |
| staff.last_name | String | Kari |
| staff.role | String | nurse |
| staff.telephone | String | 043-730-56-60 |
| staff.email | String | matilda.kari@example.com |
| staff.birth_date | Date | 1954-11-03T00:00:00.000+00: |
| staff.location | LOCATION | |

| Field | Type | Example |
|---|---|---|
| name | String | provincial health office 1 |
| type | String | provincial health office |
| department | String | epidemiology |
| telephone | String | 655-087-351 |
| location | LOCATION | |

# Appendix B

| Name | Value | Description |
|---|---|---|
| CERTIFICATES_NUMBER | 1000 | Number of certificates |
| DATE_START | 01/01/2021 | Date start of generation |
| DATE_END | 31/12/2021 | Date end of generation |
| HOUR_START | 8 | Hour when vaccination and tests can begin |
| HOUR_END | 18 | Hour when vaccination and tests end |
| EMERGENCY_CONTACTS_POOL_SIZE | 100 | Number of unique emergency contacts |
| LOT_MIN | 100 | Lot number min |
| LOT_MAX | 999 | Lot number max |
| PLACES_NUMBER | 100 | Number of unique places where to vaccinate and test |
| PLACES_TYPES | { "hospital": 0.6, "vaccination_center": 0.35, "medical_guard": 0.05, } | Place possible types and the probability of being that type |
| PRODUCTION_DELTA_MIN | 45 | Min number of days when a vaccine get produced and used |
| PRODUCTION_DELTA_MAX | 100 | Max number of days when a vaccine get produced and used |
| PERCENTAGE_VACCINATED | 0.8458 | Probability of being vaccinated |
| PERCENTAGE_BOOSTER_DOSE | 0.1 | Probability of taking a booster dose if vaccinated |
| VACCINATION_STAFF_MIN | 1 | Min number of staff members needed to vaccinate |
| VACCINATION_STAFF_MAX | 3 | Max number of staff members needed to vaccine |
| VACCINATION_SECOND_DOSE_DELTA_MIN | 25 | Min number of days to wait for the second dose |
| VACCINATION_SECOND_DOSE_DELTA_MAX | 35 | Max number of days to wait for the second dose |

| | | |
|---|---|---|
| NUMBER_OF_TESTS _MIN | 0 | Min number of tests that a person performs. If a person is not vaccinated then he performs at least one test no matter of this parameter. This number include the "initial" tests, if a test is positive and more tests are needed they will not count thowards this parameter. |
| NUMBER_OF_TESTS _MAX | 3 | Max number of tests that a person performs. This number include the "initial" tests, if a test is positive and more tests |
| POSITIVE_TEST_DEL TA_DAYS | 10 | How many days after a positive test another one is required. |
| PERCENTAGE_POSIT IVE_TESTS_VACCINA TED | 0.15 | Possibility of a test being positive if the person is vaccinated |
| PERCENTAGE_POSIT IVE_TESTS_UNVACCI NATED | 0.3 | Possibility of a test being positive if the person is not vaccinated |
| STAFF_NUMBER | 100 | Number of unique staff members, they will be doctors and |
| STAFF_DOCTORS_PE RCENTAGE | 0.25 | Percentage of doctors |
| AUTHORIZED_BODIE S_NUMBER | 10 | Number of unique authorized bodies |
| AUTHORIZED_BODIE S_TYPES | {     "provincial health office": 0.4,     "regional health office": 0.4,     "national health office": 0.2,  } | Authorized bodies type and probability of being of that type |
| AUTHORIZED_BODIE S_DEPARTMENTS | {     "immunology": 0.4,     "epidemiology": 0.4,     "covid-19 emergency": 0.2,  } | Authorized bodies departments and probability of being of that department |