

Aplicación big data sobre el efecto del coronavirus en el mundo

Bases de Datos Avanzadas

Escuela Superior Informática (UCLM)

Mario Pérez Sánchez-Montañez, David Camuñas Sánchez, Francesco Zingariello

7 de mayo de 2020

Resumen

Este proyecto tratará sobre la creación de una aplicación *big data*, para la cual se realizará un estudio sobre el efecto del coronavirus (**COVID-19**) en el mundo. Para ello se usará las diversas técnicas utilizadas para el tratamiento masivo de información

Keywords: Big Data, MongoDB, Azure Cosmos DB

Introducción

El **Big Data** (*macro datos, datos masivos, inteligencia de datos*), se conoce como el análisis masivo de datos. Esto es una cantidad de datos tan sumamente grande, que las aplicaciones que se dedicaban al procesamiento de datos que tradicionalmente se venían usando no son capaces de tratar y poner en valor en un tiempo razonable. Este término ha estado en uso desde la década de 1990.

Este término también hace referencia, a las nuevas tecnologías que hacen posible el almacenamiento y procesamiento de datos, además del uso que se hace de la información obtenida a través de ellas.

Tipos de datos que existen

A continuación, se dará una breve explicación de los tipos de datos existentes, y de los cuales hace uso esta tecnología (*Big Data*).

Datos estructurados

Este tipo de datos se suelen usar en el tratamiento de datos. Ya que sus principales características son que pueden ser almacenados en tablas, y que

tienen definido su formato.

Algunos datos de este tipo son:

- **Números.**
- **Cadenas de caracteres.**
- **Fechas.**

Datos semiestructurados

Los *datos semiestructurados* tienen una tipo de estructura, pero esta no es lo suficientemente regular, como para permitir su gestión y estructuración como si fuese similar a los datos estructurados. Este tipo de datos, posee ciertos patrones comunes que lo describen y dan información sobre las relaciones entre los mismos.

Un ejemplo de utilización, sería el lenguaje de marcas *HTML*, el cual sirve para el desarrollo de paginas web, donde estas pautas son representadas por su sistema de etiquetas.

Datos no estructurados

Este tipo de datos, se trata de datos que han sido obtenidos en su formato original. Estos no se encuentran especificados en ningún formato, el cual permita que sean almacenados como se almacenarían los anteriores tipos de datos (*estructurados* y *semiestructurados*). Ya que no disponen de una estructura definida (longitud, formato, etc) para poder desglosar la información.

Estos datos se pueden encontrar, por ejemplo: en presentaciones (*PowerPoints*), correos, documentos de texto (*Word*, *PDF*, *documentos de google*, etc).

Obtención de los datos

Los datos utilizados se han obtenido de *EU Open Data Portal* [2]. El conjunto de datos contiene los últimos datos públicos disponibles sobre COVID-19, incluida una actualización diaria de la situación, la curva epidemiológica y la distribución geográfica global (UE / EEE y el Reino Unido, en todo el mundo).

Para asegurar la precisión y confiabilidad de los datos, este proceso se refina constantemente. Esto ayuda a controlar e interpretar la dinámica de la pandemia de COVID-19 no solo en la Unión Europea (UE), el Espacio Económico Europeo (EEE), sino también en todo el mundo. Todos los días, entre las 6.00 y las 10.00 CET, un equipo de epidemiólogos examina hasta 500 fuentes relevantes para recopilar las últimas cifras. El análisis de datos es seguido por el proceso estándar de inteligencia epidémica del ECDC para el cual cada entrada de datos se valida y documenta en una base de datos del ECDC [6].

Creación de la base de datos

La base de datos se ha construido utilizando *MongoDB* [7] y *Azure Cosmos DB* [1]. De esta forma la base de datos es accesible desde cualquier parte.

MongoDB es una base de datos documental, el elemento esencial es el documento que normalmente se los agrupa en colecciones de documentos similares. Una base de datos en MongoDB es un conjunto de colecciones.

MongoDB almacena datos en documentos flexibles, similares a JSON, lo que significa que los campos pueden variar de un documento a otro y la estructura de datos se puede cambiar con el tiempo. El modelo de documento se asigna a los objetos en el código de su aplicación, lo que facilita el trabajo con los datos.

En nuestro caso, la base datos estará formada por una sola colección y el número de documentos varía con el tiempo, pero actualmente contiene más de 15000. La estructura de un documento se puede observar en la figura 1

Fields	Global Probability	Type
[covid.Worldwide]	100.0%	Collection
▶ _id	100.0%	ObjectId
▶ cases	100.0%	Int32
▶ continentExp	100.0%	String
▶ countriesAndTerritories	100.0%	String
▶ countryterritoryCode	100.0%	String
▶ dateRep	100.0%	Date
▶ day	100.0%	Int32
▶ deaths	100.0%	Int32
▶ geold	100.0%	String
▶ month	100.0%	Int32
▶ popData2018	100.0%	Double
▶ year	100.0%	Int32

Figura 1: Estructura de un documento [5]

La distribución de tipos de datos se puede ver en la Figura 2. Siendo el tipo numérico el predominante con cinco elementos del documento de este tipo, seguido del tipo cadena con cuatro elementos y el formato fecha con un elemento.

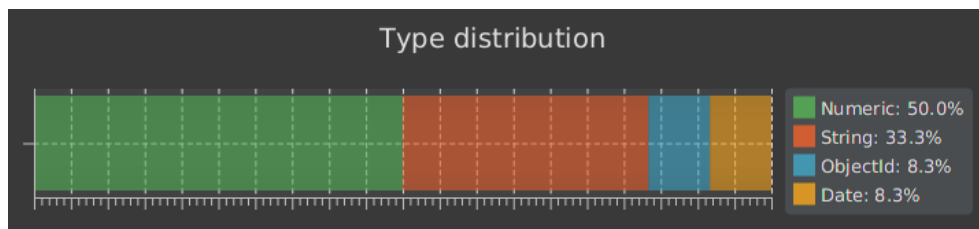


Figura 2: Distribución de tipos de datos [5]

Almacenamiento de datos

Para automatizar el proceso de almacenamiento de datos hemos creado un *script* en *Python* [8], que descarga los datos de la fuente mencionada anteriormente en formato *xlsx* y actualiza los datos en la base de datos, indicando simplemente el identificador (uri) proporcionado por *Azure Cosmos DB*.

```
1 import pandas as pd
2 import wget
3 import os
4 from pymongo import MongoClient
5
6 url = 'https://www.ecdc.europa.eu/sites/default/files' \
7 '/documents/COVID-19-geographic-disbtribution-worldwide.xlsx'
8
9
10 def excel_to_json():
11     filename = wget.download(url)
12     data = pd.read_excel(filename)
13     os.remove(filename)
14     return data.to_dict('records')
15
16
17 def upload_data(uri):
18     client = MongoClient(uri)
19     database = client['covid']
20     collection = database['Worldwide']
21     collection.drop()
22     collection.insert_many(excel_to_json())
23
24
25 upload_data(os.sys.argv[1])
```

De esta forma, cada día ejecutamos el *script* y actualizamos los datos de la base de datos en cuestión de minutos, ya que como hemos comentado anteriormente los datos de la fuente se actualizan a diario.

Ejecución de consultas

La ejecución de consultas la hemos realizado gracias a la creación de una clase *Python*, con la cual, crearemos objetos que nos permitirán conectar fácilmente con la base de datos y obtener datos específicos por país, fecha o continente. Esta clase se comporta de forma similar a un agente, utilizado en otras asignaturas para realizar la conexión con la base de datos.

```
1 from pymongo import MongoClient
2
3
4 def parse_date(date):
5     date = date.split("/")
6     day = int(date[0])
7     month = int(date[1])
8     year = int(date[2])
9
10    return day, month, year
11
12
13 class CovidClient:
14     def __init__(self, dir_bd, name_bd, collection):
15         self.__collection = MongoClient(dir_bd)[name_bd][collection]
16
17     def get_total_deaths_country(self, country):
18         try:
19             data = self.__collection.find({"countriesAndTerritories": country}, {"deaths"})
20             deaths = 0
21
22             for x in data:
23                 deaths += int(x['deaths'])
24
25             return deaths
26         except Exception as e:
27             return e
28
29     def get_total_deaths_continent(self, continent):
30         try:
31             data = self.__collection.find({"continentExp": continent}, {"deaths"})
32             deaths = 0
33
34             for x in data:
35                 deaths += int(x['deaths'])
36
37             return deaths
38         except Exception as e:
39             return e
40
41     def get_total_cases_country(self, country):
42         try:
43             data = self.__collection.find({"countriesAndTerritories": country}, {"cases"})
44             cases = 0
45
46             for x in data:
47                 cases += int(x['cases'])
48
49             return cases
50         except Exception as e:
51             return e
```

```

52
53 def get_total_cases_continent(self, continent):
54     try:
55         data = self.__collection.find({"continentExp": continent}, {"cases"})
56         cases = 0
57
58         for x in data:
59             cases += int(x['cases'])
60
61         return cases
62     except Exception as e:
63         return e
64
65 def get_worst_day_country(self, country):
66     try:
67         data = self.__collection.find({"countriesAndTerritories": country}, {"deaths", "
68         dateRep"})
69         deaths = None
70         date = None
71
72         for x in data:
73             if deaths is None:
74                 deaths = int(x['deaths'])
75                 date = x['dateRep']
76             else:
77                 if int(x['deaths']) > deaths:
78                     deaths = int(x['deaths'])
79                     date = x['dateRep']
80
81         return date, deaths
82     except Exception as e:
83         return e
84
85 def get_total_deaths(self):
86     try:
87         data = self.__collection.find({}, {"deaths"})
88         deaths = 0
89
90         for x in data:
91             deaths += int(x['deaths'])
92
93         return deaths
94     except Exception as e:
95         return e
96
97 def get_total_cases(self):
98     try:
99         data = self.__collection.find({}, {"cases"})
100         cases = 0
101
102         for x in data:
103             cases += int(x['cases'])
104
105         return cases
106     except Exception as e:
107         return e
108
109 def get_data_date_country(self, country, date):
110     try:

```

```

110     day, month, year = parse_date(date)
111     data = self.__collection.find(
112         {
113             "countriesAndTerritories": country,
114             "day": day,
115             "month": month,
116             "year": year
117         },
118         {"countriesAndTerritories", "cases", "deaths"}
119     )
120
121     return data
122 except Exception as e:
123     return e
124
125 def get_data_country(self, country):
126     try:
127         # only days with cases or deaths
128         data = self.__collection.find({"countriesAndTerritories": country,"$or": [{"cases": {
129             "$gt": 0}},{ "deaths": {"$gt": 0}}]},{ "countriesAndTerritories", "cases", "deaths", "
130             dateRep"})
131
132     return data
131 except Exception as e:
132     return e

```

Gracias a esta clase podremos realizar multitud de consultas interesantes a la base de datos, algunas de ellas las mostraremos a continuación.

Ejemplos de consultas

A continuación, se realizará una serie de consultas, las cuales hemos utilizado para llevar a cabo en el estudio de la pandemia. Cabe destacar que estos datos han sido obtenidos el día *7 de Mayo del 2020*.

Total de casos confirmados

Para ello se ha utilizado la función *get-total-cases()* [3] de la clase *Covid-Client*.

```
1 def get_total_cases(self):
2     try:
3         data = self.__collection.find({}, {"cases"})
4         cases = 0
5
6         for x in data:
7             cases += int(x['cases'])
8
9         return cases
10    except Exception as e:
11        return e
```

El **resultado** de la consulta sería el siguiente: **3623803** *casos confirmados* en todo el mundo.

Total de muertes confirmadas

Para ello se ha utilizado la función *get-total-deaths()* [3] de la clase *Covid-Client*.

```
1 def get_total_deaths(self):
2     try:
3         data = self.__collection.find({}, {"deaths"})
4         deaths = 0
5
6         for x in data:
7             deaths += int(x['deaths'])
8
9         return deaths
10    except Exception as e:
11        return e
```

El **resultado** de la consulta sería el siguiente: **256880** *muertes confirmadas* en el mundo.

Casos confirmados en España el mes de Marzo

Para ello se ha utilizado la función *get-data-country()* [3] de la clase *Covid-Client*.

```
1 data = client.get_data_country('Spain')
2 month = 3
3 cases = []
4 dates = []
5
6 # Get data
7 for d in data:
8     date = d['dateRep']
9     if date.month == month:
10         cases.append(abs(d['cases']))
11         dates.append(date.strftime("%d/%m"))
12
13 # Reverse data
14 cases = cases[::-1]
15 dates = dates[::-1]
16 y_pos = np.arange(len(dates))
17 plt.bar(y_pos, cases)
18 plt.xticks(y_pos, dates, rotation = 90)
19 plt.show()
```

El **resultado** de la consulta sería el diagrama de barras de la Figura 3 recaudando dicha información.

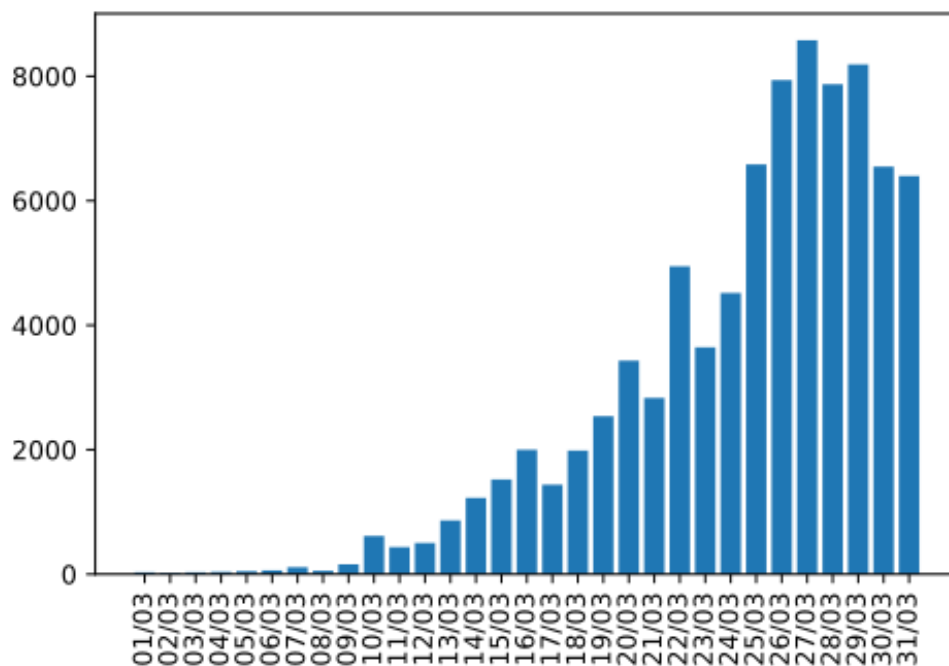


Figura 3: Casos en España mes de Marzo. [3]

Muertes en Italia en el mes de Abril

Para ello se ha utilizado la función *get-data-country()* [3] de la clase *Covid-Client*.

```
1 data = client.get_data_country('Italy')
2 month = 4
3 deaths = []
4 dates = []
5
6 for d in data:
7     date = d['dateRep']
8     if date.month == month:
9         deaths.append(abs(d['deaths']))
10        dates.append(date.strftime("%d/%m"))
11
12 # Reverse data
13 deaths = deaths[::-1]
14 dates = dates[::-1]
15 y_pos = np.arange(len(deaths))
16 plt.bar(y_pos, deaths)
17 plt.xticks(y_pos, dates, rotation = 90)
18 plt.show()
```

El **resultado** de la consulta sería el diagrama de barras de la Figura 4 recaudando dicha información.

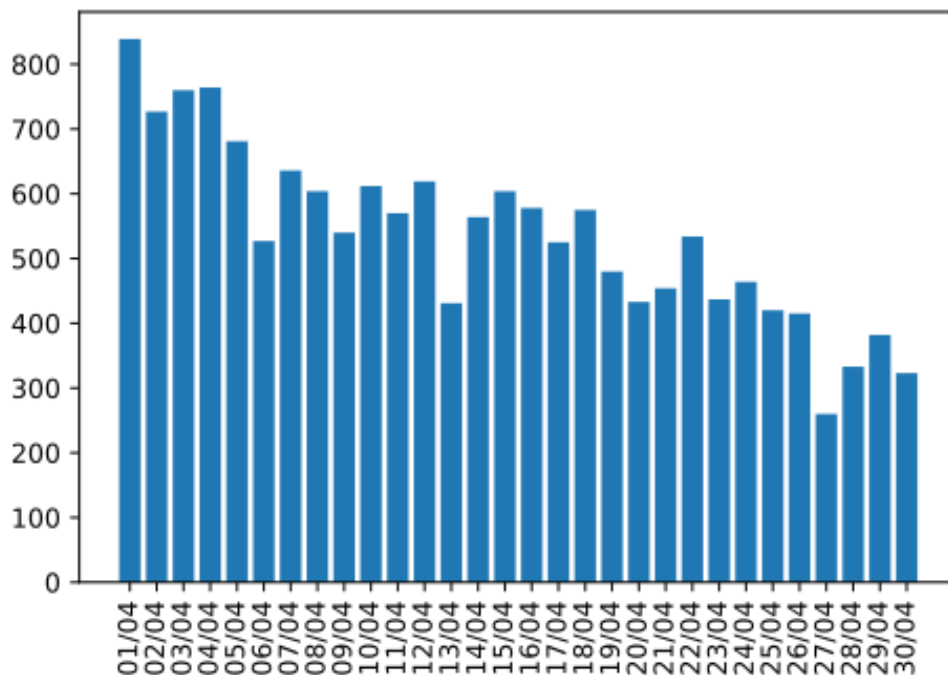


Figura 4: Muertes en Italia en el mes de Abril [4]

Casos confirmados agrupados por meses en Italia

Para ello se ha utilizado la función *get-data-country()* [3] de la clase *Covid-Client*.

```
1 data = client.get_data_country('Italy')
2 months = {1:0, 2:0, 3:0, 4:0, 5:0}
3
4 for d in data:
5     months[d['dateRep'].month] += abs(d['cases'])
6
7 y_pos = np.arange(len(months.values()))
8 plt.bar(y_pos, months.values())
9 plt.xticks(y_pos, months.keys())
10 plt.show()
```

El **resultado** de la consulta sería el diagrama de barras de la Figura 5 recaudando dicha información.

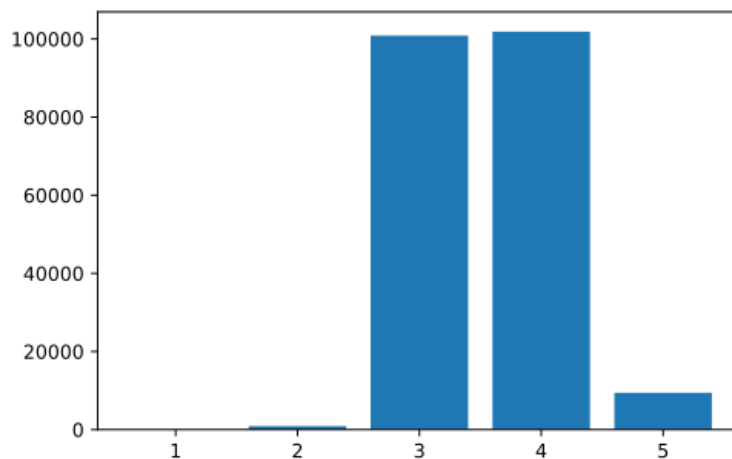


Figura 5: Casos confirmados agrupados por meses en Italia. [4]

Porcentaje de casos en Estados Unidos respecto al resto del mundo

Para ello se ha utilizado la función *get-total-cases-country()* y *get-total-cases()* [3] de la clase *CovidClient*.

```
1 country = "United_States_of_America"
2 country_cases = client.get_total_cases_country(country)
3 total_cases = client.get_total_cases() - country_cases
4
5 values = [total_cases, country_cases]
6 names = ["Rest_of_world_cases", (country.replace("_", " ") + "s_cases")]
7 plt.pie(values, labels=names, autopct="%0.2f%%")
8 plt.show()
```

El **resultado** de la consulta sería el diagrama de barras de la Figura 5 recaudando dicha información.

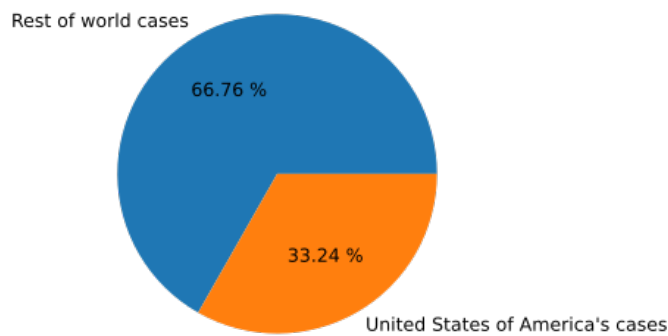


Figura 6: Porcentaje de casos en Estados Unidos respecto al resto del mundo. [4]

Porcentaje de casos de los tres países europeos más afectados

Para ello se ha utilizado la función *get-total-cases-country()* y *get-total-cases-continent()* [3] de la clase *CovidClient*.

```

1  continent = "Europe"
2  country1 = "Spain"
3  country2 = "Italy"
4  country3 = "United_Kingdom"
5
6  country1_cases = client.get_total_cases_country(country1)
7  country2_cases = client.get_total_cases_country(country2)
8  country3_cases = client.get_total_cases_country(country3)
9  continent_cases = client.get_total_cases_continent(continent) - (country1_cases -
    country2_cases - country3_cases)
10
11 values = [continent_cases, country1_cases, country2_cases, country3_cases]
12 names = [("Rest_of_" + continent.replace("_", " ") + "'s_cases"), (country1.replace("_", "
    ") + "'s_cases"), (country2.replace("_", " ") + "'s_cases"), (country3.replace("_", " ")
    ) + "'s_cases")]
13 plt.pie(values, labels=names, autopct="%0.2f%%")
14 plt.show()

```

El **resultado** de la consulta sería el diagrama de barras de la Figura 5 recaudando dicha información.

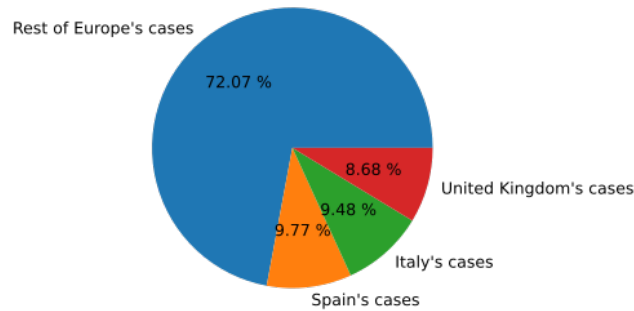


Figura 7: Porcentaje de casos de los tres países europeos más afectados. [4]

Porcentaje de muertes de los tres países más afectados en el mundo

Para ello se ha utilizado la función *get-total-deaths-country()* y *get-total-deaths()* [3] de la clase *CovidClient*.

```

1 country1 = "United_States_of_America"
2 country2 = "Spain"
3 country3 = "Italy"
4 country1_deaths = client.get_total_deaths_country(country1)
5 country2_deaths = client.get_total_deaths_country(country2)
6 country3_deaths = client.get_total_deaths_country(country3)
7 total_cases = client.get_total_deaths() - (country1_deaths - country2_deaths -
8     country3_deaths)
9 values = [total_cases, country1_deaths, country2_deaths, country3_deaths]
10 names = ["Rest_of_world_deaths", (country1.replace("_", " ") + "'s_deaths"), (country2.
11     replace("_", " ") + "'s_deaths"), (country3.replace("_", " ") + "'s_deaths")]
12 plt.pie(values, labels=names, autopct="%0.2f%%")
13 plt.show()

```

El **resultado** de la consulta sería el diagrama de barras de la Figura 5 recaudando dicha información.

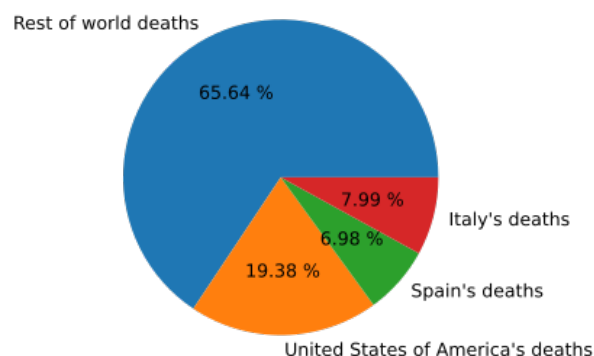


Figura 8: Porcentaje de muertes de los tres países más afectados en el mundo. [4]

Referencias

1. <https://azure.microsoft.com/es-es/services/cosmos-db/>. Azure cosmos db.
2. <https://data.europa.eu/euodp/en/data/dataset/covid-19-coronavirus-data>. Eu open data portal.
3. <https://jupyter.org/>. Jupyter.
4. <https://matplotlib.org/>. matplotlib.
5. <https://studio3t.com/>. Studio 3t.
6. <https://www.ecdc.europa.eu/en/covid-19-pandemic>. European centre for disease prevention and control.
7. <https://www.mongodb.com/>. Mongodb.
8. <https://www.python.org/>. Python software foundation.