

# COMP1204: Data Management

## Coursework Two: Representing Covid-19 data

PROKOPIS GEORGIU  
STUDENT ID: 34077499

### CORONAVIRUS DATA FROM THE EU OPEN DATA PORTAL. RELATIONAL MODEL, NORMALISATION, MODELLING, QUERYING

## Introduction

The COVID-19 pandemic has presented unprecedented challenges to global health systems and societies, prompting the need for innovative data management strategies. In their COMP1204: Data Management coursework, the students have been tasked with addressing this demand by organizing COVID-19 data sourced from the EU Open Data Portal into a SQLite database. Their objective is to establish a structured framework that facilitates insightful analysis of pandemic trends. Throughout this report, the students will outline their methodology for designing the database. This entails conceptualizing a relational model, ensuring data organization adheres to normalization principles, and constructing queries to extract meaningful insights. Each step will be thoroughly explained, from identifying relationships between different data points to optimizing the database's efficiency through normalization. Furthermore, the report will demonstrate how the students translated their plan into action. This includes setting up the database, importing the COVID-19 dataset, and executing queries to uncover valuable insights. The aim is to showcase how a well-designed database can streamline data management processes and empower informed analysis of COVID-19 statistics.

## 1 The Relational Model

**EX1: Express the relation directly represented in the dataset file. Assign relevant SQLite data types to each column.**

dataset(dateRep, day, month, year, cases, deaths, countriesAndTerritories, geoId, countryterritoryCode, popData2020, continentExp)

INTEGER	TEXT
day	dateRep
month	countriesAndTerritories
year	geoId
cases	countryterritoryCode
deaths	continentExp
popData2020	

Table 1: Attributes and types

This represents the dataset as a relational table with each column having the appropriate SQLite data type.

**EX2: List the minimal set of Functional Dependencies (FDs).**

Functional Dependency
$(\text{day, month, year}) \rightarrow \text{dateRep}$
$\text{dateRep, countriesAndTerritories} \rightarrow \text{cases}$
$\text{dateRep, countriesAndTerritories} \rightarrow \text{deaths}$
$\text{dateRep, countryterritoryCode} \rightarrow \text{cases}$
$\text{dateRep, countryterritoryCode} \rightarrow \text{deaths}$
$\text{dateRep, geoId} \rightarrow \text{cases}$
$\text{dateRep, geoId} \rightarrow \text{deaths}$

Functional Dependency
$\text{countriesAndTerritories} \rightarrow \text{countryterritoryCode}$
$\text{countriesAndTerritories} \rightarrow \text{continentExp}$
$\text{countryterritoryCode} \rightarrow \text{continentExp}$
$\text{countriesAndTerritories} \rightarrow \text{popData2020}$
$\text{geoId} \rightarrow \text{countriesAndTerritories}$
$\text{geoId} \rightarrow \text{countryterritoryCode}$
$\text{geoId} \rightarrow \text{popData2020}$

Assumptions: 1) dateRep can represent day, month and year. Thus, it would be redundant to write them separately. 2) countriesAndTerritories, geoId, countryterritoryCode, popData2020 and continentExp are all redundant since information is unnecessarily repeated.

**EX3: From your minimal set of functional dependencies, list the potential candidate keys.**

From these FDs, the potential candidate keys are: (dateRep, geoId), (dateRep, countriesAndTerritories), (dateRep, countryterritoryCode).

**EX4: Identify a suitable primary key, and justify your decision.**

The composite key (dateRep, geoId) appears to be the most suitable primary key. It provides a unique identifier for each record in the dataset and offers a high level of granularity by combining the date of reporting with the unique country or territory. This primary key ensures that each record is uniquely identifiable and allows for efficient data retrieval and manipulation.

## 2 Normalisation

**EX5: List any partial-key dependencies in the relation as it stands and any resulting additional relations you should create as part of the decomposition.**

1.  $\text{dateRep} \rightarrow \text{day}$
2.  $\text{dateRep} \rightarrow \text{month}$
3.  $\text{dateRep} \rightarrow \text{year}$
4.  $\text{geoId} \rightarrow \text{popData2020}$
5.  $\text{geoId} \rightarrow \text{continentExp}$
6.  $\text{geoId} \rightarrow \text{countryterritoryCode}$
7.  $\text{geoId} \rightarrow \text{countriesAndTerritories}$

Additional Tables:

1.  $(\text{dateRep, geoId}) \rightarrow (\text{cases, deaths})$
2.  $\text{dateRep} \rightarrow (\text{day, month, year})$
3.  $\text{geoId} \rightarrow (\text{countriesAndTerritories, countryterritoryCode, popData2020, continentExp})$

**EX6: Use decomposition and your answer to the above to achieve 2nd Normal Form, introducing appropriate new relations. List the new relations and their fields, types and keys. Explain the process you took.**

Additional Relations:

1. Condition(dateRep, geoId, cases, deaths)
2. Date(dateRep, day, month, year)
3. Country(geoId, countriesAndTerritories, popData2020, continentExp)

The primary keys for the relation Condition are (dateRep and geoId). For the relation Date it is dateRep and for relation Country, the primary key is geoId. The types are the same as in Table 1. To achieve 2nd Normal Form (2NF), I decomposed the original relation into smaller relations to eliminate partial-key dependencies. The decomposition process involved identifying attributes that were functionally dependent on only part of the primary key and separating them into new relations.

#### **Condition Relation**

This relation includes attributes that are functionally dependent on both the date of reporting (dateRep) and the unique identifier for the country or territory (geoId). The primary key for this relation is a composite key consisting of dateRep and geoId.

#### **Date Relation**

To address the partial-key dependencies related to the date of reporting (dateRep), I created a separate relation named Date. This relation includes attributes specific to the date, namely day, month, and year. The primary key for this relation is dateRep.

#### **Country Relation**

Similarly, I created a separate relation named Country to handle attributes related to countries or territories. This relation includes attributes such as countriesAndTerritories, popData2020, and continentExp, which are functionally dependent on the unique identifier for the country or territory (geoId). The primary key for this relation is geoId.

**EX7: List transitive dependencies, if any, in your new relations.**

All new relations are in 2NF as they do not contain any transitive dependencies. Therefore, no further decomposition is needed.

**EX8: Convert your relations into 3rd Normal Form using your answers to the above. List the new relations and their fields, types and keys. Explain the process you took.**

The relations Condition, Date, and Country are in 3rd Normal Form (3NF) because they meet the following criteria:

1. **Atomic Values:** All attributes contain atomic values, meaning they cannot be further divided. In each relation, the attributes represent single, indivisible pieces of information. For example, in the Condition relation, cases and deaths represent single counts of COVID-19 cases and deaths, respectively.
2. **No Transitive Dependencies:** There are no transitive dependencies in any of the relations. A transitive dependency occurs when a non-prime attribute depends on another non-prime attribute, rather than directly on the primary key. However, in all three relations, each non-prime attribute is fully functionally dependent on the primary key, and there are no indirect dependencies between non-prime attributes. For example, in the Country relation, countriesAndTerritories, popData2020, and continentExp are all directly dependent on the primary key geoId.
3. **Every Non-Prime Attribute Fully Depends on the Primary Key:** Each non-prime attribute in all three relations is fully functionally dependent on the primary key. This means that the value of any non-prime attribute is uniquely determined by the primary key alone, without any further dependencies. For instance, in the Date relation, day, month, and year are all fully dependent on the primary key dateRep.

## EX9: Convert your relations into Boyce-Codd Normal Form. Justify and explain how your relations are in BCNF.

To convert the relations into Boyce-Codd Normal Form (BCNF), we need to ensure that every determinant is a candidate key. Here's how we can justify that the relations Condition, Date, and Country are in BCNF:

### Condition Relation:

- Primary key: (dateRep, geoId)
- Non-prime attributes: cases, deaths
- Justification: The primary key (dateRep, geoId) uniquely identifies each tuple in the relation. Both cases and deaths are fully functionally dependent on the primary key. Since the primary key is a composite key, it cannot be decomposed further.

### Date Relation:

- Primary key: dateRep
- Non-prime attributes: day, month, year
- Justification: The primary key dateRep uniquely identifies each tuple in the relation. All non-prime attributes (day, month, year) are fully functionally dependent on the primary key. Since the primary key is a candidate key, and there are no dependencies on other non-prime attributes.

### Country Relation:

- Primary key: geoId
- Non-prime attributes: countriesAndTerritories, popData2020, continentExp
- Justification: The primary key geoId uniquely identifies each tuple in the relation. All non-prime attributes (countriesAndTerritories, popData2020, continentExp) are fully functionally dependent on the primary key. Since the primary key is a candidate key, and there are no dependencies on other non-prime attributes.

The relations Condition, Date, and Country are in BCNF because every determinant is a candidate key, ensuring that there are no non-trivial functional dependencies on attributes that are not candidate keys.

## 3 Modelling

### EX10

```
1 sqlite3 coronavirus.db
2 sqlite> .open coronavirus.db
3 sqlite> .mode csv
4 sqlite> .import dataset.csv dataset
5 sqlite> .output dataset.sql
6 \\Open DataGrip and set the data source and drivers
7 sqlite3 coronavirus.db .dump > database.sql
```

## EX11

```
1 CREATE TABLE Date (  
2     dateRep TEXT NOT NULL PRIMARY KEY,  
3     day     INTEGER NOT NULL,  
4     month   INTEGER NOT NULL,  
5     year    INTEGER NOT NULL  
6 );  
7 CREATE TABLE Country (  
8     geoId          TEXT NOT NULL CONSTRAINT Country_pk PRIMARY KEY,  
9     countriesAndTerritories TEXT NOT NULL,  
10    countryterritoryCode TEXT NOT NULL,  
11    popData2020     INTEGER NOT NULL,  
12    continentExp    TEXT NOT NULL  
13 );  
14 CREATE TABLE Condition (  
15     dateRep TEXT NOT NULL,  
16     geoId   TEXT NOT NULL,  
17     deaths  INTEGER,  
18     cases   INTEGER,  
19     CONSTRAINT Condition_pk PRIMARY KEY (dateRep, geoId)  
20 );
```

## EX12

```
1 -- Populate the Date table  
2 INSERT INTO Date (dateRep, day, month, year)  
3 SELECT DISTINCT dateRep, day, month, year  
4 FROM dataset;  
5  
6 -- Populate the Country table  
7 INSERT INTO Country (geoId, countriesAndTerritories, countryterritoryCode, popData2020, continentExp)  
8 SELECT DISTINCT geoId, countriesAndTerritories, countryterritoryCode, popData2020, continentExp  
9 FROM dataset;  
10  
11 -- Populate the Condition table  
12 INSERT INTO Condition (dateRep, geoId, deaths, cases)  
13 SELECT DISTINCT dateRep, geoId, deaths, cases  
14 FROM dataset;
```

## EX13

```
1 sqlite3 coronavirus.db < dataset.sql  
2 sqlite3 coronavirus.db < ex11.sql  
3 sqlite3 coronavirus.db < ex12.sql  
4 //Confirming that I have run these commands and they worked.
```

## 4 Querying

**EX14:** The worldwide total number of cases and deaths (with total cases and total deaths as columns)

```
1 SELECT  
2     SUM(cases) AS total_cases,  
3     SUM(deaths) AS total_deaths  
4 FROM  
5     Condition;
```

**EX15: The number of cases by date, in increasing date order, for the United Kingdom (with the date reported and number of cases as columns)**

```
1 SELECT
2     Date.dateRep,
3     Condition.cases
4 FROM
5     Date
6 JOIN
7     Condition ON Date.dateRep = Condition.dateRep
8 JOIN
9     Country ON Condition.geoId = Country.geoId
10 WHERE
11     Country.countriesAndTerritories = 'United_Kingdom'
12 ORDER BY
13     Date.year ASC,
14     Date.month ASC,
15     Date.day ASC;
```

**EX16: The number of cases and deaths by date, in increasing date order, for each country (with country name, date, number of cases and number of deaths as columns)**

```
1 SELECT
2     Country.countriesAndTerritories,
3     Condition.dateRep,
4     Condition.cases, Condition.deaths
5 FROM
6     Country
7 JOIN
8     Condition ON Country.geoId = Condition.geoId
9 JOIN
10    Date ON Condition.dateRep = Date.dateRep
11 GROUP BY
12     Country.countriesAndTerritories,
13     Condition.dateRep
14 ORDER BY
15     Country.countriesAndTerritories ASC,
16     Date.year ASC,
17     Date.month ASC,
18     Date.day ASC;
```

**EX17: The total number of cases and deaths as a percentage (to 2DP) of the population, for each country**

```
1 SELECT
2     Country.countriesAndTerritories,
3     ROUND((SUM(Condition.cases) * 100.0 / Country.popData2020), 2) AS "% cases",
4     ROUND((SUM(Condition.deaths) * 100.0 / Country.popData2020), 2) AS "% deaths"
5 FROM
6     Condition
7 JOIN
8     Country ON Condition.geoId = Country.geoId
9 GROUP BY
10    Country.countriesAndTerritories;
```

**EX18:** A descending list of the the top 10 countries by name, by percentage (to 2DP) total deaths out of total cases in that country.

```
1 SELECT
2     Country.countriesAndTerritories AS Country,
3     ROUND((SUM(Condition.deaths) * 100.0 / SUM(Condition.cases)), 2) AS "% deaths of country cases"
4 FROM
5     Condition
6 JOIN
7     Country ON Condition.geoId = Country.geoId
8 GROUP BY
9     Country.countriesAndTerritories
10 ORDER BY
11     "% deaths of country cases" DESC
12 LIMIT
13     10;
```

**EX19:** The date against a cumulative running total of the number of deaths by day and cases by day for the united kingdom.

```
1 SELECT Date.dateRep AS "date",
2     SUM(Condition.deaths) OVER win1 AS "cumulative UK deaths",
3     SUM(Condition.cases) OVER win1 AS "cumulative UK cases"
4 FROM Date
5 JOIN Condition ON Date.dateRep = Condition.dateRep
6 WHERE Condition.geoId = 'UK'
7 WINDOW win1 AS (ROWS BETWEEN UNBOUNDED PRECEDING AND CURRENT ROW)
8 ORDER BY Date.year, Date.month, Date.day ASC;
```