

# COMP1204: Data Management

## Coursework two

Mohammad Rayes  
35316004

May 2024

# 1 The Relational Model

## 1.1 EX1

Attribute	SQLite Data Type
dateRep	TEXT
day	INTEGER
month	INTEGER
year	INTEGER
cases	INTEGER
deaths	INTEGER
countriesAndTerritories	TEXT
geoId	TEXT
countryterritoryCode	TEXT
popData2020	INTEGER
continentExp	TEXT

Table 1: Attributes and SQLite Data Types in the Dataset

## 1.2 EX2

### 1.2.1 Minimal Set of Functional Dependencies

The minimal functional dependencies for the dataset, assuming unique identification of records by combinations of date and geographical identifiers, are as follows:

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

### Assumptions and Justifications:

- *Uniqueness of geoId:* Assumed that each **geoId** uniquely identifies a set of static attributes about a location (country, population, continent), which is common in datasets containing geographical information.
- *Stability of dateRep:* The representation of a date as **dateRep** is assumed to uniquely determine its day, month, and year components, typical for date formatting in data systems.
- *Combination keys for dynamic data:* The combination of **dateRep** and **geoId** is assumed necessary and sufficient to determine the counts of cases and deaths, reflecting the nature of reporting in epidemiological data where each entry represents a unique event/day/location.
- *Treatment of nulls and blanks:* It is assumed that all entries are complete and any nulls would be addressed prior to applying these dependencies, as incompleteness would violate the principles of functional dependency.

## 1.3 EX3

Based on the provided functional dependencies, several combinations of attributes are considered to explore their eligibility as candidate keys. Each candidate key must be minimal and able to uniquely identify every tuple in the dataset.

### 1.3.1 Explored Combinations

1. (dateRep, geoId)
2. (dateRep, countryterritoryCode)
3. (day, month, year, geoId)
4. (geoId, day, month, year, cases, deaths)

## 1.4 EX4

After evaluating potential candidate keys, (**dateRep, geoId**) is chosen as the primary key for the dataset for the following reasons:

- **Simplicity and Minimality**
- **Comprehensive Coverage**
- **Practicality in Usage**

This key not only fulfills the requirements of a candidate key but also aligns with the practical needs of database management and data integrity maintenance.

## 2 Normalisation

### 2.1 EX5

Identified partial-key dependencies from established functional dependencies illustrate that specific attributes rely solely on parts of a potential composite key. This necessitates additional relations for normalization adherence:

- $\text{dateRep} \rightarrow \text{day, month, year}$ : Suggests a separate relation for date components derived directly from the reporting date.
- $\text{geoId} \rightarrow \text{countryterritoryCode, countriesAndTerritories, continentExp, popData2020}$ : Indicates the need for a dedicated geographical information table.

Resulting relations:

1. **Date Information Table**: Includes `dateRep`, `day`, `month`, `year`.
2. **Geographic Details Table**: Contains `geoId`, `countryterritoryCode`, `countriesAndTerritories`, `continentExp`, `popData2020`.

This decomposition aligns with the principles of the Second Normal Form (2NF).

### 2.2 EX6

To achieve 2NF, we decomposed the original dataset into three tables, eliminating partial-key dependencies and ensuring full dependency on primary keys:

- **Date Information**: Holds date-related attributes, isolated to prevent partial-key dependencies. Fields include `dateRep` (TEXT), `day` (INTEGER), `month` (INTEGER), `year` (INTEGER). Primary Key: `dateRep`.
- **Geographic Information**: Consolidates geographic details, ensuring clean structure. Fields include `geoId` (TEXT), `countryterritoryCode` (TEXT), `countriesAndTerritories` (TEXT), `continentExp` (TEXT), `popData2020` (INTEGER). Primary Key: `geoId`.
- **Case and Death Records**: Captures cases and deaths linked to dates and locations. Fields include `dateRep` (TEXT), `geoId` (TEXT), `cases` (INTEGER), `deaths` (INTEGER). Primary Key: (`dateRep`, `geoId`).

These tables ensure data integrity by reducing redundancy.

### 2.3 EX7

No transitive dependencies are present in the new relations, complying with 2NF and 3NF principles.

## 2.4 EX8

Our relations inherently satisfy the requirements of 3NF as they are free from transitive dependencies and each attribute is fully functionally dependent on the primary key(s).

## 2.5 EX9

The relations conform to Boyce-Codd Normal Form (BCNF), ensuring no anomalies or redundancy:

- **Date Information Table:** All attributes are functionally dependent solely on dateRep, the primary key.
- **Geographic Information Table:** Attributes are functionally dependent solely on geoId, the primary key.
- **Case and Death Records Table:** cases and deaths depend only on the composite key (dateRep, geoId).

Thus, the relations are in BCNF, ensuring robust data integrity.

# 3 Modelling

## 3.1 EX10

To import a raw dataset into SQLite and create a dump file, follow these streamlined steps:

1. **Start SQLite:** Launch SQLite with your database.

```
sqlite3 coronavirus.db
```

2. **Create the Table:** Set up the dataset schema.

```
CREATE TABLE dataset (  
    dateRep TEXT,  
    day INTEGER,  
    month INTEGER,  
    year INTEGER,  
    cases INTEGER,  
    deaths INTEGER,  
    countriesAndTerritories TEXT,  
    geoId TEXT,  
    countryterritoryCode TEXT,  
    popData2020 INTEGER,  
    continentExp TEXT);
```

3. **Set CSV Mode:** Configure SQLite for CSV data.

```
.mode csv
```

4. **Import the Dataset:** Ensure CSV columns align with your table schema.

```
.import dataset.csv dataset
```

5. **Prepare for Dumping:** Specify the output file for the SQL dump.

```
.output dataset.sql
```

6. **Dump the Dataset:** Export the schema and data to an SQL file.

```
.dump
```

## 3.2 EX11

```
1  -- Table: Date Information
2  CREATE TABLE IF NOT EXISTS Date_Information (
3      dateRep TEXT PRIMARY KEY NOT NULL,
4      day INTEGER NOT NULL,
5      month INTEGER NOT NULL,
6      year INTEGER NOT NULL
7  );
8
9  -- Table: Geographic Information
10 CREATE TABLE IF NOT EXISTS Geographic_Information (
11     geoId TEXT PRIMARY KEY NOT NULL,
12     countryterritoryCode TEXT NOT NULL,
13     countriesAndTerritories TEXT NOT NULL,
14     continentExp TEXT NOT NULL,
15     popData2020 INTEGER
16 );
17
18 -- Table: Case and Death Records
19 CREATE TABLE IF NOT EXISTS Case_Death_Records (
20     dateRep TEXT NOT NULL,
21     geoId TEXT NOT NULL,
22     cases INTEGER,
23     deaths INTEGER,
24     PRIMARY KEY (dateRep, geoId)
25 );
```

## 3.3 EX12

```
1  -- Insert data into Date Information Table
2  INSERT OR IGNORE INTO Date_Information (dateRep, day, month, year)
3  SELECT dateRep, day, month, year
4  FROM dataset;
5
6  -- Insert data into Geographic Information Table
7  INSERT OR IGNORE INTO Geographic_Information (geoId,countriesAndTerritories
↵ ,countryterritoryCode, continentExp, popData2020)
```

```

8  SELECT geoId,countriesAndTerritories ,countryterritoryCode, continentExp,
   ↪ popData2020
9  FROM dataset;
10
11  -- Insert data into Case and Death Records Table
12  INSERT OR IGNORE INTO Case_Death_Records (dateRep, geoId, cases, deaths)
13  SELECT dateRep, geoId, cases, deaths
14  FROM dataset;

```

### 3.4 EX13

To verify the functionality of the SQL scripts, the following commands were executed:

```

1  sqlite3 coronavirus.db < dataset.sql
2  sqlite3 coronavirus.db < ex11.sql
3  sqlite3 coronavirus.db < ex12.sql

```

All of these tests were successful and executed without errors.

## 4 Querying

### 4.1 EX14

```

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

```

### 4.2 EX15

```

1      SELECT Date_Information.dateRep AS Date_Reported,
2             Case_Death_Records.cases AS Number_of_Cases
3  FROM Case_Death_Records
4  INNER JOIN Date_Information ON Case_Death_Records.dateRep =
   ↪ Date_Information.dateRep
5  INNER JOIN Geographic_Information ON Case_Death_Records.geoId =
   ↪ Geographic_Information.geoId
6  WHERE Geographic_Information.countryterritoryCode = 'GBR'
7  ORDER BY Date_Information.year, Date_Information.month, Date_Information.day ASC;

```

### 4.3 EX16

```

1  SELECT
2      g.countriesAndTerritories AS CountryName,
3      d.dateRep AS Date,
4      cd.cases AS Cases,
5      cd.deaths AS Deaths
6  FROM
7      Case_Death_Records cd
8  JOIN
9      Geographic_Information g ON cd.geoId = g.geoId

```

```

10 JOIN
11     Date_Information d ON cd.dateRep = d.dateRep
12 ORDER BY
13     d.year ASC, d.month ASC, d.day ASC, g.countriesAndTerritories ASC;

```

#### 4.4 EX17

```

1  SELECT
2      g.countriesAndTerritories AS CountryName,
3      ROUND(SUM(cd.cases) * 100.0 / g.popData2020, 2) AS Cases_Percentage,
4      ROUND(SUM(cd.deaths) * 100.0 / g.popData2020, 2) AS Deaths_Percentage
5  FROM
6      Case_Death_Records cd
7  JOIN
8      Geographic_Information g ON cd.geoId = g.geoId
9  GROUP BY
10     g.countriesAndTerritories

```

#### 4.5 EX18

```

1  SELECT
2      g.countriesAndTerritories AS "Country Name",
3      ROUND((SUM(cd.deaths) * 100.0 / SUM(cd.cases)), 2) AS "Percentage Deaths of
4      ↳ Country Cases (%)"
5  FROM
6      Case_Death_Records cd
7  JOIN
8      Geographic_Information g ON cd.geoId = g.geoId
9  GROUP BY
10     g.countriesAndTerritories
11  HAVING
12     SUM(cd.cases) > 0 -- To prevent division by zero
13  ORDER BY
14     "Percentage Deaths of Country Cases (%)" DESC
15  LIMIT 10;

```

#### 4.6 EX19

```

1  SELECT
2      d.dateRep AS "Date",
3      SUM(cd.cases) OVER (ORDER BY d.year, d.month, d.day) AS "Cumulative UK
4      ↳ Cases",
5      SUM(cd.deaths) OVER (ORDER BY d.year, d.month, d.day) AS "Cumulative UK
6      ↳ Deaths"
7  FROM
8      Case_Death_Records cd
9  JOIN
10     Geographic_Information g ON cd.geoId = g.geoId
11  JOIN
12     Date_Information d ON cd.dateRep = d.dateRep
13  WHERE
14     g.countriesAndTerritories = 'United_Kingdom'
15  ORDER BY
16     d.year, d.month, d.day;

```