

COMP1204: Data Management

Coursework Two

Zhengbin Lu
33839239

April 11, 2023

1 The Relational Model

1.1 EX1

Relation	Type
dateRep - cases	Many to many
dateRep - deaths	Many to many
dateRep - countriesAndTerritories	Many to many
dateRep - geoId	Many to many
dateRep - countryterritoryCode	Many to many
dateRep - popData2020	Many to many
day - dateRep	One to many
day - month	Many to many
day - year	Many to many
day - cases	Many to many
day - deaths	Many to many
day - countriesAndTerritories	Many to many
day - geoId	Many to many
day - countryterritoryCode	Many to many
day - popData2020	Many to many
month - dateRep	One to many
month - year	Many to many
month - cases	Many to many
month - deaths	Many to many
month - countriesAndTerritories	Many to many
month - geoId	Many to many
month - countryterritoryCode	Many to many
month - popData2020	Many to many
year - dateRep	One to many
year - cases	Many to many
year - deaths	Many to many
year - countriesAndTerritories	Many to many

(a) Relation Table 1

Relation	Type
year - geoId	Many to many
year - countryterritoryCode	Many to many
year - popData2020	Many to many
cases - deaths	Many to many
cases - countriesAndTerritories	Many to many
cases - geoId	Many to many
cases - countryterritoryCode	Many to many
cases - popData2020	Many to many
deaths - countriesAndTerritories	Many to many
deaths - geoId	Many to many
deaths - countryterritoryCode	Many to many
deaths - popData2020	Many to many
countriesAndTerritories - geoId	One to one
countriesAndTerritories - countryterritoryCode	One to one
countriesAndTerritories - popData2020	One to many
geoId - countryterritoryCode	One to one
geoId - popData2020	One to many
countryterritoryCode - popData2020	One to many
continentExp - dateRep	Many to many
continentExp - day	Many to many
continentExp - month	Many to many
continentExp - year	Many to many
continentExp - cases	Many to many
continentExp - deaths	Many to many
continentExp - countriesAndTerritories	Many to many
continentExp - geoId	Many to many
continentExp - countryterritoryCode	Many to many
continentExp - popData2020	Many to many

(b) Relation Table 2

Figure 1: Relations and Type

Attributes	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: Attribute Tab

1.2 EX2

The continentExp can't identify a continent uniquely if there is any other continent will be covered in the data in the future. The popData2020 can't uniquely identify a country when there is other countries with the same populatuion in the data. The date information, for instance, dateRep, day, month, and year are not null. The cases and deaths attributes can be a integer or a null value, but can't be a blank. The attribute countriesAndTerritories, geoId, countryterritoryCode, popData2020 can map to their countries as injections, thus these attributes can be determinant and they can uniquely identify their countries.

The FDs table is as follow.

Functional Dependencies
$dateRep \rightarrow day$
$dateRep \rightarrow month$
$dateRep \rightarrow year$
$countriesAndTerritories \rightarrow geoId$
$geoId \rightarrow countryterritoryCode$
$countryterritoryCode \rightarrow popData2020$
$countryterritoryCode \rightarrow continentExp$
$day, month, year \rightarrow dateRep$
$dateRep, geoId \rightarrow cases$
$dateRep, geoId \rightarrow deaths$
$geoId \rightarrow countriesAndTerritories$
$countryterritoryCode \rightarrow geoId$

Table 2: FDs Table

1.3 EX3

Potential Candidate Keys
dateRep, countriesAndTerritories
dateRep, geoId
dateRep, countryterritoryCode
day, month, year, countryterritoryCode
day, month, year, geoId
day, month, year, countriesAndTerritories

Table 3: Candidate Keys Table

1.4 EX4

The candidate key is composed by two part, the date attribute and the country attribute, so the primary key that I choose is

dateRep, geoId

dateRep contains the exact date of each case and death, and geoId contains the country detail. Compared with other candidate key, this one is more concise.

2 Normalisation

2.1 EX5

$dateRep \rightarrow continentExp, popData2020$
 $geoId \rightarrow continentExp, popData2020$

continentExp and popData2020 are not key attributes, they dependent on the part of candidate key geoId or dateRep, thus these two attributes partially dependent on the primary key.

The partial dependencies result in an additional relation for decomposition.

$geoId, dateRep \rightarrow continentExp, popData2020$

2.2 EX6

When decomposing, I need to introduce extra two INTEGER attributes, **data_ID** and **geo_ID** as surrogate keys. Using surrogate keys can guarantee that each row in the table can be identified uniquely. Furthermore, the surrogate keys will replace original primary key and as new primary key. Thus, we have the new relations in the tables:

$data_ID \rightarrow dateRep, day, month, year$
 $geo_ID \rightarrow geoId, countriesAndTerritories, countryterritoryCode$
 $geo_ID \rightarrow popData2020, continentExp$
 $data_ID, geo_ID \rightarrow cases, deaths$

The field of these two attributes are geo_ID and data_ID, and the type of these two attributes is INTEGER, as surrogate keys in the table.

2.3 EX7

According to the definition of the transitive dependent, the primary key in relation geo_ID –> geoId, countriesAndTerritories, countryterritoryCode is geo_ID, the rest of attributes are the part of candidate key, and dependent on the primary key, so there is no transitive dependent. For the other relation, it is obviously that there is no transitive dependent.

2.4 EX8

Depending on assumptions I have made and earlier processes, i'm already in the 3NF. In the relations I provided, there is no transitive dependency. Therefore there are no non-key attributes that are transitively dependent on any candidate key. For every FD $A \rightarrow b$, A is a super key, b is either a key or non-key. For the key attribute in the relation, it can only be implied by super key, for non-key attribute, there is no transitive dependency.

2.5 EX9

The relations above are already in BCNF, for each FD, the determinant is candidate key, and it can identify a unique row. There is only one relation $data_ID \rightarrow dateRep, day, month, year$ have more than one attribute in the candidate but the keys are disjoint. Thus, the relations are already in BCNF.

3 Modelling

3.1 EX10

1. Starts terminal and input **sqlite3 coronavirus.db** to create a empty coronavirus.db
2. Inputs **.mode csv** to shift into the csv mode and use **.import dataset.csv dataset** to import the csv file into the table dataset in coronavirus.db.
3. Uses **.output dataset.sql** to set the output file.
4. Uses the **.dump** to dump the database as the dataset.sql.
5. Runs **sqlite3 coronavirus.db < dataset.sql** to populate the dataset table.

3.2 EX11

The SQL with full normalised representation is as follow. The surrogate keys that I'm introduced is ID in the table Data and CountriesAndTerritories, these keys can increase automatically and identify each row uniquely in the table.

```
1 PRAGMA foreign_keys = ON;
2
3 CREATE TABLE Date
4 (
5     ID      INTEGER PRIMARY KEY AUTOINCREMENT,
6     dateRep TEXT    NOT NULL UNIQUE,
7     day     INTEGER NOT NULL,
8     month   INTEGER NOT NULL,
9     year    INTEGER NOT NULL
10 );
11
12 CREATE TABLE CountriesAndTerritories
13 (
14     ID                  INTEGER PRIMARY KEY AUTOINCREMENT,
15     geoId              TEXT NOT NULL UNIQUE,
16     countriesAndTerritories TEXT NOT NULL UNIQUE ,
17     countryterritoryCode TEXT NOT NULL UNIQUE
18 );
19
20 CREATE TABLE PopAndContinent
21 (
22     geo_ID      INTEGER NOT NULL ,
23     popDate2020 INTEGER NOT NULL ,
24     continentExp TEXT NOT NULL ,
25     FOREIGN KEY (geo_ID) references CountriesAndTerritories (ID)
26         ON DELETE CASCADE
27         ON UPDATE CASCADE
28 );
29
30 CREATE TABLE CasesAndDeaths
31 (
32     data_ID  INTEGER ,
33     geo_ID   INTEGER ,
34     cases    INTEGER ,
```

```

35    deaths  INTEGER,
36    PRIMARY KEY (data_ID, geo_ID),
37    FOREIGN KEY (data_ID) REFERENCES Date (ID)
38        ON DELETE CASCADE
39        ON UPDATE CASCADE,
40    FOREIGN KEY (geo_ID) REFERENCES CountriesAndTerritories (ID)
41        ON DELETE CASCADE
42        ON UPDATE CASCADE
43 );

```

3.3 EX12

```

1 INSERT INTO Date (dateRep, day, month, year)
2 SELECT DISTINCT dateRep, day, month, year
3 FROM dataset
4 LIMIT -10 OFFSET 1;
5
6 INSERT INTO CountriesAndTerritories (geoId, countriesAndTerritories, countryterritoryCode)
7 SELECT DISTINCT geoId, countriesAndTerritories, countryterritoryCode
8 FROM dataset
9 LIMIT -10 OFFSET 1;
10
11 INSERT INTO PopAndContinent (geo_ID, popDate2020, continentExp)
12 SELECT DISTINCT C.ID, popData2020, continentExp
13 FROM dataset
14     JOIN CountriesAndTerritories C ON dataset.geoId = C.geoId;
15
16 INSERT INTO CasesAndDeaths (data_ID, geo_ID, cases, deaths)
17 SELECT DISTINCT D.ID, C.ID, cases, deaths
18 FROM dataset
19     JOIN Date D ON dataset.dateRep = D.dateRep
20     JOIN CountriesAndTerritories C ON dataset.geoId = C.geoId;

```

3.4 EX13

Successfully run the command and end up with a fully populated database.

4 Querying

4.1 EX14

```

1 SELECT SUM(cases) as casesSum, SUM(deaths) as deathsSum
2 FROM CasesAndDeaths;

```

4.2 EX15

```

1 SELECT D.dateRep, cases FROM CasesAndDeaths
2 JOIN Date D on CasesAndDeaths.data_ID = D.ID
3 JOIN CountriesAndTerritories CAT on CasesAndDeaths.geo_ID = CAT.ID
4 WHERE geoId = 'UK'
5 ORDER BY year, month, day;

```

4.3 EX16

```

1 SELECT geoId,dateRep, cases, deaths
2 FROM CasesAndDeaths
3         JOIN Date D on CasesAndDeaths.data_ID = D.ID
4             JOIN CountriesAndTerritories CAT on CasesAndDeaths.geo_ID = CAT.ID
5 ORDER BY year, month, day, countriesAndTerritories;

```

4.4 EX17

```
1 SELECT CAT.geoId,
2     (SUM(CasesAndDeaths.cases) * 100.0 / PAC.popDate2020) || '%' as casesPercentage ,
3     (SUM(CasesAndDeaths.deaths) * 100.0 / PAC.popDate2020) || '%' as deathsPercentage
4 FROM CasesAndDeaths
5     JOIN CountriesAndTerritories CAT on CasesAndDeaths.geo_ID = CAT.ID
6     JOIN PopAndContinent PAC ON CAT.ID = PAC.geo_ID
7 GROUP BY CAT.geoId
8 ORDER BY CAT.geoId;
```

4.5 EX18

```
1 SELECT CAT.geoId,
2     (SUM(CasesAndDeaths.deaths) * 100.0 / SUM(CasesAndDeaths.cases)) || '%' as
3     deathsPercentage
4 FROM CasesAndDeaths
5     JOIN CountriesAndTerritories CAT on CasesAndDeaths.geo_ID = CAT.ID
6 GROUP BY CAT.geoId
7 ORDER BY deathsPercentage DESC;
```

4.6 EX19

```
1 SELECT dateRep,
2     SUM(deaths) OVER (ORDER BY year, month, day) AS cumulativeDeaths ,
3     SUM(cases) OVER (ORDER BY year, month, day) AS cumulativeCases
4 FROM CasesAndDeaths
5     JOIN Date D on D.ID = CasesAndDeaths.data_ID
6     JOIN CountriesAndTerritories CAT on CAT.ID = CasesAndDeaths.geo_ID
7 WHERE geoId = 'UK';
```

5 Extension

5.1 EX20

```
1#!/bin/bash
2
3# Create an empty array to store the temporary files
4declare -a data
5
6# Create a command which will be excuted by the gnuPlot
7plot_command="plot "
8
9# Fetch the top 10 death country and iterate them, for each loop, grab the date, sum of the
10# death and geoId from the table, and put them into a temp file after processing
11# Append the gnuPlot command
12for geoId in $(sqlite3 coronavirus.db "SELECT geoId FROM CountriesAndTerritories JOIN
13    CasesAndDeaths CAD on CountriesAndTerritories.ID = CAD.geo_ID GROUP BY geoId ORDER BY SUM(
14    deaths) DESC LIMIT 10;")
15do
16    Deaths=$(mktemp)
17    sqlite3 coronavirus.db "SELECT dateRep, SUM(deaths) OVER (ORDER BY year, month, day) AS
18        cumulativeDeaths , '$geoId' as country FROM CasesAndDeaths JOIN Date D ON D.ID =
19        CasesAndDeaths.data_ID JOIN CountriesAndTerritories CAT ON CAT.ID = CasesAndDeaths.geo_ID
20        WHERE geoId = '$geoId' ORDER BY year, month, day;" | tr '|', ' ' > "$Deaths"
21    plot_command+=" '$Deaths' using 1:2 with lines title columnheader(3),"
22    data+=("$Deaths") # Add the temporary file to the array, and delete later
23done
24
25# Remove the final comma from the gnuPlot
26plot_command="${plot_command%,}"
```

```

21
22 # Run the gnuPlot script by using a heredoc.
23 # Set the key's position on the top right of picture, out of the graph.
24 # Set the size and the font of the picture.
25 # Set the output file
26 # Set the X,Y label and the title
27 # Set the x data type as time
28 # Set the time format which is the same format as dateRep
29 # Set the display format
30 # Set the distance between two major tick marks
31 # Rotate the x coordinates
32 gnuplot << EOF
33 set key outside top right
34 set key autotitle columnheader
35 set terminal pngcairo size 1000,750 enhanced font 'Verdana,10'
36 set output 'graph.png'
37 set xlabel 'Date'
38 set ylabel 'Deaths'
39 set title 'Cumulative Deaths'
40 set xdata time
41 set timefmt "%d/%m/%Y"
42 set format x "%Y-%m-%d"
43 set xtics 60*60*24*30*2
44 set xtics mirror rotate by -60
45 $plot_command
46 EOF
47
48 # Delete all the temporary files by using a loop
49 for data in "${data[@]}"; do
50     rm "$data"
51 done

```

5.2 The graph of the top ten countries for cumulative deaths

