

# Chapter 2

## Data description

### 2.1 Overview of the database

The database used for analytics contains data on medical histories of patients between **January 2000** and **October 2018**.

It's important to note that the research work has been done on only a part of the whole database, consisting in **5 tables**. There is information available on **general practitioners**, **patients**, **diagnoses** and **prescriptions**: each macro-category is included in a separated table, so it's necessary to recognise the relationship between fields.

The 5 tables and their relative sizes are:

- *pazienti*, 1.015.618 tuples;  
Basic information about patients, identified by an encrypted UID;
- *nos\_002*, 1.015.618 tuples;  
Extension of *pazienti* with the same key, containing more detailed information and linkage with GPs;
- *cart\_pazpbl*, 15.460.199 tuples;  
Information about diagnoses and relative description;
- *cart\_terap*, 118.716.403 tuples;  
Information about therapies and prescribed medicines;
- *cart\_accert*, 151.478.456 tuples;  
Information about medical checks and examinations.

It's easy to see that the number of rows is varying a lot: there are a lot more prescriptions than diagnosis, and 150 millions of medical examinations.

Each diagnosis, prescription and examination is uniquely distinguished by the triplet patient-GP-date, which maps to *codice-userid-time\_last* in the database.

There are many other date fields, but *time\_last* (last update of the record) is the only one of type *timestamp*, making each one different from the others (it's unlikely to have a diagnosis or prescription for the same patient, by the same doctor and at the same exact moment).

Analysis is performed using dates in the *YYYY-MM-DD* format, since there's no need of unique data: there are lots of different prescription for the same patient made on the same day.

## 2.2 Description of the tables

Below is reported a brief description of the 5 tables, along with the main fields used for analytics and their meaning.

### 2.2.1 *pazienti*

The table *pazienti* includes information about patients. To maintain privacy, there are no real names: everything is **encrypted** as a 22-character string containing letters, numbers and special symbols.

Other relevant fields are:

- *data\_open*, date of beginning of the doctor-patient relationship;
- *nascita*, date of birth;
- *decesso*, eventual date of death;
- *comune\_di\_nascita*, name (and code) of the birth municipality;
- *Sesso*, birth sex;
- *pa\_convenzione*, type of convention with the Italian insurance system.

### 2.2.2 *nos\_002*

The table *nos\_002* contains the same information as *pazienti*, with additional fields which are essential to analyse data:

- *pa\_medi*, **encrypted** UID of the general practitioner of the patient;
- *pa\_cap*, zip-code of the patient (for geographical analysis);
- *pa\_pro*, province of the patient;
- *pa\_drevoca*, eventual date of termination of the doctor-patient relationship (a patient changing GP).

All the IDs of the GPs, along with all other data on GPs, are stored in an external table *users* (the research has been made considering a subset of the original DB). The latter will be used to identify active doctors, but doesn't contain any other useful information.

### 2.2.3 *cart\_pazpbl*

The table *cart\_pazpbl* comprehends the diagnoses associated to patients and relative GPs. Each diagnosis is defined by its **ICD-9** code, an international identifier for diseases maintained by the World Health Organization.

Fields summary:

- *codice*, patient ID;
- *userid*, corresponding to *pa\_medi* in *nos\_002* (general practitioner ID);
- *data\_open*, date of insertion of the diagnosis in the database;
- *time\_last*, timestamp of last edit of the tuple;
- *nome\_pbl*, a textual description of the diagnosis;
- *cp\_code*, code of the diagnosis according to the ICD-9 standards.

### 2.2.4 *cart\_terap*

The table *cart\_terap* contains the prescribed medicines for each patient. There is no linkage between diagnosis and prescriptions in the database, so additional work is required to detect correlation.

Each prescription is defined by an **ATC code**, from the Anatomical Therapeutical Chemical classification system maintained by the World Health Organization. Furthermore, there are **active principle code** and **authorisation for commerce code**.

Fields summary:

- *codice*, patient ID;
- *userid*, corresponding to *pa\_medi* in *nos\_002* (general practitioner ID);
- *data\_open*, date of insertion of the prescription in the database;
- *time\_last*, timestamp of last edit of the tuple;
- *co\_codifa*, AIC code (authorisation for commerce);
- *co\_des*, textual description of the medicine;
- *te\_npezzi*, number of boxes;
- *te\_attivo*, active principle code;
- *co\_atc*, ATC code;
- *euro*, price.

### 2.2.5 *cart\_accert*

The table *cart\_accert* includes all the medical checks and examination for each patient. The first goals of the research only included diagnoses and prescriptions, hence why the

comprehension of this table isn't as deep as the others.

Each examination is defined by an **ICD-9-CM** code, an extension of the ICD-9 database with standard procedures.

Fields summary:

- *codice*, patient ID;
- *userid*, corresponding to *pa\_medi* in *nos\_002* (general practitioner ID);
- *data\_open*, date of insertion of the examination in the database;
- *time\_last*, timestamp of last edit of the tuple;
- *ac\_nt\_code*, ICD-9-CM code;
- *ac\_des*, textual description;
- *cod\_ese*, exemption code for the examination.

## 2.3 ER model

After identifying the main fields, it's useful to build and include an **ER model**<sup>4</sup> to have a better understanding of the tables and their relationship.

