

# A Data Analytics Framework for Medical Prescription Pattern Dynamics

**Relatore:** Prof. Francesco Archetti

**Correlatore:** Prof. Antonio Candelieri

**Tutor aziendale:** Dott. Gaia Arosio

**Relazione della prova finale di:**

Ilaria Battiston  
Matricola 816339

Anno Accademico 2018-2019

## **Abstract**

Research on prescription pattern dynamics aims to highlight the issues concerning the raising antibiotic resistance among the Italian country, using a subset of data regarding the Campania region.

After a clear understanding of the problem and the available information along with its quality and possible uses, a first analytical approach allows reconstruction of diagnosis and prescriptions global trending, from which is possible to extract patient journeys and underline changes.

A deeper insight on antibiotics verifies national studies on overprescription, giving additional information on general practitioners' behaviour, brands popularity and seasonality of prescriptions, obtained with time series analysis and clustering.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Antibiotic resistance and misuse . . . . .	3
1.2	Italian National Sanitary Service . . . . .	4
1.3	Healthcare data . . . . .	6
1.4	Data classification . . . . .	6
1.5	Healthcare analytics . . . . .	9
<b>2</b>	<b>Goals definition</b>	<b>11</b>
2.1	Methodologies . . . . .	11
2.2	Considerations . . . . .	12
2.3	Practical goals . . . . .	13
2.4	Tools . . . . .	14
<b>3</b>	<b>Data description</b>	<b>15</b>
3.1	Dataset description . . . . .	15
3.2	Overview of the database . . . . .	16
3.3	Description of the tables . . . . .	17
3.4	ER model . . . . .	19
3.5	Variations through time . . . . .	19
<b>4</b>	<b>Information loss</b>	<b>20</b>
4.1	General overview . . . . .	20
4.2	Information loss on records . . . . .	21
<b>5</b>	<b>Global analytics</b>	<b>24</b>
<b>6</b>	<b>Patient journey</b>	<b>25</b>
6.1	Imposed criteria . . . . .	25
6.2	Initial data cleaning . . . . .	26
6.3	First approach . . . . .	27
6.4	Second approach . . . . .	28
6.5	Results comparison . . . . .	28
<b>7</b>	<b>Antibiotic trends analytics</b>	<b>30</b>
7.1	Identifying antibiotics . . . . .	30
7.2	Subset extraction . . . . .	31
7.3	Global statistics . . . . .	31
7.4	ATC rankings . . . . .	34
7.5	AIC rankings . . . . .	35

7.6 ATC and AIC trends . . . . .	41
<b>8 k-means approach for cluster analysis</b>	<b>46</b>
<b>9 Graph analysis</b>	<b>47</b>
9.1 Prescription coupling . . . . .	47
9.2 Graph databases . . . . .	47
9.3 Goals . . . . .	47
9.4 Practical approach . . . . .	48
9.5 Visualisation and analytics . . . . .	50
9.6 Considerations . . . . .	56
<b>10 Assessments of results and future directions</b>	<b>57</b>

# Chapter 1

## Introduction

### 1.1 Antibiotic resistance and misuse

Antimicrobial resistance is a rising global problem which threatens the effective prevention and treatment of an ever-increasing range of infections caused by bacteria, parasites, viruses and fungi.<sup>15</sup>

Microorganisms exposed to antimicrobial drugs develop the ability to defeat substances designed to kill them, making infections persist in the body due to the unsuccessful action of agents.

This issue threatens public health causing higher healthcare costs to treat patients, and potentially compromising surgeries and chemotherapy results due to the ineffectiveness of antibiotics. No one can completely avoid the risk of resistant infections, but some people are at greater risk than others (for example, people with chronic illnesses).<sup>16</sup>

Antimicrobial resistance occurs naturally over time, usually through genetic changes. However, **the misuse and overuse of antimicrobials is accelerating this process**. In many places, antibiotics are overused and misused in people and animals, and often given without professional oversight.<sup>15</sup>

Infections such as the common cold or sore throats are often countered with antibiotics, which have no effect against viruses and could as well put patients at the risk of suffering adverse reaction.<sup>18</sup>

Those critical issues are worsened by the fact that at the moment there are no antibiotic drugs in development, and no trials in the past 30 years led to discoveries of new antimicrobial medicines.<sup>40</sup>

Therefore, to minimise the development of resistance, contributing factors must be reduced, optimising the use of drugs. This work requires effective surveillance and follow-up of consumption, at a local and national level.

To be effective in the long run, work to optimise use of antibiotics must influence the prescribing practices of individual physicians. The goal is “rational use”, i.e. the correct patient receives the correct antibiotic at the correct dose and for the correct duration of treatment, in accordance with evidence-based guidelines. Over-prescribing should be

avoided without resulting in under-prescribing.<sup>17</sup>

Such work should be carried out close to the prescriber, something which also requires high resolution prescription data, down at the level of **individual prescribers**.

## **1.2 Italian National Sanitary Service**

The Italian National Sanitary Service (SSN) consists the complex of functions, activities and healthcare services offered by the State. It is based on subsidiarity, a general principle of the European Union law, stating that a central authority should have a subsidiary function, performing only those tasks which cannot be performed at a more local level.<sup>19</sup>

SSN is articulated in different responsibility levels divided among the State, Regions, institutions and organisations, along with private structures and the Health Ministry, which coordinates the national sanitary plan.

Citizens benefit of healthcare services paying a related ticket,<sup>20</sup> which represents the established way to contribute to expenses. It is used for:

- Specialist examinations;
- First-aid help in non-emergency situations;
- Thermal care.

Sanitary assistance on the territory is free, in fact general practitioners' visits are exempted from payment of tickets (while additional services such as certificates may require a fee). A general practitioner (GP) is a way for the citizen to access SSN in terms of global care and health education.

GPs manage types of illness that present in an undifferentiated way at an early stage of development, which may require urgent intervention. Their duties are not confined to specific organs of the body, and they have particular skills in treating people with multiple health issues.<sup>11</sup>

According to WONCA (World Organization of Family Doctors), they are responsible of supplying integrated and continuative care. Some fundamental skills and activities<sup>11</sup> to pursue this goal are:

- Communication with patients;
- Management of the practice;
- Clinical tasks;
- Problem solving;
- Holistic modelling.

In Italy, GPs have a crucial role in preventing diseases, understanding the symptoms, introducing patients to therapeutical approaches and monitoring the development or regression of illnesses.<sup>24</sup> Primary aid is guaranteed through diagnoses, prescription, therapy and basic levels of assistance.

Visits take place at the medical office, according to methodologies established by the doctor, generically on appointment or at the patient's domicile. After a visit, a common task of the general practitioner in agreement with SSN is prescribing drugs or further medical checks through a prescription.

Prescriptions are healthcare documents that govern the plan of care for an individual patient,<sup>25</sup> consisting in written authorization to purchase a specific medicine from a pharmacy.

Drugs are dispensed according to the national guidelines,<sup>22</sup> with the following regimes:

- OTC (Over The Counter), not subject to medical prescription;
- RR (Repeatable Prescription), to be sold after presenting a prescription;
- RNR (Non-Repeatable Prescription), to be sold after presenting a prescription which has to be renewed each time;
- RL (Limitative Prescription), used in hospitals and clinics;
- RMR (Ministerial Tracing Prescription), for narcotics and psychotropic substances.

### **1.2.1 Drugs and prescriptions**

The Italian Pharmacy Agency (AIFA) is the public institution for regulatory activity of drugs in Italy. Its main duty consists in all the activities related to the regulatory process of drugs, registering and authorising them to commercialization with a negotiated price.

Before a drug can be sold in pharmacies among the Italian territory, it must have received the authorisation by AIFA: each medicine is subject to checks regarding chemical, pharmaceutical, biological, toxicological, and clinical aspects and researches, to see if it satisfies security and efficacy standards.<sup>30</sup>

After passing all the quality controls, a product is assigned an unique AIC code for it to be identified with its specific details.

AIFA guarantees uniformity and equity of the pharmaceutical system, coordinating national and local authorities such as Regions. Procedures are based on safeness, innovation and accessible healthcare: pharmaceutical costs are regulated in a context of financial compatibility with industry competitiveness, while pursuing goals of economic balance and population' safeguard.

From the year 2000 every form of participation to sanitary expenses from citizens has been abolished,<sup>20</sup> yet most Regions have introduced special drugs classes with a fixed quota for each medical prescription or package, to remedy the profit deficit.

Currently existing categories according to reimbursement are:<sup>21</sup>

- A, entirely at the expense of SSN, comprehending essential medicines and the ones for chronic diseases;
- H, at the expense of SSN only in a hospital environment;
- C, fully paid by citizens according to the brand price.

Complete lists for classes A and H are publicly available, while each single drug or active principle can be individually looked up on the official AIFA database. Prescriptions can fall into any category, while OTCs follow a C regime.

Generic drugs have reduced costs, and to have a brand product the citizen must explicitly ask for it and pay the additional price.<sup>20</sup>

### 1.2.2 Antibiotic resistance in Italy

Italy is the European nation with the highest antibiotic resistance mortality ( 10.000 deaths every year), in terms of infection caused by resistant bacteria.<sup>32</sup>

Resistance to methicillin goes up to 38%, which is a non-negligible emergency; furthermore, people affected by MRSA (methicillin-resistant *Staphylococcus aureus*) have 64% more probability of death compared to individuals who didn't develop an antibiotic resistant infection.<sup>40</sup>

The problem is aggravated by the fact antibiotics are drugs sold over the counter in Italy, therefore individuals are able to get medicines according to their own judgement without a prescription by an expert — even if their diagnosis doesn't involve an antibiotic cure.

AIFA recommends actions to raise awareness among the population:<sup>4</sup>

- Using antibiotics only if prescribed by a doctor;
- Completing the therapy without interrupting it;
- Avoiding taking several antibiotics in short spans of time.

Patients are not the only ones contributing to antibiotic resistance: another aspect of it is the lack of ethical prescriptions from general practitioners.

## 1.3 Healthcare data

Electronic Health Records is a widespread application of big data in medicine. Patients have their own digital record which includes demographics, medical history, allergies, laboratory test results etc. Records are shared via secure information systems and are available for providers from both public and private sector.<sup>38</sup>

## 1.4 Data classification

Healthcare data needs to be classified according to national or international standards, in a way that information cannot be misinterpreted and fields can be mapped to wider categories.

### 1.4.1 ICD

ICD (International Classification of Diseases) is the foundation for global identification of health trends and statistics, and the international standard for reporting diseases and health conditions. It is the diagnostic classification standard for all clinical and research purposes,<sup>15</sup> maintained by the World Health Organization (WHO).

ICD defines the universe of diseases, disorders, injuries and other related health conditions, listed in a comprehensive, hierarchical fashion that allows for:

- Easy storage, retrieval and analysis of health information for evidenced-based decision-making;
- Sharing and comparing health information between hospitals, regions, settings and countries;
- Data comparisons in the same location across different time periods.<sup>26</sup>

The ICD is revised periodically and is currently in its 10th version, while in Italy the latter has only been adopted to classify death causes.<sup>27</sup> Until a complete upgrade, the diagnostic system is standardised using ICD-9.

#### ICD-9

ICD-9, officialized in 1978, is the 9th version of the International Classification of Diseases, ordering diseases and traumas in groups according to defined criteria, allowing a common language to code information related to morbidity and mortality for comparisons and statistics.<sup>27</sup>

ICD9-CM is an adaption to ICD-9 used in Italy to assign diagnostic and procedure identifiers, providing additional morbidity detail. Diagnoses are extended with codes for surgical, diagnostic and therapeutical procedures.

It is composed by a group of three digits followed by up to two optional ones adding further details, separated with a dot:

1. The first number (001-999) represents the macro-category based on the type of the disease or the injury they describe;
2. The second group provides more specific information about the type, location, and severity of the disease or injury.

There are also two sets of alphanumeric codes in ICD-9-CM. E-codes describe external causes of injury, while V-codes describe factors that influence health status and/or describe interactions with health services.<sup>28</sup>

Example: 414.12, falling in diseases of the circulatory system (390-459), specifically into ischaemic heart disease (410-414).

- 414: other forms of chronic ischaemic heart disease;
- 1: aneurysm and dissection of heart;
- 2: dissection of coronary artery.

### 1.4.2 ATC

The Anatomical Therapeutic Chemical Classification System (ATC) is a drug coding system adopted worldwide, controlled by World Health Organisation.

Medicines are divided into different groups according to the organ or system on which they act, their therapeutic intent or nature, and the drug's chemical characteristics. Different brands share the same code if they have the same active principle and indications.

One single drug can have multiple codes, since ATC also comprehends instructions regarding administration or use, and a code can represent more than one active ingredient.

The ATC classification is composed by seven alphanumeric symbols split into five adjacent hierarchical sets, defined levels and having the following structure:<sup>31</sup>

1. One letter indicating the anatomical/pharmacological main group among 14;
2. Two digits indicating the therapeutic subgroup;
3. One letter indicating the pharmacological subgroup;
4. One letter indicating the chemical subgroup;
5. Two digits indicating the chemical substance.

The ATC system also includes many defined daily doses (DDDs). This is a measurement of the assumed average maintenance dose per day for a drug used for its main indication in adults.

Alterations in ATC classification can be made when the main use of a product has clearly changed, and when new groups are required to accommodate new substances or to achieve better specificity in the groupings. Changes are twice annually submitted to the WHO official database.

Example: C03BA12.

- C: cardiovascular system;
- 03: diuretics;
- B: low-ceiling diuretics, excluding thiazides;
- A: sulfonamides, plain;
- 12: Clorexolone.

### 1.4.3 AIC

AIC represents authorisation to admission to commerce of a medicine, and is a 9-digit code conceded by AIFA after a careful check of safeness and efficacy. It is a sort of "identity card" of the drug, since it contains the essential characteristics defining it.<sup>29</sup> Different brands are identified by different AIC.

AIC establishes:

- The drug name;

- Its composition (active principle);
- Description of the fabrication method;
- Therapeutical instructions, contraindications and adverse reactions;
- Dosage and way of administration;
- Conservation measures;
- Characteristics of the product and its packaging;
- Brochure;
- Risks evaluation for the environment.

Every possible modification of those characteristics involves a further request for authorization to AIFA. Official databases such as Fedefarma contain information related to every product, along with its eventual expire of authorisation.

## 1.5 Healthcare analytics

**Healthcare analytics** is a field of growing importance which helps understanding the statistical perspective of results collected in the healthcare area.

Extracting insights can be a complex challenge: health big data gives a huge *volume* and *variety* of information, therefore accessing the resources in a quick way is necessary.

Other issues to deal with are *veracity*, *validity* and *viability*, fundamental characteristics to ensure reliable and relevant analytics. Checking for integrity and quality can be difficult to verify without domain knowledge.<sup>1</sup>

One of the possible applications, given the concerning issue of antibiotic resistance, is explaining the situation through statistics and trends, obtaining practical results alongside theoretical scientific research.

Big data can assess the appropriateness of prescribing through the existing classification systems, comparing their patterns within time ranges. A general view is essential to identify specific unusual changes, extending national studies with a perspective centred on products and reduced geographical areas.

There are five main necessary fields for analysis:<sup>5</sup>

1. **Spatial data**, in different granularity levels;
2. **Personal data** of patients;
3. **Temporal data**, for time-series analysis;
4. **Pharmacotherapeutic data**, classifiable according to different identification codes;
5. **Diagnostic data**, for cross-validation of diagnoses.

The two main risks encountered while doing analysis are information loss and inappropriate prescribing, that compromise the quality of statistics. Data may be incomplete,

biased or filled with noise: another goal of analytics is to contrast incompleteness and incorrectness, obtaining coherent and clear results.

# Chapter 2

## Goals definition

### 2.1 Methodologies

This research is specifically focussed on the doctor-patient relationship: person-centred care concerning diagnoses and prescriptions, analysing their changes according to habits, physiological aspects, time and geographical area of both interested parts.

There are several external factors influencing market trends of medicines, for instance the advent of generic drugs, having the same active principle and bioequivalence with brand medicines, but lower prices thanks to their dissociation from pharmaceutical companies.

The concept of generic drug has been introduced in Italy in 1995, to be legally formalised in 1996,<sup>3</sup> yet it has initially been perceived by general practitioners and pharmacists as a mere instrument to save money at the expense of quality.<sup>10</sup>

Only in the past 10 years, advertisement efforts have been made to make the population aware of the strict quality checks and the reliability of generic drugs, starting a slow switching process. If sales of a brand product decrease, then, it might have been replaced with its equivalent.

This is only an instance of event which cannot be predicted using the raw data: information about market withdrawals, advertisement or economic availability need to be cross-checked with analytical results.

Considering a wide time span is essential to have an overall idea of data quality, information loss and potentiality of available resources. All those factors contribute to progressive rescaling of the dataset:

1. Detailed analytics must be done according to restricted areas (pathologies, products);
2. Different time spans can be compared, going deeper into the general view;
3. Incorrect or unclear information has to be removed, causing retention of total records whose amount gets narrower while cleaning is in progress (funnel).

The dataset is modelled in a relational schema, which allows organization of records in structures (tables) to maintain integrity and compatibility between different kind of fields.

This allows flexibility using dynamic views and query optimisation based on set theory. For a better fruition of the workload, table names have been changed to standard English keywords.

## 2.2 Considerations

Since health big data can provide a high variety of information, it is required to have some clear *objectives* in mind to keep on track and avoid losing focus.

The research is centred on the **changes of prescription patterns of antibiotics**: this is a wide goal, and more information is needed to achieve results. It's necessary to narrow the field down, only concentrating on some classes of medicines, and define subgroups of GPs and patients in a restricted amount of time.

To obtain constraints the more objective as possible, some further analysis can be useful. For instance, focussing on chronic patients is a relatively fast way to reduce the huge amount of rows to elaborate, but causes the loss of most information.

The first step to take, having a deep understanding of the data, is recognising the extent and impact of the **progressive information loss**, to define the final amount of clear records. Trying to fix mistakes is a risk, since the outcome could be incorrect, so deleting is the most practical option.

Removing unclear and futile data may not be enough to have consistent results: 18 years of data is a wide range, and *splitting the dataset* or deciding to only consider a smaller scope can be beneficial for the analysis quality.

Some constraints can be imposed on general practitioners as well: since the results have to be coherent and accurate, it's best practice to only consider **active GPs** with a **constant number of patients** (to be defined).

This would partially remedy the fact that doctors may have different approaches to the same disease.

The creation of cohorts (chronic patients) among the statistical population is an example of **cluster sampling**.

After the initial parsing, there will be a rough draft of the final result which then will be subject of the following steps:

1. Further analysis on data correctness (record linkage);
2. Elaboration of the statistics and time series clustering.

Another relevant instance for analysis is the **subset** of diseases to consider: choices have to be made according to *external studies, marketing researches* and further *discoveries on the provided data*. Focussing on the **most common ones** is a guideline to start.

Having an idea of which illnesses and prescription have unstable patterns might give a better vision, and can be done through statistics on the whole database.

Some examples of analytics are:

- Most common diseases through the years;
- Most common *chronic* diseases through the years;
- Changes of the number of prescriptions for diseases in the same area;
- Changes of prescriptions based on the patient phenotype or market trends.

An obstacle to perceiving the meaning of results is the restricted domain knowledge: to compensate, confronting some experts in the field is required. The team comprehends computer scientists, statisticians, biologists and healthcare workers.

## 2.3 Practical goals

Having a better vision of the medical domain and the rising issues in the Italian system, before approaching analytical procedures it is necessary to underline practical objectives and expected outcomes.

The first essential statistic to extract is about progressive information loss: this makes possible to estimate how much data can give reliable and complete results, compared to the total. It is defined *progressive* since the more iterations of cleaning are being made, the more data is going to be lost.

Parsing remaining data is essential to have a correct functioning of database interrogations. Fields need to be checked for correctness using lookup tables, record linkage or regular expression matching; text has to be cast to numbers to apply mathematical operations, and dates can be divided into months and years.

Performances require a high level of optimisation due to the huge workload of information, obtained through targeted expressions, strict constraints and better memory management of records.

A global overview of the data allows to detect anomalies or trends of specific areas to focus on: potential analysis pertains dividing patients according to distinctive characteristics (sex, age group) to observe variations of diagnoses and prescriptions.

After collecting the first batch of results and checking them using external knowledge, information with unusual patterns is outlined, and further examination is made on a subset of features.

Final outcomes are then subject of more advanced techniques, which include:

- Time series analysis;
- Clustering of trajectories;
- Graph algorithms.

## 2.4 Tools

Modern technologies make possible to process big data with reduced costs: there are plenty of data stores, development and integration tools for each research purpose.

Since the project requires a relational structure along with statistical computing and machine learning, analysing and elaborating the health data is made through:

- **PostgreSQL**,<sup>12</sup> an open source object-relational database management system known for its robustness and reliability:
  - The development platform to interface with the web server containing the dataset is **PgAdmin 4**;
  - Indexes and Common Table Expressions are useful to avoid huge computational times.
- **Neo4j**,<sup>13</sup> a native graph database which gives data a different representation, processing entities as nodes while highlighting their connections:
  - Queries are expressed using **Cypher**;
  - The additional plugin Graph Algorithms returns implemented, parallel version of common network problems.
- **R**,<sup>14</sup> a free software environment for statistics and graphics computing, offering a wide range of techniques and formulae:
  - *ggplot2* is a package to create and visualise plots;
  - Clustering is performed using embedded functions.

Reports and slide-shares have been created and accessed using the **Google Suite**.

Due to the amount of sensitive data, detailed results are going to be omitted: the final conclusions will be a product of aggregation and schematisation.

# Chapter 3

## Data description

### 3.1 Dataset description

The available dataset is provided by Dedalus, market leader of the clinical software area, supporting doctors and their processes through its society Millennium.

Dedalus offers a wide range of solutions, such as surgical journey management, drug management systems, tracking for healthcare and enterprise resource planning.<sup>37</sup>

Millenium is a specific national infrastructure focussing on primary healthcare, using Dedalus products to develop national and regional level projects.

Research is conducted using data collected by an interface used by general practitioners to track interactions between them and the population benefiting of the national healthcare system.

The database contains recorded medical history of patients using healthcare services in the region Campania, focussing on the doctor-patient relationship.

Since the database is not regulated by local laws, it encounters a greater risk of inaccuracy, unlike pharmacies or tax registers.

General practitioners are the only responsible of filling values, therefore there is no assurance of completeness and correctness of data: mistakes are common, as well as missing information. A part of patients journey happens in hospitals or specialised medical offices, and those records aren't present since belonging to external sources.

Only a part of the whole amount of drugs (and examinations) require a written prescription: most medicines are given over the counter, and there is no certainty that a patient is going to buy that specific drug or the generic equivalent. Linkage between prescriptions and actual purchase is missing.

Furthermore, not all prescriptions are ethical: antibiotic resistance is an ascertained issue, and general practitioners can be influenced by pharmaceutical companies, resulting in lack of objectiveness.

Data is highly sensitive: despite encryption of all names, a considerable amount of geographical information is available. To avoid cross-checking using location, for results

to be published patients or doctors must be aggregated in groups whose numerosness exceeds a fixed value (rule of thumb states 3).

Other sensitive information such as email addresses and passwords to log in the system is irrelevant for analytics, therefore it is safe to remove anything not strictly related to the research purposes.

## 3.2 Overview of the database

The database used for analytics contains data on medical histories of patients between **January 2000** and **October 2018**. This leads to some observations:

1. The year 2018 is present only up to June, so it cannot be used while making time series within years (there is going to be a drop of values due to incompleteness which may lead to wrong conclusions);
2. A timespan of 20 years is too wide to make consistent analytics;
3. Early dated records might contain outdated or incomplete information.

Global inferences have been made with the entire dataset, while the need of detailed recent reports leads to the decision of using a limited range of years for prescription pattern changes and patient journey.

The research work has been done on only a part of the original Millewin database, consisting in **4 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 identify the relationship between fields.

The 4 tables with their sizes are:

- *patients*, 1 015 618 tuples;  
Basic information about patients, identified by an encrypted UID;
- *patients\_doctors*, 1 015 618 tuples;  
Extension of *patients* with the same key, containing more detailed information about patient-doctor relationships and linkage with GPs identifiers;
- *diagnoses*, 15 460 199 tuples;  
Information about diagnoses and relative description;
- *prescriptions*, 118 716 403 tuples;  
Information about therapies and prescribed medicines.

It is noticeable that the number of rows is varying: there are more prescriptions than diagnoses, since the first tend to happen more often.

Each diagnosis and prescription is uniquely distinguished by the triplet patient, doctor, date. Dates are at level of timestamp, making each one different from the others (it's improbable 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 non-unique data still allows to aggregate results and identify patterns. There are several different prescriptions for the same patient made on the same day.

## 3.3 Description of the tables

Before being able to work with the data, it is essential to understand its structure, functioning and trending within time.

Below is reported a brief description of the 4 tables, along with the main fields used for analytics, statistics and machine learning.

### 3.3.1 *patients*

The table *patients* includes information about patients. To ensure privacy dealing with sensitive data, there are no full names: everything is **encrypted** as a 22-character string containing letters, numbers and special symbols.

Other relevant fields are:

- *birthdate*, date of birth;
- *death*, eventual date of death;
- *birth\_municipality*, name (and code) of the birth municipality;
- *sex*, birth sex;
- *convention*, type of convention with the Italian insurance system.

### 3.3.2 *patients\_doctors*

The table *patients\_doctors* contains information similar to *pazienti*, with additional fields focussing on their relationship with the general practitioners, which are essential to link and analyse data:

- *userid*, **encrypted** UID of the general practitioner of the patient;
- *date*, date of beginning of the doctor-patient relationship;
- *postcode*, zip-code of the patient (for geographical analysis);
- *province*, province of the patient;
- *revocation*, 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 does not contain any other information relevant for analysis, since active doctors can be extracted from other tables.

### 3.3.3 *diagnoses*

The table *diagnoses* 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.

Summary of most important features:

- *id*, patient ID from *patients*;
- *userid*, corresponding to *doctor* in *patient\_doctor* (general practitioner ID);
- *date*, date of insertion of the diagnosis in the database;
- *last\_update*, timestamp of last edit of the tuple;
- *description*, a textual description of the diagnosis;
- *IDC9*, code of the diagnosis according to the ICD-9 standards.

### 3.3.4 *prescriptions*

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

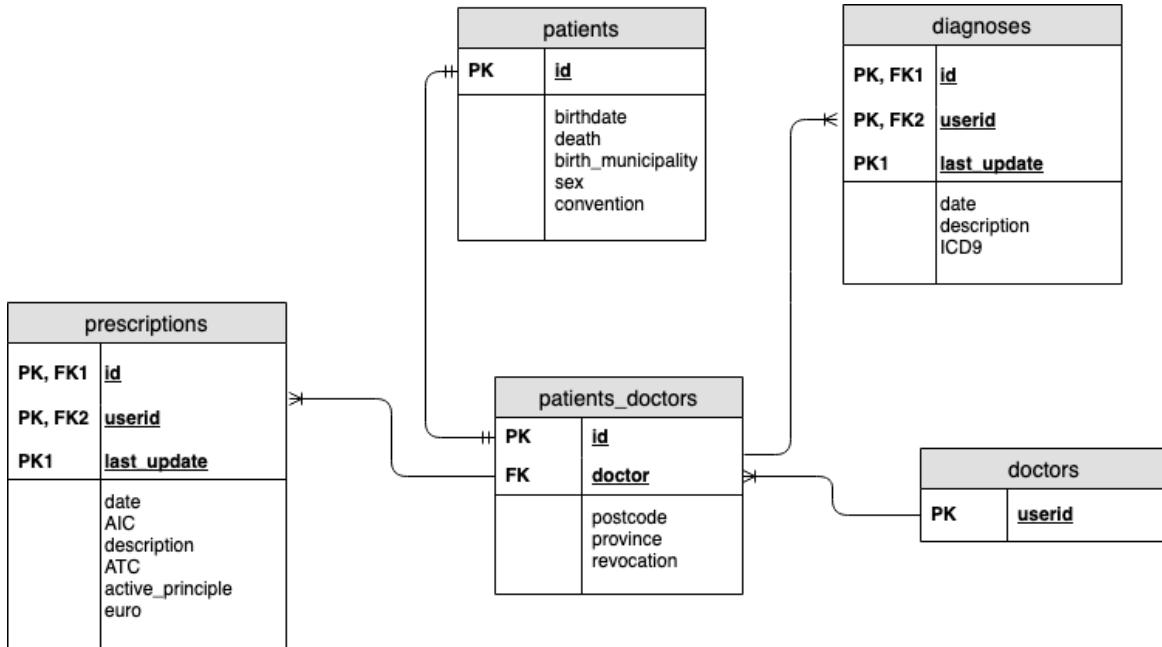
Each prescription is defined by an **ATC code**, from the Anatomical Therapeutical Chemical classification system maintained by the World Health Organization. Furthermore, there are **active principle code** and **authorisation for commerce code(AIC)**.

Fields summary:

- *id*, patient ID;
- *userid*, corresponding to *pa\_medi* in *nos\_002* (general practitioner ID);
- *date*, date of insertion of the prescription in the database;
- *last\_update*, timestamp of last edit of the tuple;
- *AIC*, AIC code (authorisation for commerce);
- *description*, textual description of the medicine with dosage;
- *pieces*, number of boxes;
- *active\_principle*, active principle code;
- *ATC*, ATC code;
- *euro*, price.

## 3.4 ER model

After identifying the main fields, it is possible to build and include an **ER model**,<sup>6</sup> to provide immediate comprehension of how entities are related, visualising information to identify keys and unique fields, essential for joining data.



## 3.5 Variations through time

To give a general idea of variations of patients and magnitude orders, a snapshot of 2010-2017 is used to show patterns and differences, counting the number of patients getting at least a diagnosis for each year:

	2010	2011	2012	2013	2014	2015	2016	2017
Number	329 715	320 253	316 431	320 948	324 920	323 441	330 641	326 987
Age mean	47,43	47,89	48,44	48,68	48,98	49,78	50,12	50,54
Age st. dev.	20,7	20,81	20,84	20,86	20,87	20,84	20,85	20,79
Women	185 035	179 700	178 248	180 304	182 286	181 313	185 729	183 419
Men	144 016	140 330	137 514	139 780	141 727	141 158	143 931	142 363

Another example can be obtained counting the number of prescriptions each year:

2010	2011	2012	2013	2014	2015	2016	2017
7 326 923	7 108 288	7 269 855	7 541 539	7 777 753	7 847 313	7 856 246	7 785 325

# Chapter 4

## Information loss

The first analysis aims to give a **qualitative assessment** of the whole database, giving an overall idea on how impactful is the progressive information loss.

After understanding the correctness and completeness of records, some fields will be eventually excluded from the analysis, while others that cannot be removed will have a major impact on the results.

Having missing fields, particularly in the process of joining tables, can lead to a cumulative augmentation of the information loss: empty data such as the patient date of birth or gender will cause the deletion of the entire patient, in case analytics is centred on pathologies by age or gender.

Joining is in fact an operation which requires all fields of reference to be present, combining entries of the selected tables.<sup>8</sup>

### 4.1 General overview

Before starting running queries, there are some factors to consider and issues which sometimes cannot be addressed just with database interrogations:

1. The geographic information is sometimes imprecise and hard to comprehend, since it consists in text fields;
2. Some diagnoses descriptions don't match with the corresponding ICD-9 code;
3. A consistent amount of missing data is originated while a general practitioner doesn't prescribe anything but makes other operations (medical certificates, examinations and such);
4. Prescriptions in hospitals are missing;
5. Some medicines are given over the counter without requiring a prescription, therefore there is no entry in the DB;
6. General practitioners might prescribe a medicine to a different patient than the one who has the pathology (relatives, friends, ...);

7. There is inappropriate prescribing, antibiotic resistance, misuse or over-use of medicines;
8. A patient can change doctor, so the approach to the same disease may vary.

Furthermore, there have been noticed some common instances of incorrect data:

- Some dates don't fall in the acceptable range (e.g. 1999, 2034, ...),  
a date is considered wrong if it's before 01-01-2000 or after 10-01-2018;
- A lot of fields are empty or null.

Overall, the loss on single records isn't a relevant issue since the total amount of rows allows safe removal, yet the cumulative augmentation (funnel analysis) gives a progressive deletion of information.

## 4.2 Information loss on records

### 4.2.1 *patients* and *patient\_doctors* (1 million of tuples)

#### Summary

Both tables contain similar data related to patients, with a 1 : 1 correspondence between primary keys, so joining is an elementary operation and there is no information loss.

- Patients with null or empty sex: 2 752 → 0, 27%;
- Patients with sex different from M and F: 54 → 0.005%;
- Patients with beginning of the patient-doctor relationship outside the accepted range: 226 858 → 22, 33%;
- Patients with null province: 99 981 → 9, 84%.

A subset of patients is created through an auxiliary view to highlight the final progressive data loss compared to the original tables, joining patients and patients\_doctors according to those constraints

1. Join on equal patient code (*id*);
2. Not null date of birth;
3. Dates between 2000 and 2018;
4. Existing and not null sex;
5. Existing and not null province.

The total rows respecting all those constraints are 713 352, so approximately 300 000 tuples (patients) have been deleted.

This result implies that during analysis there will be at least 1/3 of the data which is going to be removed due to incompleteness and inaccuracy: not considering patients will imply deleting their diagnoses and prescriptions as well.

### 4.2.2 *diagnoses* (15 millions of tuples)

#### ICD-9

An ICD-9 code is correct if it's in the official ICD-9 database.<sup>7</sup> General practitioners may use different formats and notations, so before removing codes each one has been subject to some parsing.

An ICD-9 code is considered wrong if there is no match in the official DB after any of those transformations:

1. Removal of the dot;
2. Addition of 0 at the beginning;
3. Addition of 0 at the end;
4. Removal of the last 0.

There are 76 incorrect ICD-9 codes, a very small amount considering the total records.

#### Summary

- Incorrect ICD-9 codes: 76 → 0,0005%;
- Null or empty ICD-9 codes: 2.103.169 → 13.02%;
- Null or empty descriptions: 656.067 → 4,24%;
- Dates out of range: 807.985 → 5,23%.

### 4.2.3 *prescriptions* (118 millions of tuples)

#### ATC code

The ATC code, unlike ICD-9, has an univocal format (a numeric string of 7 digits), and no parsing is needed. All the codes have been checked with the ontologies in a well-known biology portal,<sup>9</sup> and an ATC is considered incorrect if there is no match.

There are 1 750 292 non-existent codes, which are caused by:

1. Alteration of the codes within the years without updating the database;
2. Single prescriptions indexed using the superclass code;
3. Codes only recognised by local pharmacies.

The latter is the most common reason: there are up to 90 000 occurrences of a single unofficial code. Those numbers, despite being high, aren't really impacting results considering the total amount of 118 millions of records.

## Summary

- Incorrect ATC codes: 1 750 292 → 0, 1, 47%;
- Null or empty ATC codes: 1 577 749 → 1, 33%;
- Null or empty AC codes: 1 310 719 → 1, 1%;
- Null or empty descriptions: 101 248 410 → 85, 28%;
- Dates out of range: 3 472 119 → 2, 92%.

Clearly, the prescription description is not a field which can be used for reliable analytics since empty values prevail.

The field *pieces* (number of boxes) might be useful to check for appropriate prescribing, but since the field is a string it requires casting and parsing to integer, and it's more relevant to focus on prescription patterns analysis.

# Chapter 5

## Global analytics

# Chapter 6

## Patient journey

After the preliminary analysis and knowing the main focus area, there is enough information to begin reconstructing the **patient journey**, a complete set of data representing a patient history and relationships with primary healthcare, to then identify patterns and changes.

An objective definition of patient journey can be created using the following guidelines:

1. Patients with **complete medical history** for a fixed amount of years;
2. Records with patient, prescribing GP, diagnosis and prescription on the **same date**;
3. Only **first-time** diagnoses and prescriptions considered.

The imposed criteria is strict: taking diagnoses and prescriptions on the same day means removing *all prescriptions* following the first diagnosis. In other words, all the instances of patients coming back to their GP to renew a prescription for a chronic disease have been deleted.

This can be useful to extract a cohort of patients beginning their treatment, and analyse the variations of first-time prescriptions, especially for chronic illnesses. It's important to notice that not all doctors may have patients getting new diagnoses.

### 6.1 Imposed criteria

Aside from the completeness and correctness of the data, there are more restrictions to maintain the consistency:

- The prescribing general practitioner mustn't change in the time range;
- The patient mustn't be deceased;
- There must be a sanitary convention;
- The general practitioner must be active.

All those constraints can be checked using the related fields in the database: *revocation* for interruption of the relationship, *death* for death and *convention* for the sanitary convention.

The table *users* contains all the IDs of active general practitioners, so joining it with other tables is enough to remove all the rows with an inactive GP. Data is up to date, but since the patient journey includes 2018 and requires consistency of history (the focus is on the most recent information) there is no need to check for active GPs in the previous years.

The biggest risk is again the **loss of information**: the impact of data cleansing is heavy, and the obtained results might not give an insightful enough prospective.

## 6.2 Initial data cleaning

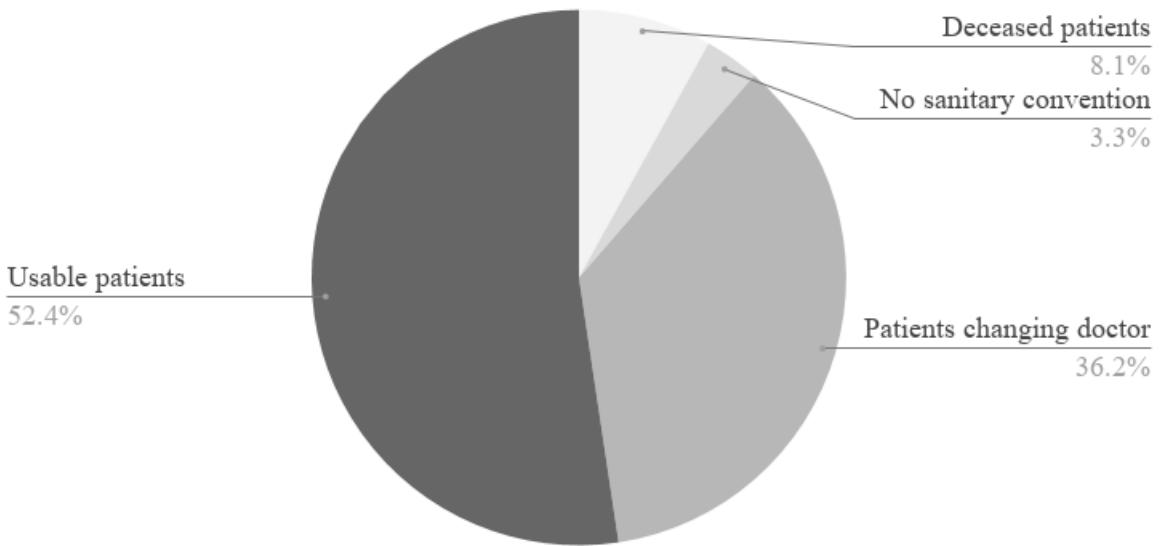
An initial data cleaning has been made on the whole database to have a first understanding of the potential information loss.

In this case, having such a big amount of tuples is useful: it's possible to remove a considerable percentage of them without losing generality and still having numerous samples.

Information on the general loss is already available thanks to the specific analysis on each single field, so the shown data cleaning will only consider patients and GPs.

The active general practitioners are **432**: this result has been retrieved counting the different IDs in *users* (438) and removing the ones not present in *patients\_doctors* (6).

The following *pie chart* illustrates the data loss on patients according to the criteria defined in the previous section.



About half of patients is going to be lost, due to not respecting the consistency criteria. Starting from a million of records, concrete results are still obtainable.

## 6.3 First approach

The first approach consists in testing with an **arbitrary range constraint**: all dates must fall in the span between 2010 and 2018.

The first analysis has pure research and testing purposes, to understand the impact of a cutting the dataset in terms of information loss. All the previously introduced criteria must be considered as well, so there must be a continuous doctor-patient relationship between active GPs and non-deceased patients with sanitary conventions.

The outcome is a patient journey table containing data from 2000 to 2018, with a total amount of **144 618** tuples: this means that there are roughly 150k of first-time diagnosis and prescriptions to patients.

### 6.3.1 Results breakdown

Seeing that the starting tables had number of rows in the order of millions, some deeper analysis is necessary to figure out the causes of this significant loss.

The 144 618 complete tuples are composed by:

- 27 733 patients;
- 422 general practitioners;
- 1 381 unique diagnoses;
- 904 unique prescriptions.

Further causes for those low values can be found counting how many dates would not fall in the considered range. The percentage of records with date earlier than 2010 in each table is:

- Patients: 75%;
- Diagnosis: 54,6%;
- Prescriptions: 44,7%.

There is an enormous loss on patients: the possible reason might be that most patients started treatment earlier than 2010, plus diagnosis and relative prescription are in different dates.

8 years is a too wide range to obtain a consistent patient journey, and such information loss isn't negligible: the conclusion of the first approach is that introducing time boundaries is something which needs an accurate control, to avoid missing out most of the data.

## 6.4 Second approach

Since just removing according to the date is an abrupt approach, it's necessary to tune parameters and introduce more detailed constraint, to have a bigger amount of information.

The focus is on the number of prescriptions: given a small range of years, only patients with at least one new prescription are going to be considered. All the previous criteria must be respected, so there must be a continuous doctor-patient relationship between active GPs and non-deceased patients with sanitary conventions.

This methodology allows cluster sampling without having to remove half of the dates: criteria based on the number of prescriptions creates another patient cohort which can be used to accurately select rows from the other tables.

The proposed time range is 2016-2018: pharmaceutical companies generally use the last two years of sales, so picking the last three years gives additional information without compromising the consistency of data.

The outcome is a patient journey consisting of 1 465 005 tuples: almost 10 times the previous result. This leads to two important statements:

1. The time range is appropriate, since the number is large enough to make analysis without loss of generality;
2. The new imposed criterion gives more consistent data and the possibility to build time series.

More cleaning is required to link diagnoses and prescriptions, since there is not a 1 : 1 correspondence: multiple diagnosis and prescriptions may be associated to the same date. This can be done using a lookup table.

### 6.4.1 Results breakdown

The 1 665 005 complete tuples are composed by:

- 230 381 patients;
- 422 general practitioners;
- 1 381 unique diagnoses;
- 904 unique prescriptions.

Only 7% of the total prescription has been taken in consideration, yet a million and half is still a satisfying amount.

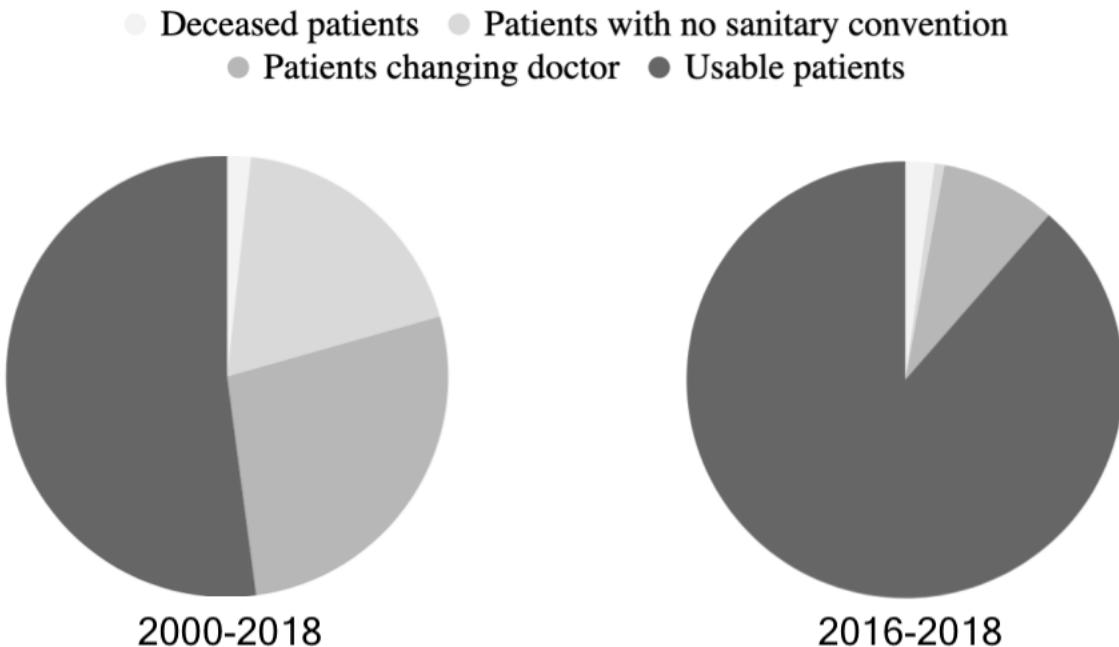
## 6.5 Results comparison

Comparing the two patient journey outcomes through graphs is a good way to visualize changes and improvements.

### 6.5.1 Changes in data composition

### 6.5.2 Improvement on patients information loss

A pie chart for data loss on patients in the range 2016-2018 has been made to compare with the previous one.



It's easy to see that the number of usable patients has noticeably increased: using a smaller time span reduces the chances of death and change of GP.

### 6.5.3 Funnel graph of patients

A funnel graph is useful to see how imposing every restriction made the patients number decrease: starting from a million, in the end there only is about  $\frac{1}{4}$  of it.

# Chapter 7

## Antibiotic trends analytics

Having obtained a general analytical view on couples of first-time diagnoses and prescriptions, and cross-checking those results with global trends, the research focusses on antibiotic prescriptions to identify changes of patterns.

Defined time frames in limited temporal windows are selected to make a detailed trajectory analysis related to prescriptive appropriateness compared to antibiotic resistance, without going into details of pathologies.

Such analytics are useful for research, highlighting rising health issues alongside the AIFA reports, and to support pharmaceutical companies using the potentiality of healthcare data. Having information about AICs allows to give a new insight not only on ethical matters, but also on the global market trends.

### 7.1 Identifying antibiotics

Antibiotics, also known as antibacterials, are medications produced by microorganisms that destroy or slow down the growth of bacteria. They include a range of powerful drugs and are used to treat diseases caused by bacteria.

A doctor can prescribe a broad-spectrum antibiotic to treat a wide range of infections. A narrow-spectrum antibiotic is only effective against a few types of bacteria. In some cases, a healthcare professional may prescribe to prevent rather than treat an infection, as might be the case before surgery.

ATC codes related to antibiotics divided by class are:

- A01AB, A02BD, A07A;
- D01, D06, D07C, D09AA, D10AF;
- G01;
- J, the whole class;
- R02AB;
- S01A, S02A, S03A (removing hormones).

## 7.2 Subset extraction

The first criteria to impose, considering analytics is going to be made on antibiotic prescriptions, is the actual presence of such. This can be reworded by removing all data related to patients who never had a prescription of an ATC code among the ones previously listed.

There are 794 267 patient with at least one antibiotic prescription in the whole time range (2000-2018), whose total prescriptions consist in 114 044 470 tuples. This implies the remaining 200 000~ patients only have 4m~ prescriptions, which is probably justified by a healthier physical state.

Since antibiotics patterns are subject to major changes during years, a big amount of data can be dispersive: pharmaceutical companies make analysis only considering 2-3 years, yet for research purposes an intermediate range is the most informative. A good time span is 10 years, considering most recent ones: since 2018 has to be excluded due to incompleteness, 2008-2017 is the final choice.

Progressive data loss has to be considered, imposing additional constraints of completeness and accuracy. The latest version of the used dataset is a subset of the table prescriptions, with the following restrictions:

- AIC corresponding to an antibiotic;
- Prescription date between 2008-01-01 and 2017-12-31;
- Active general practitioners;
- Patients with usable information about sex, date of birth and location.

The obtained record set is composed by 8 386 057 tuples, each representing a single prescription.

## 7.3 Global statistics

Global statistics help to have a general overview of the data, aggregating values to identify most impactful changes and the broad behaviour.

### 7.3.1 Average prescriptions

The average number of prescriptions has been calculated only considering the number of patients with at least one antibiotic prescription from 2008 onwards (670 634).

Values are calculated using the R functions `mean` and `sd`, only considering the number of patients who received at least one prescription in the determined year.

Year	Mean	Standard deviation
2008	2,72	2,73
2009	2,75	2,78
2010	2,72	2,76
2011	2,69	2,73
2012	2,68	2,76
2013	2,74	2,84
2014	2,79	2,88
2015	2,79	2,79
2016	2,77	2,91
2017	2,72	2,84

Values show the mean tends to slowly increase, and although there are patients getting up to 30 antibiotic prescriptions each month, the standard deviation is low: the mode both for year and month is 1, so lower values weigh more.

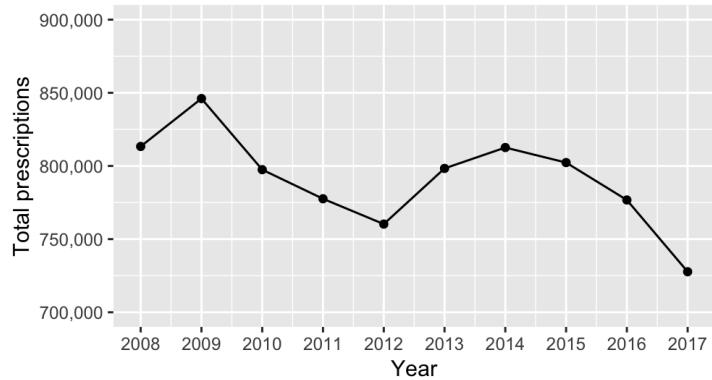
Statistics related to the number of patients with only one antibiotic prescription, making statistics comparing time spans between 2008-2017:

Interval	Mean	Standard deviation
Year	118 301,4	2 831,65
Month	40 130,83	6 004,2

### 7.3.2 Yearly prescriptions

The number of yearly antibiotic prescriptions is:

Year	Antibiotics
2008	813 337
2009	846 067
2010	797 472
2011	777 582
2012	760 280
2013	798 320
2014	812 594
2015	802 323
2016	776 756
2017	727 723

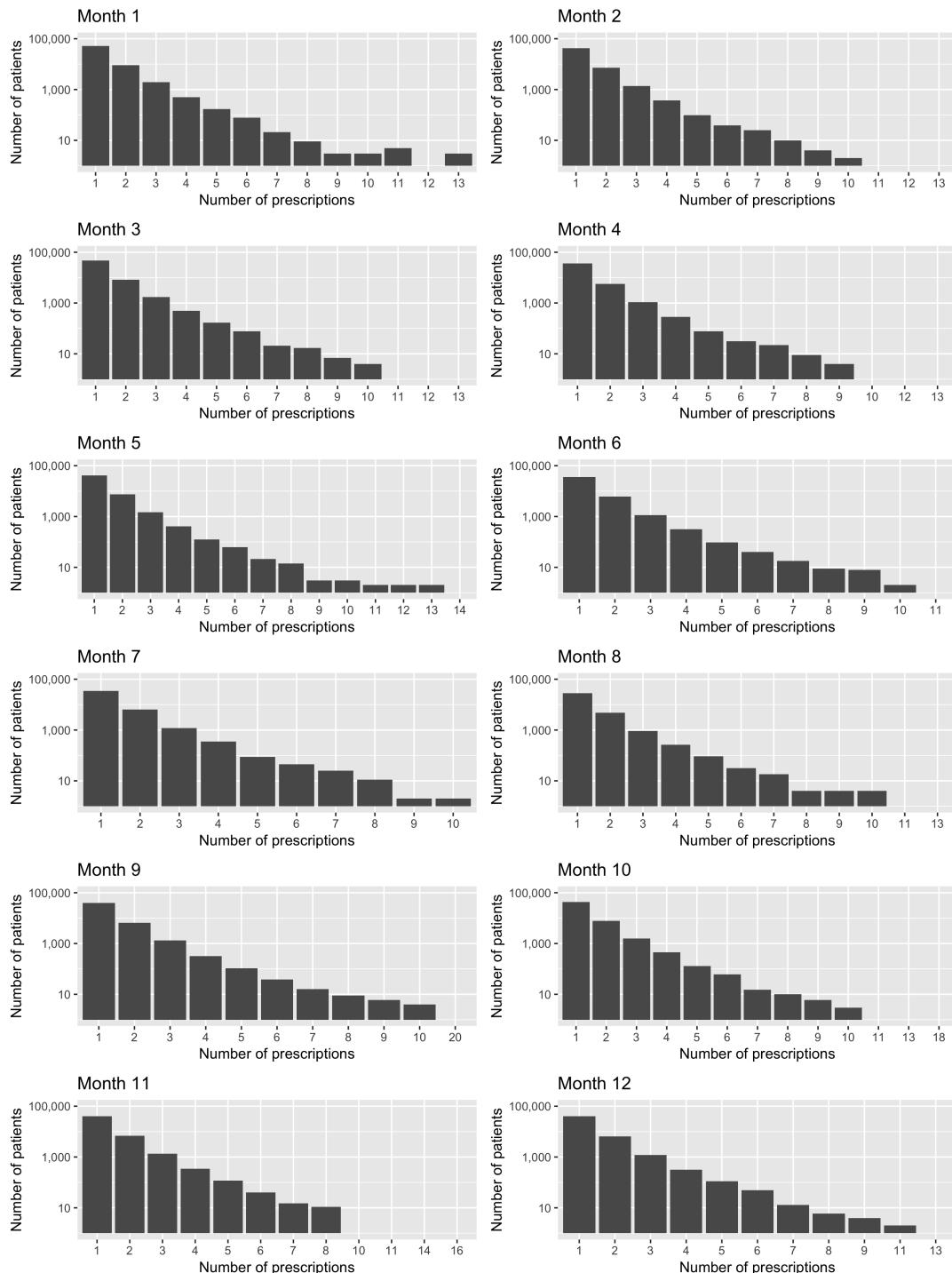


The 2012 fall might be caused by the substantial decrease of the approved antibiotics by AIFA,<sup>34</sup> but is noticeable that the number rises starting from 2013 due to antibiotic resistance.

### 7.3.3 Number of prescriptions trends

A barplot is the most informative way to visualise trending of number of prescriptions and highlight seasonality. For each month, the amount of patients having the number of prescriptions stated in the x-axis is plotted, using a logarithmic scale.

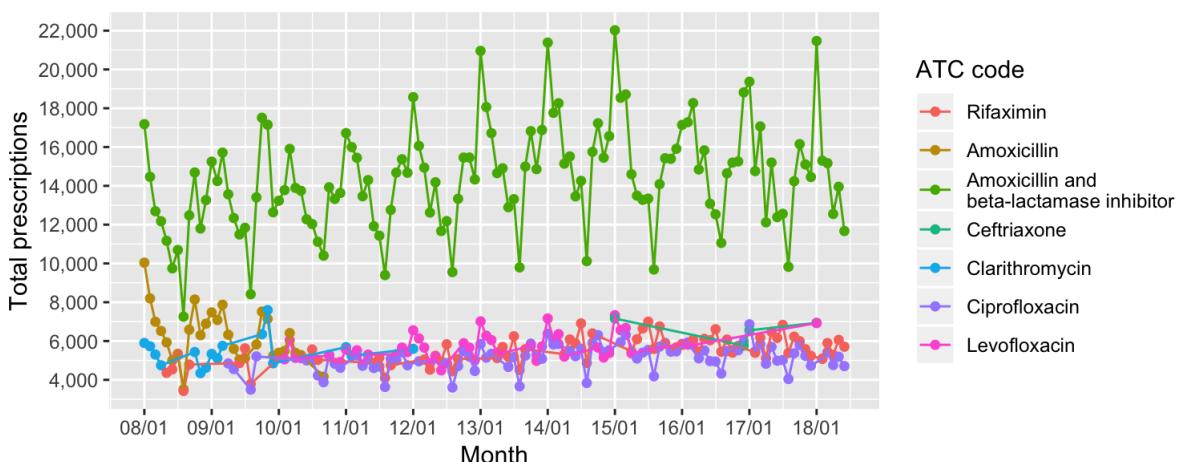
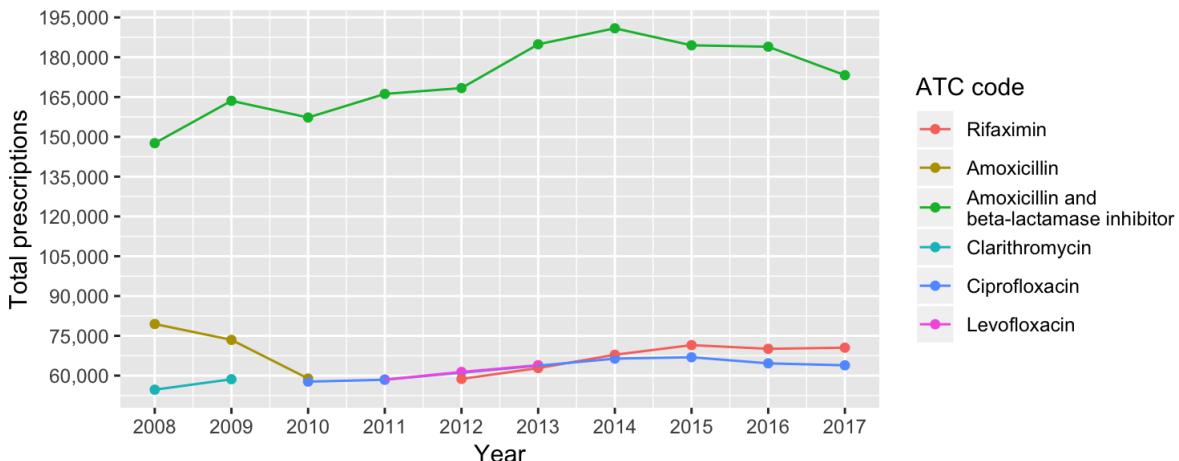
The scale has been adopted to standardise the quantity of patients having a smaller number of prescriptions (1-2), giving higher bars to larger values.



In winter there are more patients having a larger number of prescriptions, while during the summer higher values tend to disappear and even patients with one prescription decrease.

## 7.4 ATC rankings

For each year and each month, the 3 most common antibiotics are extracted, according to their ATC code. Top 3 changes during the years, hence why there are more than 3 labels in the plots.



Most ATCs follow a similar pattern: the amount decreases in 2010-2012 to then have a small peak in 2016 and get lower again in 2017, similarly to the global amount of prescriptions.

The most noticeable unusual trend corresponds to amoxicillin and beta-lactose inhibitors (ATC J01CR02), which considerably detaches from the others. Amoxicillin is one of the most popular group of drugs, confirmed by global analytics on the database and USA trends.<sup>35</sup>

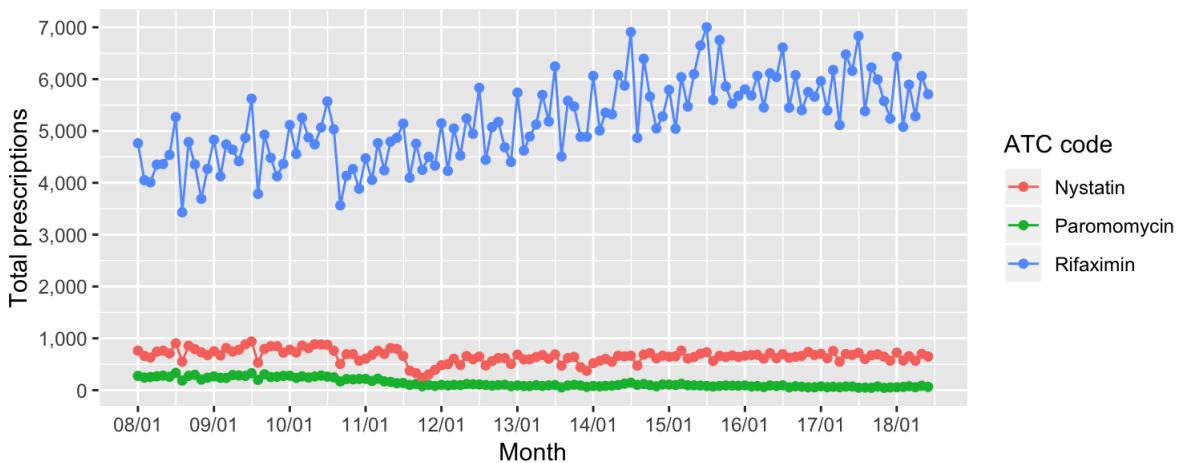
It is evident that the monthly prescriptions follow a seasonal trend: the number of prescription rises during the winter and falls during the summer.

### 7.4.1 Comparison with ICD-9

Since ICD-9 have some unusual trends, comparing them with antibiotics might provide additional information. The interested area comprehends ICD-9 class A, diseases of the digestive tract whose diagnoses have doubled in the span of 10 years. Antibiotics to treat this subset are:

- Nyastatin;
- Rifaximin;
- Paromomycin.

Cross-checking related codes with the antibiotics dataset, rifaximin is indeed one of the most prescribed drugs, while its numbers are still considerably inferior to amoxicillin (the most popular antibiotic). Seasonality is still present, yet less pronounced.



## 7.5 AIC rankings

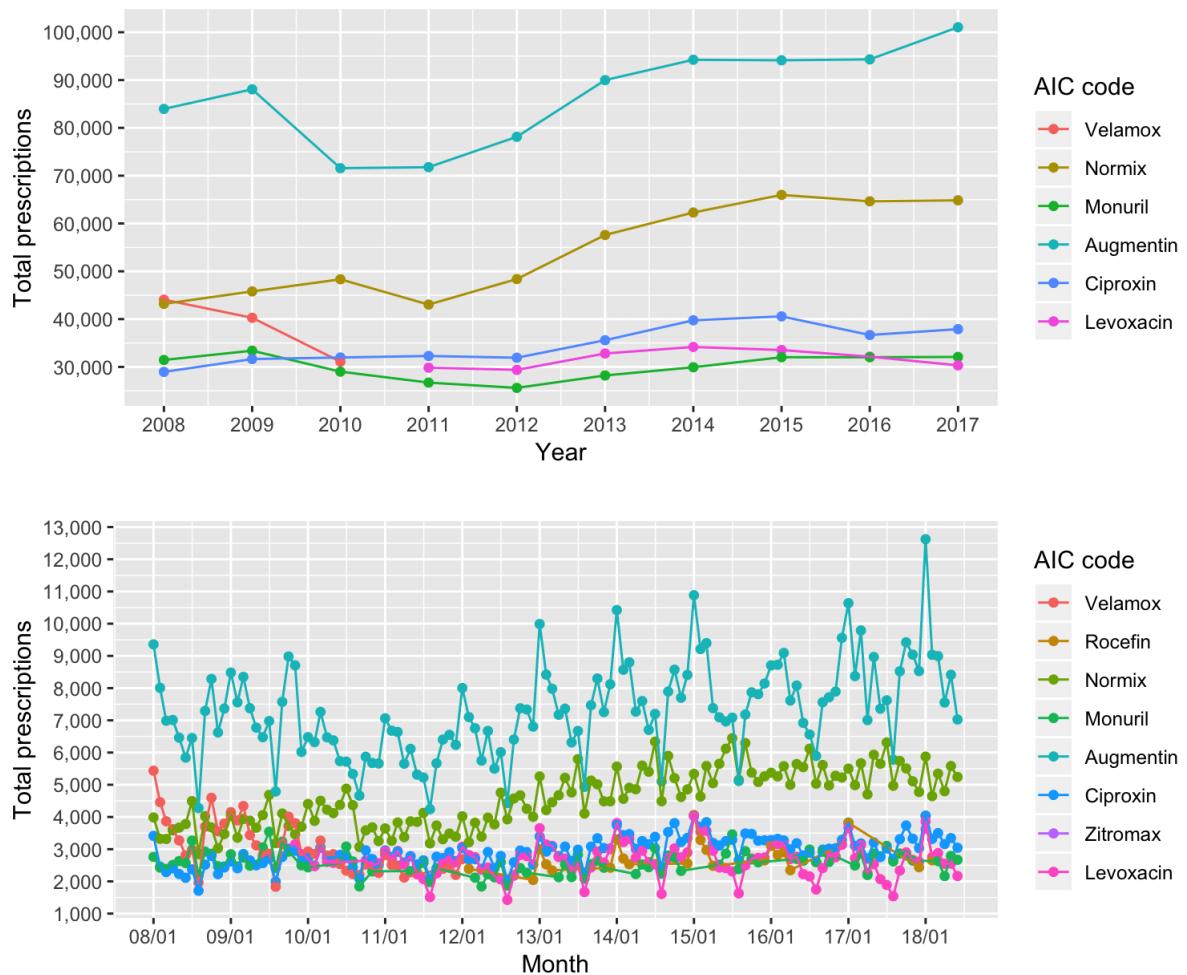
Prescriptions with AICs don't have a public lookup table, yet drugs can be singularly searched on a public database by AIFA. Analytics have been made on the antibiotics subset of data, using the 5 most prescribed medicines for months and years.

Augmentin is the top ranked prescription, detaching from the others and with a progressive increase starting from 2010: values go up to 100 000 prescriptions in 2017, and January 2017 is the month with the maximum number.

Having a difference of 30 000 between 2010 and 2017 without a correspondent increase of diagnoses shows frequent episodes of antibiotic resistance and overprescribing.

The drug Velamox has an opposite behaviour compared to the others: while most medicines are increasing their number of prescriptions, while Velamox is subject to progressive decrease until disappearance from the most prescribed drugs.

A better understanding between orders of magnitudes could be obtained counting the total prescriptions in 10 years of the most common AICs:



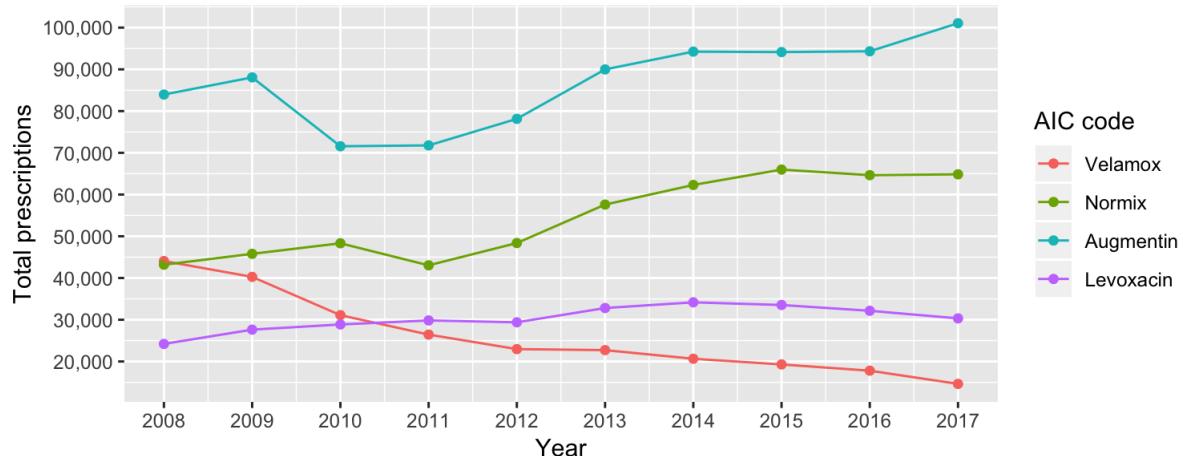
Code	Prescriptions	Antibiotic
026089019	934 942	Augmentin 875 mg + 125 mg 12 coated tablets
025300029	586 782	Normix 200 mg 12 coated tablets
026664021	373 720	Ciproxin 500 mg 6 coated tablets
033940038	264 604	Levoxacin 500 mg 5 coated tablets
025680024	229 816	Monuril 3 g 2 sachets
023097102	132 787	Velamox 1 g 12 dispersible tablets

### 7.5.1 An insight on Velamox

Velamox, according to AIFA, is an amoxicillin-based drug used for respiratory tract, ear and genital infections. It is sold in different packages, whose the most popular is the 1 g dispersible tablets with 12 tablets.

Since the general trending of most prescribed antibiotics only gave a partial vision of the comparisons between products, a complete chart is extracted using Velamox and three other popular drugs: Augmentin, Normix and Levoxacin. The difference is evident:

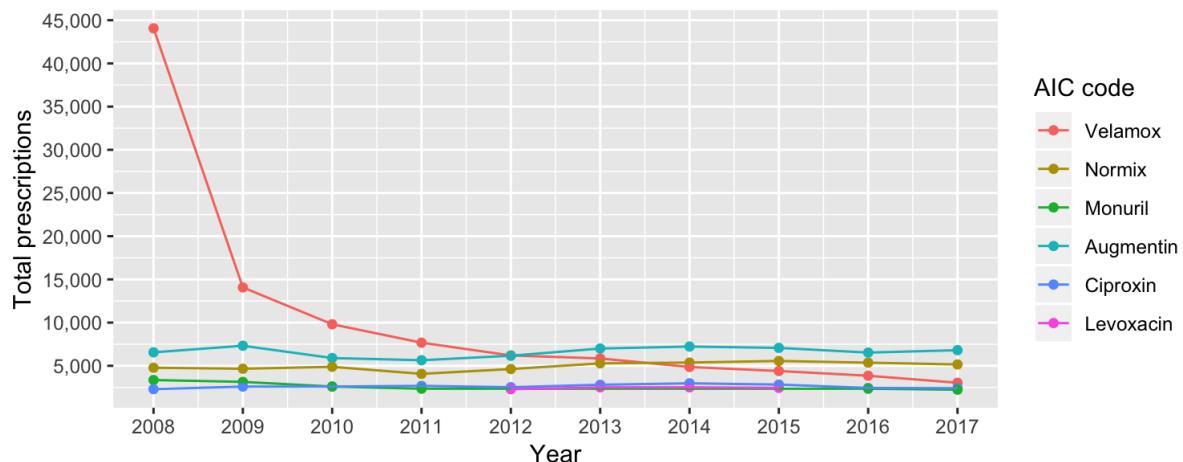
The first hypothesis of such a drastic drop of a product is its switch with another one, changing prescriptive habits of doctors. If that was the case, graphs related to the specific subset of patients would have a correspondent rise of values belonging to another



antibiotic.

Velamox in 2008, the year of maximum amount of prescriptions, has been given to approximately 32 000 patients. Extracting only the top 3 antibiotic prescriptions of those, in the last 10 years, along with Velamox, shows no substitution by another drug.

The fall starts in 2009, having a drop of 30 000 units, and gradually continues until 2017 with only 3 037 prescriptions:



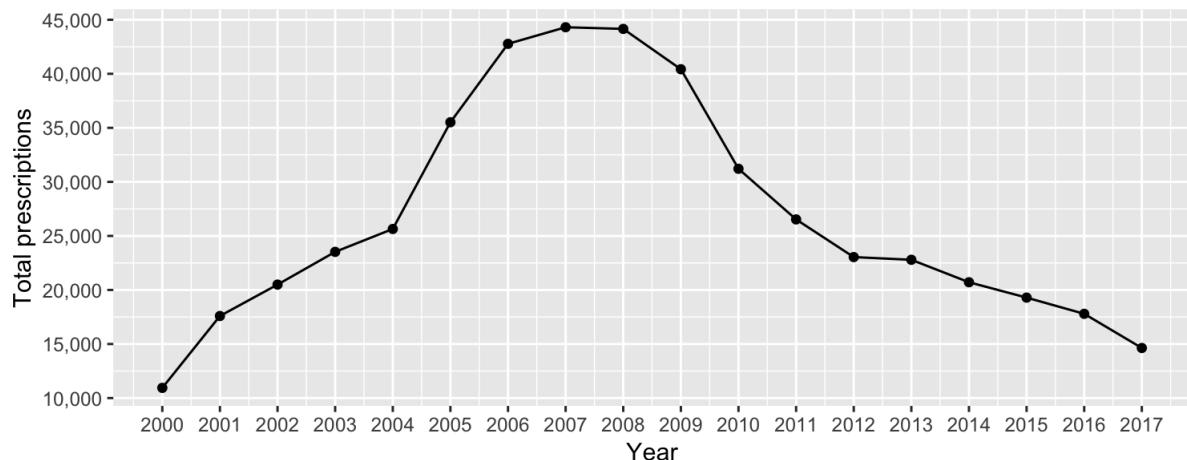
ATC trends related to Velamox's AIC show usual patterns: since the same ATC is linked with more products, the curve is still leaning upwards, yet there is an offset equal to the decrease in Velamox.

Since analysis on a subset of patients gives little insight, a wider time range can give additional information on a complete picture of global trends.

Further knowledge is obtained counting the total number of Velamox prescriptions among all patients from 2000 to 2017:

Velamox has acquired popularity from 2005 to 2010, when antibiotic resistance still hadn't been recognised as a problem and before the advent of generic drugs, yet this does not provide an explanation for the decrease.

Additional analysis is made on the subset of 32 000 patients:



- 1 500 are dead;
- 8 000 have changed GP;
- 18 000 are females;
- 14 000 are males;
- 16 000 are born before 1960.

General practitioners assigned to the subset of patients have generally decreased their prescriptions quota, roughly halving it comparing 2008 to 2017.

AIFA has also revoked authorisation for most instances of Velamox packages, only leaving three out of ten. The company which produces the drug, Mediolanum Farmacy, has been acquired by Neopharmmed Gentili in 2018.

Summarising, changes in Velamox prescriptions are partially caused by:

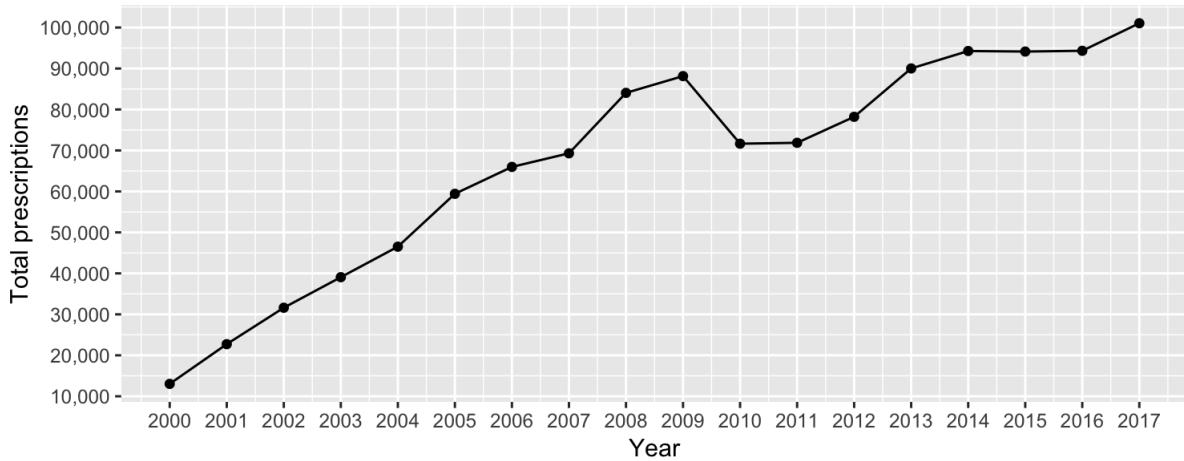
- Doctors gradually abandoning it after authorisation revoking;
- Patients dying or switching doctor;
- Substitutions with other drugs, none of whose is prescribed enough to have a significant trend, such as generics;
- Different advertisement campaigns after the company acquisition.

### 7.5.2 An insight on Augmentin

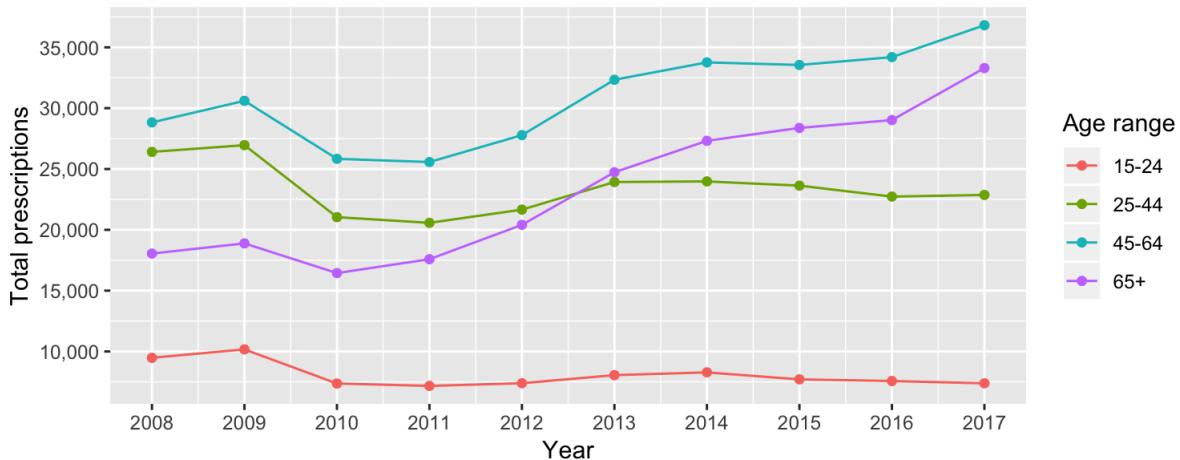
After analysing a drug with an unusual decrease, the focus can shift to the opposite side, considering the one with a noticeable increase.

Augmentin is the most prescribed antibiotic, with an almost constant positive trend:

2010 and 2011 have lower values, possibly explained by the Italian economic crisis or the diffusion of generic drugs. Starting from 2013 numbers continue to rise, reaching their highest peak in 2017.



Most Augmentin consumers are adults in the range of 45-64 years of age, following global trends previously obtained.



### 7.5.3 According to geographical area

There is information loss on geographical area of patients, therefore it is necessary to extract a further subset of antibiotic prescriptions, applying a level of cleansing.

Patients with complete and correct province of residence fields are selected, obtaining:

- 590 472 patients, 74,34%;
- 7 458 312 prescriptions, 88,94%.

Campania has five provinces: Naples, Avellino, Caserta, Benevento and Salerno. Analytics are expected to show an uniform distribution of prescriptions among those, yet outcomes are skewed towards Naples (significantly bigger) and other cities.

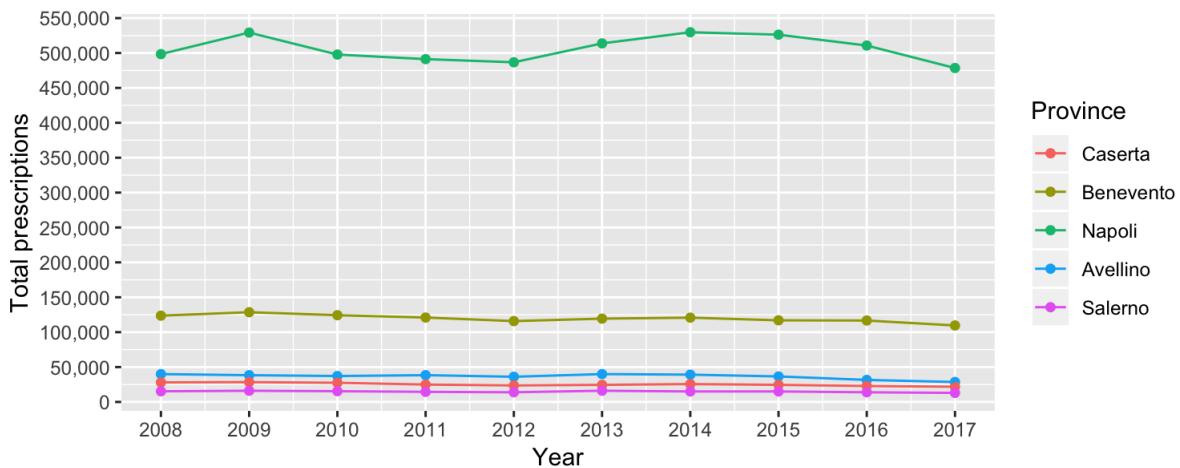
Top 5 municipalities for total prescriptions in 2008-2017:

Postcode	Province	Prescriptions
063049	Naples	2 854 056
063002	Afragola	618 519
063024	Castellammare di Stabia	497 309
063005	Arzano	314 640
063035	Gragnano	210.960

Prescriptions come from 932 different cities, while there are 550 distinct cities in Campania. This implies a slice of patients having residence in a different place outside the region.

Results might be altered by the non-uniform provenance of doctors and patients (mostly from the Naples area), and other big municipalities have large population although not qualifying as provinces.

Discrepancies among provinces are also shown by the following plot, identifying trends:



## Pharmacies

The number of pharmacies divided by province with related population (ISTAT 2018) confirms the gap between Naples and the others:

1. 857 pharmacies in the Naples area, 3 101 002 inhabitants;
2. 375 pharmacies in the Salerno area, 1 101 763 inhabitants;
3. 254 pharmacies in the Caserta area, 923 445 inhabitants;
4. 161 pharmacies in the Avellino area, 421 523 inhabitants;
5. 107 pharmacies in the Benevento area, 279 127 inhabitants.

## 7.6 ATC and AIC trends

### 7.6.1 By sex

### 7.6.2 By age

Analytics by age are deployed dividing the subset in the 5 official age ranges, however since patients in the youngest group should receive prescriptions by paediatricians and there is no related information in the database, the range 1-14 has to be considered incorrect and therefore removed.

This causes a difference of approximately 150 000 tuples, and the remaining ones are arranged by age of patient counting the amount of prescriptions in 2008-2017:

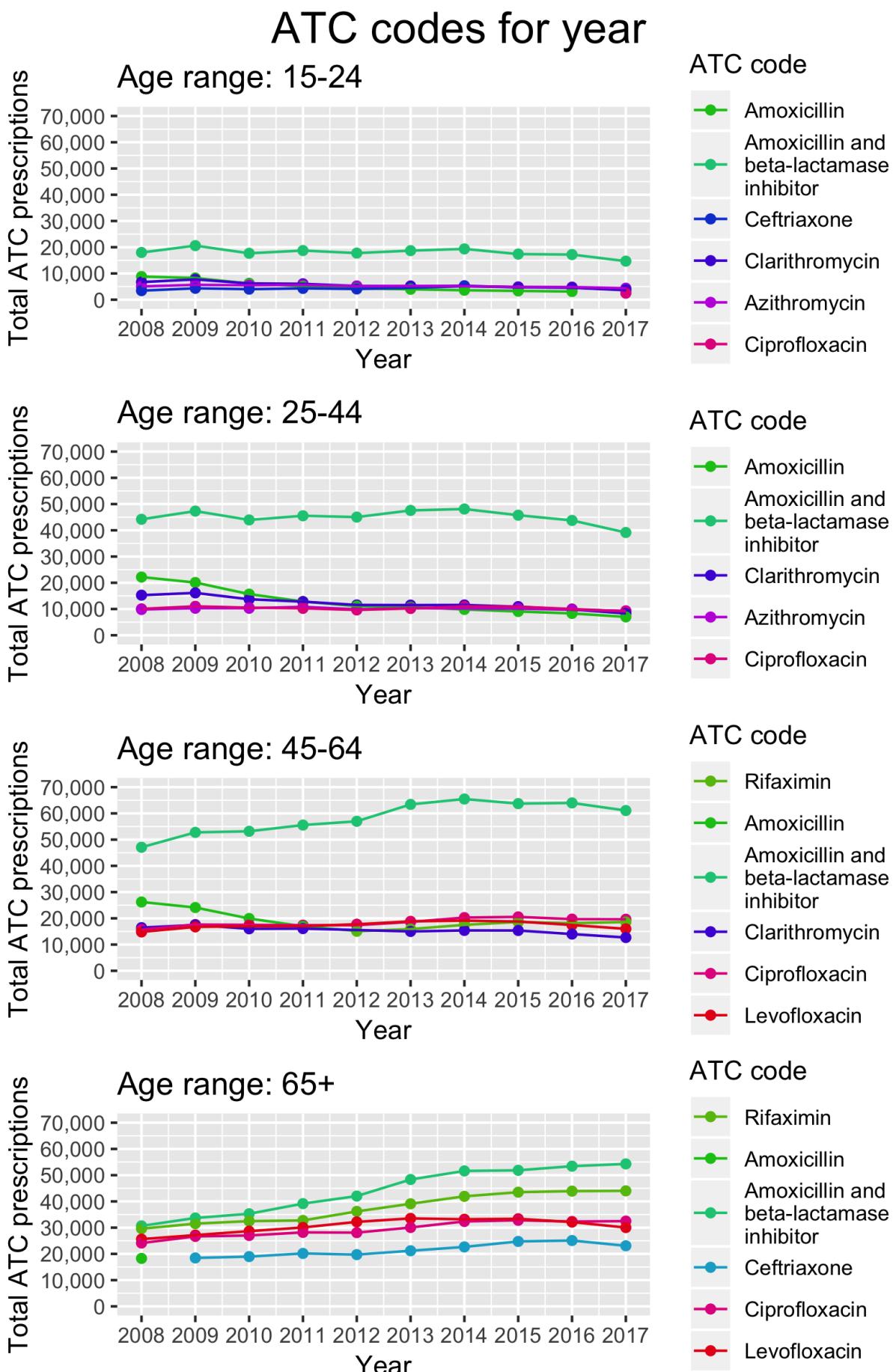
Age range	Prescriptions	Mean
15-24	643 540	5,25
25-44	1 776 228	7,29
45-64	2 674 494	10,81
65+	3 115 334	16,16

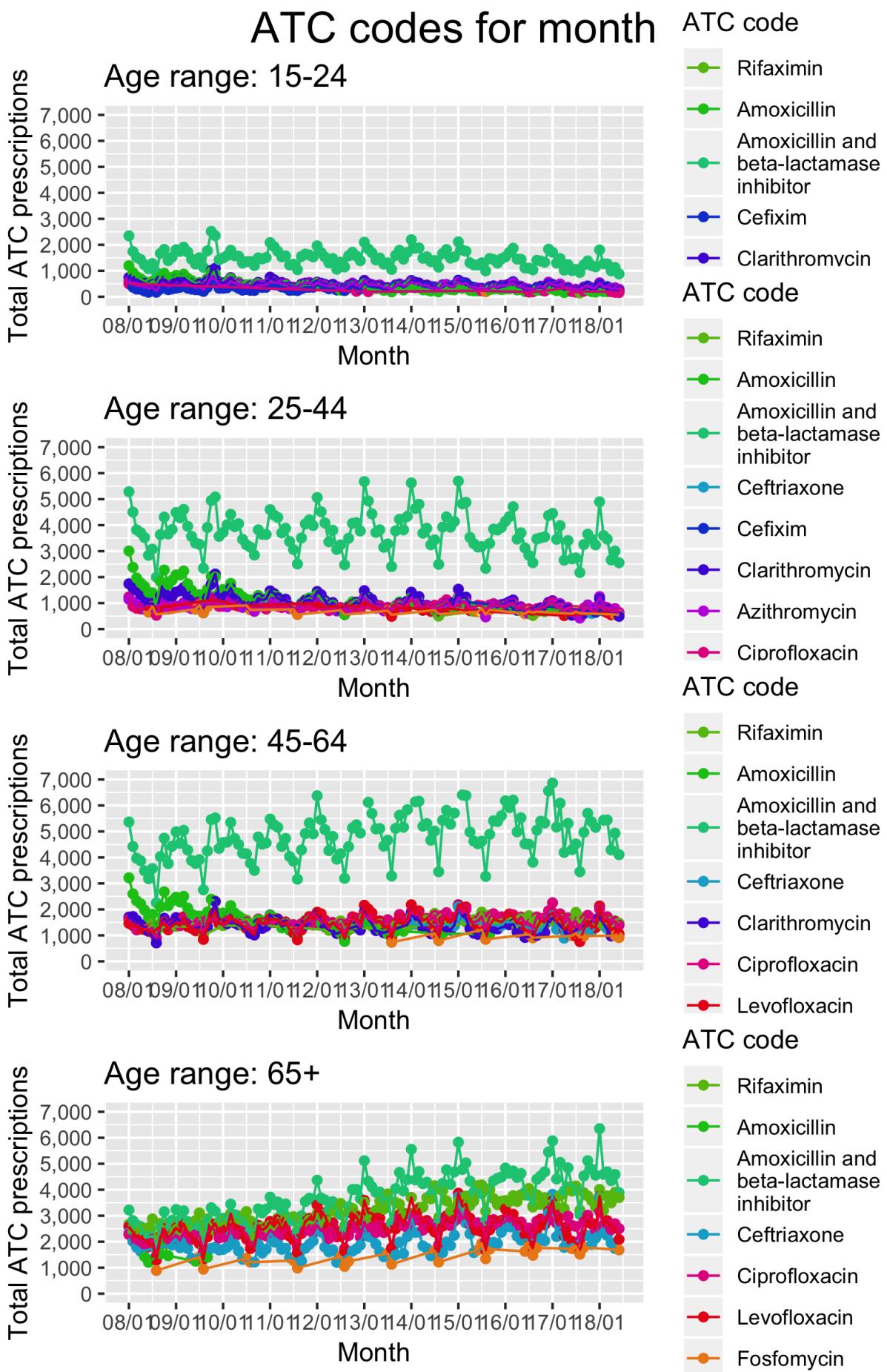
A patient on average receives 9,8 prescriptions in 10 years, but there are big differences between older and younger people: the number progressively increases.

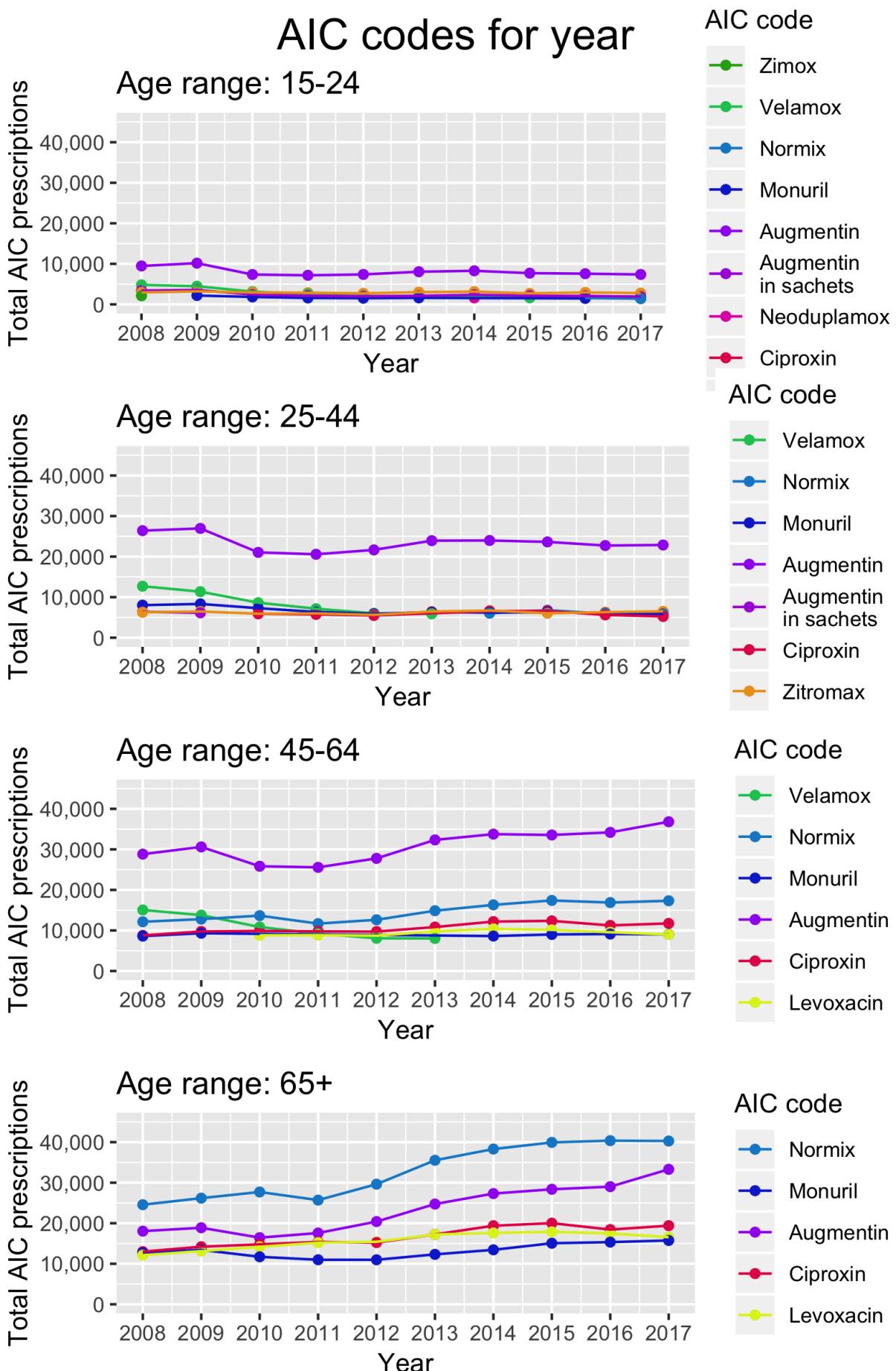
ATC and AIC popular codes are related, observing that some of the most common ones gets prescribed starting from an early age while others are prevalent between higher ranges.

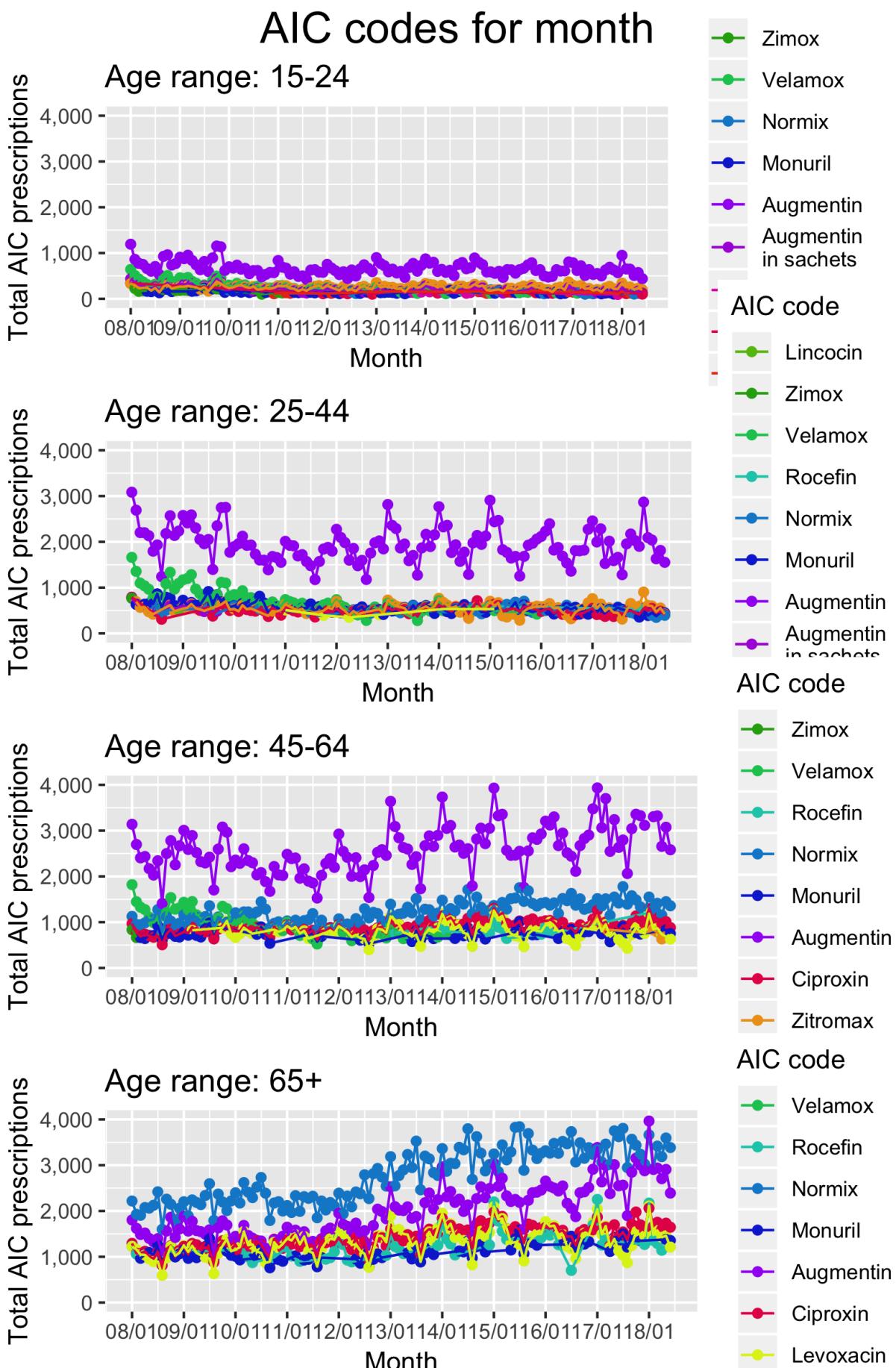
5 top AICs and ATCs for age, yearly and monthly:

### 7.6.3 ATC and AIC correlation









# Chapter 8

## k-means approach for cluster analysis

# Chapter 9

## Graph analysis

### 9.1 Prescription coupling

An excessive usage of antibiotics causes death of microorganisms in the human body which provide to maintaining immune cells and killing certain oral infections.<sup>36</sup>

To equilibrate the intestinal flora, lactic ferments are often taken together with antibiotics, so that new “good” bacteria can restore the probiotic action.

If this hypothesis is correct, the dataset will show antibiotic prescriptions paired with other drugs, on the same date — or it will highlight potential linkings between infections and other pathologies receiving a specific prescription.

### 9.2 Graph databases

Graph databases are data management systems allowing persistent representation of entity and relationship in a graph structure, implementing the Property Graph Model efficiently down to the storage level.

A graph  $G = \langle E, V \rangle$  is an abstract data type showing connections (edges  $E$ ) between pairs of vertices ( $V$ ). Nodes identify entities and their properties, while relationships are joining attributes between tables with eventual additional characteristics.

Queries allow to match pattern of nodes and relationships in a graph, providing ACID transaction compliance without specifying details on how to implement operations. Graph-crossing and related algorithms are highly efficient.

### 9.3 Goals

Goals of analytics through graphs is completion of antibiotic patterns changes and patient journey, providing a different point of view on those two important aspects altogether.

This part of the research aims to focus on:

- Co-prescriptions, understanding whether specified couples of drugs are often prescribed together;
- Clustering, to identify similar kinds of patients according to their prescription history;
- Centrality measures of nodes, to highlight particularly important entities in the graph.

## 9.4 Practical approach

### 9.4.1 Relational database structure

The available data comprehends patients, general practitioners, and their prescriptions in the time span from 2000 to 2018. Summarising the amount of records for each entity:

- 888 219 patients;
- 2 486 doctors;
- 118 716 403 prescriptions;
- 33 523 drugs.

Due to the amount and veracity of data, identifying a subset of records is useful to have detailed and targeted results, removing dispersive information and leaving a restricted pool of prescriptions, setting acceptability conditions.

Since analytics are aimed to identify antibiotic prescription patterns, similarly to previous work a new dataset has been extracted, imposing the following constraints:

1. AIC corresponding to an antibiotic;
2. Prescription date between 2008-01-01 and 2017-12-31;
3. Active general practitioners;
4. Patients with usable information about sex, date of birth and location.

This leads to obtaining a new relationship, composed by:

- 670 634 patients;
- 1 377 doctors;
- 8 386 057 prescriptions;
- 2 802 antibiotics.

To allow analytics on patient journey and co-prescriptions, it is necessary to access all the prescriptions assigned to all patients belonging in the subset. A major extraction is performed from the main table, comprehending:

1. Identifier of patients who received at least one other antibiotic prescription;
2. Prescription date between 2008-01-01 and 2017-12-31.

This reduces the number of other prescriptions, adding drugs not belonging to the antibiotic class. Duplicates, mistakes and empty fields are removed. The final composition of data is:

- 670 634 patients;
- 1 377 doctors;
- 8 328 272 prescriptions of antibiotics;
- 2 465 antibiotics;
- 7 587 009 other prescriptions;
- 21 248 other medicines.

#### 9.4.2 Migration of the database and graph modelling

The database has to be structured following the SQL to Cypher practices and guidelines, assigning nodes and relationships in an appropriate way considering the existing dataset and the related goals.

After having a final version of the data to import, the entity-relationship model translates with the following nodes and attributes:

- Patient;
  - ID, birthdate, sex;
- Doctor;
  - ID;
- Antibiotic;
  - AIC code, ATC code, active principle;
- Medicine (anything not Antibiotic);
  - AIC code, ATC code, active principle;
- Prescription;
  - patient, doctor, date, drug;
- OtherPrescription (not Antibiotic Prescription);
  - patient, doctor, date, drug.

All nodes are imported, and main indexes are created for optimisation of queries speed. Relationships are then created according to IDs and AIC codes:

- Prescription – TO → Patient;
- Prescription – FROM → Doctor;
- OtherPrescription – OF → Antibiotic;
- OtherPrescription – TO → Patient;

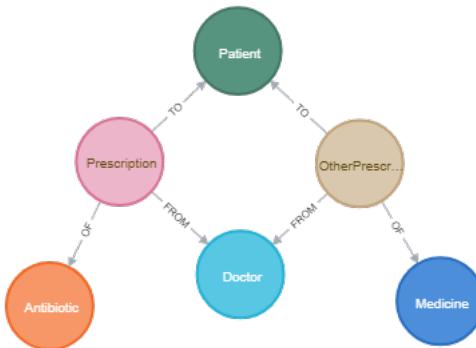
- OtherPrescription – FROM → Doctor;
- OtherPrescription – OF → Medicine.

After migrating all elements of the graph database, additional features of Patient representing age is added. Since age changes within years, and the considered timespan is 2008-2017, 10 node properties are created for each patient (age2008, age2009, ...).

## 9.5 Visualisation and analytics

### 9.5.1 Sample graph

A sample graph is obtained using the `apoc` functions.



### 9.5.2 Examples

An example of graph subset can be obtained extracting one of the individuals with the most antibiotic prescriptions (the 10th in descending order), a male patient born in 1943, and his associated drugs and doctors.

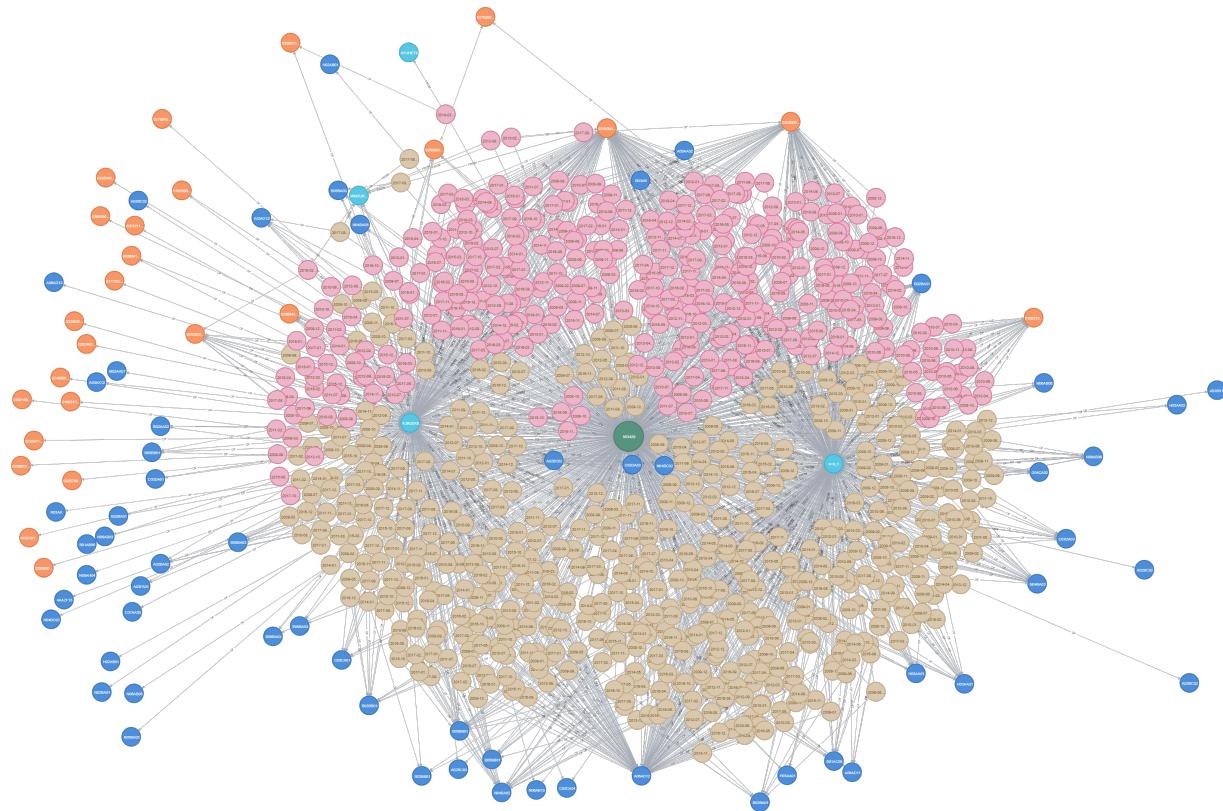
The graph displays:

1. 1 patient, the green node in the middle;
2. 360 antibiotic prescriptions, the pink nodes;
3. 524 other prescriptions, the brown nodes;
4. 4 doctors, the light blue nodes;
5. 25 antibiotics, the orange nodes;
6. 59 other medicines, the blue nodes.

Having a view focussed on prescriptions, the same procedure is applied extracting the 500th antibiotic in descending order according to number of prescriptions, corresponding to Locabiotol spray bottle 15 ml (50 mg / 5 ml).

The graph displays:

1. 1 antibiotic, the orange node in the middle;



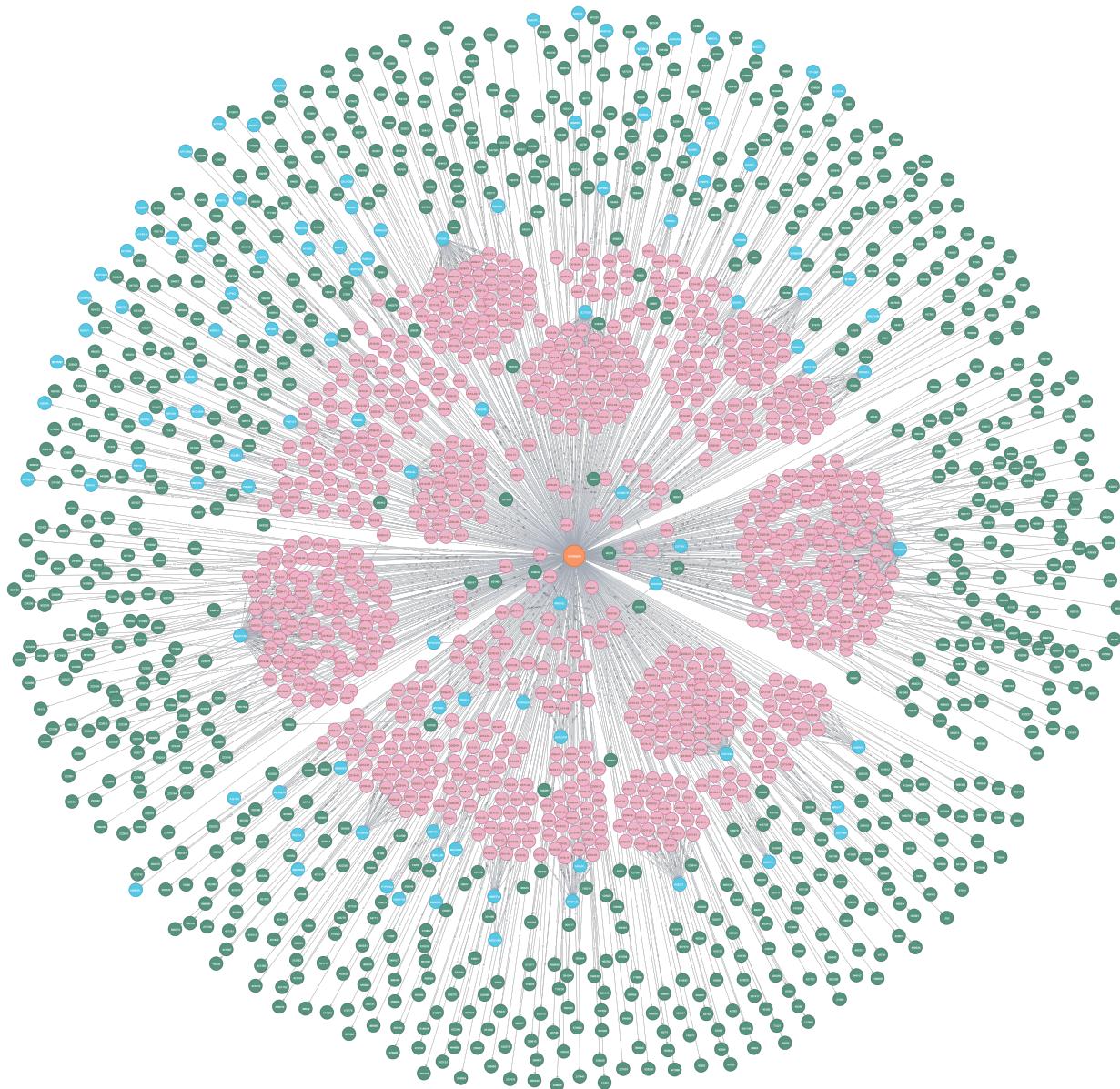
2. 896 antibiotic prescriptions, the pink nodes;
3. 107 doctors, the light blue nodes;
4. 817 patients, the green nodes.

Those two examples allow to have a general idea of how nodes interact with each other, grouping in clusters.

### 9.5.3 Graph statistics

A first set of global and local statistics are used to get the first insight on the graph and its components.

Number of nodes	16 611 005
Number of relationships	47 745 841
Average prescriptions per doctor	11 557,93
Standard deviation	15 457,6
Maximum per doctor	96 841
Minimum per doctor	1
Average prescriptions per patient	23,73
Standard deviation	40,46
Maximum per patient	1 567
Minimum per patient	1



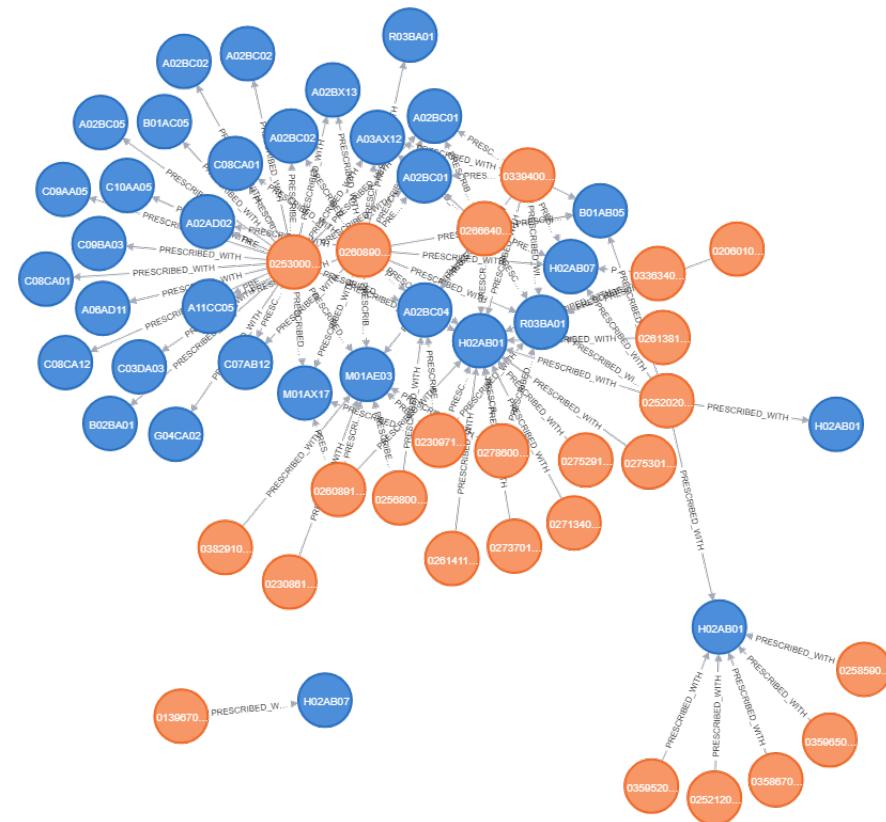
#### 9.5.4 Projecting a co-prescription graph

Since the restriction on generic prescriptions involves having the same date and patient of another antibiotic prescription, couples are analysed adding a relationship between Antibiotic and Medicines in the main graph.

After counting the number of repetitions for each couple Antibiotic-Medicine, the first 100 most popular ones are used to couple nodes, with the amount as property of the relationship PRESCRIBED\_WITH.

Visualising the newly created relationships, two connected components are highlighted (orange antibiotics, blue other medicines).

The component with only one couple corresponds to Plaquinil - Deltacortene, prescribed together approximately 5 000 times. Plaquinil can be used as antibiotic (antimalarial), but it is mostly given to treat arthritis, while Deltacortene is a corticosteroid for rheumatisms.



The range of values varies from 4 145 to 38 153 in the timespan of 10 years. The 5 most popular co-prescriptions are:

1. Augmentin - Oki;
2. Rocefin - Bentelan;
3. Augmentin - Bentelan;
4. Normix - Cardioaspirin;
5. Augmentin - Aulin.

Bentelan is a corticosteroid which cannot be used without an antibiotic in presence of systemic (concerning the whole organism) infections,<sup>39</sup> since it is an immunosuppressive drug, and this would explain the frequent co-prescriptions.

All the antibiotics are among the most prescribed ones, which justifies their presence in the co-prescriptions as well.

### 9.5.5 Projecting prescriptive habits

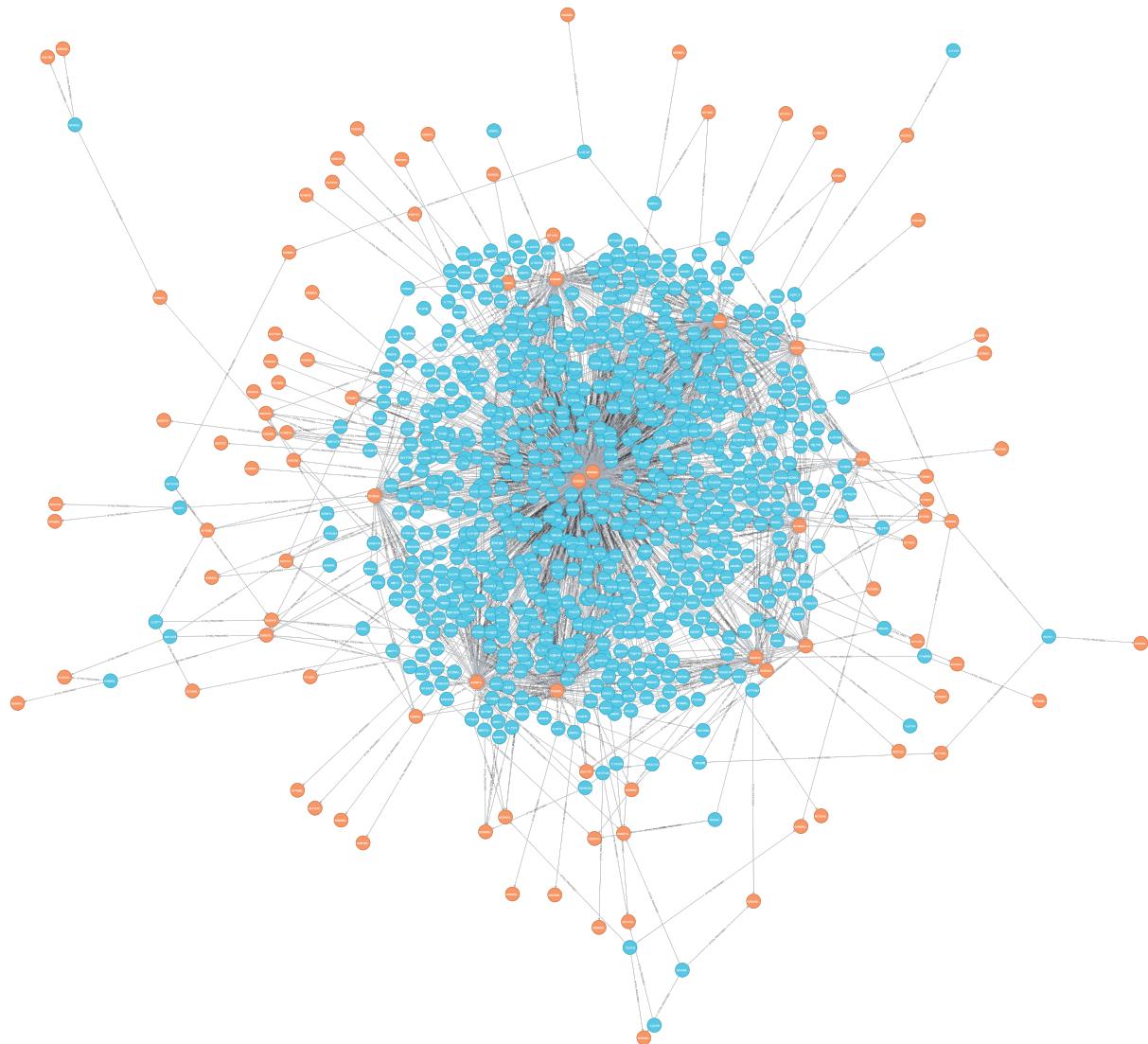
Having a detailed view of doctors' most common prescriptions gives another insight on how antibiotics are related.

The three most prescribed antibiotics for each doctor are taken, along with their count, and a new relationship OFTEN\_PRESCRIBED between Doctor and Antibiotic is created. Only amounts of prescriptions greater or equal to 50 are taken into account, for

consistency of results, resulting in 2 353 links.

The obtained graph displays:

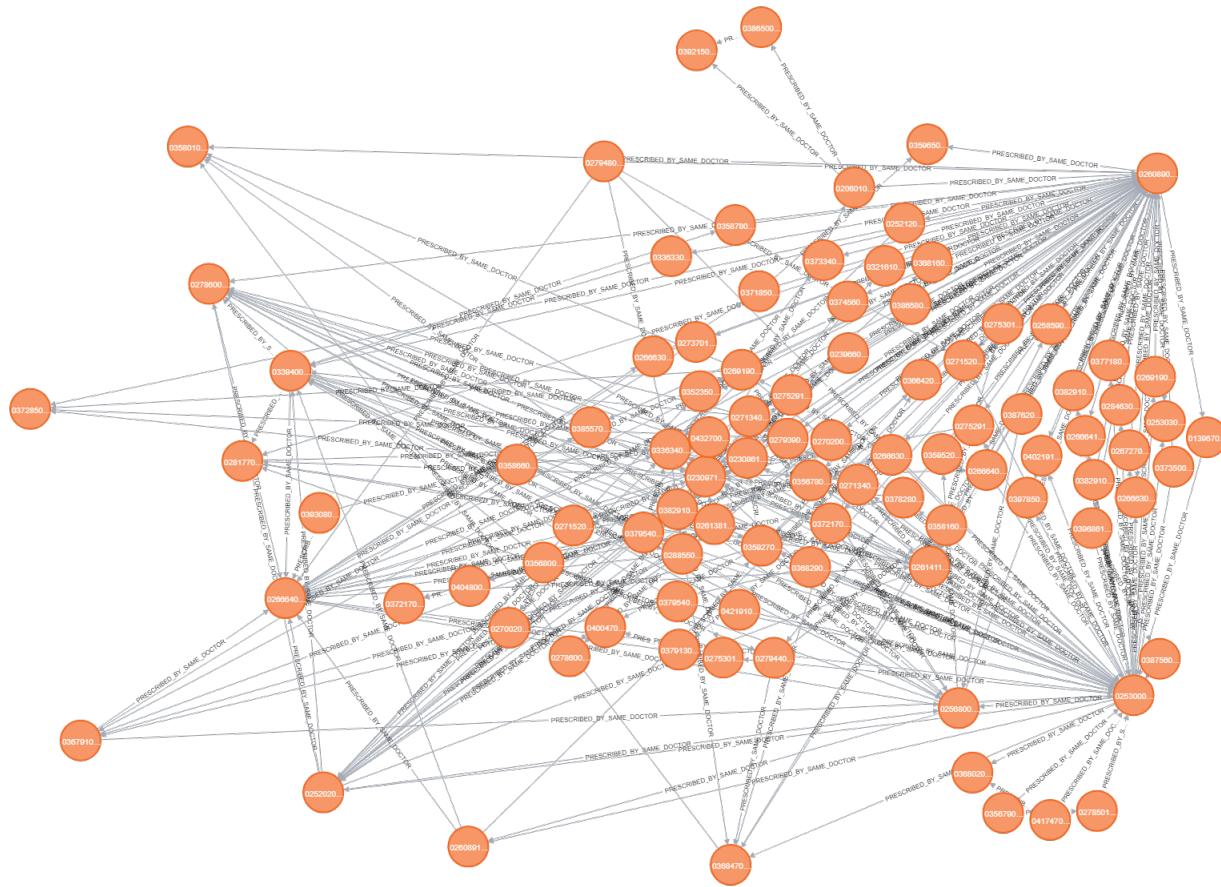
1. 809 doctors, the light blue nodes;
2. 95 antibiotics, the orange nodes.



Comparing the quantity of antibiotics and doctors and their relationships, it can be seen that most doctors prescribe a very restricted set of antibiotics.

To have a better understanding of the popularity of specific drugs, an auxiliary relationship PRESCRIBED\_BY\_SAME\_DOCTOR is created between Antibiotics.

Antibiotics are interlinked with 294 new relationships:



### 9.5.6 Centrality algorithms

#### Betweenness

The Betweenness Centrality algorithm calculates the shortest (weighted) path between every pair of nodes in a connected graph, using the breadth-first search algorithm.

Each node receives a score, based on the number of these shortest paths that pass through the node. Nodes that most frequently lie on these shortest paths will have a higher betweenness centrality score.

Betweenness among the top 5 antibiotics:

Antibiotic	Centrality
Augmentin (tablets)	2 562
Normix	1 852, 74
Velamox	683, 78
Ciproxin	562, 71
Augmentin (bottles)	489,93

#### Degree

Degree Centrality is the simplest of all the centrality algorithms. It measures the number of incoming and outgoing relationships from a node, analysing its influence.

Degree among the top 5 antibiotics, according to co-prescriptions:

Antibiotic	Degree
Augmentin (tablets)	32
Normix	22
Ciproxin	14
Augmentin (bottles)	12
Velamox	10

Calculating degree among doctors is useful to determine whether the top antibiotics prescribers are also the top prescribers, using the top 5 doctors.

Doctor	Degree (antibiotics)	Degree (medicine)	Total degree	Patients
1	45 625	51 221	96 846	4 466
2	42 554	41 594	84 148	2 022
3	35 645	22 383	58 028	1 816
4	34 587	40 315	74 902	2 028
5	34 380	28 763	63 143	1 874

Seeing the obtained results, the overprescribing of some general practitioners is clear: most of them makes a number of antibiotic prescriptions equal to others having double the amount of patients.

Doctors prescribing too much and not strictly when needed is one of the main causes of antibiotic resistance in Italy.

### 9.5.7 Community detection

Communities are identified according to patients getting the same antibiotic or other medicine. To ease computational time, a new relationship GETS is created, with an attribute representing the number of time a patient received that determined prescription.



Since data is in the magnitude order of millions, only the top 10 most frequent antibiotics and medicines are selected, creating roughly 2 millions of relationships between 580 282 patients and 20 total drugs.

This means 85,5% of patients received a prescription among a subset representing 0,08% of antibiotics and other medicines.

### 9.5.8 Similarity detection

## 9.6 Considerations

# Chapter 10

## Assessments of results and future directions

# Bibliography

- [1] <https://healthitanalytics.com/news/understanding-the-many-vs-of-healthcare-big-data-analytics>
- [2] <https://www.millewin.it/>
- [3] [https://www.normattiva.it/uri-res/N2Ls?urn:nir:stato:legge:1995-12-29;549\\_art3!vig=1](https://www.normattiva.it/uri-res/N2Ls?urn:nir:stato:legge:1995-12-29;549_art3!vig=1)
- [4] <http://www.agenziafarmaco.gov.it/content/la-resistenza-agli-antibiotici-emergenza-mondiale-il-primo-punto-di-contatto>
- [5] Davide Castaldi, *Richiesta Dati CNCM per Analisi Appropriatezza Prescrittiva*, Consorzio Milano Ricerche, 2018.
- [6] Made with draw.io.
- [7] Manuale ICD-9-CM versione italiana 2007.  
[http://www.salute.gov.it/portale/documentazione/p6\\_2\\_2\\_1.jsp?lingua=italiano&id=2251](http://www.salute.gov.it/portale/documentazione/p6_2_2_1.jsp?lingua=italiano&id=2251)
- [8] Davide Castaldi, *Allegato Tech DB Campania*, Consorzio Milano Ricerche, 2018.
- [9] <https://bioportal.bioontology.org/ontologies/ATC>
- [10] [http://www.agenziafarmaco.gov.it/sites/default/files/medicinali-equivalenti-qualita\\_sicurezza\\_efficacia.pdf](http://www.agenziafarmaco.gov.it/sites/default/files/medicinali-equivalenti-qualita_sicurezza_efficacia.pdf)
- [11] <https://www.woncaeurope.org/sites/default/files/documents/Definizione%20WONCA%202011%20ita-A4.pdf>
- [12] <https://www.postgresql.org/>
- [13] <https://neo4j.com/>
- [14] <https://www.r-project.org/>
- [15] <https://www.who.int/en/news-room/fact-sheets/detail/antimicrobial-resistance>
- [16] <https://www.cdc.gov/drugresistance/about.html>
- [17] <https://www.folkhalsomyndigheten.se/contentassets/dae82c7afd424a57b57ec81818793346/swedish-workplace-drug-resistance-report-2018.pdf>
- [18] [bmj.com/cgi/pmidlookup?view=long&pmed=9270458](https://bmj.com/cgi/pmidlookup?view=long&pmed=9270458)
- [19] <https://en.oxforddictionaries.com/definition/subsidiarity>
- [20] <http://www.salute.gov.it/portale/esenzioni/detttaglioContenutiEsenzioni.jsp?lingua=italiano&id=460>
- [21] <http://www.fcr.re.it/classificazione-dei-farmaci-ai-fini-della-imborsabilita>
- [22] <https://web.archive.org/web/20111129162006/http://www.farmaciadicello.it/ricetta-01.htm>
- [23] <https://web.archive.org/web/20140611065109/http://www.woncaeurope.org/sites/default/files/documents/Definizione%20WONCA%202011%20ita-A4.pdf>
- [24] [http://www.salute.gov.it/portale/temi/p2\\_6.jsp?lingua=italiano&id=1698&area=tumori&menu=percorso](http://www.salute.gov.it/portale/temi/p2_6.jsp?lingua=italiano&id=1698&area=tumori&menu=percorso)
- [25] <https://ascpt.onlinelibrary.wiley.com/doi/full/10.1038/clpt.2008.24>

- [26] <https://www.who.int/classifications/icd/en/>
- [27] [http://www.salute.gov.it/portale/temi/p2\\_6.jsp?lingua=italiano&id=1982&area=statisticheSSN&menu=1](http://www.salute.gov.it/portale/temi/p2_6.jsp?lingua=italiano&id=1982&area=statisticheSSN&menu=1)
- [28] <https://www.medicalbillingandcodingonline.com/icd-cm-codes/>
- [29] <http://www.agenziafarmaco.gov.it/glossary/term/1432>
- [30] <http://www.agenziafarmaco.gov.it/content/1%E2%80%99autorizzazione-all%E2%80%99immissione-commerciale>
- [31] [https://www.whocc.no/filearchive/publications/2019\\_guidelines\\_web.pdf](https://www.whocc.no/filearchive/publications/2019_guidelines_web.pdf)
- [32] [https://www.repubblica.it/salute/medicina-e-ricerca/2019/03/13/news/antibioticoresistenza-in-italia\\_il\\_primo\\_europeo\\_di\\_decessi-221467306/](https://www.repubblica.it/salute/medicina-e-ricerca/2019/03/13/news/antibioticoresistenza-in-italia_il_primo_europeo_di_decessi-221467306/)
- [33] <https://www.medicalnewstoday.com/articles/10278.php>
- [34] <https://www.aboutpharma.com/blog/2019/01/10/antibiotici-continua-il-calo-della-ricerca-e-sviluppo>
- [35] <https://clincalc.com/DrugStats/Top300Drugs.aspx>
- [36] <https://www.infectioncontroltoday.com/antibiotics-antimicrobials/study-shows-antibiotics-destroy>
- [37] <https://www.dedalus.eu/>
- [38] <https://www.datapine.com/blog/big-data-examples-in-healthcare/>
- [39] <https://www.my-personaltrainer.it/Foglietti-illustrativi/Bentelan.html>
- [40] <http://www.agenziafarmaco.gov.it/content/la-resistenza-agli-antibiotici-emergenza-mondiale-il-primo-europeo-di-decessi-221467306/>