

Design of a Data Warehouse

Perna Mirco
Pizzuto Salvatore Leonardo

A.A. 2024/2025

Index

1	Introduction	3
2	Creation and population of the schema	3
2.1	Analysis of the provided data	3
2.2	Queries	7
3	Cleaning Data and ETL Process	7
3.1	Patients	7
3.2	Organizations	9
4	Data Warehouse design	18
4.1	Design choices	18
4.2	Data mart	19
4.3	Administer	19
4.4	Affected_By	22
4.5	Carry_Out	26
4.6	Conduct	28
4.7	Plan	31
4.8	Require	34
4.9	Take_Place_In	37
4.10	Glossary of Measures	40
5	Use of the Data Warehouse	41
5.1	Administer	41
5.2	Affected_By	41
5.3	Plan	44

1 Introduction

The project involves designing a Data Warehouse for the analysis of healthcare-related information in the state of Massachusetts, United States of America, to meet the needs of decision-making processes. Each chapter describes the steps taken to design and develop the required Data Warehouse.

2 Creation and population of the schema

2.1 Analysis of the provided data

The data were provided in CSV format, which consists of simple text files containing lists of values. To enable a more accurate analysis, a preliminary database was created to reflect the original structure of the data. DBeaver, an SQL client, was used to facilitate the import of the files into the MySQL database. DBeaver offers a dedicated tool for importing data from CSV files. It is represented in Figure 1.

<div><div>conditions</div><div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: description</div></div></div>	<div><div>immunizations</div><div><div>date</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: base_cost</div></div></div>	<div><div>supplies</div><div><div>date</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: quantity</div></div></div>	<div><div>devices</div><div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: udi</div></div></div>	<div><div>payer_transitions</div><div><div>A: patient</div><div>A: memberid</div><div>start_year</div><div>end_year</div><div>A: payer</div><div>A: secondary_payer</div><div>A: ownership</div><div>A: ownername</div></div></div>
<div><div>claims</div><div><div>A: id</div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: reasoncode</div><div>A: reasondescription</div></div></div>	<div><div>observations</div><div><div>date</div><div>A: patient</div><div>A: encounter</div><div>A: category</div><div>A: code</div><div>A: description</div><div>A: value</div><div>A: units</div><div>A: type</div></div></div>	<div><div>procedures</div><div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: base_cost</div><div>A: reasoncode</div><div>A: reasondescription</div></div></div>	<div><div>organizations</div><div><div>A: id</div><div>A: name</div><div>A: address</div><div>A: city</div><div>A: state</div><div>A: zip</div><div>A: lon</div><div>A: lat</div><div>A: phone</div><div>A: revenue</div><div>A: utilization</div></div></div>	<div><div>providers</div><div><div>A: id</div><div>A: organization</div><div>A: name</div><div>A: gender</div><div>A: specialty</div><div>A: address</div><div>A: city</div><div>A: state</div><div>A: zip</div><div>A: lon</div><div>A: utilization</div></div></div>
<div><div>imaging_studies</div><div><div>A: id</div><div>date</div><div>A: patient</div><div>A: encounter</div><div>A: series_id</div><div>A: bodysite_code</div><div>A: bodysite_description</div><div>A: modality_code</div><div>A: modality_description</div><div>A: instance_id</div><div>A: sop_code</div><div>A: sop_description</div><div>A: procedure_code</div></div></div>	<div><div>medications</div><div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: base_cost</div><div>A: payer_coverage</div><div>A: dispenses</div><div>A: totalcost</div><div>A: reasoncode</div><div>A: reasondescription</div></div></div>	<div><div>allergies</div><div><div>start</div><div>stop</div><div>A: patient</div><div>A: encounter</div><div>A: code</div><div>A: description</div><div>A: system</div><div>A: type</div><div>A: category</div><div>A: reaction1</div><div>A: description1</div><div>A: severity1</div><div>A: reaction2</div><div>A: description2</div><div>A: severity2</div></div></div>	<div><div>encounters</div><div><div>A: id</div><div>start</div><div>stop</div><div>A: patient</div><div>A: organization</div><div>A: provider</div><div>A: payer</div><div>A: encounterclass</div><div>A: code</div><div>A: description</div><div>A: base_encounter_cost</div><div>A: total_claim_cost</div><div>A: payer_coverage</div><div>A: reasoncode</div><div>A: reasondescription</div></div></div>	<div><div>payers</div><div><div>A: id</div><div>A: name</div><div>A: address</div><div>A: city</div><div>A: state_headquartered</div><div>A: zip</div><div>A: phone</div><div>A: amount_covered</div><div>A: amount_uncovered</div><div>A: revenue</div><div>A: covered_encounters</div><div>A: uncovered_encounters</div><div>A: covered_medications</div><div>A: uncovered_medications</div><div>A: covered_procedures</div><div>A: uncovered_procedures</div><div>A: covered_immunizations</div><div>A: uncovered_immunizations</div><div>A: unique_customers</div><div>A: qqls_avg</div><div>A: member_months</div></div></div>
<div><div>patients</div><div><div>A: id</div><div>birthdate</div><div>deathdate</div><div>A: son</div><div>A: drivers</div><div>A: passport</div><div>A: prefix</div><div>A: first</div><div>A: last</div><div>A: suffix</div><div>A: maiden</div><div>A: marital</div><div>A: race</div><div>A: ethnicity</div><div>A: gender</div><div>A: birthplace</div><div>A: address</div><div>A: city</div><div>A: state</div><div>A: county</div><div>A: zip</div><div>A: lat</div><div>A: lon</div><div>A: healthcare_expenses</div><div>A: healthcare_coverage</div></div></div>	<div><div>claims</div><div><div>A: id</div><div>A: patientid</div><div>A: providerid</div><div>A: primarypatientinsuranceid</div><div>A: secondarypatientinsuranceid</div><div>A: departmentid</div><div>A: patientdepartmentid</div><div>A: diagnosis1</div><div>A: diagnosis2</div><div>A: diagnosis3</div><div>A: diagnosis4</div><div>A: diagnosis5</div><div>A: diagnosis6</div><div>A: diagnosis7</div><div>A: diagnosis8</div><div>A: referringproviderid</div><div>A: appointmentid</div><div>A: currentthessdate</div><div>A: servicedate</div><div>A: supervisingproviderid</div><div>A: status1</div><div>A: status2</div><div>A: statusup</div><div>A: outstanding1</div><div>A: outstanding2</div><div>A: outstandingp</div><div>A: lastbilleddate1</div><div>A: lastbilleddate2</div><div>A: lastbilleddatep</div><div>A: healthcareclaimtypeid1</div><div>A: healthcareclaimtypeid2</div></div></div>	<div><div>claimsTransaction</div><div><div>A: id</div><div>A: claimid</div><div>A: charged</div><div>A: patientid</div><div>A: type</div><div>A: amount</div><div>A: method</div><div>A: fromdate</div><div>A: todate</div><div>A: placeofservice</div><div>A: procedurecode</div><div>A: modifier1</div><div>A: modifier2</div><div>A: diagnosisref1</div><div>A: diagnosisref2</div><div>A: diagnosisref3</div><div>A: diagnosisref4</div><div>A: units</div><div>A: departmentid</div><div>A: notes</div><div>A: unitamount</div><div>A: transferorid</div><div>A: transfereeid</div><div>A: payments</div><div>A: adjustments</div><div>A: transfers</div><div>A: outstanding</div><div>A: appointmentid</div><div>A: linenote</div><div>A: patientinsuranceid</div><div>A: feescheduleid</div><div>A: providerid</div><div>A: supervisingproviderid</div></div></div>		

Figure 1: CSV files represented as tables.

However, the dates contained in the various CSV files do not match the format required by the DBeaver tool. To solve this issue, we decided to implement a Python program that converts the date format to the one required by the tool. The following is an example of the implemented code.

```

1 from datetime import datetime
2 import csv
3
4 with open('procedures.csv', newline='') as infile, open('dati_csv_cleaned/
   procedures_cleaned.csv', 'w', newline='') as outfile:
5     reader = csv.DictReader(infile)
6     fieldnames = reader.fieldnames
7     writer = csv.DictWriter(outfile, fieldnames=fieldnames)
8     writer.writeheader()
9
10    for row in reader:
11        if row['START']:
12            dt = datetime.fromisoformat(row['START'].replace('Z', '+00:00'))
13            row['START'] = dt.strftime('%Y-%m-%d %H:%M:%S')
14        if row['STOP']:
15            dt = datetime.fromisoformat(row['STOP'].replace('Z', '+00:00'))
16            row['STOP'] = dt.strftime('%Y-%m-%d %H:%M:%S')
17        writer.writerow(row)

```

Once the data was loaded into the database, a series of queries was executed to gain a deeper understanding of the information contained in the files, to identify primary keys and referential constraints among the data and to evaluate the data in terms of accuracy, completeness and consistency. The analysis showed that the majority of entities are related to the entity 'Encounter'. This will be more clearly observable following the normalization process, that contains all the information about a patient's visit to a specific hospital, including the reason for the visit and the doctor in charge. After completion of this initial analysis, it was decided to normalize the database to ensure data integrity, reduce redundancy, and improve the efficiency of the analytical operations required. The results obtained are shown in Figure 2.



To conclude this first step, the E-R diagram was created, which plays a significant role in the following steps, especially in the definition of facts. It is shown in Figure 3. Note that the attributes are not shown for readability reasons, but they can be read from the normalized database shown in Figure 2.

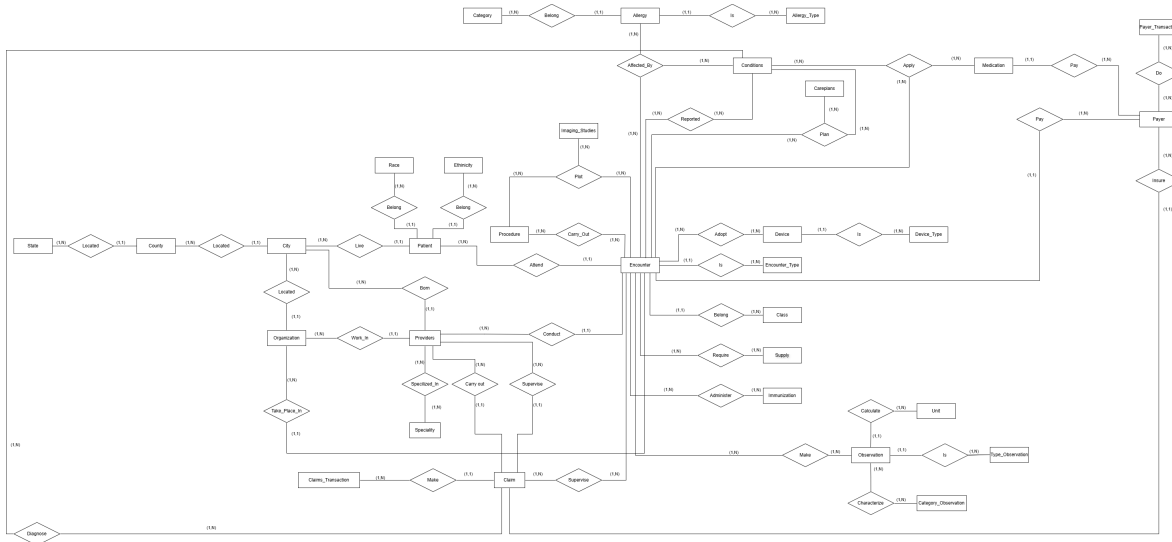


Figure 3: E/R diagram

2.2 Queries

1. What is the number of visits made to each hospital facility in the last five years, broken down by gender and age group?
2. What are the most used supplies each year, broken down by county and hospital facility?
3. What is the number of respiratory therapies recorded in various hospitals, broken down by gender and age group?
4. What is the average age of patients affected by acute bronchitis or prediabetes in each county recorded in the last five years?
5. What is the average number of treatment days required to complete a therapy plan for various diseases in each hospital over the years?
6. What is the cost of medical procedures broken down by year for the various counties and hospital facilities?
7. What are the most common allergies, grouped by patient age group, and county where they were recorded?
8. What is the incidence of the various side effects associated with each allergy, broken down by patient age group?
9. How many individuals, in the age group between 10 and 50 years, were vaccinated against HPV between the years 2015 and 2019, in the counties of Bristol, Essex, and Middlesex, grouped by sex, county, and year of administration?
10. What is the number of visits performed by each doctor, broken down by patient age group, gender, race, and, in addition, county and year of the conducted visit?

3 Cleaning Data and ETL Process

In order to use the previously discussed normalized database, an ETL service has been implemented, i.e. a process of data extraction, transformation, and loading. During this process, data cleaning was also performed, which involved handling null values and inconsistent data.

Below are some stages of the ETL process, implemented using the MySQL dialect.

3.1 Patients

During the processing of the 'Patient' entity, an issue was encountered in handling null values in the fields 'drivers', 'passport', 'prefix', 'suffix', 'maiden' and 'marital'. To address this issue, it was decided to replace the null values with standard default values, as shown in the table below.

Value	Fields
-	<ul style="list-style-type: none">• marital
NO	<ul style="list-style-type: none">• prefix

UNKNOWN	<ul style="list-style-type: none"> • drivers • maiden • passport • suffix
---------	---

Table 1: Values used to replace null values in the fields.

From the 'Patients' table of the non-normalized database, we also extracted information related to the 'States', 'Counties', 'Cities', 'Races', and 'Ethnicities' tables. To avoid repeated queries to the non-normalized database, it was decided to use the temporary table 'patients_not_cleaned'.

```

1 CREATE TEMPORARY TABLE patients_not_cleaned AS
2 SELECT * FROM BigData.patients;
3
4 INSERT INTO BigData_CleanedData.States (Name)
5 SELECT State
6 FROM patients_not_cleaned
7 group by State;
8
9 INSERT INTO BigData_CleanedData.Counties (Name, State)
10 SELECT REPLACE(p.County, ' County', ''), s.id
11 FROM patients_not_cleaned p
12 INNER JOIN BigData_CleanedData.States s
13     ON s.Name = p.State
14 group by s.id, p.county;
15
16 INSERT INTO BigData_CleanedData.Cities (Name, County)
17 SELECT p.city, c.id
18 FROM patients_not_cleaned p
19 INNER JOIN BigData_CleanedData.States s
20     ON s.Name = p.State
21 INNER JOIN BigData_CleanedData.Counties c
22     ON c.Name = REPLACE(p.County, ' County', '')
23 group by c.id, p.city;
24
25 INSERT INTO BigData_CleanedData.Races (Name)
26 SELECT p.race
27 FROM patients_not_cleaned p
28 GROUP BY p.race;
29
30 INSERT INTO BigData_CleanedData.Ethnicities (Name)
31 SELECT p.ethnicity
32 FROM patients_not_cleaned p
33 GROUP BY p.ethnicity;
34
35 INSERT INTO BigData_CleanedData.Patients
36 (id, birthdate, deathdate, ssn, drivers, passport, prefix, first_name,
37  last_name, suffix, maiden, marital, race, ethnicity, gender, birthplace,
38  address, city, zip, lat, lon, healthcare_expenses, healthcare_coverage)
39 SELECT p.id,
40        birthdate,
41        deathdate,
42        ssn,
43        IFNULL(NULLIF(drivers, ''), 'UNKNOWN') as drivers,
44        IFNULL(NULLIF(passport, ''), 'UNKNOWN') as passport,
45        IFNULL(NULLIF(prefix, ''), 'NO') as prefix, 'first', 'last',

```



```

44     IFNULL(NULLIF(suffix, ''), 'UNKNOWN') as suffix,
45     IFNULL(NULLIF( maiden, ''), 'UNKNOWN') as maiden,
46     IFNULL(NULLIF(marital, ''), '-') as marital,
47     r.id as race,
48     e.id as ethnicity,
49     gender,
50     birthplace,
51     address,
52     ci.id as city,
53     zip,
54     lat,
55     lon,
56     healthcare_expenses,
57     healthcare_coverage
58 FROM patients_not_cleaned p
59 INNER JOIN BigData_CleanedData.States s
60     ON s.Name = p.State
61 INNER JOIN BigData_CleanedData.Counties c
62     ON c.Name = REPLACE(p.County, ' County', '') and c.state = s.id
63 INNER JOIN BigData_CleanedData.Cities ci
64     on ci.Name = p.city and ci.county = c.id
65 INNER JOIN BigData_CleanedData.Races r
66     ON r.Name = p.race
67 INNER JOIN BigData_CleanedData.Ethnicities e
68     ON e.Name = p.Ethnicity;

```

3.2 Organizations

During the processing of the 'Organizations' entity, inconsistencies were found in the 'city' field. Many entries were written in uppercase, some referred to neighborhoods of specific cities, and others contained abbreviations. To address these issues, we performed the following steps:

1. Created temporary table 'organizations_not_cleaned'
2. Created the table 'temp_machusetts_cities' to map cities to counties.
3. Filtering of cities and counties already inserted during the 'Patients' processing.
4. Insertion of counties, cities, and organizations.

```

1 CREATE TEMPORARY TABLE organizations_not_cleaned AS
2 SELECT * FROM BigData.Organizations ;
3
4 CREATE TEMPORARY TABLE temp_machusetts_cities (
5     city VARCHAR(100),
6     county VARCHAR(100)
7 );
8
9 INSERT INTO temp_machusetts_cities (city, county) VALUES
10 ('Alford', 'Berkshire'),
11 ('Aquinnah', 'Dukes'),
12 ('Abington', 'Plymouth'),
13 ('Acton', 'Middlesex'),
14 ('Acushnet', 'Bristol'),
15 ('Adams', 'Berkshire'),
16 ('Agawam', 'Hampden'),
17 ('Amesbury', 'Essex'),
18 ('Amherst', 'Hampshire'),
19 ('Andover', 'Essex'),
20 ('Arlington', 'Middlesex'),
21 ('Ashburnham', 'Worcester'),

```

22 ('Ashby', 'Middlesex'),
 23 ('Ashfield', 'Franklin'),
 24 ('Ashland', 'Middlesex'),
 25 ('Athol', 'Worcester'),
 26 ('Attleboro', 'Bristol'),
 27 ('Auburn', 'Worcester'),
 28 ('Avon', 'Norfolk'),
 29 ('Ayer', 'Middlesex'),
 30 ('Barnstable', 'Barnstable'),
 31 ('Barre', 'Worcester'),
 32 ('Becket', 'Berkshire'),
 33 ('Bedford', 'Middlesex'),
 34 ('Belchertown', 'Hampshire'),
 35 ('Bellingham', 'Norfolk'),
 36 ('Belmont', 'Middlesex'),
 37 ('Berkley', 'Bristol'),
 38 ('Berlin', 'Worcester'),
 39 ('Bernardston', 'Franklin'),
 40 ('Beverly', 'Essex'),
 41 ('Billerica', 'Middlesex'),
 42 ('Blackstone', 'Worcester'),
 43 ('Blandford', 'Hampden'),
 44 ('Bolton', 'Worcester'),
 45 ('Boston', 'Suffolk'),
 46 ('Boxborough', 'Middlesex'),
 47 ('Boxford', 'Essex'),
 48 ('Boylston', 'Worcester'),
 49 ('Braintree', 'Norfolk'),
 50 ('Brewster', 'Barnstable'),
 51 ('Bridgewater', 'Plymouth'),
 52 ('Brimfield', 'Hampden'),
 53 ('Brockton', 'Plymouth'),
 54 ('Brookfield', 'Worcester'),
 55 ('Brookline', 'Norfolk'),
 56 ('Buckland', 'Franklin'),
 57 ('Burlington', 'Middlesex'),
 58 ('Cambridge', 'Middlesex'),
 59 ('Canton', 'Norfolk'),
 60 ('Carlisle', 'Middlesex'),
 61 ('Carver', 'Plymouth'),
 62 ('Charlemont', 'Franklin'),
 63 ('Charlton', 'Worcester'),
 64 ('Chatham', 'Barnstable'),
 65 ('Chelmsford', 'Middlesex'),
 66 ('Chelsea', 'Suffolk'),
 67 ('Cheshire', 'Berkshire'),
 68 ('Chester', 'Hampden'),
 69 ('Chesterfield', 'Hampshire'),
 70 ('Chicopee', 'Hampden'),
 71 ('Chilmark', 'Dukes'),
 72 ('Clarksburg', 'Berkshire'),
 73 ('Clinton', 'Worcester'),
 74 ('Cohasset', 'Norfolk'),
 75 ('Colrain', 'Franklin'),
 76 ('Concord', 'Middlesex'),
 77 ('Conway', 'Franklin'),
 78 ('Cummington', 'Hampshire'),
 79 ('Dalton', 'Berkshire'),
 80 ('Danvers', 'Essex'),
 81 ('Dartmouth', 'Bristol'),
 82 ('Dedham', 'Norfolk'),
 83 ('Deerfield', 'Franklin'),

84 ('Dennis', 'Barnstable'),
 85 ('Dighton', 'Bristol'),
 86 ('Douglas', 'Worcester'),
 87 ('Dover', 'Norfolk'),
 88 ('Dracut', 'Middlesex'),
 89 ('Dudley', 'Worcester'),
 90 ('Dunstable', 'Middlesex'),
 91 ('Duxbury', 'Plymouth'),
 92 ('East Bridgewater', 'Plymouth'),
 93 ('East Brookfield', 'Worcester'),
 94 ('East Longmeadow', 'Hampden'),
 95 ('Eastham', 'Barnstable'),
 96 ('Easthampton', 'Hampshire'),
 97 ('Easton', 'Bristol'),
 98 ('Edgartown', 'Dukes'),
 99 ('Egremont', 'Berkshire'),
 100 ('Erving', 'Franklin'),
 101 ('Essex', 'Essex'),
 102 ('Everett', 'Middlesex'),
 103 ('Fairhaven', 'Bristol'),
 104 ('Fall River', 'Bristol'),
 105 ('Falmouth', 'Barnstable'),
 106 ('Fitchburg', 'Worcester'),
 107 ('Florida', 'Berkshire'),
 108 ('Foxborough', 'Norfolk'),
 109 ('Framingham', 'Middlesex'),
 110 ('Franklin', 'Norfolk'),
 111 ('Freetown', 'Bristol'),
 112 ('Gardner', 'Worcester'),
 113 ('Georgetown', 'Essex'),
 114 ('Gill', 'Franklin'),
 115 ('Gloucester', 'Essex'),
 116 ('Goshen', 'Hampshire'),
 117 ('Gosnold', 'Dukes'),
 118 ('Grafton', 'Worcester'),
 119 ('Granby', 'Hampshire'),
 120 ('Granville', 'Hampden'),
 121 ('Great Barrington', 'Berkshire'),
 122 ('Greenfield', 'Franklin'),
 123 ('Groton', 'Middlesex'),
 124 ('Groveland', 'Essex'),
 125 ('Hadley', 'Hampshire'),
 126 ('Halifax', 'Plymouth'),
 127 ('Hamilton', 'Essex'),
 128 ('Hampden', 'Hampden'),
 129 ('Hancock', 'Berkshire'),
 130 ('Hanover', 'Plymouth'),
 131 ('Hanson', 'Plymouth'),
 132 ('Hardwick', 'Worcester'),
 133 ('Harvard', 'Worcester'),
 134 ('Harwich', 'Barnstable'),
 135 ('Hatfield', 'Hampshire'),
 136 ('Haverhill', 'Essex'),
 137 ('Hawley', 'Franklin'),
 138 ('Heath', 'Franklin'),
 139 ('Hingham', 'Plymouth'),
 140 ('Hinsdale', 'Berkshire'),
 141 ('Holbrook', 'Norfolk'),
 142 ('Holden', 'Worcester'),
 143 ('Holland', 'Hampden'),
 144 ('Holliston', 'Middlesex'),
 145 ('Holyoke', 'Hampden'),

```

146 ('Hopedale', 'Worcester'),
147 ('Hopkinton', 'Middlesex'),
148 ('Hubbardston', 'Worcester'),
149 ('Hudson', 'Middlesex'),
150 ('Hull', 'Plymouth'),
151 ('Huntington', 'Hampshire'),
152 ('Ipswich', 'Essex'),
153 ('Indian Orchard', 'Hampden'),
154 ('Kingston', 'Plymouth'),
155 ('Lakeville', 'Plymouth'),
156 ('Lancaster', 'Worcester'),
157 ('Lanesborough', 'Berkshire'),
158 ('Lawrence', 'Essex'),
159 ('Lee', 'Berkshire'),
160 ('Leicester', 'Worcester'),
161 ('Lenox', 'Berkshire'),
162 ('Leominster', 'Worcester'),
163 ('Leverett', 'Franklin'),
164 ('Lexington', 'Middlesex'),
165 ('Leyden', 'Franklin'),
166 ('Lincoln', 'Middlesex'),
167 ('Littleton', 'Middlesex'),
168 ('Longmeadow', 'Hampden'),
169 ('Lowell', 'Middlesex'),
170 ('Ludlow', 'Hampden'),
171 ('Lunenburg', 'Worcester'),
172 ('Lynn', 'Essex'),
173 ('Lynnfield', 'Essex'),
174 ('Malden', 'Middlesex'),
175 ('Manchester-by-the-Sea', 'Essex'),
176 ('Mansfield', 'Bristol'),
177 ('Marblehead', 'Essex'),
178 ('Marion', 'Plymouth'),
179 ('Marlborough', 'Middlesex'),
180 ('Marshfield', 'Plymouth'),
181 ('Mashpee', 'Barnstable'),
182 ('Mattapoisett', 'Plymouth'),
183 ('Maynard', 'Middlesex'),
184 ('Medfield', 'Norfolk'),
185 ('Medford', 'Middlesex'),
186 ('Medway', 'Norfolk'),
187 ('Melrose', 'Middlesex'),
188 ('Mendon', 'Worcester'),
189 ('Merrimac', 'Essex'),
190 ('Methuen', 'Essex'),
191 ('Middleborough', 'Plymouth'),
192 ('Middlefield', 'Hampshire'),
193 ('Middleton', 'Essex'),
194 ('Milford', 'Worcester'),
195 ('Millbury', 'Worcester'),
196 ('Millis', 'Norfolk'),
197 ('Millville', 'Worcester'),
198 ('Milton', 'Norfolk'),
199 ('Monroe', 'Franklin'),
200 ('Monson', 'Hampden'),
201 ('Montague', 'Franklin'),
202 ('Monterey', 'Berkshire'),
203 ('Montgomery', 'Hampden'),
204 ('Mount Washington', 'Berkshire'),
205 ('Nahant', 'Essex'),
206 ('Nantucket', 'Nantucket'),
207 ('Natick', 'Middlesex'),

```

208 ('Needham', 'Norfolk'),
 209 ('New Ashford', 'Berkshire'),
 210 ('New Bedford', 'Bristol'),
 211 ('New Braintree', 'Worcester'),
 212 ('New Marlborough', 'Berkshire'),
 213 ('New Salem', 'Franklin'),
 214 ('Newbury', 'Essex'),
 215 ('Newburyport', 'Essex'),
 216 ('Newton', 'Middlesex'),
 217 ('Norfolk', 'Norfolk'),
 218 ('North Adams', 'Berkshire'),
 219 ('Northampton', 'Hampshire'),
 220 ('North Andover', 'Essex'),
 221 ('North Attleborough', 'Bristol'),
 222 ('Northborough', 'Worcester'),
 223 ('Northbridge', 'Worcester'),
 224 ('North Brookfield', 'Worcester'),
 225 ('Northfield', 'Franklin'),
 226 ('Norton', 'Bristol'),
 227 ('Norwell', 'Plymouth'),
 228 ('Norwood', 'Norfolk'),
 229 ('Oak Bluffs', 'Dukes'),
 230 ('Oakham', 'Worcester'),
 231 ('Orange', 'Franklin'),
 232 ('Orleans', 'Barnstable'),
 233 ('Otis', 'Berkshire'),
 234 ('Oxford', 'Worcester'),
 235 ('Palmer', 'Hampden'),
 236 ('Paxton', 'Worcester'),
 237 ('Peabody', 'Essex'),
 238 ('Pelham', 'Hampshire'),
 239 ('Pembroke', 'Plymouth'),
 240 ('Pepperell', 'Middlesex'),
 241 ('Peru', 'Berkshire'),
 242 ('Petersham', 'Worcester'),
 243 ('Phillipston', 'Worcester'),
 244 ('Pittsfield', 'Berkshire'),
 245 ('Plainfield', 'Hampshire'),
 246 ('Plainville', 'Norfolk'),
 247 ('Plymouth', 'Plymouth'),
 248 ('Plympton', 'Plymouth'),
 249 ('Princeton', 'Worcester'),
 250 ('Provincetown', 'Barnstable'),
 251 ('Quincy', 'Norfolk'),
 252 ('Randolph', 'Norfolk'),
 253 ('Raynham', 'Bristol'),
 254 ('Reading', 'Middlesex'),
 255 ('Rehoboth', 'Bristol'),
 256 ('Revere', 'Suffolk'),
 257 ('Richmond', 'Berkshire'),
 258 ('Rochester', 'Plymouth'),
 259 ('Rockland', 'Plymouth'),
 260 ('Rockport', 'Essex'),
 261 ('Rowe', 'Franklin'),
 262 ('Rowley', 'Essex'),
 263 ('Royalston', 'Worcester'),
 264 ('Russell', 'Hampden'),
 265 ('Rutland', 'Worcester'),
 266 ('Salem', 'Essex'),
 267 ('Salisbury', 'Essex'),
 268 ('Sandisfield', 'Berkshire'),
 269 ('Sandwich', 'Barnstable'),

```

270 ('Saugus', 'Essex'),
271 ('Savoy', 'Berkshire'),
272 ('Scituate', 'Plymouth'),
273 ('Seekonk', 'Bristol'),
274 ('Sharon', 'Norfolk'),
275 ('Sheffield', 'Berkshire'),
276 ('Shelburne', 'Franklin'),
277 ('Sherborn', 'Middlesex'),
278 ('Shirley', 'Middlesex'),
279 ('Shrewsbury', 'Worcester'),
280 ('Shutesbury', 'Franklin'),
281 ('Somerset', 'Bristol'),
282 ('Somerville', 'Middlesex'),
283 ('Southampton', 'Hampshire'),
284 ('Southborough', 'Worcester'),
285 ('Southbridge', 'Worcester'),
286 ('South Hadley', 'Hampshire'),
287 ('Southwick', 'Hampden'),
288 ('Spencer', 'Worcester'),
289 ('Springfield', 'Hampden'),
290 ('Sterling', 'Worcester'),
291 ('Stockbridge', 'Berkshire'),
292 ('Stoneham', 'Middlesex'),
293 ('Stoughton', 'Norfolk'),
294 ('Stow', 'Middlesex'),
295 ('Sturbridge', 'Worcester'),
296 ('Sudbury', 'Middlesex'),
297 ('Sunderland', 'Franklin'),
298 ('Sutton', 'Worcester'),
299 ('Swampscott', 'Essex'),
300 ('Swansea', 'Bristol'),
301 ('Taunton', 'Bristol'),
302 ('Templeton', 'Worcester'),
303 ('Tewksbury', 'Middlesex'),
304 ('Tisbury', 'Dukes'),
305 ('Tolland', 'Hampden'),
306 ('Topsfield', 'Essex'),
307 ('Townsend', 'Middlesex'),
308 ('Truro', 'Barnstable'),
309 ('Tyngsborough', 'Middlesex'),
310 ('Tyringham', 'Berkshire'),
311 ('Upton', 'Worcester'),
312 ('Uxbridge', 'Worcester'),
313 ('Wakefield', 'Middlesex'),
314 ('Wales', 'Hampden'),
315 ('Walpole', 'Norfolk'),
316 ('Waltham', 'Middlesex'),
317 ('Ware', 'Hampshire'),
318 ('Wareham', 'Plymouth'),
319 ('Warren', 'Worcester'),
320 ('Warwick', 'Franklin'),
321 ('Washington', 'Berkshire'),
322 ('Watertown', 'Middlesex'),
323 ('Wayland', 'Middlesex'),
324 ('Webster', 'Worcester'),
325 ('Wellesley', 'Norfolk'),
326 ('Wellfleet', 'Barnstable'),
327 ('Wendell', 'Franklin'),
328 ('Wenham', 'Essex'),
329 ('Westborough', 'Worcester'),
330 ('West Boylston', 'Worcester'),
331 ('West Bridgewater', 'Plymouth'),

```

```

332 ('West Brookfield', 'Worcester'),
333 ('Westfield', 'Hampden'),
334 ('Westford', 'Middlesex'),
335 ('Westhampton', 'Hampshire'),
336 ('Westminster', 'Worcester'),
337 ('West Newbury', 'Essex'),
338 ('Westport', 'Bristol'),
339 ('West Springfield', 'Hampden'),
340 ('West Stockbridge', 'Berkshire'),
341 ('West Tisbury', 'Dukes'),
342 ('Westwood', 'Norfolk'),
343 ('Weymouth', 'Norfolk'),
344 ('Whately', 'Franklin'),
345 ('Whitman', 'Plymouth'),
346 ('Wilbraham', 'Hampden'),
347 ('Williamsburg', 'Hampshire'),
348 ('Williamstown', 'Berkshire'),
349 ('Wilmington', 'Middlesex'),
350 ('Winchendon', 'Worcester'),
351 ('Winchester', 'Middlesex'),
352 ('Windsor', 'Berkshire'),
353 ('Winthrop', 'Suffolk'),
354 ('Woburn', 'Middlesex'),
355 ('Worcester', 'Worcester'),
356 ('Worthington', 'Hampshire'),
357 ('Wrentham', 'Norfolk'),
358 ('Yarmouth', 'Barnstable'),
359 ('Shelburne Falls', 'Franklin'),
360 ('North Falmouth', 'Barnstable'),
361 ('East Wareham', 'Plymouth'),
362 ('South Attleboro', 'Bristol'),
363 ('West Wareham', 'Plymouth');
364
365 CREATE TEMPORARY TABLE temp_city_abbreviations (
366     short_name VARCHAR(100),
367     full_name VARCHAR(100)
368 );
369
370 INSERT INTO temp_city_abbreviations (short_name, full_name) VALUES
371 ('W CONCORD', 'Concord'),
372 ('INDIAN ORCHARD', 'Springfield'),
373 ('MARSTONS MILLS', 'Barnstable'),
374 ('NORTH DARTMOUTH', 'Dartmouth'),
375 ('SOUTH DARTMOUTH', 'Dartmouth'),
376 ('EAST BOSTON', 'Boston'),
377 ('FEEDING HILLS', 'Agawam'),
378 ('JAMAICA PLAIN', 'Boston'),
379 ('SOUTH WEYMOUTH', 'Weymouth'),
380 ('BRIGHTON', 'Boston'),
381 ('E SANDWICH', 'Sandwich'),
382 ('EAST WAREHAM', 'Wareham'),
383 ('DORCHESTER', 'Boston'),
384 ('SOUTH ATTLEBORO', 'Attleboro'),
385 ('MIDDLEBORO', 'Middleborough'),
386 ('WABAN', 'Newton'),
387 ('HYANNIS', 'Barnstable'),
388 ('BRADFORD', 'Haverhill'),
389 ('WELLESLEY HILLS', 'Wellesley'),
390 ('FOXBORO', 'Foxborough'),
391 ('POCASSET', 'Bourne'),
392 ('WESTBORO', 'Westborough'),
393 ('SAGAMORE BEACH', 'Bourne'),

```

```

394 ('SHELBURNE FLS', 'Shelburne Falls'),
395 ('ROSLINDALE', 'Boston'),
396 ('WEST WAREHAM', 'Wareham'),
397 ('NO FALMOUTH', 'Falmouth'),
398 ('W SPRINGFIELD', 'Springfield'),
399 ('N ADAMS', 'Adams'),
400 ('S ATTLEBORO', 'Attleboro'),
401 ('FLORENCE', 'Northampton'),
402 ('TEATICKET', 'Falmouth'),
403 ('PITTSFILED', 'Pittsfield'),
404 ('CEDARVILLE', 'Plymouth'),
405 ('WEST ROXBURY', 'Boston'),
406 ('FISKDALE', 'Sturbridge'),
407 ('CENTERVILLE', 'Barnstable'),
408 ('BALDWINVILLE', 'Templeton'),
409 ('HARWICH PORT', 'Harwich'),
410 ('Leeds', 'Northampton'),
411 ('NORTH ATTLEBORO', 'North Attleborough'),
412 ('S YARMOUTH', 'South Yarmouth'),
413 ('SOUTHBORO', 'Southborough'),
414 ('HYDE PARK', 'Boston'),
415 ('DENNIS PORT', 'Dennis'),
416 ('EAST HARWICH', 'Harwich'),
417 ('MARLBORO', 'Marlborough'),
418 ('W WAREHAM', 'Wareham'),
419 ('E WAREHAM', 'Wareham'),
420 ('SHELBURNE FALLS', 'Shelburne Falls');
421
422 CREATE TEMPORARY TABLE cities_to_insert AS
423     SELECT CONCAT(UCASE(SUBSTRING(n.city, 1, 1)), LOWER(SUBSTRING(n.city,
424         2))) AS city
425     FROM(
426     SELECT CASE
427         WHEN t.short_name is null THEN o.city
428         ELSE t.full_name
429     END as city
430     FROM organizations_not_cleaned o
431     LEFT JOIN temp_city_abbreviations t
432     ON o.city = t.short_name
433     ) as n
434     where CONCAT(UCASE(SUBSTRING(n.city, 1, 1)), LOWER(SUBSTRING(n.city,
435         2))) not in
436         (
437         SELECT c.name FROM BigData_CleanedData.Cities c
438         );
439
440 CREATE TEMPORARY TABLE cities_and_county_to_insert AS
441 SELECT c.city, mc.county
442 FROM cities_to_insert c
443 INNER JOIN temp_massachusetts_cities mc
444     on mc.city = c.city
445 group by c.city, mc.county
446
447
448 INSERT INTO BigData_CleanedData.Counties (Name, state)
449 SELECT t.county , 1
450 FROM cities_and_county_to_insert t
451 LEFT JOIN BigData_CleanedData.Counties c
452     ON t.county = c.name
453 WHERE c.name is null

```



```

454 GROUP BY t.county;
455
456
457 INSERT INTO BigData_CleanedData.Cities (Name, county)
458 SELECT t.city , c.id
459 FROM cities_and_county_to_insert t
460 INNER JOIN BigData_CleanedData.Counties c
461     ON t.county = c.name
462 LEFT JOIN BigData_CleanedData.Cities ci
463     ON t.city = ci.name
464 where ci.name is null
465 GROUP BY t.city,c.id;
466
467 INSERT INTO BigData_CleanedData.Organizations
468 (id, name, address, city, zip_code, lat, lon, phone, revenue, utilization)
469 SELECT o.id,
470     o.name,
471     address,
472     c.id as city,
473     zip,
474     lat,
475     lon,
476     phone,
477     revenue,
478     utilization
479 FROM BigData.organizations o
480 LEFT JOIN temp_city_abbreviations t
481     ON o.city = t.short_name
482 INNER JOIN BigData_CleanedData.Cities c
483     ON c.Name = o.city or c.Name = t.full_name
484 GROUP BY o.id, o.name, address, c.id, zip, lat, lon, phone, revenue,
    utilization;

```

4 Data Warehouse design

4.1 Design choices

For the design of the data marts, the analysis of the E-R schema (Section 2.1) and the query list (Section 2.2) has led to the identification of 7 facts, which divide the query list as represented in Table 1.

Identified Fact	Query List
Administer	<ul style="list-style-type: none">• How many individuals, in the age group between 10 and 50 years, were vaccinated against HPV between the years 2015 and 2019, in the counties of Bristol, Essex, and Middlesex, grouped by sex, county, and year of administration?
Affected_By	<ul style="list-style-type: none">• What are the most common allergies, grouped by patient age group, and county where they were recorded?• What is the incidence of the various side effects associated with each allergy, broken down by patient age group?
Carry_Out	<ul style="list-style-type: none">• What is the cost of medical procedures broken down by year for the various counties and hospital facilities?
Conduct	<ul style="list-style-type: none">• What is the number of visits performed by each doctor, broken down by patient age group, gender, race, and, in addition, county and year of the conducted visit?
Plan	<ul style="list-style-type: none">• What is the number of respiratory therapies recorded in various hospitals, broken down by gender and age group?• What is the average age of patients affected by acute bronchitis or prediabetes in each county recorded in the last five years?• What is the average number of treatment days required to complete a therapy plan for various diseases in each hospital over the years?
Require	<ul style="list-style-type: none">• What are the most used supplies each year, broken down by county and hospital facility?
Take_Place_In	<ul style="list-style-type: none">• What is the number of visits made to each hospital facility in the last five years, broken down by gender and age group?

Table 2: Identified facts and associated queries.

4.2 Data mart

For each previously identified fact, the corresponding data mart design schemes follow. To obtain a fact schema, we need to perform the following steps:

1. Construction of the attribute tree
2. Pruning and grafting of the attribute tree
3. Definition of dimensions
4. Definition of measures
5. Creation of the fact schema

The following paragraphs illustrate how the steps are applied to obtain a fact schema.

4.3 Administer

The N-to-N relationship between Immunization and Encounter was reified by creating the entity Administer, Figure 4. The attribute tree was then constructed, as shown in Figure 5.

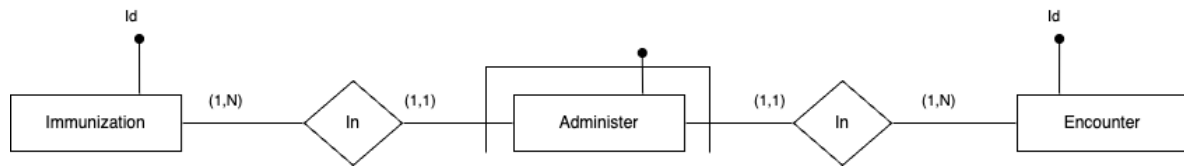


Figure 4: Reification of the "Administer" relationship

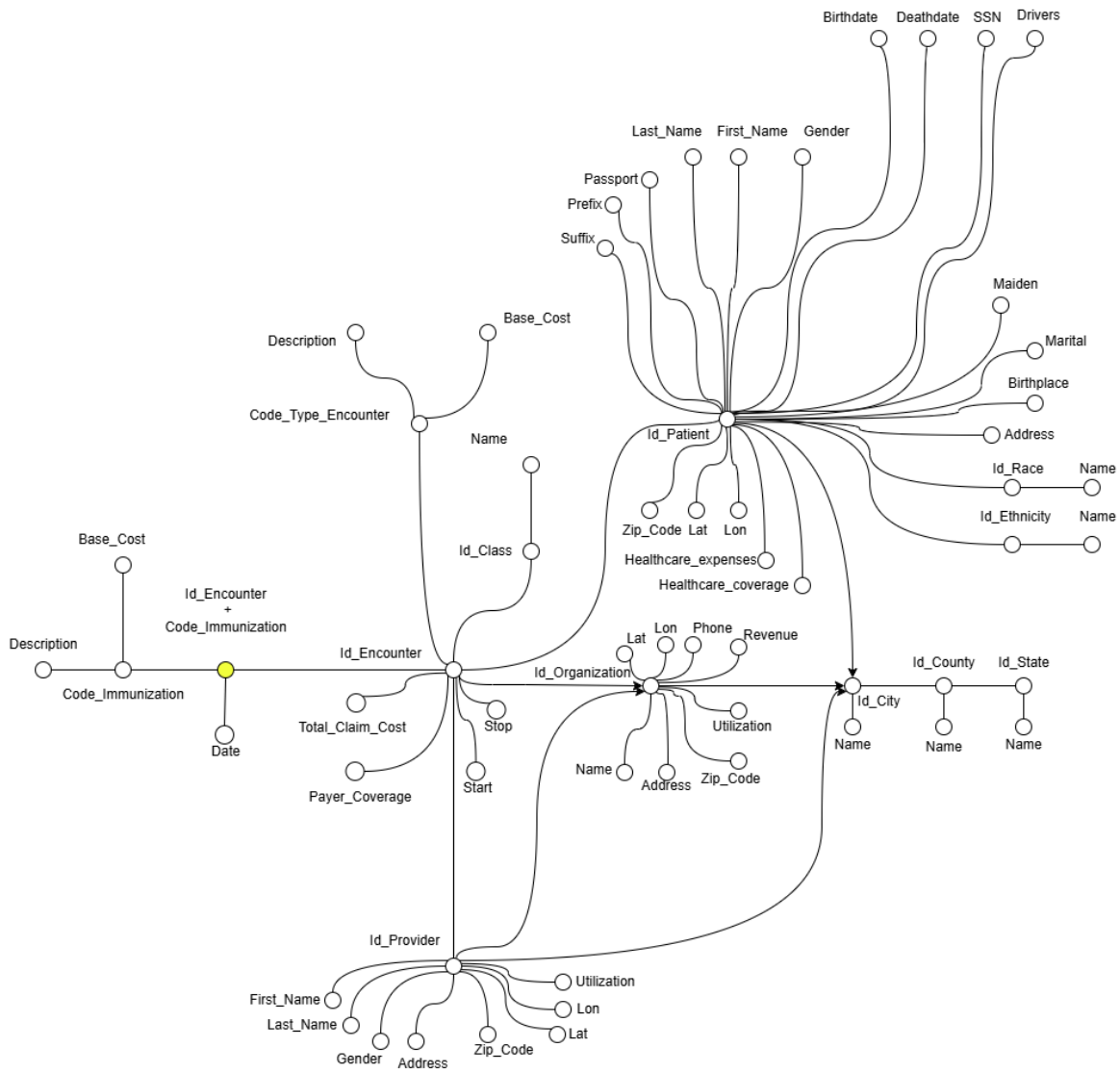


Figure 5: Complete attribute tree of the "Administer" fact

For the construction of the final attribute tree of the fact "Administer", we pruned:

- "Base_Cost" linked to "Code_Immunization"
- "Code_Type_Encounter" and its subtree
- "Id.Provider" and its subtree
- "Id.Class" and its subtree
- "Total_Claim_Cost", "Payer_Coverage", "Start" and "Stop" from "Id_Encounter"
- All attributes associated with "Id.Patient", except for "Gender", "Birthdate", and "Id.Race" along with its subtree
- All attributes associated with "Id.Organization", except for "Id.City" and its subtree

We then applied the grafting operation to the following:

- "Gender", "Birthdate", and "Id.Race" of "Id.Patient"
- "Id.City" of "Id.Organization"
- "Name" attributes, from "Id.County" and "Id.State" have been renamed in "County_Name" and "State_Name", before being grafted

The result is shown in the following figure:

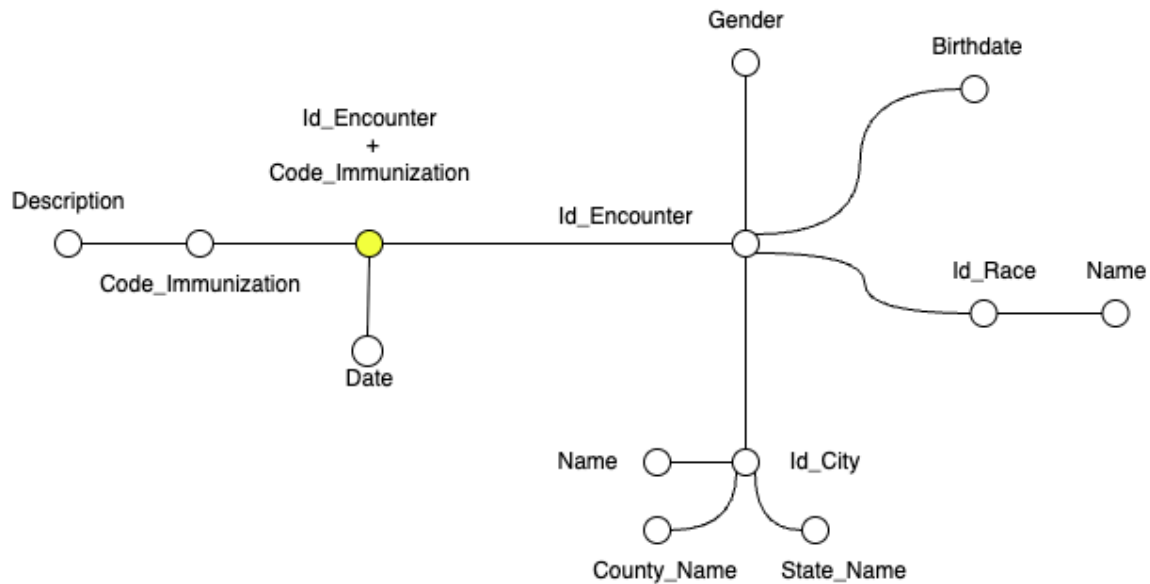


Figure 6: Pruned and grafted attribute tree of the "Administer" fact

Based on the final attribute tree, the fact schema was then constructed (Figure 7). The selected dimensions are: Immunization, Encounter and Date. The chosen measure is Amount_Immunization, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, a date hierarchy was added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 8).

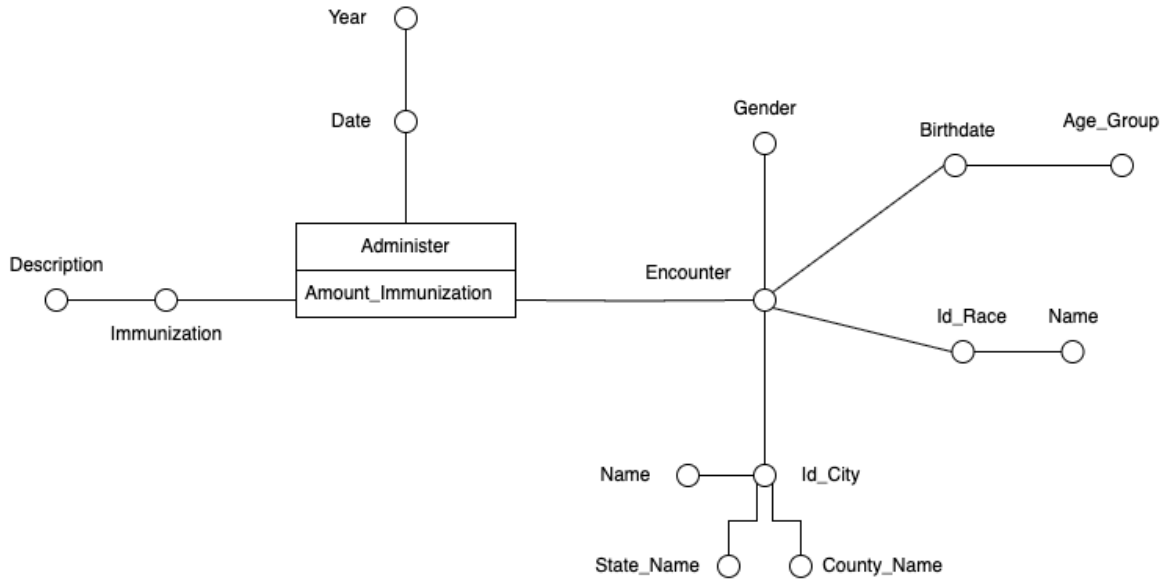


Figure 7: Fact schema of the "Administer" fact

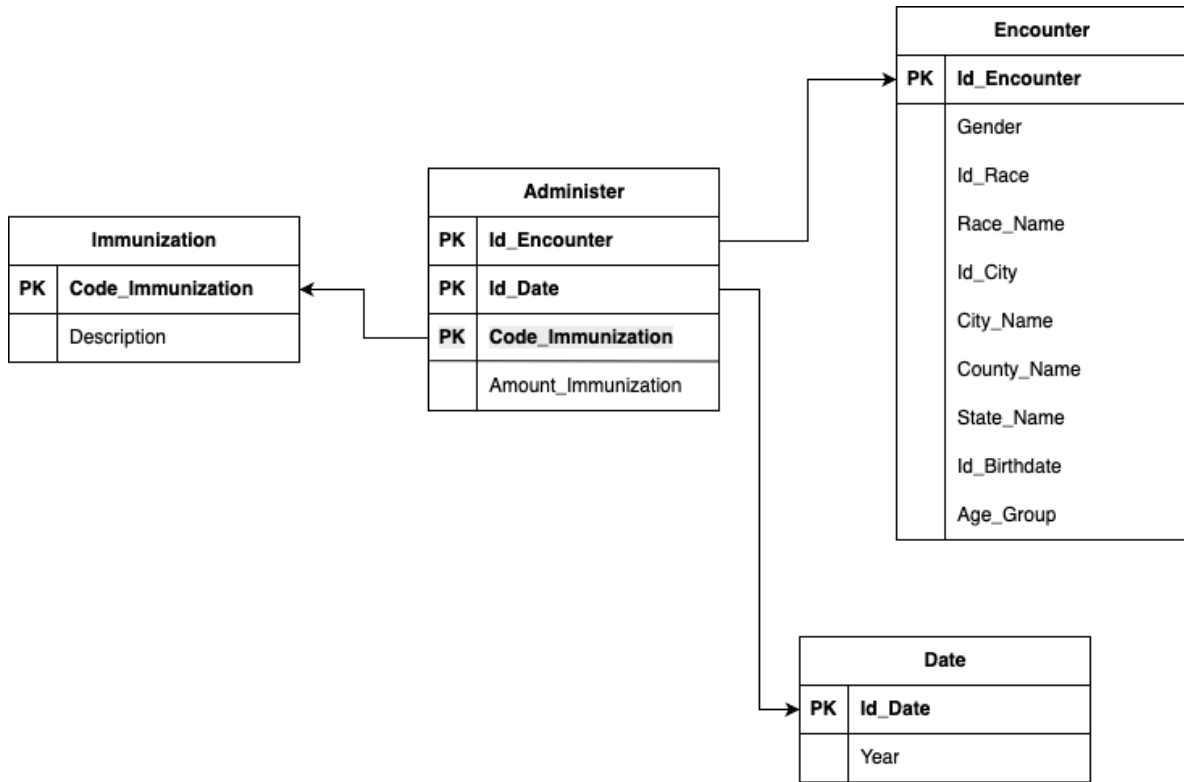


Figure 8: Star schema of the "Administer" fact

4.4 Affected_By

The N-to-N relationship between Conditions, Allergy and Encounter was reified by creating the entity Affected_By, Figure 9. The attribute tree was then constructed, as shown in Figure 10.

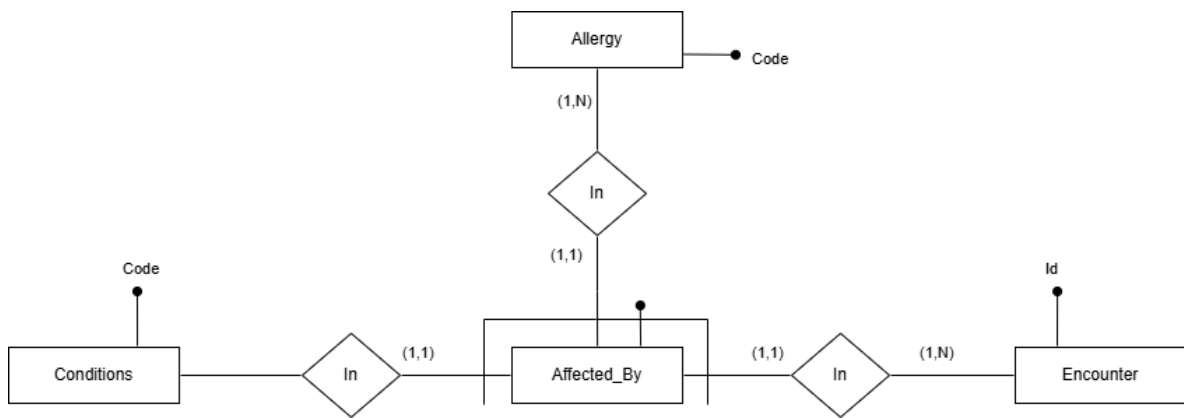


Figure 9: Reification of the "Affected_by" relationship

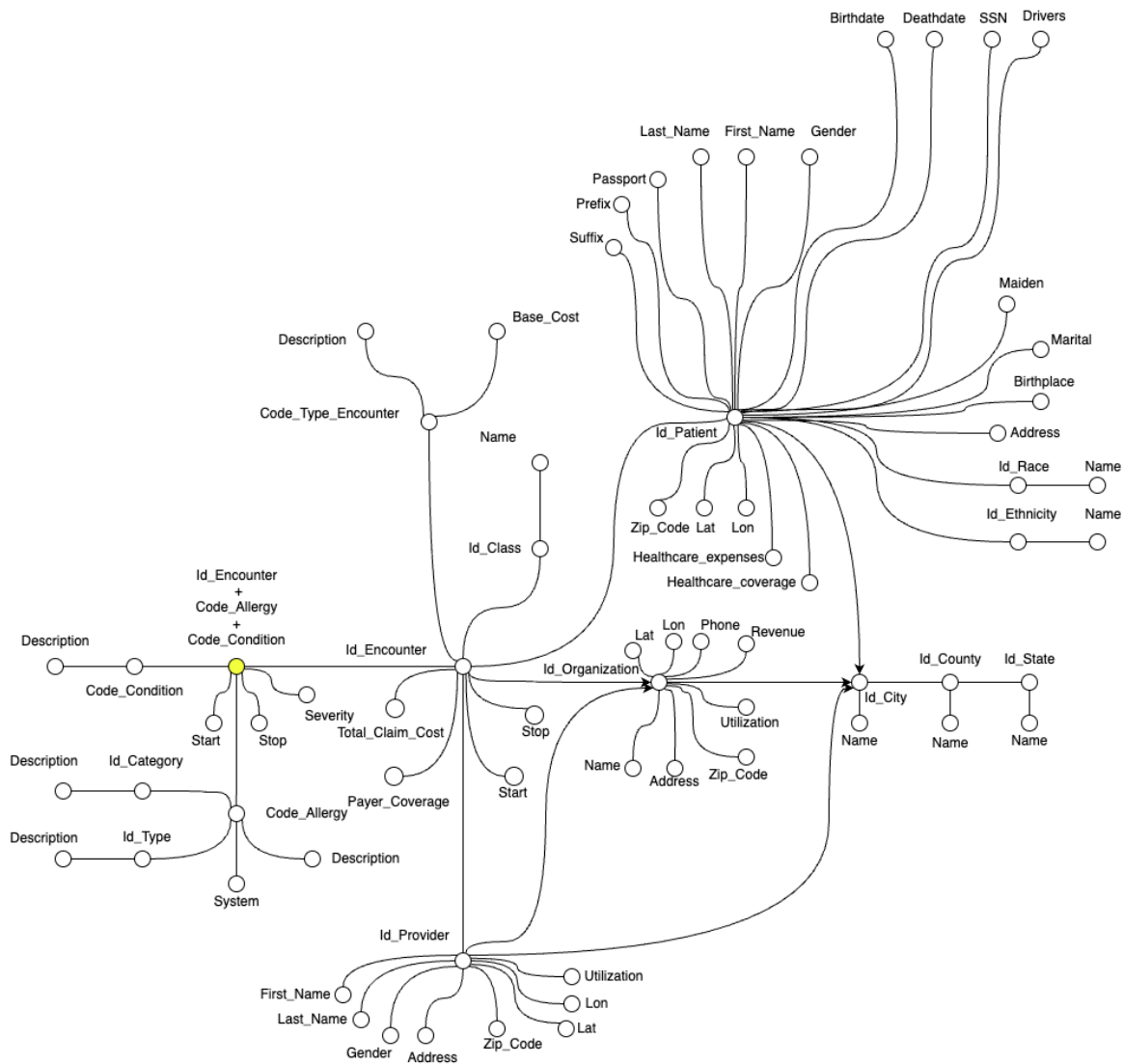


Figure 10: Complete attribute tree of the "Affected_By" fact

For the construction of the final attribute tree of the fact "Affected.By", we pruned:

- "Start", "Stop" and "Severity" from the starting node
- All attributes linked to "Code_Allergy", excluding "Description"
- "Code.Type.Encounter" and its subtree
- "Id.Class" and its subtree
- All attributes associated with "Id.Patient", except for "Gender", "Birthdate", and "Id.Race" along with its subtree
- All attributes associated with "Id.Organization", except for "Id.City" and its subtree

We then applied the grafting operation to the following:

- "Gender", "Birthdate", and "Id.Race" of "Id.Patient"
- "Id.City" of "Id.Organization"
- "Name" attributes, from "Id.County" and "Id.State" have been renamed in "County_Name" and "State_Name", before being grafted

The result is shown in the following figure:

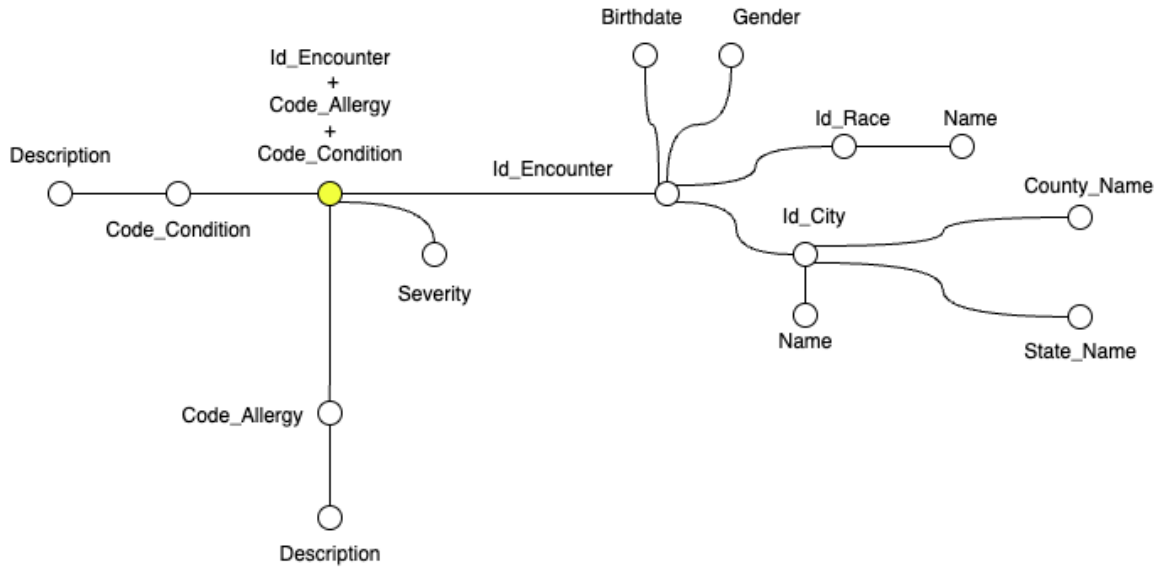


Figure 11: Pruned and grafted attribute tree of the "Affected.By" fact

Based on the final attribute tree, the fact schema was constructed (Figure 12). The selected dimensions are: Condition, Encounter and Allergy. The chosen measure are Amount.Condition and Amount.Allergy, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). In the end, the star schema was built (Figure 13).

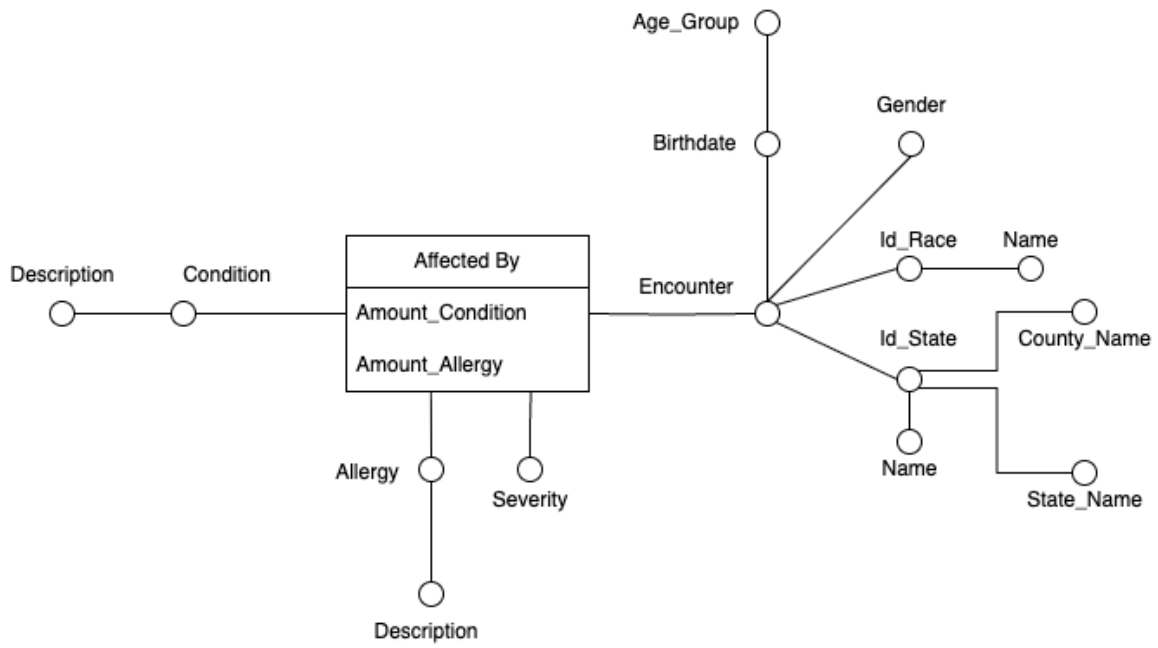


Figure 12: Fact schema of the "Affected_By" fact

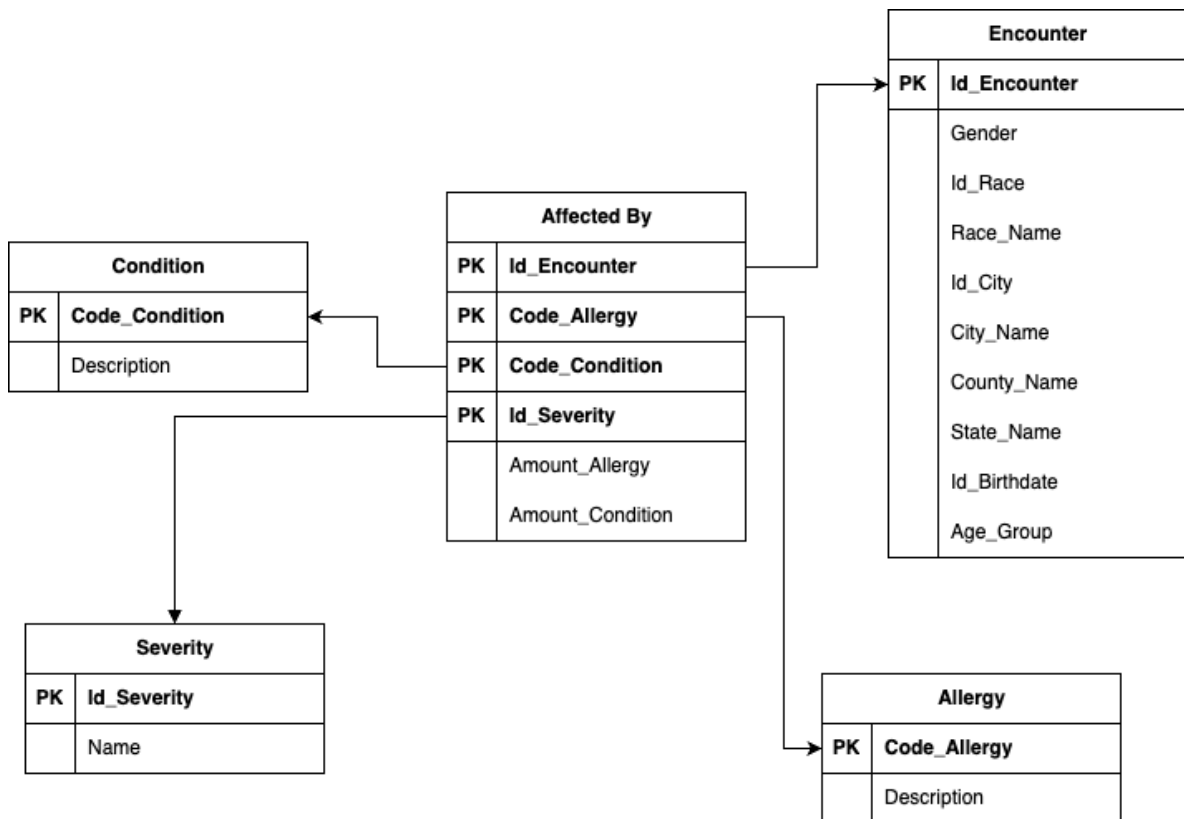


Figure 13: Star schema of the "Affected_By" fact

4.5 Carry_Out

The N-to-N relationship between Procedure and Encounter was reified by creating the entity Carry_Out, Figure 14. The attribute tree was then constructed, as shown in Figure 15.

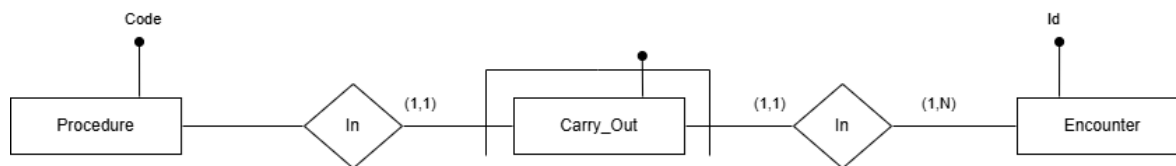


Figure 14: Reification of the "Carry_out" relationship

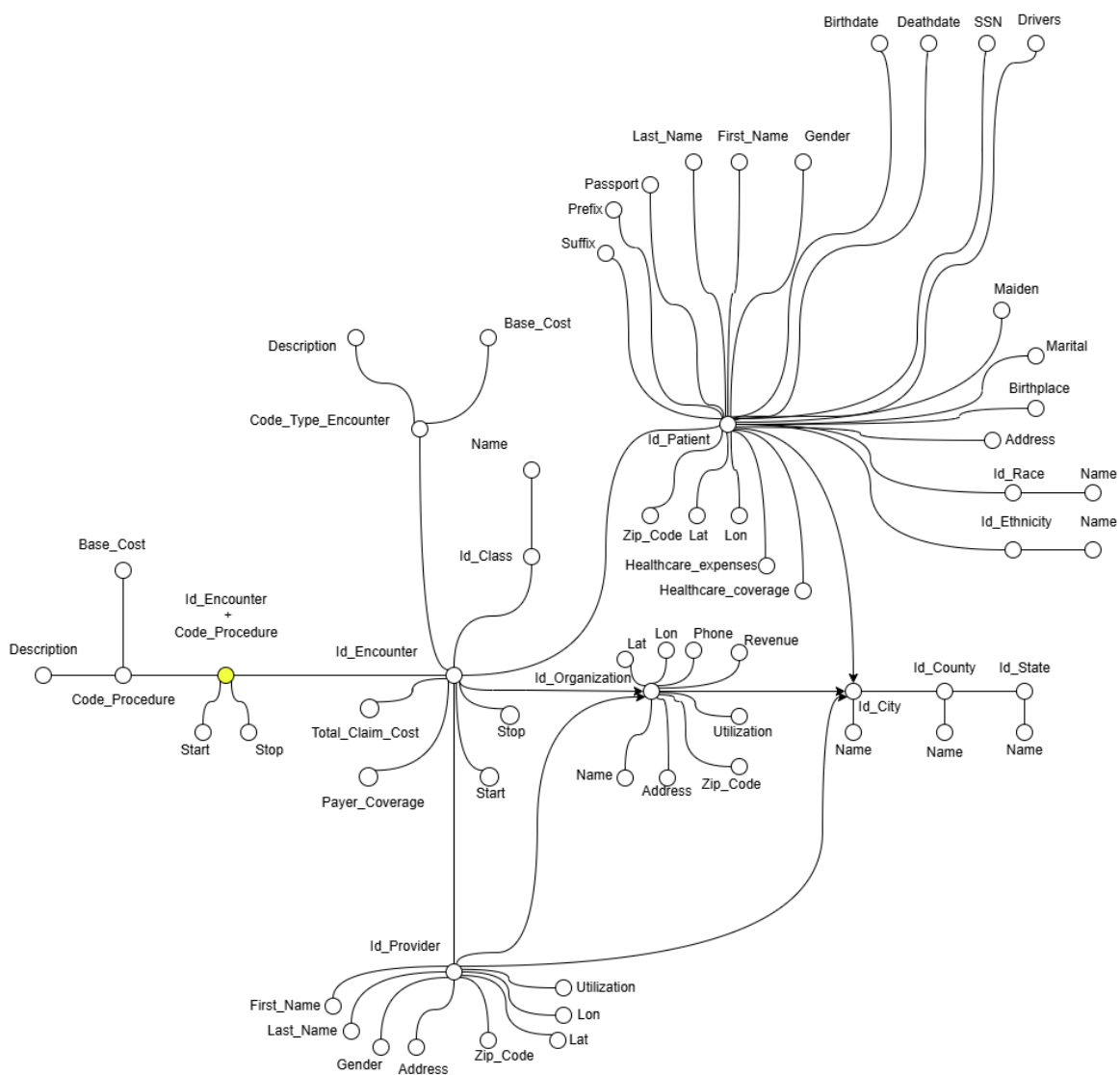


Figure 15: Complete attribute tree of the "Carry_Out" fact

For the construction of the final attribute tree of the fact "Carry_Out", we pruned:

- "Stop" from the starting node
- All attributes linked to "Id_Encounter", excluding "Id.Organization" and its subtree
- All attributes linked to "Id.Organization", excluding "Name" and "Id.City"
- "Name" from "Id.City"
- "Id.State" and its subtree

We then applied the grafting operation to the following:

- "Id.County" of "Id.City"

The result is shown in the following figure:

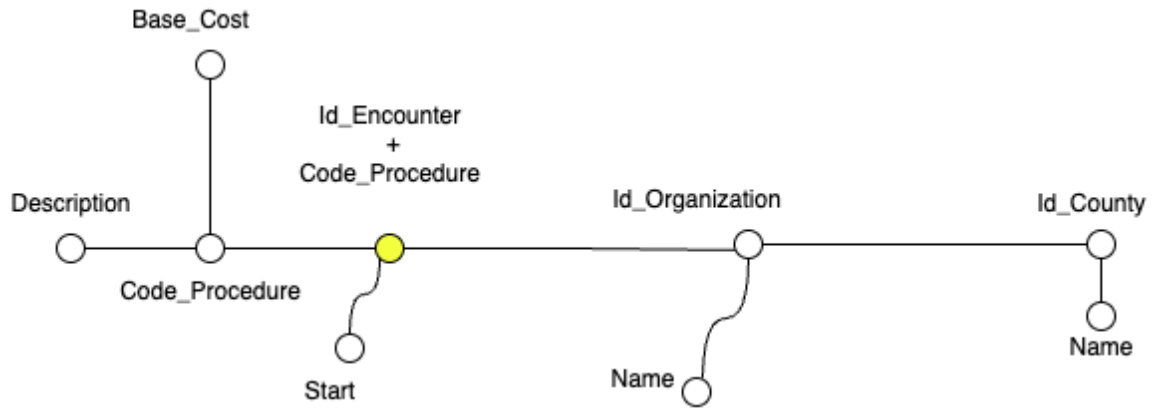


Figure 16: Pruned and grafted attribute tree of the "Carry_out" fact

Based on the final attribute tree, the fact schema was constructed (Figure 17). The selected dimensions are: Organization, Procedure and Start. The chosen measure is Total_Cost, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, a start hierarchy was added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 18).

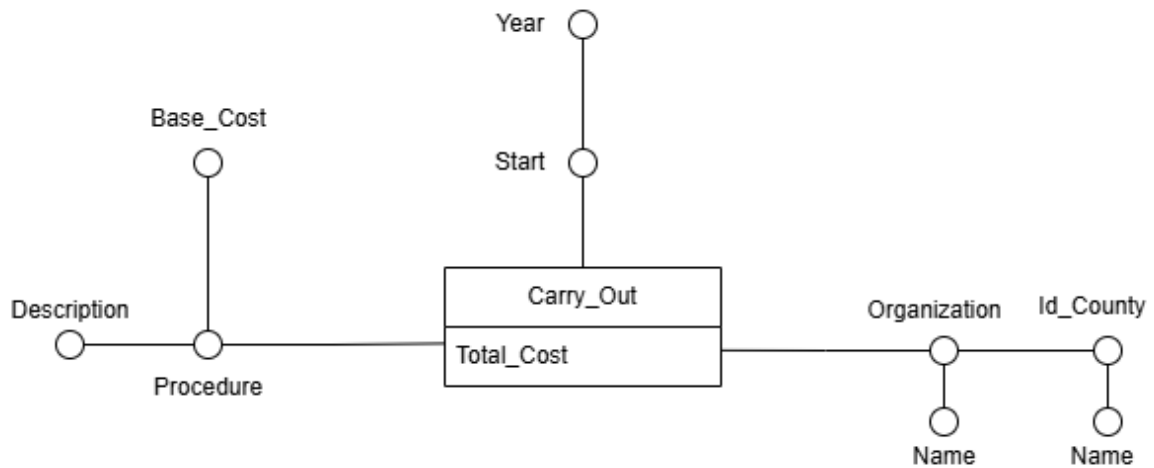


Figure 17: Fact schema of the "Carry_Out" fact

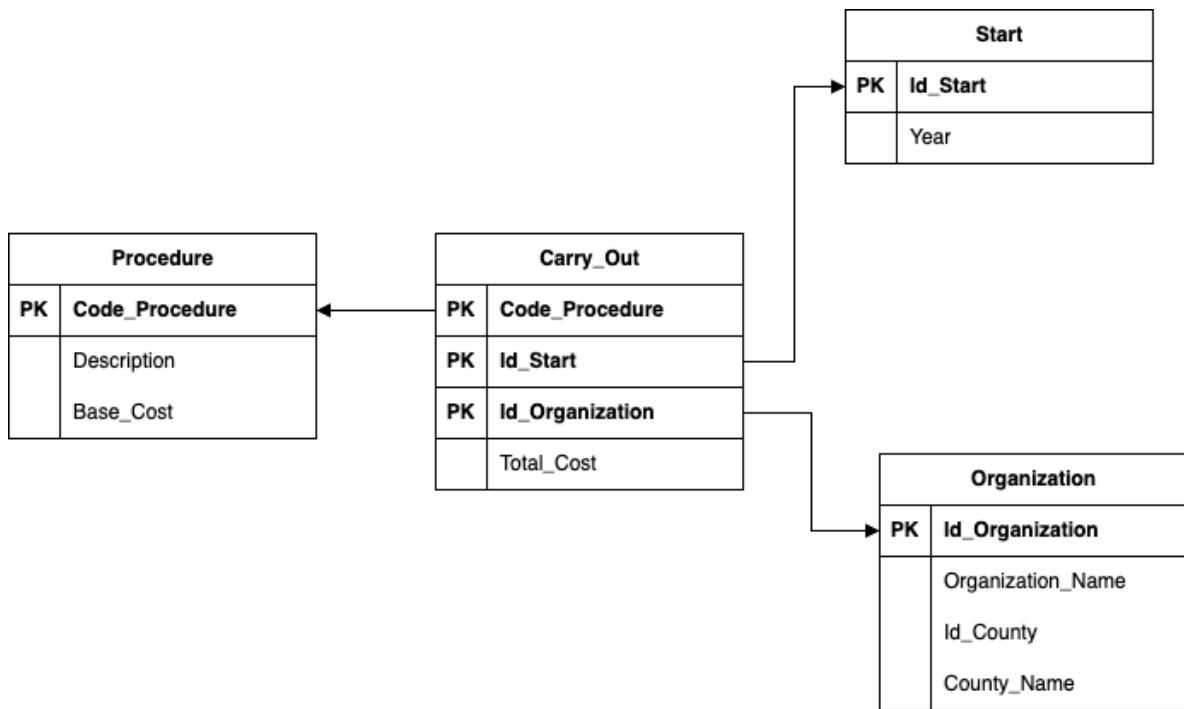


Figure 18: Star schema of the "Carry_Out" fact

4.6 Conduct

The N-to-N relationship between Providers and Encounter was reified by creating the entity Conduct, Figure 19. The attribute tree was then constructed, as shown in Figure 20.

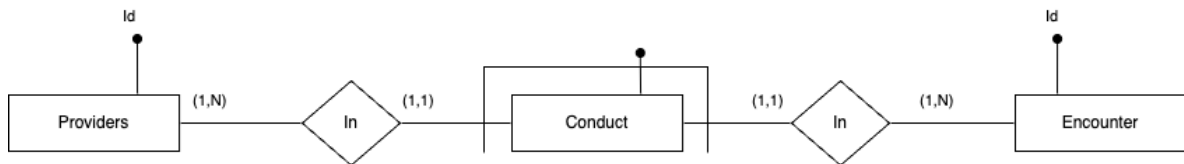


Figure 19: Reification of the "Conduct" relationship

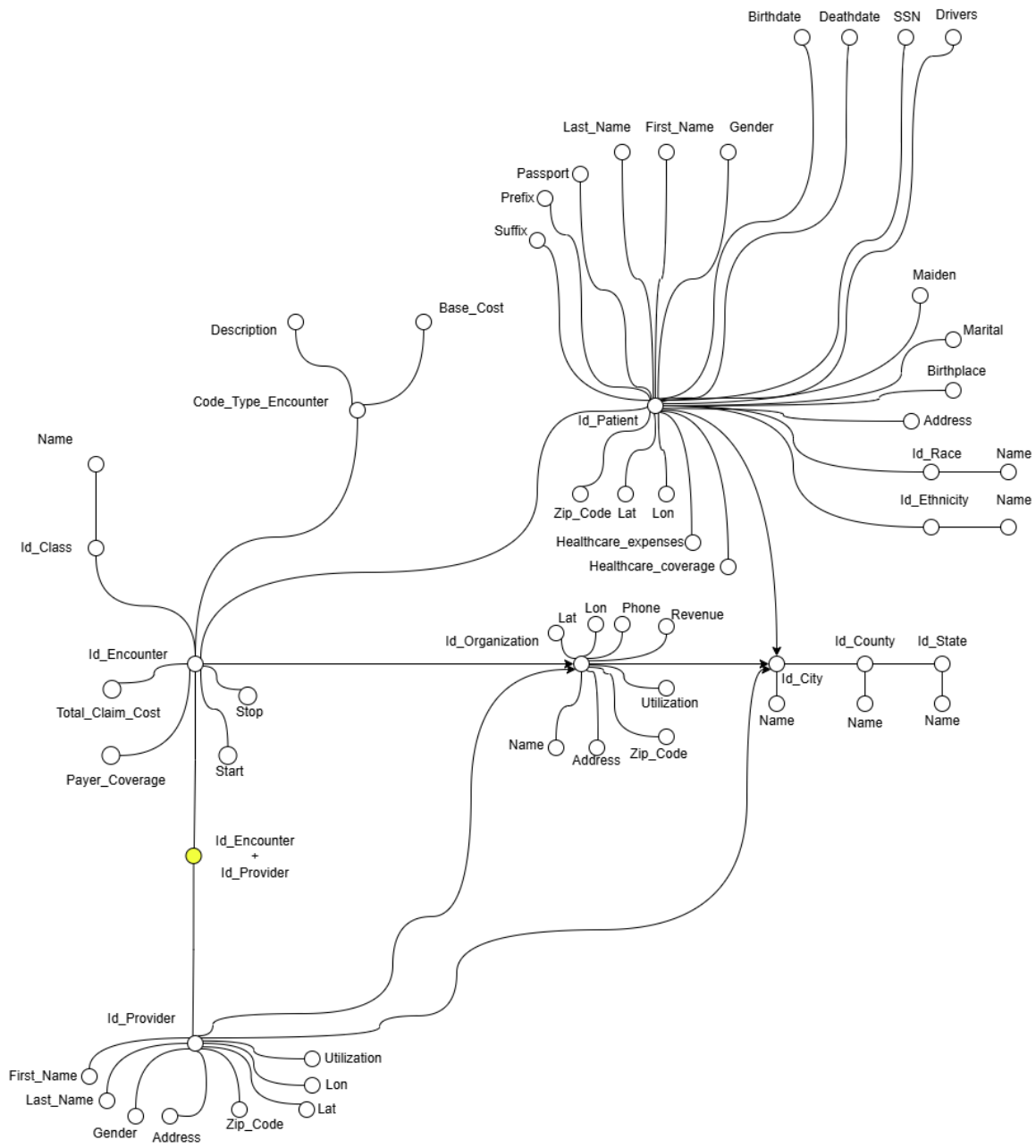


Figure 20: Complete attribute tree of the "Conduct" fact

For the construction of the final attribute tree of the fact "Conduct", we pruned:

- All attributes linked to "Id.Provider", except for "First_Name" and "Last_Name"
- All attributes linked to "Id.Encounter", excluding "Start", "Id.Organization" and its subtree and "Id.Patient" with its subtree
- All attributes of "Id.Patient", except for "Gender", "Birthdate", "Id.Race" and its subtree
- All attributes of "Id.Organization", excluding "Id.City" and its subtree
- "Name" from "Id.City"
- "Id.State" and its subtree

We then applied the grafting operation to the following:

- "Id.County" of "Id.City" and then of "Id.Organization"

The result is shown in the following figure:

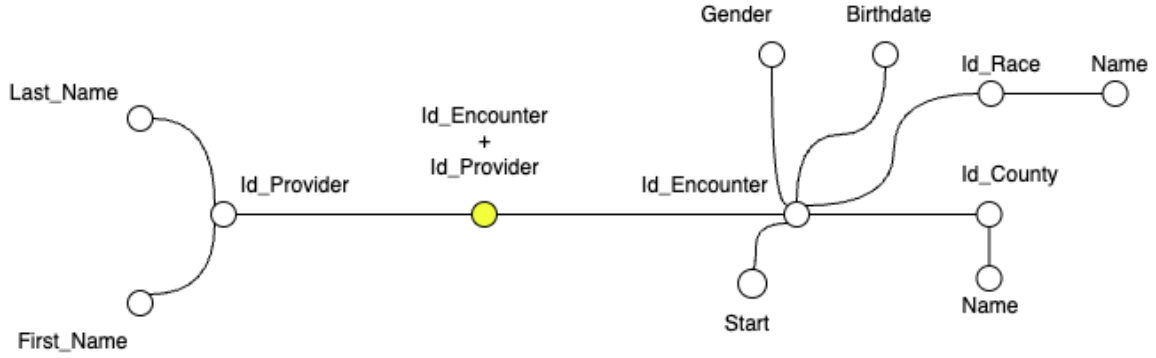


Figure 21: Pruned and grafted attribute tree of the "Conduct" fact

Based on the final attribute tree, the fact schema was constructed (Figure 21). The selected dimensions are: Provider and Encounter. The chosen measure is Amount_Visit, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, a start hierarchy and a birthdate hierarchy were added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 22).

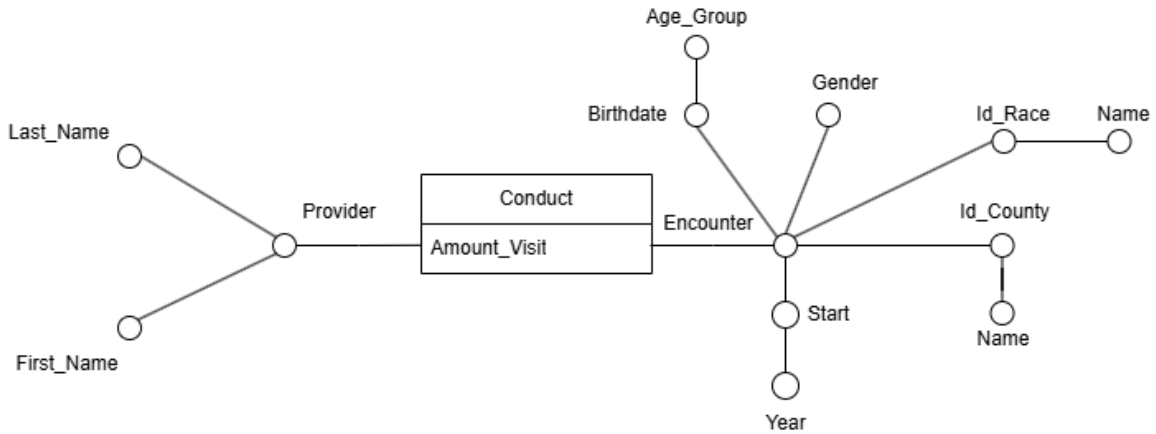


Figure 22: Fact schema of the "Conduct" fact

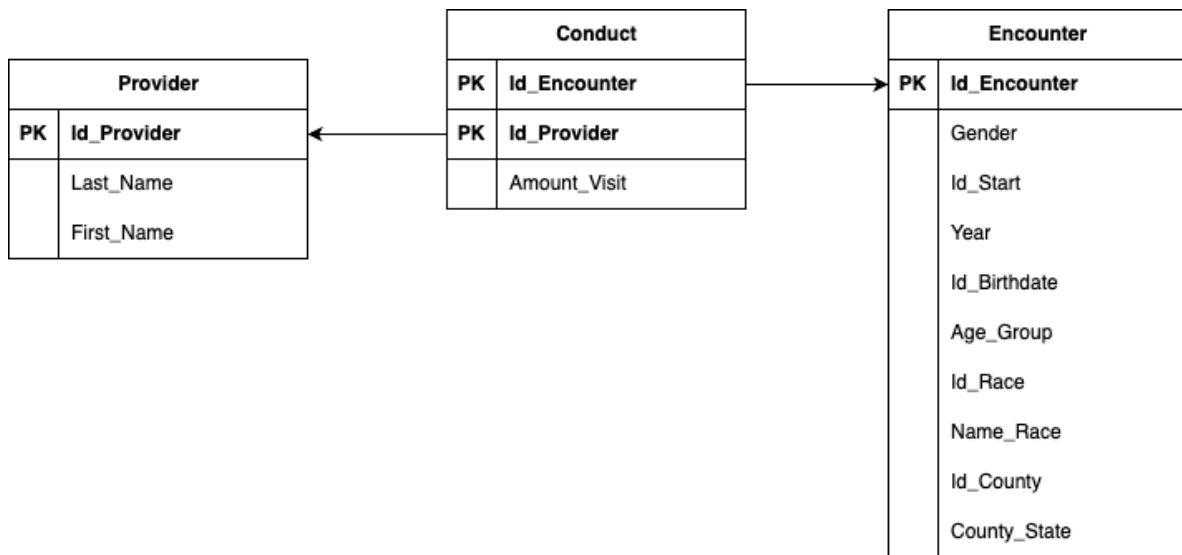


Figure 23: Star schema of the "Conduct" fact

4.7 Plan

The N-to-N relationship between Conditions, Careplans, and Encounter was reified by creating the entity Plan, Figure 24. The attribute tree was then constructed, as shown in Figure 25.

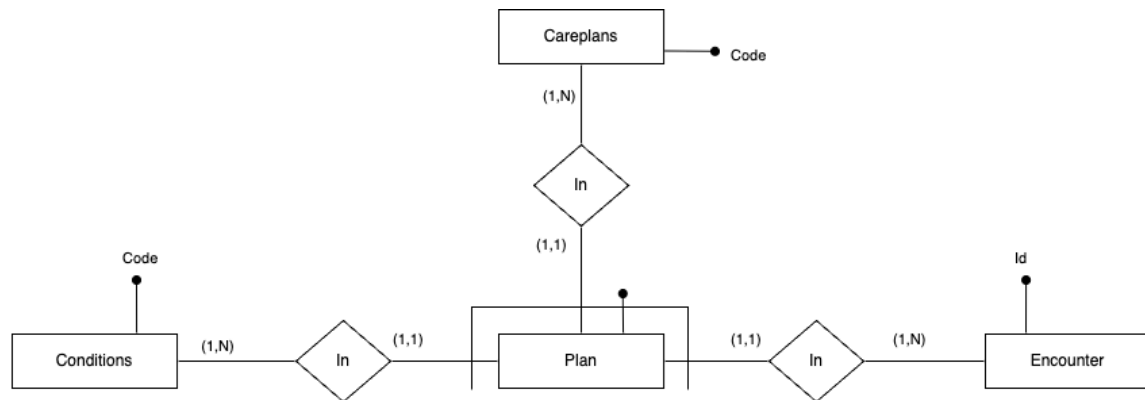


Figure 24: Reification of the "Plan" relationship

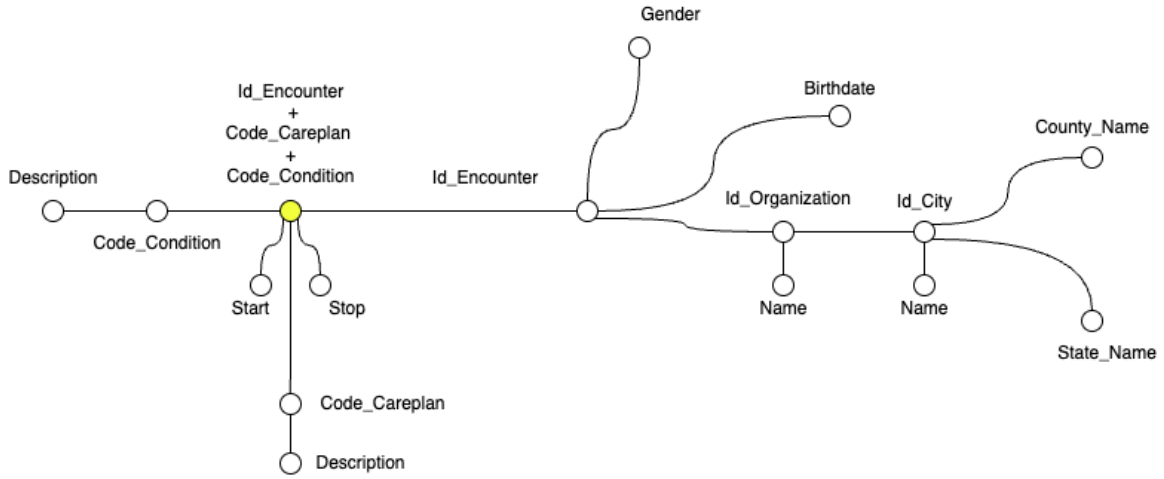


Figure 26: Pruned and grafted attribute tree of the "Plan" fact

Based on the final attribute tree, the fact schema was constructed (Figure 27). The selected dimensions are: Careplan, Condition, Encounter and Start. The chosen measure are Average_Age, Average_Days_Of_Hospitalization and Amount_Therapies, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, several hierarchies were added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 28).

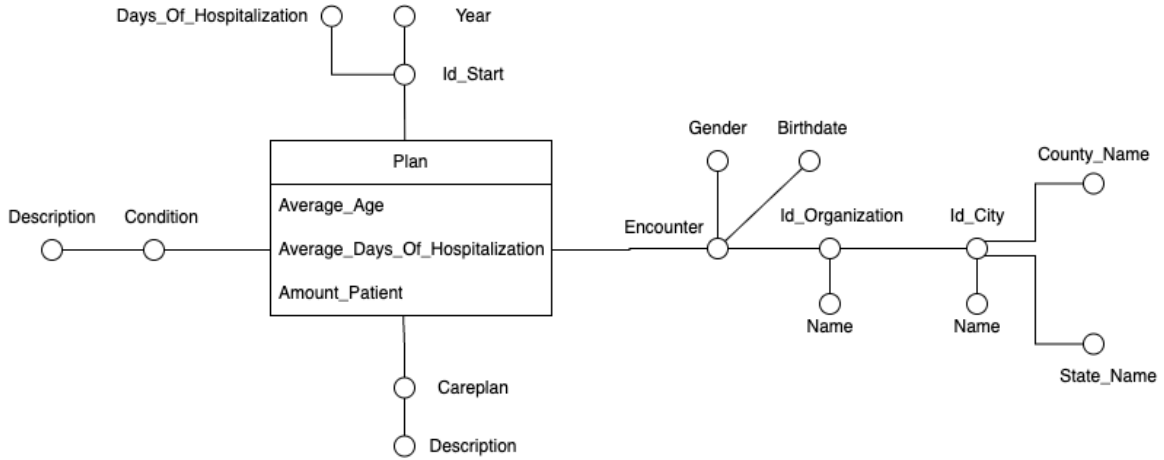


Figure 27: Fact schema of the "Plan" fact

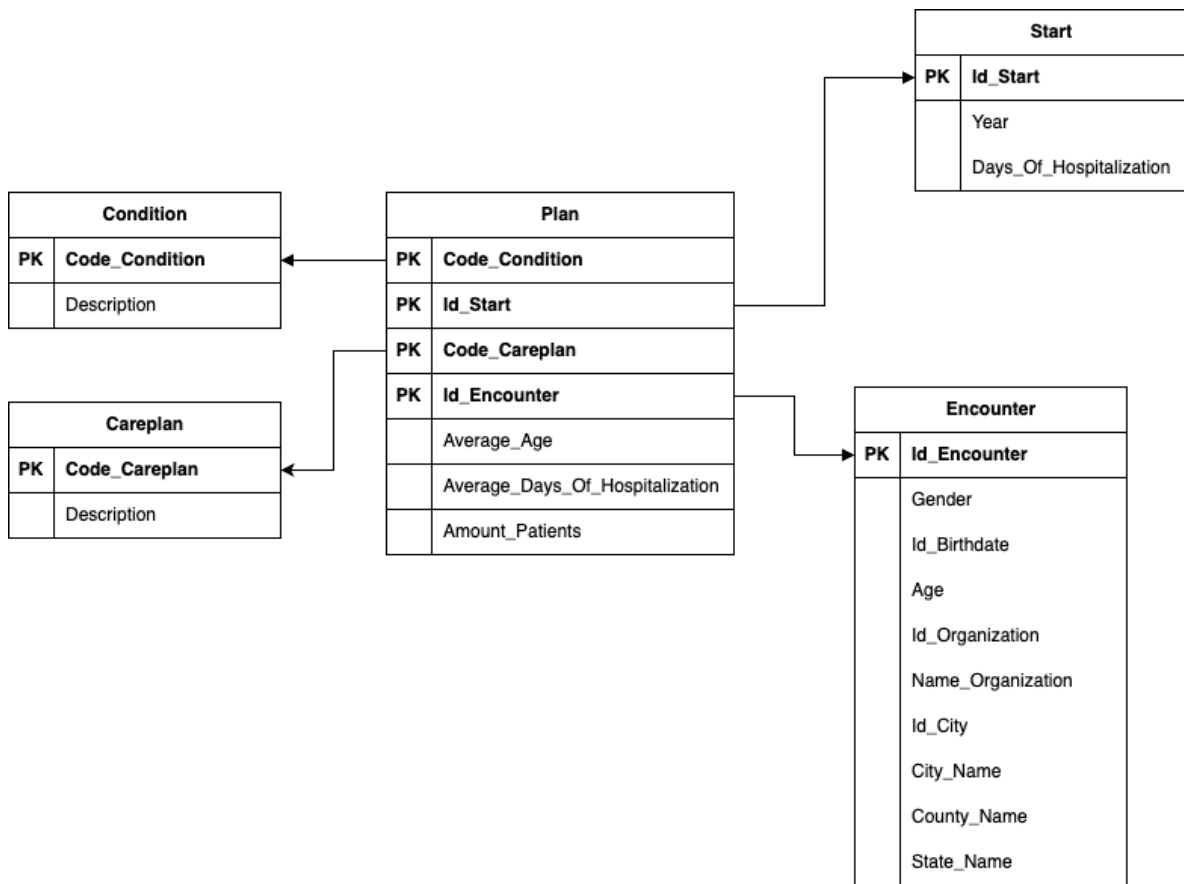


Figure 28: Star schema of the "Plan" fact

4.8 Require

The N-to-N relationship between Supply and Encounter was reified by creating the entity Require, Figure 29. The attribute tree was then constructed, as shown in Figure 30.

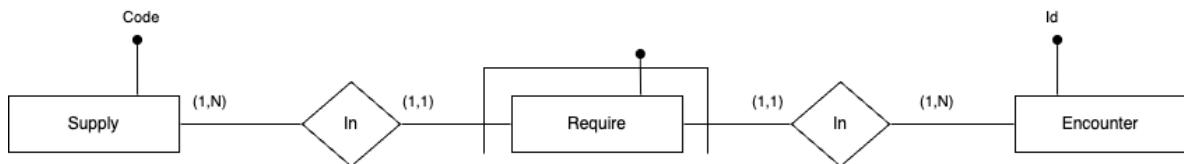


Figure 29: Reification of the "Require" relationship

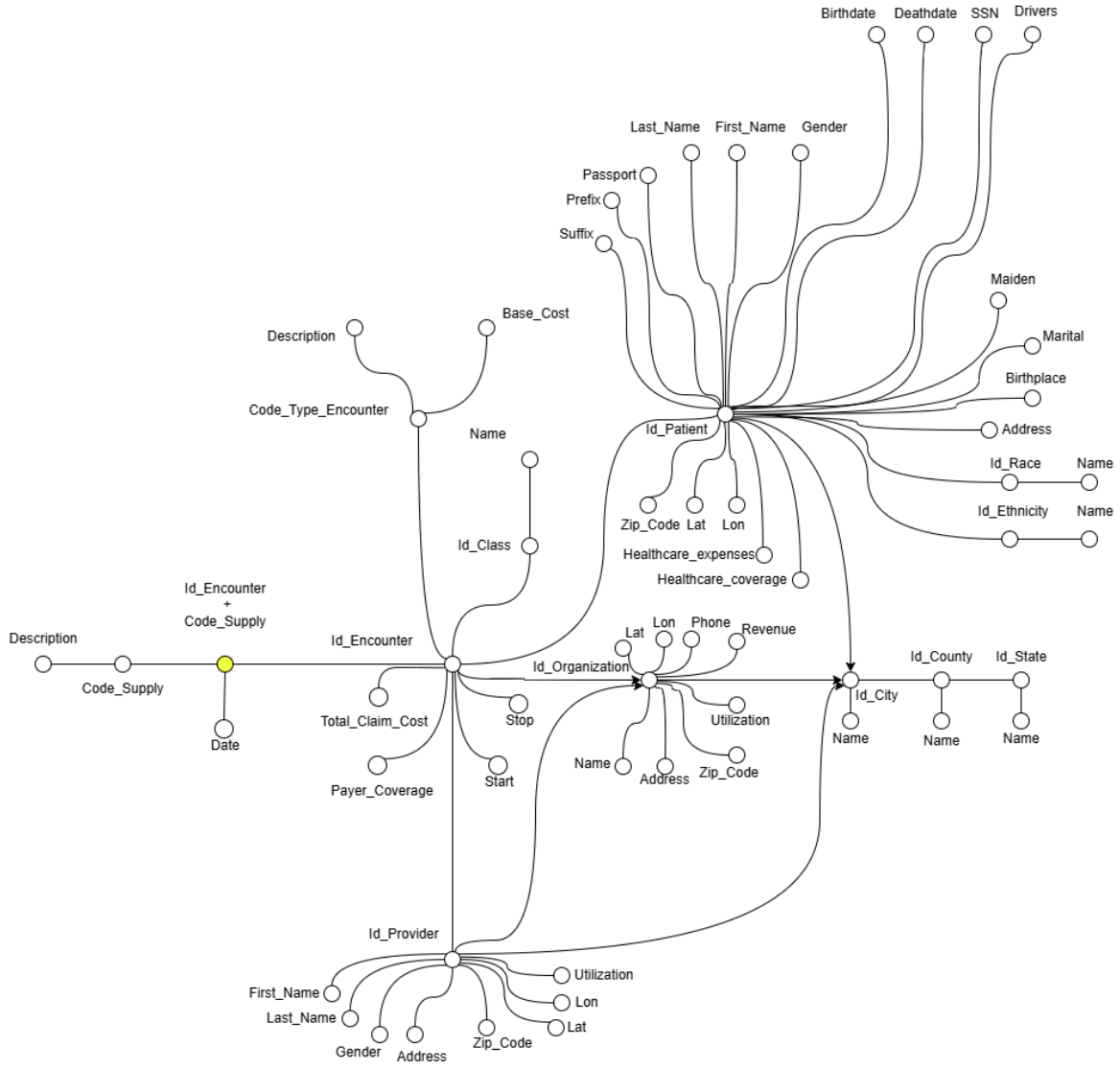


Figure 30: Complete attribute tree of the "Require" fact

For the construction of the final attribute tree of the fact "Require", we pruned:

- All attributes linked to "Id_Encounter", except for "Id_Organization"
- All attributes of "Id_Organization", except for "Id_City"
- "Name" from "Id_City"
- "Id_State" and its subtree

We then applied the grafting operation to the following:

- "Id_County" of "Id_City"

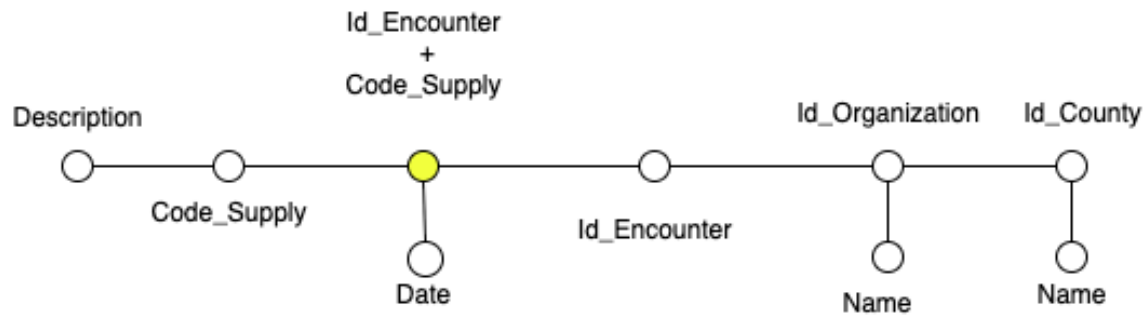


Figure 31: Pruned and grafted attribute tree of the "Require" fact

Based on the final attribute tree, the fact schema was constructed (Figure 32). The selected dimensions are: Date, Encounter and Supply. The chosen measure is Amount_Supply, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, a Date hierarchy was added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 33).

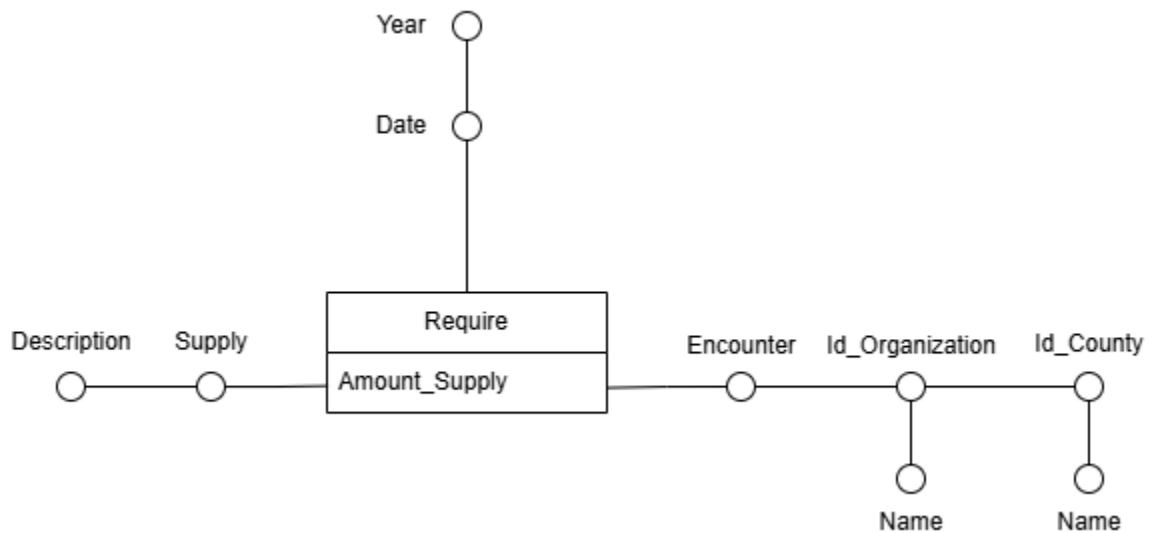


Figure 32: Fact schema of the "Require" fact

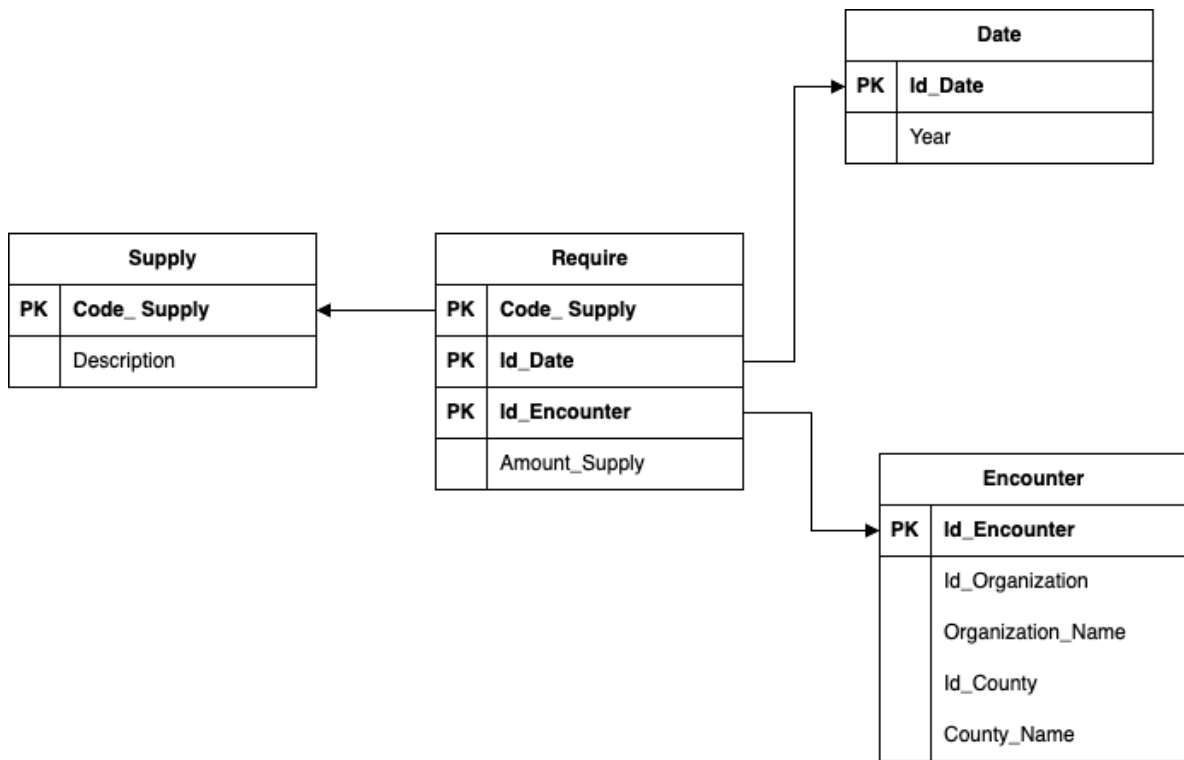


Figure 33: Star schema of the "Require" fact

4.9 Take_Place_In

The N-to-N relationship between Organization and Encounter was reified by creating the entity Take_Place_In, Figure 34. The attribute tree was then constructed, as shown in Figure 35.

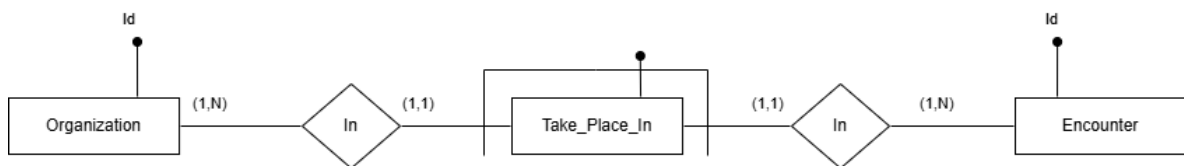


Figure 34: Reification of the "Take_Place_In" relationship

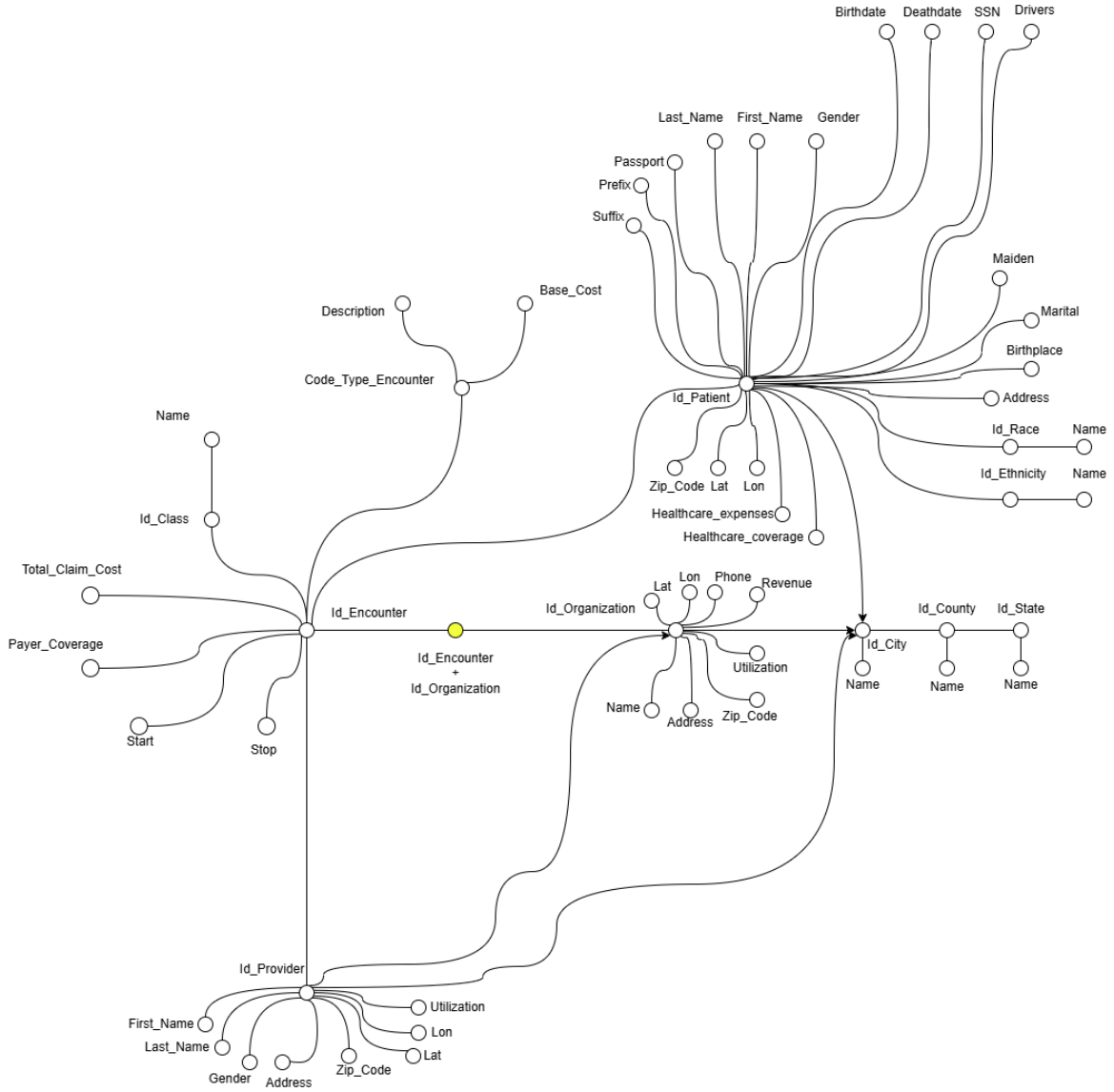


Figure 35: Complete attribute tree of the "Take_Place_In" fact

For the construction of the final attribute tree of the fact "Take_Place_In", we pruned:

- All attributes linked to "Id_Encounter", except for "Id_Patient" and "Start"
- All attributes of "Id_Organization", except for "Name"

We then applied the grafting operation to the following:

- "Birthdate" and "Gender" of "Id_Patient"

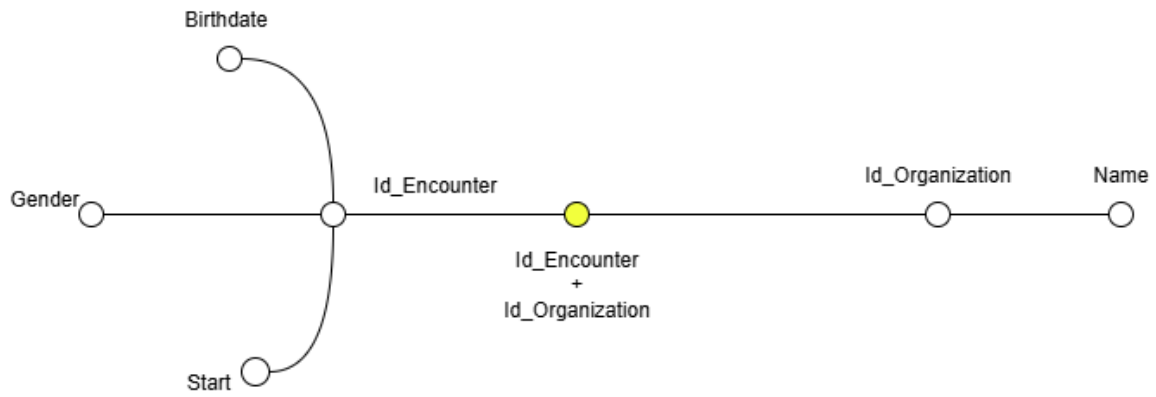


Figure 36: Pruned and grafted attribute tree of the "Take.Place.In" fact

Based on the final attribute tree, the fact schema was constructed (Figure 37). The selected dimensions are: Encounter and Organization. The chosen measure is Amount_Visit, obtained by aggregating the attributes from the tree (see Glossary of Measures 4.10). Finally, several hierarchy was added to facilitate data aggregation during query operations. In the end, the star schema was built (Figure 38).

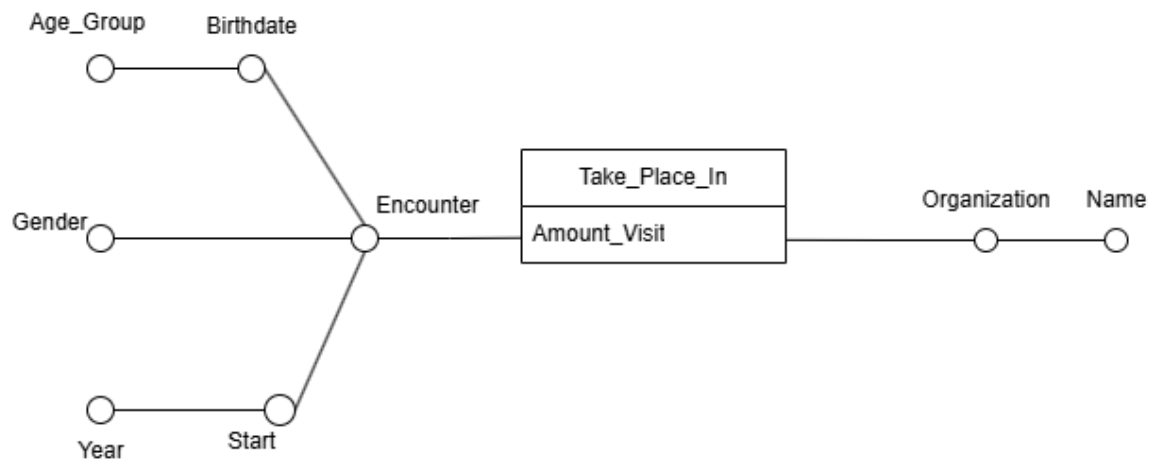


Figure 37: Fact schema of the "Take.Place.In" fact

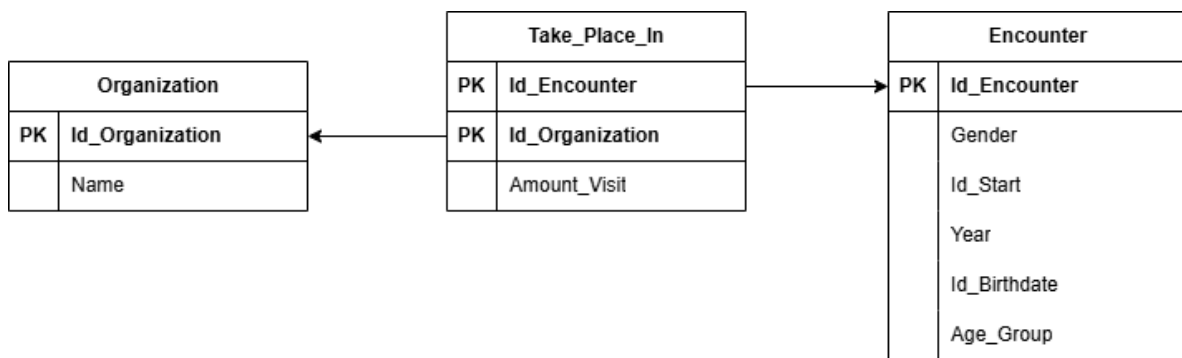


Figure 38: Star schema of the "Take_Place_In" fact

4.10 Glossary of Measures

Data mart	Measures
Administer	<ul style="list-style-type: none"> Amount_Immunization \leftarrow COUNT(Encounter.Id_Encounter)
Affected_By	<ul style="list-style-type: none"> Amount_Allergy \leftarrow COUNT(Encounter.Id_Encounter) GROUP BY Allergy.Code_Allergy Amount_Condition \leftarrow COUNT(Encounter.Id_Encounter) GROUP BY Condition.Code_Condition
Carry_Out	<ul style="list-style-type: none"> Total_Cost \leftarrow SUM(Procedure.Base_Cost)
Conduct	<ul style="list-style-type: none"> Amount_Visit \leftarrow COUNT(Encounter.Id_Encounter)
Plan	<ul style="list-style-type: none"> Amount_Patient \leftarrow COUNT(Encounter.Id_Encounter) Average_Days_Of_Hospitalization \leftarrow AVG(Start.Days_Of_Hospitalization) Average_Age \leftarrow AVG(Encounter.Age)
Require	<ul style="list-style-type: none"> Amount_Supply \leftarrow COUNT(Encounter.Id_Encounter)
Take_Place_In	<ul style="list-style-type: none"> Amount_Visit \leftarrow COUNT(Encounter.Id_Encounter)

Table 3: Facts identified for design and associated queries.

5 Use of the Data Warehouse

The Data Marts identified during the conceptual design phase were implemented. The software used for graphical representation is Tableau Software, a visual analytics platform that empowers individuals and organizations to understand and utilize data for problem-solving.

5.1 Administer

- How many individuals, in the age group between 10 and 50 years, were vaccinated against HPV between the years 2015 and 2019, in the counties of Bristol, Essex, and Middlesex, grouped by sex, county, and year of administration?

The selected dimensions are: Id_Date, County_Name, and Gender. The chosen measure is Amount_Immunization. Figure 39 shows the histogram resulting from the query.

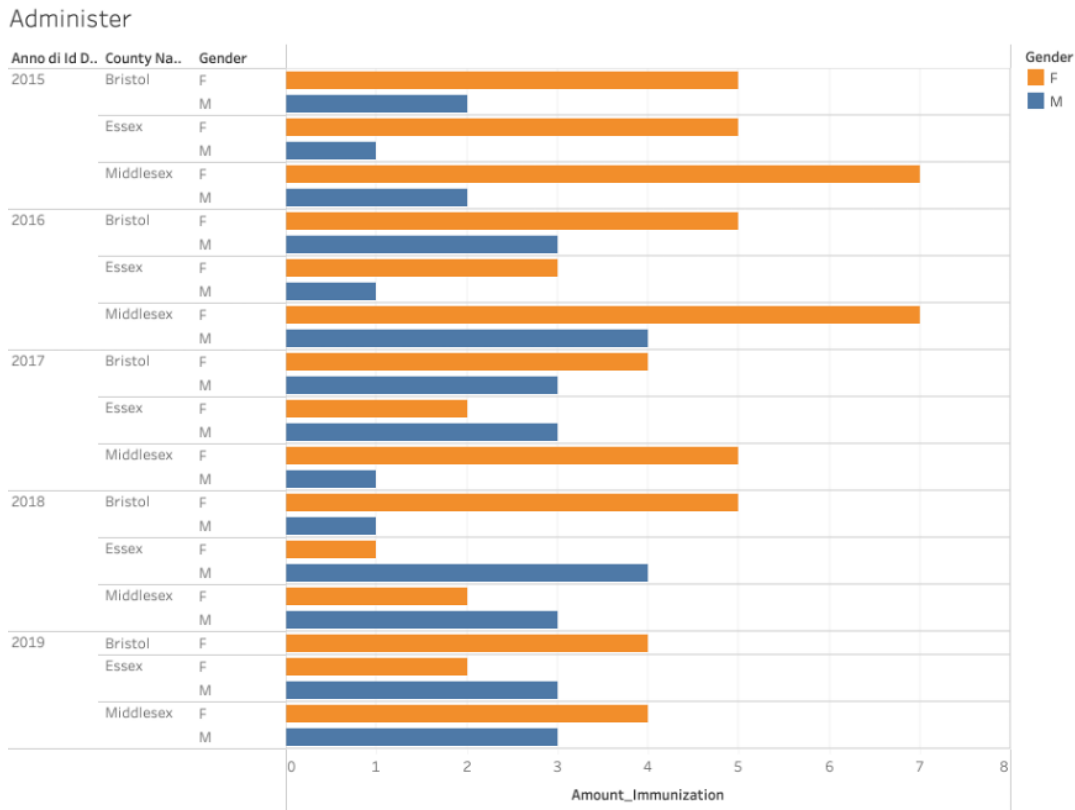


Figure 39: "Administer" query analyzed on tableau

5.2 Affected_By

- What are the most common allergies, grouped by patient age group, and county where they were recorded?

The selected dimensions are: Description of the allergy, County_Name and Age_Group. The chosen measure is Amount_Allergy. Figure 40 shows the histogram resulting from the query.

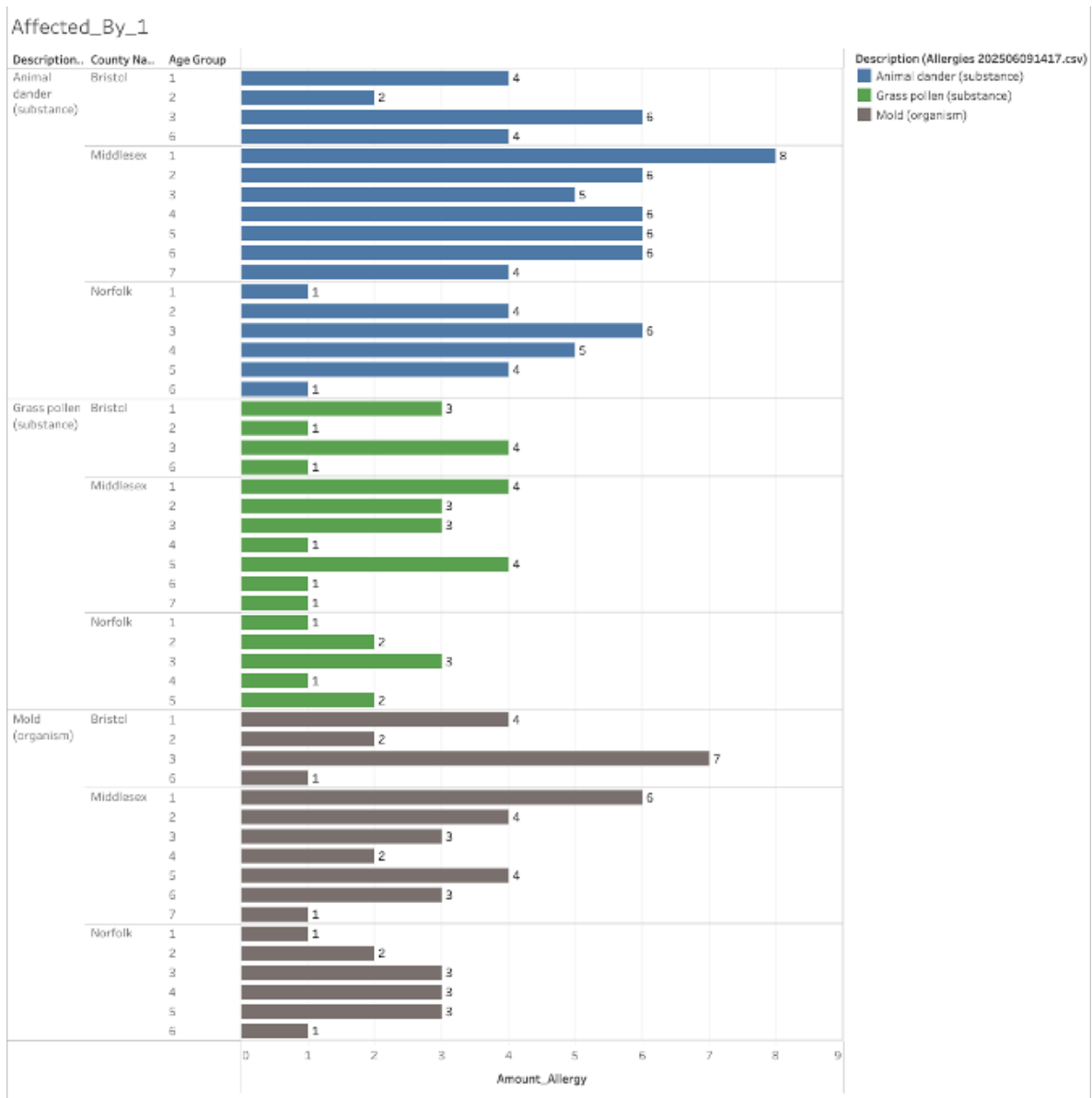


Figure 40: First "Affected_By" query analyzed on tableau

- What is the incidence of the various side effects associated with each allergy, broken down by patient age group?

The selected dimensions are: Description (of the side effect), Description (of the allergen), and Age_Group. The chosen measure is Amount_Conditions. Figure 41 shows the histogram resulting from the query. Note that only two allergens are plotted, because of readability issues.

Affected_By_2

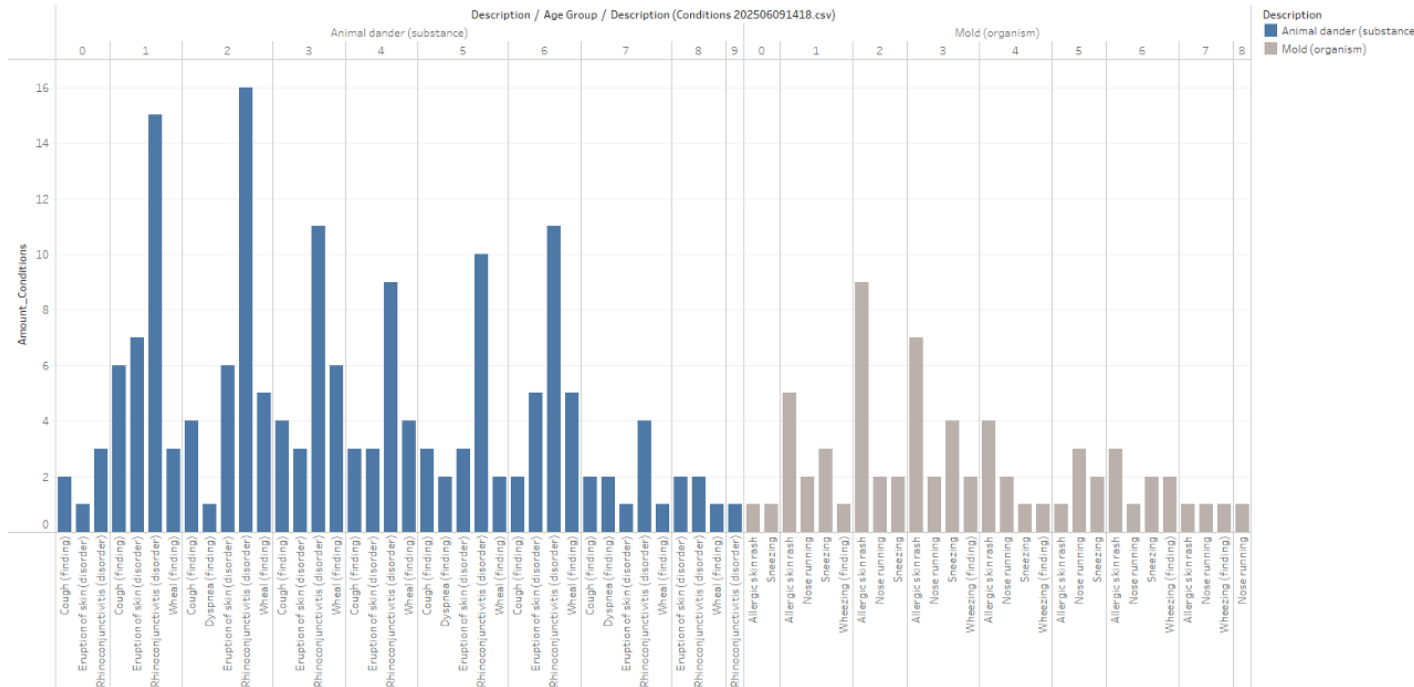


Figure 41: Second "Affected_By" query analyzed on tableau

5.3 Plan

- What is the number of respiratory therapies recorded in various hospitals, broken down by gender and age group?

The selected dimensions are: Organization_Name, Gender and Age_Group. The chosen measure is Amount_Patients. Figure 42 shows the histogram resulting from the query. Note that only four Organizations are plotted, because of readability issues.

Plan_1

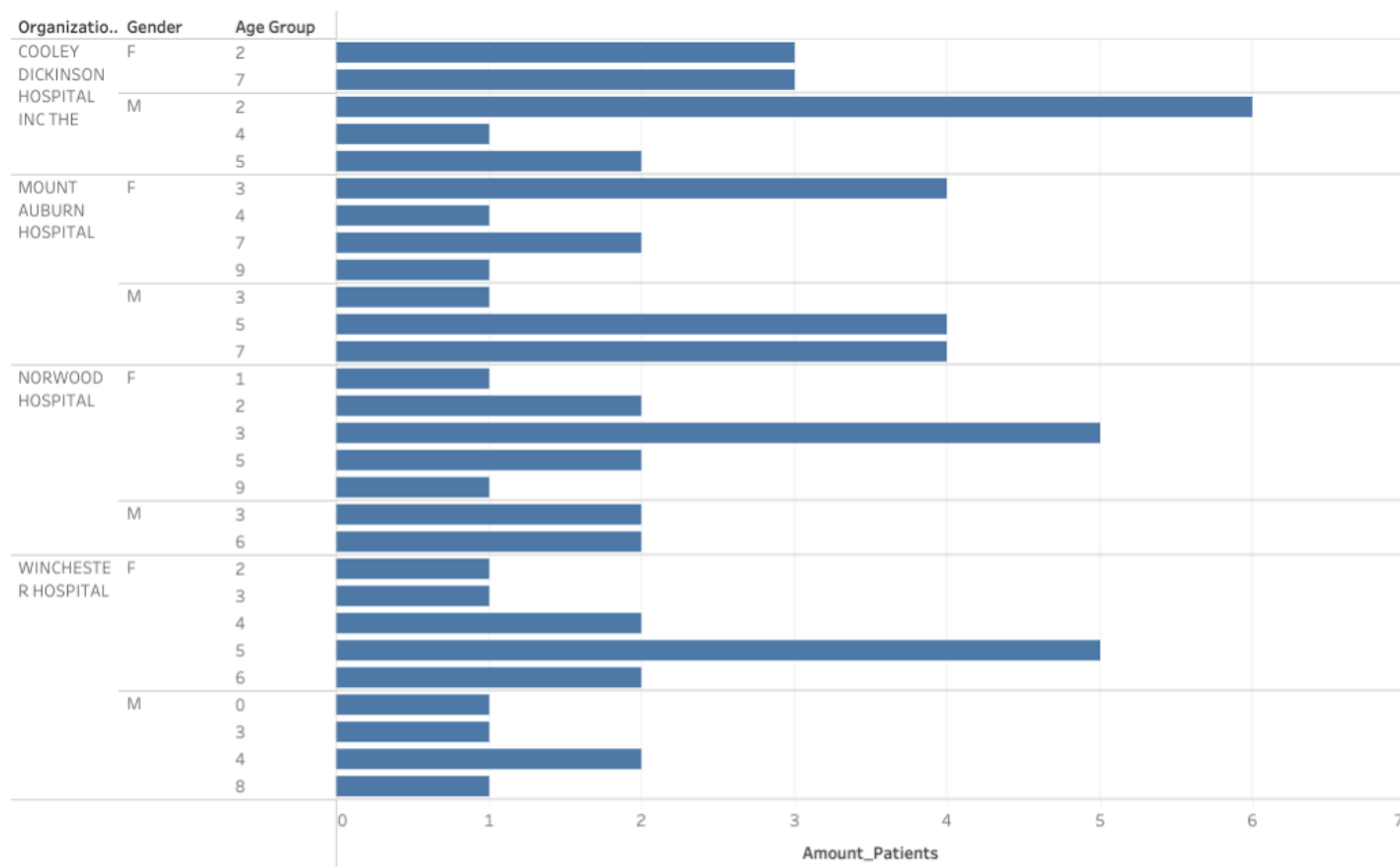


Figure 42: First "Plan" query analyzed on tableau

- What is the average age of patients affected by acute bronchitis or prediabetes in each county recorded in the last five years?

The selected dimensions are: Id_Start, Description of the disease and County_Name. The chosen measure is Average_Age_Of_Patients. Figure 42 shows the histogram resulting from the query.

Plan_2

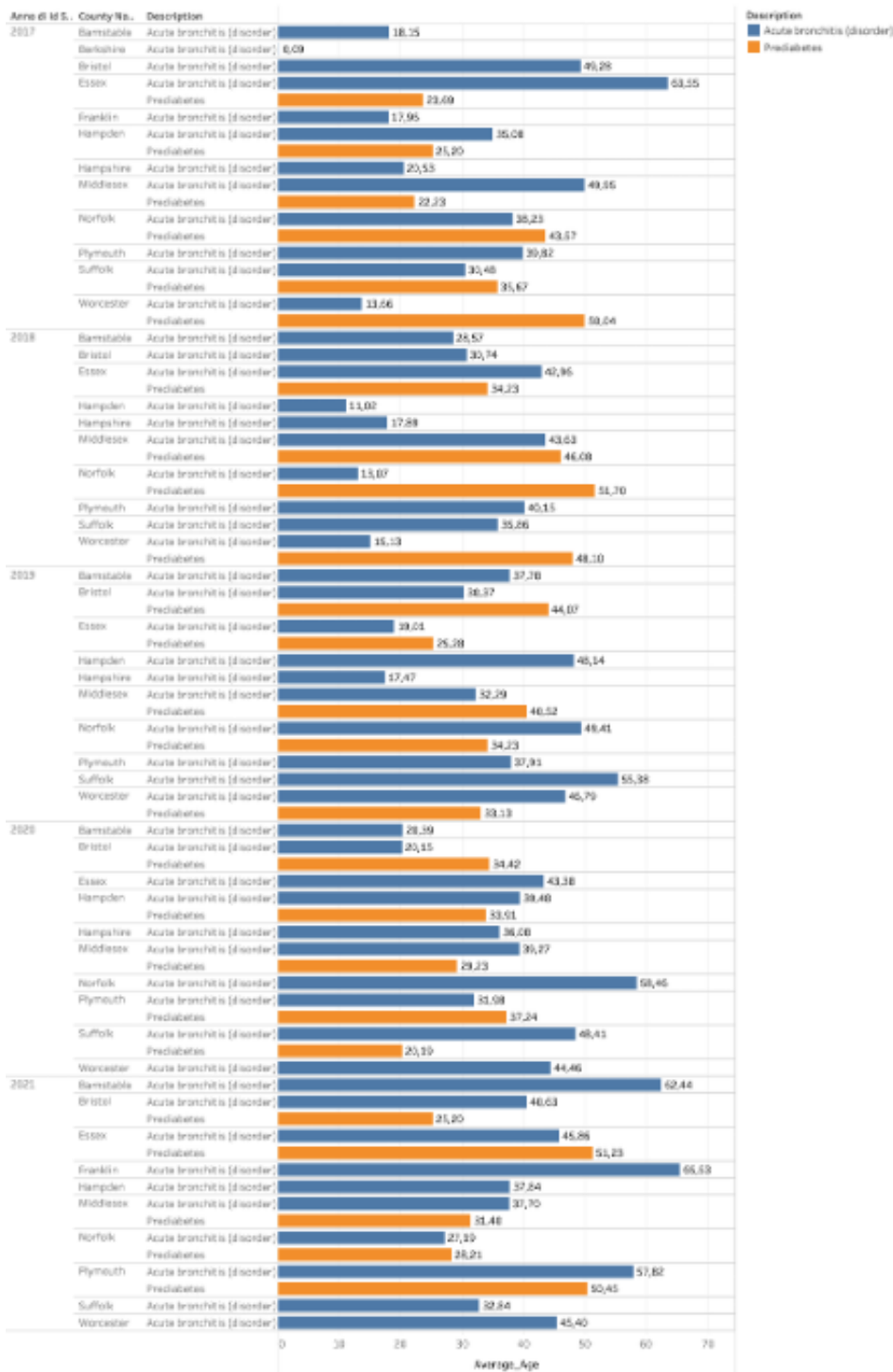


Figure 43: Second "Plan" query analyzed on tableau

- What is the average number of treatment days required to complete a therapy plan for various diseases in each hospital over the years?

The selected dimensions are: Description of the careplan, Description of the disease and Year of start of the plan. The chosen measure is Average_Days_Of_Hospitalization. Figure 43 shows the histogram resulting from the query.

Plan_3

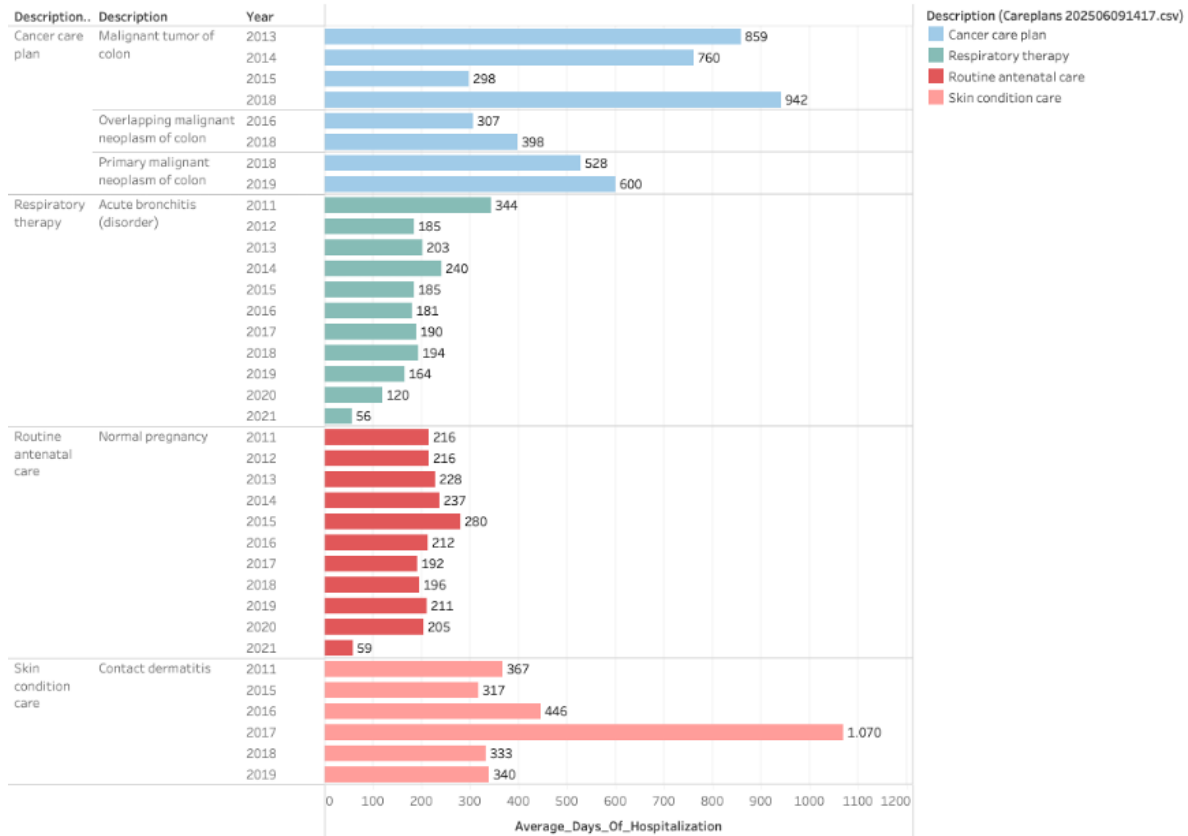


Figure 44: Third "Plan" query analyzed on tableau