

# **Tidy R programming with the OMOP common data model**

Edward Burn, Adam Black, Marti Catala, Berta Raventós

2023-02-05T00:00:00+00:00

# Table of contents

<b>Preface</b>	<b>4</b>
<b>1 Getting started</b>	<b>5</b>
1.1 A first data analysis in R with a database . . . . .	5
1.2 Getting set up . . . . .	5
1.3 Taking a peek at the data . . . . .	6
1.4 Inserting data into a database . . . . .	7
1.5 Translation from R to SQL . . . . .	8
1.6 Example analysis . . . . .	9
1.7 Further reading . . . . .	13
<b>2 Creating a reference to a database using the OMOP common data model</b>	<b>14</b>
2.1 Connecting to a database from R using DBI . . . . .	14
2.2 Creating a reference to the OMOP common data model . . . . .	14
2.3 Database snapshot . . . . .	17
2.4 Further reading . . . . .	17
<b>3 Exploring the CDM</b>	<b>18</b>
3.0.1 tally() . . . . .	18
3.0.2 summarise() . . . . .	19
3.0.3 group_by() . . . . .	20
3.0.4 filter() . . . . .	21
<b>4 Adding a cohort</b>	<b>24</b>
4.1 Finding codes with CodelistGenerator . . . . .	24
4.2 Defining a cohort with capr . . . . .	24
4.3 Adding a cohort to your cdm reference . . . . .	24
<b>5 Working with databases from R</b>	<b>25</b>
5.1 show_query() . . . . .	26
5.2 filter(), select(), mutate() . . . . .	26
5.3 right_join(), left_join(), inner_join(), and anti_join() . . . . .	26
5.4 summarise() . . . . .	26
5.5 collect() and compute() . . . . .	26
5.6 working with dates . . . . .	26
5.7 working with strings . . . . .	26

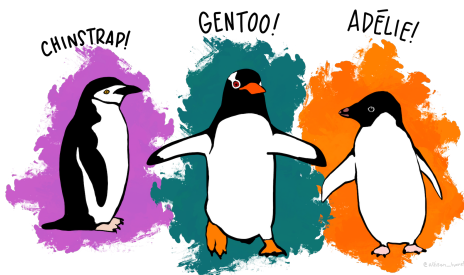
5.8	bespoke sql . . . . .	26
<b>8</b>	<b>right_join(), left_join(), inner_join(), and anti_join()</b>	<b>29</b>
8.0.1	and union() and union_all() . . . . .	29
<b>9</b>	<b>Getting to tidy data</b>	<b>30</b>
9.0.1	compute() . . . . .	30
9.0.2	collect() . . . . .	30
9.0.3	pull() . . . . .	30
<b>10</b>	<b>Analysis in R</b>	<b>31</b>
<b>11</b>	<b>Organising data analyses with projects and renv</b>	<b>32</b>
	<b>References</b>	<b>33</b>

# Preface

The source code for the book can be found at this [Github repository](#). Please open an issue there if you have a question or suggestion. Pull requests with suggested changes and additions are also most welcome.

# 1 Getting started

## 1.1 A first data analysis in R with a database



*Artwork by @allison\_horst*

Before we start thinking about working with health care data spread across the OMOP common data model, let's first do a quick data analysis using a simpler dataset. For this we'll use data from [palmerpenguins package](#), which contains data on penguins collected from the [Palmer Station](#) in Antarctica.

## 1.2 Getting set up

Assuming that you have R and RStudio already set up, first we need to install a few packages not included in base R if we don't already have them.

```
install.packages("dplyr")
install.packages("ggplot2")
install.packages("DBI")
install.packages("duckdb")
install.packages("palmerpenguins")
```

Once installed, we can load them like so.

```
library(dplyr)
library(ggplot2)
library(DBI)
library(duckdb)
library(palmerpenguins)
```

## 1.3 Taking a peek at the data

We can get an overview of the data using the `glimpse()` command.

```
glimpse(penguins)
```

```
Rows: 344
Columns: 8
$ species      <fct> Adelie, Adelie, Adelie, Adelie, Adelie, Adelie, Adel~
$ island       <fct> Torgersen, Torgersen, Torgersen, Torgersen, Torgerse~
$ bill_length_mm <dbl> 39.1, 39.5, 40.3, NA, 36.7, 39.3, 38.9, 39.2, 34.1, ~
$ bill_depth_mm <dbl> 18.7, 17.4, 18.0, NA, 19.3, 20.6, 17.8, 19.6, 18.1, ~
$ flipper_length_mm <int> 181, 186, 195, NA, 193, 190, 181, 195, 193, 190, 186~
$ body_mass_g   <int> 3750, 3800, 3250, NA, 3450, 3650, 3625, 4675, 3475, ~
$ sex          <fct> male, female, female, NA, female, male, female, male~
$ year         <int> 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007, 2007~
```

Or we could take a look at the first rows of the data using `head()`

```
head(penguins, 5)
```

```
# A tibble: 5 x 8
  species island bill_length_mm bill_depth_mm flipper_l~1 body_~2 sex year
  <fct>   <fct>         <dbl>         <dbl>         <int>   <int> <fct> <int>
1 Adelie Torgersen      39.1           18.7           181     3750 male  2007
2 Adelie Torgersen      39.5           17.4           186     3800 fema~  2007
3 Adelie Torgersen      40.3            18            195     3250 fema~  2007
4 Adelie Torgersen      NA              NA              NA        NA <NA>   2007
5 Adelie Torgersen      36.7           19.3           193     3450 fema~  2007
# ... with abbreviated variable names 1: flipper_length_mm, 2: body_mass_g
```

## 1.4 Inserting data into a database

Let's put our penguins data into a duckdb database. We create the duckdb database, add the penguins data, and then create a reference to the table containing the data.

```
db<-dbConnect(duckdb::duckdb(), dbdir=":memory:")
dbWriteTable(db, "penguins", penguins)
penguins_db<-tbl(db, "penguins")
```

Now the data is in a database we could use SQL to get the first rows that we saw before

```
dbGetQuery(db, "SELECT * FROM penguins LIMIT 5")
```

	species	island	bill_length_mm	bill_depth_mm	flipper_length_mm	body_mass_g
1	Adelie	Torgersen	39.1	18.7	181	3750
2	Adelie	Torgersen	39.5	17.4	186	3800
3	Adelie	Torgersen	40.3	18.0	195	3250
4	Adelie	Torgersen	NA	NA	NA	NA
5	Adelie	Torgersen	36.7	19.3	193	3450

	sex	year
1	male	2007
2	female	2007
3	female	2007
4	<NA>	2007
5	female	2007

But we could also use the same R code as before

```
head(penguins_db, 5)
```

```
# Source:   SQL [5 x 8]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/:memory:]
  species island  bill_length_mm bill_depth_mm flipper_l~1 body_~2 sex   year
  <fct>   <fct>         <dbl>         <dbl>         <int>   <int> <fct> <int>
1 Adelie  Torgersen      39.1           18.7           181     3750 male   2007
2 Adelie  Torgersen      39.5           17.4           186     3800 fema~  2007
3 Adelie  Torgersen      40.3            18            195     3250 fema~  2007
4 Adelie  Torgersen      NA              NA              NA        NA <NA>   2007
5 Adelie  Torgersen      36.7           19.3           193     3450 fema~  2007
# ... with abbreviated variable names 1: flipper_length_mm, 2: body_mass_g
```

## 1.5 Translation from R to SQL

The magic here is provided by dbplyr which takes the R code and converts it into SQL, which in this case looks like

```
head(penguins_db, 1) %>%  
  show_query()
```

```
<SQL>  
SELECT *  
FROM penguins  
LIMIT 1
```

More complicated SQL can also be written in what might be familiar dplyr code, for example

```
penguins_db %>%  
  group_by(species) %>%  
  summarise(min_bill_length_mm=min(bill_length_mm),  
            median_bill_length_mm=median(bill_length_mm),  
            max_bill_length_mm=max(bill_length_mm)) %>%  
  mutate(min_max_bill_length_mm=paste0(min_bill_length_mm,  
                                       " to ",  
                                       max_bill_length_mm)) %>%  
  select("species",  
        "median_bill_length_mm",  
        "min_max_bill_length_mm")
```

```
# Source:   SQL [3 x 3]  
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/:memory:]  
  species   median_bill_length_mm min_max_bill_length_mm  
  <fct>                <dbl> <chr>  
1 Adelie                38.8 32.1 to 46.0  
2 Gentoo                47.3 40.9 to 59.6  
3 Chinstrap            49.6 40.9 to 58.0
```

with the corresponding SQL looking like

```
penguins_db %>%  
  group_by(species) %>%  
  summarise(min_bill_length_mm=min(bill_length_mm),
```



```

        median_bill_length_mm=median(bill_length_mm),
        max_bill_length_mm=max(bill_length_mm)) %>%
mutate(min_max_bill_length_mm=paste0(min, " to ", max)) %>%
select("species",
       "median_bill_length_mm",
       "min_max_bill_length_mm") %>%
show_query()

```

<SQL>

```

SELECT
  species,
  median_bill_length_mm,
  CONCAT_WS(' ', .Primitive("min"), ' to ', .Primitive("max")) AS min_max_bill_length_mm
FROM (
  SELECT
    species,
    MIN(bill_length_mm) AS min_bill_length_mm,
    PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY bill_length_mm) AS median_bill_length_mm,
    MAX(bill_length_mm) AS max_bill_length_mm
  FROM penguins
  GROUP BY species
) q01

```

## 1.6 Example analysis

Let's start by getting a count by species

```

penguins_db %>%
  group_by(species) %>%
  count()

```

```

# Source:   SQL [3 x 2]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/:memory:]
  species      n
  <fct>      <dbl>
1 Adelie     152
2 Gentoo     124
3 Chinstrap   68

```

Now suppose we are particularly interested in the body mass variable. We can first notice that there are a couple of missing records for this.

```
penguins_db %>%  
  mutate(missing_body_mass_g = if_else(  
    is.na(body_mass_g), 1, 0  
  )) %>%  
  group_by(species, missing_body_mass_g) %>%  
  tally()
```

```
# Source:   SQL [5 x 3]  
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/:memory:]  
# Groups:   species  
  species  missing_body_mass_g      n  
  <fct>                <dbl> <dbl>  
1 Adelie                0    151  
2 Adelie                1      1  
3 Gentoo                0    123  
4 Gentoo                1      1  
5 Chinstrap            0     68
```

We can get the mean for each of the species (dropping those two missing records).

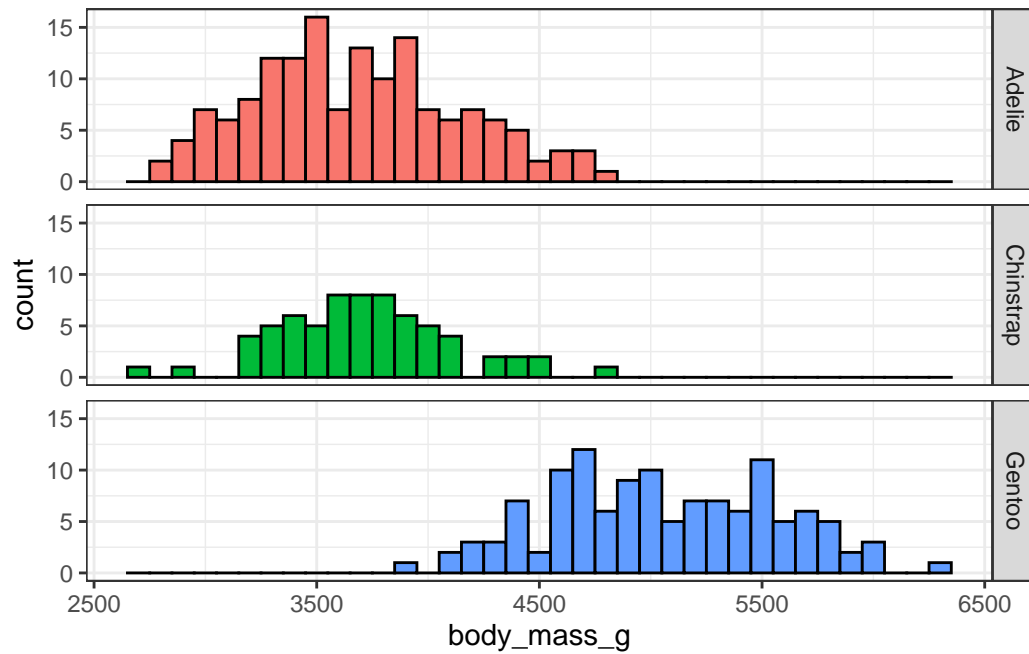
```
penguins_db %>%  
  group_by(species) %>%  
  summarise(mean_body_mass_g = round(mean(body_mass_g, na.rm=TRUE), 0))
```

```
# Source:   SQL [3 x 2]  
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/:memory:]  
  species  mean_body_mass_g  
  <fct>                <dbl>  
1 Adelie                3701  
2 Gentoo                5076  
3 Chinstrap            3733
```

We can then also do a histogram for each of the species. For this we need to bring the data into R so that we can work with `ggplot()`, and we use `collect()` to do this.

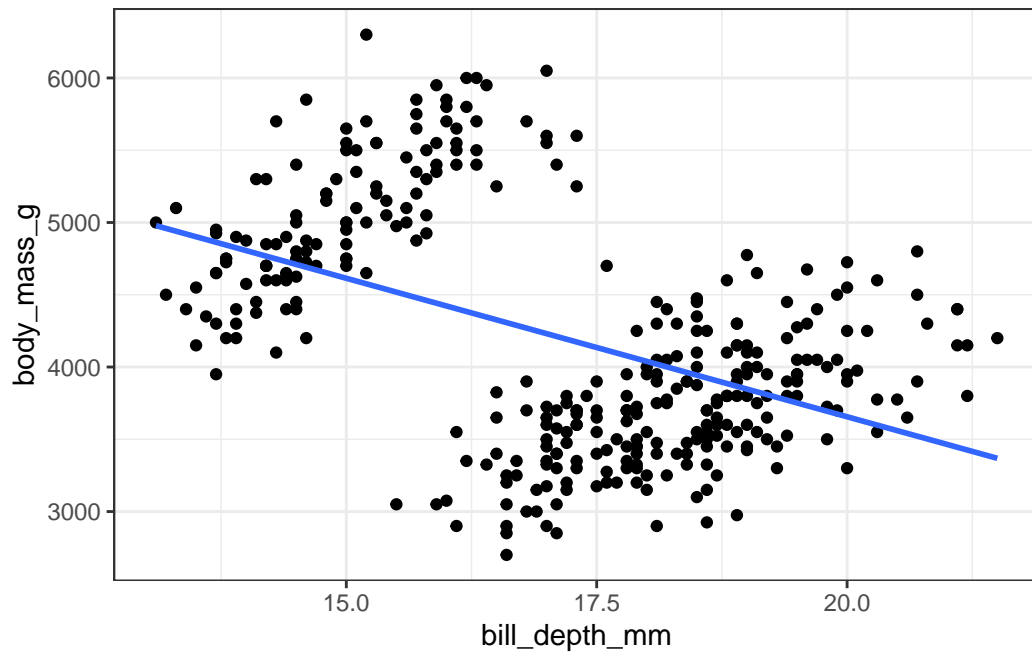
```
penguins_db %>%  
  collect() %>%
```

```
ggplot(aes(group=species, fill=species))+
  facet_grid(species~ .) +
  geom_histogram(aes(body_mass_g), colour="black", binwidth = 100)+
  theme_bw()+
  theme(legend.position = "none")
```



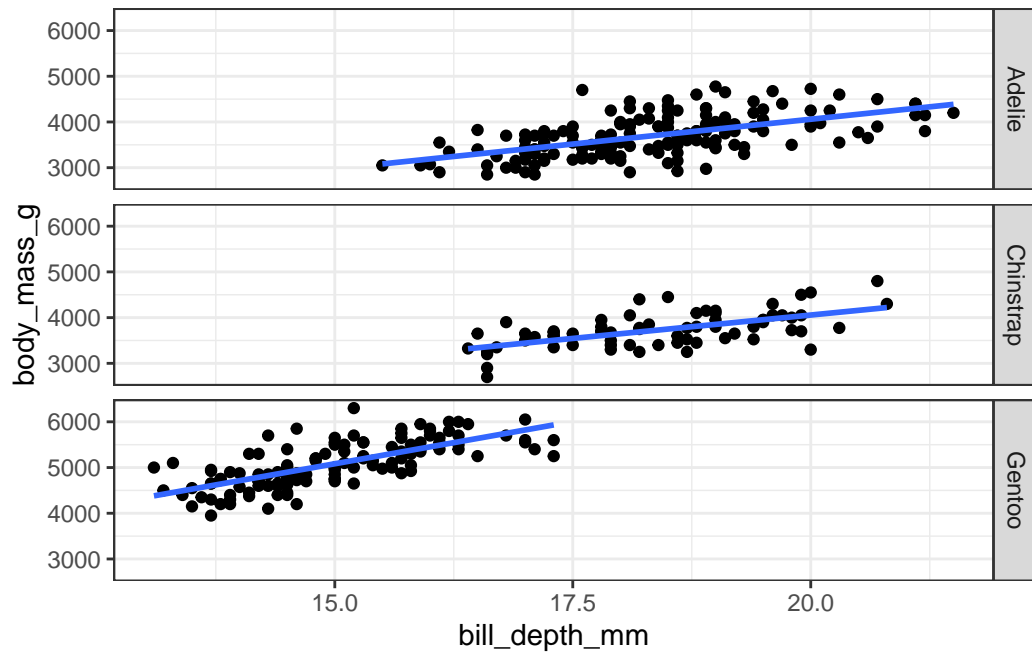
How about the relationship between body mass and bill depth?

```
penguins %>%
  collect() %>%
  ggplot(aes(x=bill_depth_mm,y=body_mass_g))+
  geom_point()+
  geom_smooth(method="lm",se=FALSE )+
  theme_bw()+
  theme(legend.position = "none")
```



But what about by species?

```
penguins %>%  
  collect() %>%  
  ggplot(aes(x=bill_depth_mm,y=body_mass_g))+  
  facet_grid(species~ .) +  
  geom_point()+  
  geom_smooth(method="lm",se=FALSE )+  
  theme_bw()+  
  theme(legend.position = "none")
```



As well as having an example of working with data in database from R, you also have an example of [Simpson's paradox](#)! And now we've reached the end of this example, we can close the database like so

## 1.7 Further reading

- [R for Data Science \(Chapter 13: Relational data\)](#)
- [Writing SQL with dbplyr](#)
- [Data Carpentry: SQL databases and R](#)

## 2 Creating a reference to a database using the OMOP common data model

### 2.1 Connecting to a database from R using DBI

Database connections from R can be made using the [DBI package](#). The back-end for DBI is facilitated by database specific driver packages. As an example, lets say we want to work with a local duckdb from R. In this case the we can use the duckdb R package as the driver.

```
library(DBI)
db<-dbConnect(duckdb::duckdb(), dbdir=":memory:")
```

If we instead wanted to connect to other database management systems, these connections would be supported by the associated back-end packages and could look something like the below example for Postgres:

```
# Postgres
db <- DBI::dbConnect(RPostgres::Postgres(),
                     dbname = Sys.getenv("CDM5_POSTGRESQL_DBNAME"),
                     host = Sys.getenv("CDM5_POSTGRESQL_HOST"),
                     user = Sys.getenv("CDM5_POSTGRESQL_USER"),
                     password = Sys.getenv("CDM5_POSTGRESQL_PASSWORD"))
```

### 2.2 Creating a reference to the OMOP common data model

As seen in the previous chapter, once a connection to the database has been created then we could create references to the various tables in the database and build queries using in a familiar dplyr style. However, as we already know what the structure of the OMOP CDM looks like, we can avoid the overhead of building *ad hoc* references by instead using the CDMConnector package to quickly create a reference to the OMOP CDM data as a whole.

If you don't already have it installed, the first step would be to install CDMConnector from CRAN.

```
install.packages("CDMConnector")
```

Once we have it installed, we can then load it as with other R packages.

```
library(CDMConnector)
```

For this example, we'll use the Eunomia example data contained in a duckdb database. First we need to download the data. And once downloaded, make sure to add the path to your Renviron.

```
# change pathToData to the location you want to save the data
CDMConnector::downloadEunomiaData(
  pathToData = here::here(),
  overwrite = TRUE
)
# once downloaded, save your pathToData to your Renviron (and then restart R)
# EUNOMIA_DATA_FOLDER="....."
```

```
db <- DBI::dbConnect(duckdb::duckdb(),
  dbdir = CDMConnector::eunomia_dir())
cdm <- CDMConnector::cdm_from_con(con = db,
  cdm_schema = "main")

cdm
```

```
# OMOP CDM reference (tbl_duckdb_connection)
```

Tables: person, observation\_period, visit\_occurrence, visit\_detail, condition\_occurrence, drug\_exposure

Once we have created the our reference to the overall OMOP CDM, we can reference specific tables using the “\$” operator or [[ “”]].

```
cdm$observation_period
```

```
# Source:   table<main.observation_period> [?? x 5]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp]
  observation_period_id person_id observation_period_start~1 observat~2 perio~3
      <dbl>         <dbl> <date>                                <date>         <dbl>
1           6           6 1963-12-31                    2007-02-06  4.48e7
2          13          13 2009-04-26                    2019-04-14  4.48e7
3          27          27 2002-01-30                    2018-11-21  4.48e7
```

```

4          16          16 1971-10-14          2017-11-02 4.48e7
5          55          55 2009-05-30          2019-03-23 4.48e7
6          60          60 1990-11-21          2019-01-23 4.48e7
7          42          42 1909-11-03          2019-03-13 4.48e7
8          33          33 1986-05-12          2018-09-10 4.48e7
9          18          18 1965-11-17          2018-11-07 4.48e7
10         25          25 2007-03-18          2019-04-07 4.48e7
# ... with more rows, and abbreviated variable names
#   1: observation_period_start_date, 2: observation_period_end_date,
#   3: period_type_concept_id

```

```
cdm[["observation_period"]]
```

```

# Source:   table<main.observation_period> [?? x 5]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
  observation_period_id person_id observation_period_start~1 observat~2 perio~3
                <dbl>      <dbl> <date>                <date>      <dbl>
1                 6         6 1963-12-31          2007-02-06 4.48e7
2                13        13 2009-04-26          2019-04-14 4.48e7
3                27        27 2002-01-30          2018-11-21 4.48e7
4                16        16 1971-10-14          2017-11-02 4.48e7
5                55        55 2009-05-30          2019-03-23 4.48e7
6                60        60 1990-11-21          2019-01-23 4.48e7
7                42        42 1909-11-03          2019-03-13 4.48e7
8                33        33 1986-05-12          2018-09-10 4.48e7
9                18        18 1965-11-17          2018-11-07 4.48e7
10               25        25 2007-03-18          2019-04-07 4.48e7
# ... with more rows, and abbreviated variable names
#   1: observation_period_start_date, 2: observation_period_end_date,
#   3: period_type_concept_id

```

When we created our reference we could have also specified a subset of cdm tables that we want to read:

```

cdm <- CDMConnector::cdm_from_con(db,
                                cdm_tables = c("person","observation_period"))
cdm

```

```
# OMOP CDM reference (tbl_duckdb_connection)
```

```
Tables: person, observation_period
```



Moreover, we can also specify a write schema and the tables that we are interested in it when creating our reference. For example, if we wanted to create a reference to the person and observation period tables in the common data model along with cohort tables in a schema we have write access to, we could do this like so:

```
cdm <- CDMConnector::cdm_from_con(db,
  cdm_schema = "main",
  cdm_tables = c("person", "observation_period"),
  write_schema = "results",
  cohort_tables = c("exposure_cohort", "outcome_cohort"))
```

## 2.3 Database snapshot

We can also use `CDMConnector` to provide a summary of the metadata for the OMOP CDM data we have connected to

```
cdm_from_con(con = db,
  cdm_schema = "main") %>%
  snapshot() %>%
  glimpse()
```

List of 7

```
$ cdm_source_name      : chr "Synthea synthetic health database"
$ cdm_version          : chr "v5.3.1"
$ cdm_holder           : chr "OHDSI Community"
$ cdm_release_date     : Date[1:1], format: "2019-05-25"
$ vocabulary_version   : chr "v5.0 18-JAN-19"
$ person_cnt           : num 2694
$ observation_period_cnt: num 5343
- attr(*, "class")= chr "cdm_snapshot"
```

## 2.4 Further reading

- [CDMConnector package](#)

## 3 Exploring the CDM

Let's first connect again to our Eunomia data and create the reference to the common data model.

```
library(DBI)
library(dplyr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

`filter`, `lag`

The following objects are masked from 'package:base':

`intersect`, `setdiff`, `setequal`, `union`

### 3.0.1 tally()

Let's say we want to get a count of the people in the person table. For this we can use the `tally` or `count` verbs from `dbplyr`

```
cdm$person %>%
  count()
```

```
# Source:   SQL [1 x 1]
```

```
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
```

```
      n
```

```
<dbl>
```

```
1  2694
```

This count was done on the database side, with the code we wrote in `dplyr` style translated into `sql`.

```
cdm$person %>%
  count() %>%
  show_query()
```

```
<SQL>
SELECT COUNT(*) AS n
FROM main.person
```

### 3.0.2 summarise()

Another way to get the same count would be to use the summarise verb

```
cdm$person %>%
  summarise(n = n())
```

```
# Source:   SQL [1 x 1]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
      n
<dbl>
1  2694
```

```
cdm$person %>%
  summarise(n = n())%>%
  show_query()
```

```
<SQL>
SELECT COUNT(*) AS n
FROM main.person
```

We can also use summarise for various other calculations

```
cdm$person %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))
```

```
# Source:   SQL [1 x 1]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
      median
<dbl>
1    1961
```

```
cdm$person %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))%>%
  show_query()
```

<SQL>

```
SELECT PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY year_of_birth) AS median
FROM main.person
```

### 3.0.3 group\_by()

What if we want to get a count of people in the person table by gender concept id? In this case we can use group\_by

```
cdm$person %>%
  group_by(gender_concept_id) %>%
  count()
```

# Source: SQL [2 x 2]

# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp

	gender_concept_id	n
	<dbl>	<dbl>
1	8532	1373
2	8507	1321

```
cdm$person %>%
  group_by(gender_concept_id) %>%
  count() %>%
  show_query()
```

<SQL>

```
SELECT gender_concept_id, COUNT(*) AS n
FROM main.person
GROUP BY gender_concept_id
```

Similarly we could use group\_by to calculate median year of birth by gender concept id.

```
cdm$person %>%
  group_by(gender_concept_id) %>%
```

```
summarise(median = median(year_of_birth, na.rm=TRUE))
```

```
# Source:   SQL [2 x 2]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
  gender_concept_id median
          <dbl>   <dbl>
1             8532   1961
2             8507   1961
```

```
cdm$person %>%
  group_by(gender_concept_id) %>%
  summarise(median = median(year_of_birth, na.rm=TRUE)) %>%
  show_query()
```

```
<SQL>
SELECT
  gender_concept_id,
  PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY year_of_birth) AS median
FROM main.person
GROUP BY gender_concept_id
```

### 3.0.4 filter()

Or if we wanted a count within only for those with a specific gender concept id we can use the filter verb to subset the data before summarising it

```
cdm$person %>%
  filter(gender_concept_id == "8532") %>%
  count()
```

```
# Source:   SQL [1 x 1]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
      n
  <dbl>
1  1373
```

```
cdm$person %>%
  filter(gender_concept_id == "8532") %>%
```

```
count() %>%
show_query()
```

```
<SQL>
SELECT COUNT(*) AS n
FROM main.person
WHERE (gender_concept_id = '8532')
```

Similarly we could have

```
cdm$person %>%
  filter(year_of_birth < 1970) %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))
```

```
# Source:   SQL [1 x 1]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
median
<dbl>
1    1955
```

```
cdm$person %>%
  filter(year_of_birth < 1970) %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))%>%
  show_query()
```

```
<SQL>
SELECT PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY year_of_birth) AS median
FROM main.person
WHERE (year_of_birth < 1970.0)
```

We can combine the above, with a filter, followed by a group\_by, and then followed by a summarise

```
cdm$person %>%
  filter(year_of_birth < 1970) %>%
  group_by(gender_concept_id) %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))
```

```
# Source:   SQL [2 x 2]
# Database: DuckDB 0.5.0 [eburn@Windows 10 x64:R 4.2.1/C:\Users\eburn\AppData\Local\Temp\Rtmp
  gender_concept_id median
          <dbl>   <dbl>
1             8532   1955
2             8507   1956
```

```
cdm$person %>%
  filter(year_of_birth < 1970) %>%
  group_by(gender_concept_id) %>%
  summarise(median = median(year_of_birth, na.rm=TRUE))%>%
  show_query()
```

```
<SQL>
SELECT
  gender_concept_id,
  PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY year_of_birth) AS median
FROM main.person
WHERE (year_of_birth < 1970.0)
GROUP BY gender_concept_id
```

## **4 Adding a cohort**

### **4.1 Finding codes with CodelistGenerator**

### **4.2 Defining a cohort with capr**

### **4.3 Adding a cohort to your cdm reference**



## 5 Working with databases from R

Let's start by taking some data and putting it in a database. Here we'll use an in-memory duckdb database, but the same code should work for other databases with only the connection details and the package used to connect to the database changing.

For this example let's use data on Darwin's finches as that seems rather appropriate ([link to wiki article on darwin finches](#))

```
# install packages
# commented out as you might already have them
# but if not then uncomment and run
# install.packages("DBI")
# install.packages("SQLite")
# install.packages("dbplyr")
# install.packages("dplyr")
#
# # load packages
# library(DBI)
# library(SQLite)
# library(dbplyr)
# library(dplyr)

# get data

# move into a database
```

## **5.1 show\_query()**

## **5.2 filter(), select(), mutate()**

## **5.3 right\_join(), left\_join(), inner\_join(), and anti\_join()**

## **5.4 summarise()**

## **5.5 collect() and compute()**

## **5.6 working with dates**

Here be dragons

## **5.7 working with strings**

## **5.8 bespoke sql**

Alternative approaches

1) Where to do computation

- Database side vs in local memory vs R

2) Scope of a package

3) Scope of analysis code All in one vs one at a time

# 6

strings

**7**

Dates

## 8 `right_join()`, `left_join()`, `inner_join()`, and `anti_join()`

### 8.0.1 `and`, `union()` and `union_all()`

## 9 Getting to tidy data

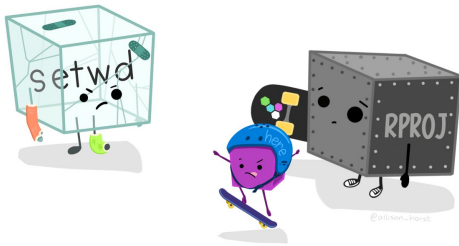
9.0.1 `compute()`

9.0.2 `collect()`

9.0.3 `pull()`

## 10 Analysis in R

# 11 Organising data analyses with projects and renv



*Artwork by @allison\_horst*



## References

# 12

Learning R